

iThermTroj: Exploiting Intermittent Thermal Trojans in Multi-Processor System-on-Chips

Mehdi Elahi¹, Mohamed R. Elshamy², Abdel-Hameed Badawy², Ahmad Patooghy¹

¹ Department of Computer Systems Technology, North Carolina A&T State University, NC, 27411

² New Mexico State University, United states

Abstract—Thermal Trojan attacks present a pressing concern for the security and reliability of System-on-Chips (SoCs), especially in mobile applications. The situation becomes more complicated when such attacks are more evasive and operate sporadically to stay hidden from detection mechanisms. In this paper, we introduce Intermittent Thermal Trojans (iThermTroj) that exploit the chips’ thermal information in a random time-triggered manner. According to our experiments, iThermTroj attack can easily bypass available threshold-based thermal Trojan detection solutions. We investigate SoC vulnerabilities to variations of iThermTroj through an in-depth analysis of Trojan activation and duration scenarios. We also propose a set of tiny Machine Learning classifiers for run-time anomaly detection to protect SoCs against such intermittent thermal Trojan attacks. Compared to existing methods, our approach improves the attack detection rate by 29.4%, 17.2%, and 14.3% in scenarios where iThermTroj manipulates up to 80%, 60%, and 40% of SoC’s thermal data, respectively. Additionally, our method increases the full protection resolution to 0.8 degrees Celsius, meaning that any temperature manipulations exceeding ± 0.8 degrees will be detected with 100% accuracy.

Index Terms—Intermittent Thermal Trojan Attacks, Machine Learning Anomaly Detection, System-on-Chips (SoCs)

I. INTRODUCTION

Demand for high-performance mobile systems equipped with multi-core processors is continually on the rise [1]–[3]. Mobile System-on-Chips (SoCs) have emerged as the backbone of these systems, catering to a myriad of tasks integral to the human’s daily life. However, despite the ever-increasing performance advancements, mobile SoCs face a significant constraint in the form of a limited thermal-power budget [4]. This limitation is primarily attributed to the challenges posed by maintaining compact form factors and limited cooling capabilities essential for handheld devices like smartphones [5], [6]. Ensuring user comfort by mitigating skin temperature to prevent discomfort or burning sensations adds another layer of complexity to the thermal management puzzle [7]. Hence, Dynamic Thermal Management (DTM) techniques have been widely utilized among mobile systems to effectively regulate and minimize high operating temperatures [8].

The accuracy and reliability of thermal sensors (that DTM methods rely on) are crucial for effective thermal management

strategies [9]. Thermal sensor malfunctions can occur due to a variety of factors, ranging from unintentional faults to deliberate tampering. Unintentional faults may include issues arising from fabrication defects, aging effects, susceptibility to noise, and process variations [10]–[12]. These factors can lead to inaccuracies in thermal measurements, which can, in turn, impact the overall performance and reliability of the chip. Additionally, there are significant security risks associated with deliberate tampering, such as the insertion of Hardware Trojans (HT), which are malicious modifications trying to alter the behavior of thermal sensors, leading to falsified temperature values [13]. Erroneous temperature readings can trigger unnecessary frequency reduction or core throttling by DTM, which may reduce the performance of the chip. Furthermore, incorrect thermal data can accelerate the aging process of the chip, thereby reducing its lifespan. The insertion of thermal HTs thus poses a significant threat to the security and reliability of thermal management systems.

To address these challenges, several strategies have been proposed in the literature. Authors in [14] introduce the Blind Identification Countermeasure (BIC), a technique derived from the Blind Power Identification (BPI) algorithm [15]. BIC aims to detect, contain, and isolate malicious sensors, thereby offering accurate temperature estimations and safeguarding chip integrity. Through extensive testing across multiple processor architectures, the efficacy and accuracy of BIC are demonstrated, marking a significant advancement in enhancing the reliability and performance of mobile SoCs amidst evolving security threats and thermal management challenges. The contributions of this paper are as follows.

- We introduce a new thermal Trojan that sporadically tampers with SoC’s thermal data. Through conducting extensive experiments, we demonstrate that the BIC method fails when confronted with the novel Trojan.
- We assess SoC vulnerability to thermal Trojans with thermal footprints ranging from 0.1 to 15 degrees Celsius. Existing methods fail to detect thermal attacks with smaller temperature modifications.
- We study the capability of tiny machine learning classifiers in detecting thermal Trojans by learning from the steady-state temperature readings.

The work is partially supported by the National Science Foundation under grants number 2219679 and 2219680.

II. THREAT MODEL AND IMPACTS

This work investigates adversaries capable of compromising DTM systems in mobile SoCs through the insertion or exploitation of Hardware Trojans (HTs). The adversaries can operate across multiple stages of the semiconductor supply chain, including malicious insiders in design or fabrication processes [16], external attackers leveraging untrusted third-party intellectual property (3PIP) cores, and those employing compromised electronic design automation (EDA) tools [17]–[19]. HTs are strategically embedded in thermally sensitive regions, such as near CPU/GPU cores, power management units, or within analog/digital interfaces of thermal sensors [14]. By under-reporting temperatures to disable performance throttling or over-reporting to trigger unnecessary power and clock adjustments, these HTs subvert DTM decision-making logic. Often, such HTs remain dormant until activated by specific thermal profiles (e.g., sustained high-temperature workloads) or external triggers, enabling stealthy, context-aware attacks. Beyond these foundational threats, adversaries may deploy increasingly sophisticated HTs with adaptive capabilities, posing unique challenges to detection and mitigation.

Building on these hardware-level threats, adversaries may deploy adaptive HTs that dynamically alter their activation patterns in response to runtime conditions or defensive measures, enhancing their ability to evade detection. These advanced HTs may leverage machine learning or analog circuitry to mimic benign sensor noise, bypassing static detection mechanisms such as threshold checks or Blind Identification Countermeasures (BIC) [14]. Particularly challenging are analog-based HTs embedded in sensor interfaces, which operate outside digital scan chains and exploit subtle signal distortions to manipulate temperature readings [20].

Such attacks can severely impact SoC performance, reliability, and security. Performance suffers as the DTM system makes incorrect decisions about frequency and core throttling, reducing speed and efficiency, while reliability is compromised by inadequate overheating protection, risking system failures and accelerated aging. Additionally, these manipulations create security vulnerabilities, enabling attackers to exploit the SoC for unauthorized access or operational disruption, threatening user data and device integrity. To counter these risks, implementing solutions like BIC is vital to detect and mitigate attacks, preserving the SoC’s functionality.

III. iTHERMTROJ: AN INTERMITTENT THERMAL TROJAN

According to the literature, the persistence of thermal anomalies injected by an attacker results in a continuous alteration of the temperature readings within the victim core of the target SoC. This sustained manipulation fundamentally challenges the efficacy of existing detection and mitigation mechanisms, including BIC. The consistent nature of thermal

Algorithm 1 iThermTroj Attack

```

1: Attack_Scenario = Lowering, Elevation, or Fluctuation
2: Attack_Rate = Pick from {100%, 80%, 60%, 40%}
3: IDX = Randomly Chosen Core Index
4: for steady_state_values do
5:   random_value = generate_random_number(0, 1)
6:   if random_value ≤ Attack_Rate then
7:     if T_error ≠ 0 then
8:       Switch(Attack_Scenario)
9:       {
10:        CASE = Lowering :
11:           $T(\text{IDX}) = T(\text{IDX}) - T_{\text{error}}$ 
12:        CASE = Elevation :
13:           $T(\text{IDX}) = T(\text{IDX}) + T_{\text{error}}$ 
14:        CASE = Fluctuation :
15:           $T(\text{IDX}) = T(\text{IDX}) \mp T_{\text{error}}$ 
16:       }
17:     end if
18:   end if
19: end for
```

attacks allows them to evade detection methods, making the SoC more susceptible to nuanced and stealthy threats. However, thermal attacks can become more sophisticated and stealthy through intermittent thermal alterations.

In this section, we introduce a novel and more evasive attack called Intermittent Thermal Trojans (iThermTroj) that works by taking into account SoC’s actual thermal traces. This can be done in multiple ways, for example, the attacker can feed the layout information to the HotSpot 6.0 thermal simulator to generate the corresponding thermal traces. The attacker then selects the victim core (can be selected randomly) and injects Δt_{error} , which is done by altering the HotSpot’s reported temperature. The Δt_{error} might be injected to the chosen core sporadically according to *Temperature Lowering*, *Temperature Elevation* or *Temperature Fluctuation* scenarios.

The key point of the proposed iThermTroj attack is that unlike traditional thermal Trojan attack scenarios, it does not involve permanent temperature manipulation. Instead, it follows an intermittent injection process in which some sensor readings are tampered with and the rest of the data points are left untouched. This will help iThermTroj to stay hidden and bypass the most recent detection methods presented in the literature [14]. In this regard, we have considered iThermTroj impacting different percentages of the SoC’s thermal traces i.e., scenarios at which iThermTroj is active for 80%, 60%, and 40% of the victim core. The detailed structure of the proposed attack is provided in Algorithm 1. In the first two lines of the algorithm, the attacker selects the attacking scenario as well as the attack rate. The attack rate determines the percentage of the

thermal data which will be manipulated by iThermTroj. Then, in line 3, the victim core is selected and a random value is generated to be used to manipulate the actual thermal reading according to the selected scenario (lines 10-15).

To highlight this vulnerability, we have performed a detailed analysis comparing the failure rates of BIC under a conventional persistent attack scenario with those observed during our proposed iThermTroj attack. Figure 1.a and Figure 1.b compare the efficacy of BIC once encountered persistent and intermittent thermal Trojans respectively. Each plot reports the number of times BIC failed to detect a thermal attack (the lower the better) as a function of various Δt_{error} values. The experiment illustrates that under typical persistent attack conditions, BIC demonstrates robust performance, maintaining system integrity for Δt_{error} values above 3 and below -2 . However, this effectiveness significantly diminishes when the SoC is subjected to the iThermTroj attack, as illustrated in Figure 1.b. This stark contrast highlights a critical vulnerability, emphasizing the urgent necessity for developing a more resilient countermeasure capable of defending against sophisticated threats like the iThermTroj attack. Enhanced security measures must be prioritized to ensure the SoC's protection and maintain system reliability in the face of evolving attack methodologies.

IV. PROPOSED MACHINE LEARNING DETECTION

This research addresses limitations in traditional thermal monitoring systems which struggle to detect sophisticated thermal attacks like iThermTroj. TinyML classifiers emerge as a promising solution due to their ability to operate efficiently within mobile devices' resource constraints while processing thermal data in realtime. These classifiers can be deployed directly on SoCs to minimize detection latency and enable continuous monitoring without network dependency. The study evaluates five different classifiers-SVM, Logistic Regression, Random Forest, Decision Tree, and two variants of Naive Bayes-for anomaly detection in thermal data. The methodology involved collecting steady-state temperature readings,

deliberately introducing adversarial attacks on 80% of the data to simulate malicious activity, and using this combined dataset to train and test the models. This approach aims to assess each classifier's effectiveness in identifying thermal anomalies that could indicate attacks, ultimately enhancing thermal management systems while extending device longevity by preventing excessive thermal stress.

V. EVALUATION RESULTS & DISCUSSIONS

We utilize a heterogeneous 6-core mobile processor layout as detailed by [21] for our evaluations. For each scenario, we undertake a series of methodical steps. Initially, we employ the HotSpot 6.0 Thermal Simulator [22] to process the layout information along with the power traces, thereby generating the corresponding thermal traces necessary for our analysis. Subsequently, we act as an attacker by selecting one of the six cores at random to apply a thermal error, denoted as Δt_{error} , to the targeted core (intermittently and sporadically). This attack is repeated across different portions of the thermal traces, i.e., 80%, 60% and 40% of the generated thermal traces to assess its impact under varying conditions. We used this data to train used ML classifiers with the distribution of 70% and 30% for training and inference.

Figure 2 shows the performance of the used ML classifiers in terms of accuracy, recall, precision and F1-score under an iThermTroj attack affecting 80% of thermal data. Results consistently exhibit the high performance of these classifiers across all evaluated metrics, indicating their robustness and reliability in detecting and counteracting the iThermTroj attack. To extend our analysis we repeated the simulations for 60% and 40% thermal data corruption caused by the iThermTroj attack. To account for variations, average precision data across all Δt_{error} and injection rates were calculated and visualized in Figure 3. This figure offers a two-fold interpretation of the classifiers' performance under attack:

Accuracy Analysis: The charts indicate that detection failure rates (i.e., the cases at which ML classifiers failed in detecting iThermTroj attack) increase as the thermal injection becomes smaller, with the BIC method demonstrating greater vulnerability compared to our proposed iThermTroj at various injection rates. Specifically, BIC shows a detection failure rate of 52.38% under certain conditions. However, when applying an 80% Trojan injection by iThermTroj attack, the proposed ML detection countermeasure significantly reduces the detection failure rate to approximately 23%, representing a reduction of around 29% compared to the highest failure rate observed with the BIC method. Additionally, as attacks become more sophisticated and evasive, the detection failure rate correspondingly rises. For instance, with fault injection rates of 60% and 40%, the proposed ML detection countermeasure fails to detect 35.16% and 38.10% of attack cases,

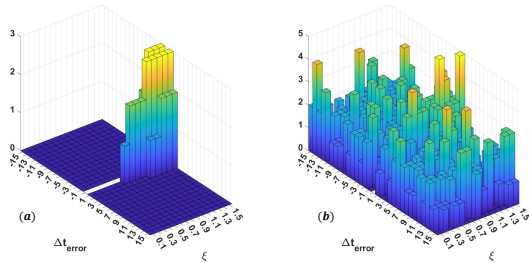


Fig. 1: The number of detection failures relative to Δt_{error} for a heterogeneous core layout a) with normal attack b) with iThermTroj attack

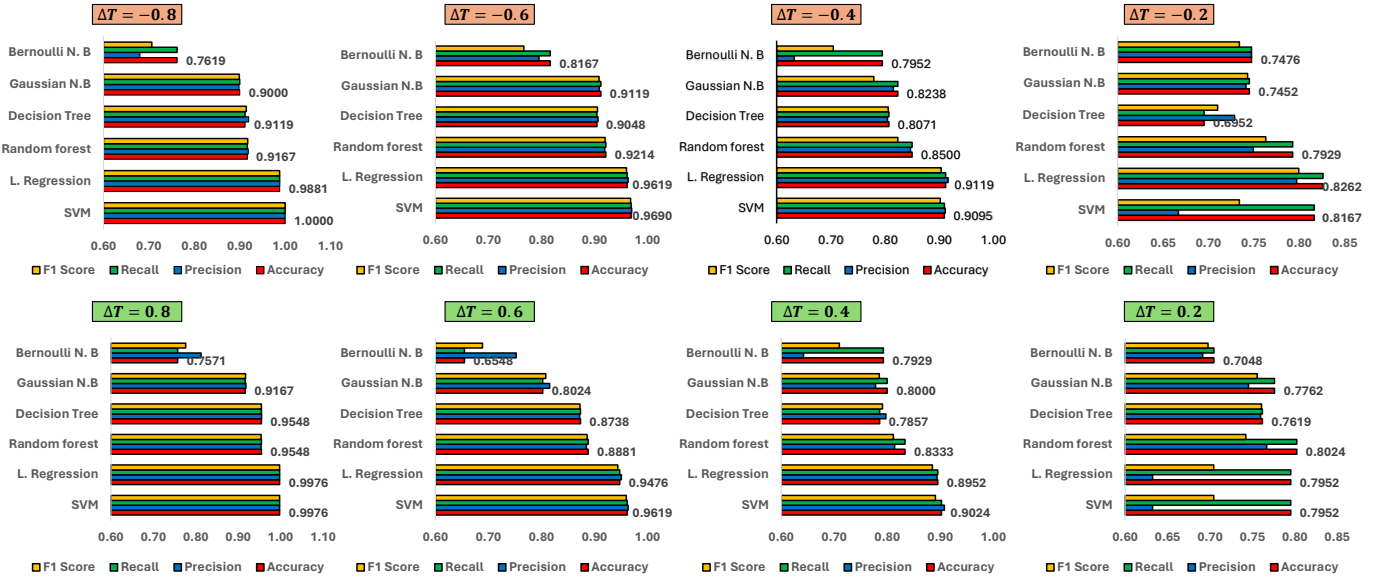


Fig. 2: Output parameters of six machine learning classifiers applied to 80% of the thermal data, tested across Δt_{error} intervals from -0.8 to 0.8 with a step size of 0.2.

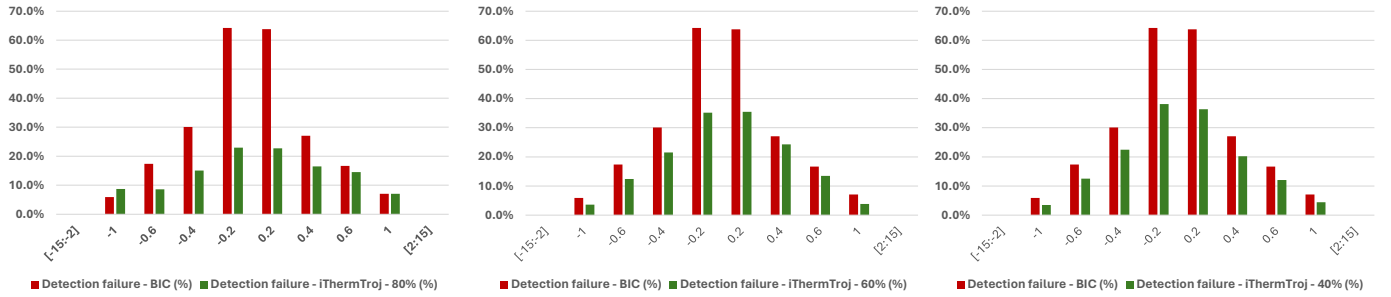


Fig. 3: Detection Failure Rates of the ML Countermeasure Across Different Trojan Injection Rates and Δt_{error} Thresholds

respectively. Moreover, the distribution form of the chart data appears to follow a normal distribution, which suggests that the variations in detection failure rates are systematically related to the evasiveness of the attacks.

Resolution Analysis: The reported results indicate that iThermTroj's ML countermeasure offers a compelling solution, demonstrating not only high accuracy but also superior resolution in terms of Δt_{error} . Specifically, the BIC method can fully secure the system for Δt_{error} values below -2 and above 3. In contrast, even in the worst-case scenario for iThermTroj's ML countermeasure, with a 40% fault injection rate, we observe complete detection coverage for Δt_{error} values below -2 and above 2. This resolution is further enhanced at 60% and 80% injection rates, where we achieve complete detection for Δt_{error} values below -1.2 and -1, and above 1.2 and 1, respectively. This demonstrates the superior detection capabilities of iThermTroj's ML countermeasure, particularly at higher injection rates, ensuring robust system security across a broader range of Δt_{error} values.

VI. CONCLUSIONS

This study presents a significant advancement in the field of thermal management and security for mobile System-on-Chips (SoCs) through the introduction of Intermittent Thermal Trojans (iThermTroj) and the application of machine learning-based anomaly detection techniques. By identifying the vulnerabilities posed by intermittent thermal Trojan attacks and proposing the integration of tiny Machine Learning classifiers, the research offers a novel approach to enhance the resilience and effectiveness of SoC security measures. The findings demonstrate the potential of machine learning models to adapt to complex thermal behaviors and detect anomalies in real-time, thereby improving the responsiveness and longevity of mobile devices. This research contributes valuable insights into mitigating the risks associated with thermal attacks on SoCs and sets a foundation for further advancements in securing modern processors against evolving security threats.

REFERENCES

- [1] J. S. Jeong, J. Lee, D. Kim, C. Jeon, C. Jeong, Y. Lee, and B.-G. Chun, "Band: coordinated multi-dnn inference on heterogeneous mobile processors," in *Proceedings of the 20th Annual International Conference on Mobile Systems, Applications and Services*, 2022, pp. 235–247.
- [2] A. Patooghy, M. Elahi, M. F. Torkaman, S. S. Dokhtfaroughi, and R. Rajaei, "Addressing benign and malicious crosstalk in modern system-on-chips," *IEEE Access*, vol. 11, pp. 142 263–142 275, 2023.
- [3] M. R. Elshamy, M. Elahi, A. Patooghy, and A.-H. A. Badawy, "Fine-grained clustering-based power identification for multicores," in *2024 IEEE 15th International Green and Sustainable Computing Conference (IGSC)*, 2024, pp. 165–170.
- [4] C. Wang, Q. Xu, C. Nie, H. Cao, J. Liu, and Z. Li, "An efficient thermal model of chiplet heterogeneous integration system for steady-state temperature prediction," *Microelectronics Reliability*, vol. 146, p. 115006, 2023.
- [5] S. Pagani, P. D. S. Manoj, A. Jantsch, and J. Henkel, "Machine learning for power, energy, and thermal management on multicore processors: A survey," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 39, no. 1, pp. 101–116, 2020.
- [6] A. Noori, P. Devkota, S. D. Mohanty, and P. Manda, "Llms in action: Robust metrics for evaluating automated ontology annotation systems," *Information*, vol. 16, no. 3, 2025. [Online]. Available: <https://www.mdpi.com/2078-2489/16/3/225>
- [7] T. Gong, L. Li, M. Shi, L. Kang, L. Gao, and J. Li, "Performance assessment and optimization of a thin-film thermoelectric cooler for on-chip transient thermal management," *Applied Thermal Engineering*, vol. 224, p. 120079, 2023.
- [8] J. Yang, X. Zhou, M. Chrobak, Y. Zhang, and L. Jin, "Dynamic thermal management through task scheduling," in *ISPASS 2008 - IEEE International Symposium on Performance Analysis of Systems and software*, 2008, pp. 191–201.
- [9] B. Ding, Z.-H. Zhang, L. Gong, M.-H. Xu, and Z.-Q. Huang, "A novel thermal management scheme for 3d-ic chips with multi-cores and high power density," *Applied thermal engineering*, vol. 168, p. 114832, 2020.
- [10] A. Oukaira, D. E. Touati, A. Hassan, M. Ali, Y. Savaria, and A. Lakhssassi, "Fem-based thermal profile prediction for thermal management of system-on-chips," *2022 8th International Conference on Optimization and Applications (ICOA)*, pp. 1–4, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:253423797>
- [11] R. Eitan and A. Cohen, "Untrimmed low-power thermal sensor for soc in 22 nm digital fabrication technology," *Journal of Low Power Electronics and Applications*, vol. 4, no. 4, pp. 304–316, 2014. [Online]. Available: <https://www.mdpi.com/2079-9268/4/4/304>
- [12] A. Patooghy, M. F. Torkaman, and M. Elahi, "Your hardware is all wired up! attacking network-on-chips via crosstalk channel," in *Proceedings of the 12th International Workshop on Network on Chip Architectures*, 2019, pp. 1–6.
- [13] K. Hasegawa, K. Yamashita, S. Hidano, K. Fukushima, K. Hashimoto, and N. Togawa, "Node-wise hardware trojan detection based on graph learning," *IEEE Transactions on Computers*, 2023.
- [14] M. Abdelrehim, A. Patooghy, A. Malekmohammadi, and A.-H. A. Badawy, "Bic: Blind identification countermeasure for malicious thermal sensor attacks in mobile socs," in *2022 23rd International Symposium on Quality Electronic Design (ISQED)*. IEEE, 2022, pp. 1–6.
- [15] M. Said, S. Chetoui, A. Belouchrani, and S. Reda, "Understanding the sources of power consumption in mobile socs," in *2018 Ninth International Green and Sustainable Computing Conference (IGSC)*. IEEE, 2018, pp. 1–7.
- [16] G. Bhat, S. Gumussoy, and U. Y. Ogras, "Power-temperature stability and safety analysis for multiprocessor systems," *ACM Transactions on Embedded Computing Systems (TECS)*, vol. 16, no. 5s, pp. 1–19, 2017.
- [17] S. M. Sebt, A. Patooghy, H. Beitollahi, and M. Kinsy, "Circuit enclaves susceptible to hardware trojans insertion at gate-level designs," *IET Computers & Digital Techniques*, vol. 12, no. 6, pp. 251–257, 2018.
- [18] A. Mohammadi, R. Ahmari, V. Hemmati, F. Owusu-Ambrose, M. N. Mahmoud, P. Kebria, and A. Homaifar, "Detection of multiple small biased gps spoofing attacks on autonomous vehicles using time series analysis," *IEEE Open Journal of Vehicular Technology*, pp. 1–13, 2025.
- [19] R. Ahmari, V. Hemmati, A. Mohammadi, M. Mynuddin, P. Kebria, M. Mahmoud, and A. Homaifar, "Evaluating trojan attack vulnerabilities in autonomous landing systems for urban air mobility," *Proceedings of the Automation, Robotics & Communications for Industry*, vol. 4, no. 5.0, p. 80, 2025.
- [20] M. Elahi, M. R. Elshamy, A.-H. Badawy, M. Fazeli, and A. Patooghy, "Matter: Multi-stage adaptive thermal trojan for efficiency & resilience degradation," 2024. [Online]. Available: <https://arxiv.org/abs/2412.00226>
- [21] Y.-H. Gong, J. J. Yoo, and S. W. Chung, "Thermal modeling and validation of a real-world mobile ap," *IEEE Design & Test*, vol. 35, no. 1, pp. 55–62, 2017.
- [22] W. Huang, S. Ghosh, S. Velusamy, K. Sankaranarayanan, K. Skadron, and M. R. Stan, "Hotspot: A compact thermal modeling methodology for early-stage vlsi design," *IEEE Transactions on very large scale integration (VLSI) systems*, vol. 14, no. 5, pp. 501–513, 2006.