

LLMs for Hardware Security: Boon or Bane?

Rahul Kande, Vasudev Gohil, Matthew DeLorenzo, Chen Chen, Jeyavijayan Rajendran

Texas A&M University, USA

{rahulkande, gohil.vasudev, matthewdelorenzo, chenc, jv.rajendran}@tamu.edu

Abstract—Large language models (LLMs) have emerged as transformative tools within the hardware design and verification lifecycle, offering numerous capabilities in accelerating design processes. Recent research has showcased the efficacy of LLMs in translating design specifications into source code through hardware description languages. Researchers are also using LLMs to generate test cases and write assertion rules to bolster the detection of hardware vulnerabilities. Thus, the semiconductor industry is swiftly integrating LLMs into its design workflows. However, this adoption is not without its challenges.

While LLMs offer remarkable benefits, they concurrently introduce security concerns that demand a thorough examination. These concerns manifest as potential vulnerabilities indirectly introduced into the designs while generating the design code, or by directly equipping the attackers with novel avenues for exploitation. In this paper, we discuss the emerging security implications due to the capabilities introduced by LLMs in the context of hardware design verification, evaluate the capabilities of existing security detection and mitigation techniques, and highlight the possible future security attacks that use LLMs.

Index Terms—LLM, hardware, verification, security

I. INTRODUCTION

A. Hardware Design Verification

Hardware verification techniques can be primarily classified into regression-based testing and formal verification. Hardware regression-based testing involves repeatedly generating and simulating input tests with the target design (e.g., random regression [1]). Most regression-based techniques use a golden reference model or assertions inserted into the hardware to detect vulnerabilities [2]. Formal verification techniques are widely applied in the industry [3], [4], which explore design spaces exhaustively and prove if a design-under-test (DUT) satisfies specified properties [5]. Based on security specifications, designers write assertions and use commercial formal tools to prove if the design violates any requirements. While regression-based techniques scale to large designs, they suffer from their slow design coverage in detecting vulnerabilities. On the other hand, formal techniques can fully verify the given functional and security properties, but do not scale well to large and complex modern hardware. Also, writing properties requires expert knowledge of the DUT, which is error-prone and time-consuming [6], [7].

Recently, hardware fuzzing has shown its effectiveness in detecting vulnerabilities in large-scale designs such as modern processors [8]–[11]. Hardware fuzzing techniques generate test cases dynamically and analyze the execution logs of a DUT to see if the design violates any assertions [12] or triggers vulnerabilities [8]–[10], [13]–[15]. However, fuzzers require

time and compute resources to execute the hardware, which is precious given the market demand to shrink verification time.

B. Supply-Chain Security

The globalization of the integrated circuit (IC) supply chain has brought significant benefits to the production of electronic products, such as cost reductions, access to specialized expertise, and quicker market entry. However, it has also introduced complex hardware security issues that undermine trust in the IC ecosystem, with the threat of hardware Trojans (HTs) during IC fabrication being a particularly alarming concern. These malicious modifications, covertly inserted during the manufacturing process, can severely compromise the integrity and functionality of ICs upon activation, facilitating a variety of cyberattacks. Such attacks can alter an IC's operations, degrade performance, expose sensitive data, or even disable entire systems, necessitating sophisticated security strategies for their detection and neutralization [16].

IC security issues persist even after a chip has been manufactured, exposing it to multiple forms of attack. Among these, fault injection methods, including voltage and clock glitching, have become a significant threat [17]. These techniques can severely impact the operational reliability and security of ICs, leading to dire outcomes. In recent years, reinforcement learning (RL) has emerged as an effective strategy to address these hardware security concerns. Researchers have devised various RL formulations to detect hardware Trojans [18]–[20] and evaluate defenses [21], assess vulnerability of ciphers to fault injection attacks [22], and even to evaluate the robustness of machine learning based HT detection techniques against attacks [23]. However, these RL-based approaches require (i) extensive training processes, which are time- and compute-intensive, and (ii) a deep understanding of the target problems with steep learning curves.

C. Large Language Models

Large Language Models (LLMs) have revolutionized various fields of natural language processing and code-related tasks. In natural language processing, LLMs excel in areas including language translation [24], sentiment analysis [25], summarization, question answering, and more, owing to their ability to understand and generate human-like text. Moreover, LLMs can generate code snippets or even entire programs based on natural language descriptions, facilitating rapid prototyping and code design. In code completion tasks, they provide intelligent suggestions for completing code segments, improving developer productivity, and reducing errors [26].

Overall, LLMs represent a powerful tool for enhancing natural language processing and code-related tasks, offering versatility and efficiency across various applications.

D. Hardware Security in the Context of LLMs

In the realm of hardware security, LLMs are emerging as pivotal assets, shaping academic research and industry practices in hardware design and verification workflow [27]–[29]. LLM applications in hardware design span from code completion [28] to code correction [30]. Furthermore, LLMs are undergoing continuous fine-tuning to enhance their performance in comprehending and generating hardware description language (HDL) code, amplifying their role in aiding hardware design [28], [31], [32]. These capabilities of LLMs are evident through the shifting trends in semiconductor industries that are adopting LLMs in their electronic design automation (EDA) flow [33]. Beyond hardware design, LLMs are also starting to play a major role in hardware verification to detect functional and security vulnerabilities, including detecting hardware trojans and faults (see Sec. II).

However, as LLMs augment design and verification tasks, they are also a cause for concern, as they can potentially introduce new security vulnerabilities and attack surfaces into hardware. Existing hardware verification techniques can address some concerns, like vulnerabilities inserted by LLMs during code completion. However, the capabilities of LLMs can result in new attacks, necessitating novel security defenses to safeguard against emerging threats. Therefore, while leveraging LLMs' potential for hardware design and verification, it is imperative to concurrently address the evolving security landscape and potential attack vectors introduced by these powerful models (see Sec. III).

II. LLMs FOR HARDWARE VERIFICATION

Other than generating hardware code, LLMs are also aiding with the verification of hardware as we discuss below.

A. Analyzing Code to Aid Debugging

Given the ability of LLMs to analyze and generate text (Sec. I), LLMs can be utilized as a method for detecting code vulnerabilities. One strategy includes utilizing in-context learning, in which an LLM is provided a few ideal examples (input-output pairs) before being prompted with a test input (i.e. code to analyze). This strategy enabled ChatGPT-3.5 to perform competitively with the fine-tuned CodeBERT in determining if tested codes contain a vulnerability, with an accuracy of up to 62.7% [34]. LLMs have also demonstrated the ability to debug hardware (HDL). Through fine-tuning an LLM with a dataset of defective hardware designs, the ability of LLMs to effectively identify and correct defects within Verilog programs increased [35]. However, even advanced language models (GPT-4, PaLM2), are often unreliable when prompted to identify specific vulnerabilities within source codes, as the LLM may provide responses that are based upon incorrect reasoning, non-deterministic, or not robust [36].

LLMs can also aid the development process at the IC level, such as in writing architecture specifications. To minimize the

human error and time associated with this task, specific LLM prompting strategies can result in the successful generation of effective specifications, ranging from high- to low-level architectural descriptions [37]. Furthermore, RTL code can be directly utilized in this prompting framework, providing a confined scope for the specification. Additionally, the text analysis capabilities of LLMs enable the ability to generate an effective overview of a provided specification, even containing additional unique insights [37].

B. Generating Test Cases Using LLMs

Detecting vulnerabilities entails generating test inputs that trigger activity in various hardware components, thereby revealing the vulnerabilities present. While techniques like fuzzing use mutation strategies to explore the design space, they still use randomly generated inputs as the starting points [8]. LLMs, however, can understand a given hardware design and generate test inputs targeting various regions of hardware [38]. These LLM-generated test cases, coupled with mutations created by the fuzzer, can increase the overall coverage and accelerate vulnerability detection.

C. Generating Assertions Using LLMs

Both regression-based and formal verification techniques use assertions, i.e., functional and security properties, to detect vulnerabilities [39]. These properties define the expected behavior of the hardware designs, from which any deviation is flagged as a vulnerability. Assertions are effective in verifying vulnerabilities at the early stages of design, unlike reference model-based detection techniques, since assertions can be verified on individual design modules and do not require the output of the entire design. These properties are developed either manually by referring to specifications or mined from the existing hardware designs. However, manual creation of the assertions consumes time and requires design expertise, while mining the properties requires existing hardware designs [6], [7]. This limitation in the generation of assertions limits the adaption of assertions into hardware verification workflows.

Using LLMs to generate hardware assertions resolves this limitation, as LLMs enable faster and more automated generation of assertions to detect various hardware vulnerabilities [40]–[44]. However, the ability of LLMs to generate correct assertions is largely dependent on the input prompt [41]. LLMs can generate security assertions to detect vulnerabilities belonging to various common weakness enumeration (CWE) [45] types using the hardware design specification and description of the hardware CWEs [40].

D. Fixing Code Using LLMs

LLMs have been utilized to not only detect code vulnerabilities but also generate the corrected (vulnerability repaired) code. One such application is through a zero-shot approach, in which prompts are directly engineered to result in an ideal response. This can be implemented on vulnerable software and hardware programs through constructing a prompt that contains the original vulnerable code, some comments describing the vulnerability, and an instruction to generate a

fixed version [30]. Although this process resulted in a high success rate in tested synthetic and handcrafted buggy codes, the reliability of the tested LLMs (OpenAI's Codex and others) was determined insufficient to replace automatic program repairs in real-world scenarios [30].

A method to more effectively utilize LLMs for fixing code is fine-tuning. This process updates the weights of a pre-trained language model through training on a dataset related to an intended task, such as vulnerability repair. In particular, researchers used a training dataset of vulnerable C codes [46]. After fine-tuning various LLMs (Llama and Minstral) on this dataset, the test found a significant increase in the accuracy and adaptability of the fine-tuned LLMs' ability to repair vulnerabilities when compared to previous repair techniques [46].

E. Detecting Hardware Trojans Using LLMs

As mentioned in Sec. I, existing techniques for detecting HTs are either ineffective or require extensive training processes and an in-depth understanding of the problem. Off-the-shelf LLMs can aid in overcoming these limitations, as there is no need for additional training or extensive formulations. Researchers have performed a preliminary investigation into the capabilities of LLMs for detecting HTs inserted in the AES benchmark [27]. The investigation reveals that when no context about HT-insertion methods is provided, GPT-3.5 does not detect any HTs, but GPT-4 detects an impressive 81.01% of the HTs. On the other hand, when context about potential HT insertion methods is provided, the performances of GPT-3.5 and GPT-4 jump to 16.36% and 91.67%. GPT-4's performance demonstrates the tremendous power and advantage of LLMs in zero-shot detection of HTs as opposed to traditional methods, which require extensive formulations, coding, and training (for ML-based techniques). Finally, note that the research on LLM-based HT detection is still in its infancy, with many unanswered questions and scope for future research. For instance, [27] only tests on HT-inserted designs, so an understanding of the LLMs' false-positive rates is missing. Moreover, locating HTs in a complex design is another challenging task on which LLMs can be tested.

F. Analyzing Fault Injection Vulnerability Using LLMs

Researchers have also used LLMs in the context of fault injection attacks. In particular, [27] tested the ability of LLMs to compute the fault injection feasibility in finite-state machine transitions. They showed that, when using a sequential multi-step reasoning process, GPT-4 can correctly compute the fault injection feasibility.

III. FUTURE RESEARCH OPPORTUNITIES

LLMs have revolutionized natural language processing, and researchers have explored their potential for various hardware security applications. However, this also introduces security and privacy challenges.

Data Leakage Vulnerabilities. LLMs are trained on vast datasets sourced from the internet and carry inherent risks related to data leakage and misuse. For instance, LLMs can

inadvertently memorize and regurgitate sensitive information seen during training, potentially leading to privacy breaches. [47] details this and demonstrates the capability of attackers to extract specific pieces of information from a model, thereby highlighting the risk of sensitive data leakage.

Malicious Content Generation Attacks leverage LLMs to produce harmful or deceptive content, exploiting the models' ability to generate coherent and contextually relevant text. Sophisticated attackers can utilize LLMs to generate or refine malicious code and software exploits [48]. By providing a context or a set of requirements, LLMs can produce code snippets or entire programs designed to breach security systems, exploit vulnerabilities, or perform unauthorized actions on targeted systems.

IV. CONCLUSION

The advent and integration of LLMs into the semiconductor industry signify a revolutionary shift in the way hardware design and verification processes are approached. LLMs have not only streamlined these processes by automating complex tasks such as translating design specifications into source code and generating robust test cases, but have also enhanced the overall efficiency and reliability of hardware systems. However, as we embrace these advancements, the security implications introduced by the deployment of LLMs within this domain cannot be overlooked. The dual-edged nature of LLMs, capable of both augmenting the design process and potentially introducing vulnerabilities, calls for a balanced approach towards their adoption. It necessitates ongoing research to understand and mitigate the security risks associated with LLMs. As we move forward, the development of advanced security detection and mitigation techniques specifically tailored to counter LLM-induced vulnerabilities will be paramount. Furthermore, the hardware design and semiconductor industry must foster a culture of security-aware LLM usage, ensuring that the benefits of these models are harnessed without compromising the integrity and security of hardware systems. Ultimately, the journey towards fully realizing the potential of LLMs in hardware security is just beginning. With concerted efforts from researchers, industry professionals, and security experts, the future holds the promise of not only innovative but also secure hardware design methodologies, safeguarding the foundational elements of our increasingly digital world.

ACKNOWLEDGMENT

Our research work was partially funded by Intel's Scalable Assurance Program, the US Office of Naval Research (ONR Award #N00014-18-1-2058), and the Lockheed Martin Corporation. This work does not in any way constitute an Intel endorsement of a product or supplier. Any opinions, findings, conclusions, or recommendations expressed herein are those of the authors and do not necessarily reflect those of Intel, the US Government, or the Lockheed Martin Corporation.

REFERENCES

- [1] Y. Naveh, M. Rimon, I. Jaeger, Y. Katz, M. Vinov, E. s Marcu, and G. Shurek, "Constraint-Based Random Stimuli Generation for Hardware Verification," *AI magazine*, vol. 28, no. 3, pp. 13–13, 2007.