



**QUEEN'S
UNIVERSITY
BELFAST**

Ten years of hardware Trojans: a survey from the attacker's perspective

Xue, M., Gu, C., Liu, W., Yu, S., & O'Neill, M. (2020). Ten years of hardware Trojans: a survey from the attacker's perspective. *IET Computers And Digital Techniques*. Advance online publication. <https://doi.org/10.1049/iet-cdt.2020.0041>

Published in:
IET Computers And Digital Techniques

Document Version:
Peer reviewed version

Queen's University Belfast - Research Portal:
[Link to publication record in Queen's University Belfast Research Portal](#)

Publisher rights
Copyright 2020 Institution of Engineering and Technology. This work is made available online in accordance with the publisher's policies. Please refer to any applicable terms of use of the publisher.

General rights
Copyright for the publications made accessible via the Queen's University Belfast Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy
The Research Portal is Queen's institutional repository that provides access to Queen's research output. Every effort has been made to ensure that content in the Research Portal does not infringe any person's rights, or applicable UK laws. If you discover content in the Research Portal that you believe breaches copyright or violates any law, please contact openaccess@qub.ac.uk.

Open Access
This research has been made openly available by Queen's academics and its Open Research team. We would love to hear how access to this research benefits you. – Share your feedback with us: <http://go.qub.ac.uk/oa-feedback>

Ten years of hardware Trojans: A survey from the attacker's perspective

Mingfu Xue, *Member, IEEE*, Chongyan Gu, *Member, IEEE*, Weiqiang Liu, *Senior Member, IEEE*, Shichao Yu, and Máire O'Neill, *Senior Member, IEEE*

Abstract—In the last decade, hardware Trojan has emerged as a serious concern in integrated circuit (IC) industry. As such, hardware Trojan detection techniques have been studied extensively. However, in order to develop reliable and effective defenses, it is important to figure out how hardware Trojans are implemented in practical scenarios. In this paper, we attempt to make a review of the hardware Trojan design and implementations in the last decade and also provide an outlook. Unlike all previous surveys that discuss Trojans from the defender's perspective, for the first time, we study the Trojans from the attacker's perspective, focusing on the attacker's methods, capabilities and challenges when he designs and implements a hardware Trojan. Particularly, the following questions are explored. What are the current methods and capabilities of attackers after ten years of research and development? By considering more and more sophisticated hardware Trojan detection techniques, what challenges do the attackers face, and vice versa? First, we present adversarial models in terms of the adversary's methods, adversary's capabilities and adversary's challenges in seven practical hardware Trojan implementation scenarios: in-house design team attacks, third-party intellectual property (3PIP) vendor attacks, computer-aided design (CAD) tools attacks, fabrication stage attacks, testing stage attacks, distribution stage attacks, and field programmable gate array (FPGA) Trojan attacks. Second, we analyze the hardware Trojan implementation methods under each adversarial model in terms of seven aspects/metrics: hardware Trojan attack scenarios, the attacker's motivation, feasibility (the practicality of the attacks), detectability (anti-detection capability of that kind of Trojan), protection and prevention suggestions for the designer, overhead analysis, and case studies of Trojan implementations. Finally, future directions on hardware Trojan attacks and defenses are discussed. This paper also presents several new insights and assumptions for the first time, including considering the Trojans not only from the copyright owner's perspective, but also from the users' perspective, and discussing the hardware Trojan attacks in the testing phase and in the distribution phase. This paper can hopefully provide a reference for future works on building effective hardware Trojan defenses.

I. INTRODUCTION

In the last decade, the malicious modifications of integrated circuits (ICs), also referred to as hardware Trojans (HTs), have become emerging security concerns in the IC industry. ICs that contain HTs can cause malfunction, leakage of confidential

information, or lead to other disastrous consequences. Therefore, HT has been a matter of concern for industry, academia, government and military (1; 2; 3; 4).

Since the first research on HT published in 2007 by Agrawal *et al.* (5), HTs have been developed for more than ten years. A lot of research has been conducted on detecting HTs. However, there have been very little research on the implementation of HTs in practice. In order to develop reliable HT detection and defense techniques, it is necessary to understand the feasibility of inserting HTs in practical implementations (6). More specifically, how stealthy HTs can be inserted into a target circuit, how feasible is HT for a specific application model, and what are the challenges to implement such HTs (7). This remains a field that has received relatively little attention in the research community where most HTs referred to in the literature are small or medium-sized circuits added at register transfer level (RTL) during the IC design flow (7).

In this paper, we attempt to make a review of the HT implementations in the last decade and also make an outlook. In particular, unlike all previous surveys that discuss Trojans from the defender's perspective, for the first time, we will discuss the Trojans from the attacker's perspective, focusing on the attacker's methods, feasibility, anti-detection capability, and challenges when designing and implementing a HT. As Chinese strategist Sun Tzu said, "If you know yourself and the enemy, you will never lose a battle", (Sun Tzu, *The Art of War*, ancient Chinese military philosophy book). Discussing Trojans from an attacker's perspective can give readers a clear understanding of an attacker's considerations when implementing a HT, including advantages and deficiencies of the attacker, trade-offs, and the methods that the attacker can take. This can hopefully help designers better understand the Trojan insertion, and provide guidelines for the defenders to design better detection and defense techniques.

Particularly, in this survey, we want to explore the following questions:

- Q1. HT techniques have been developed for more than ten years. What attack methods do the attackers at various stages have?
- Q2. What are the attackers' capabilities at various abstraction levels to launch the attacks?
- Q3. With the continuous development of Trojan detection techniques, what are the challenges and difficulties faced by the attackers?
- Q4. How attackers could overcome these state-of-the-art Trojan detection techniques and what kind of new detection technique is required?

M. Xue is with College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, People's Republic of China e-mail: (mingfu.xue@nuaa.edu.cn).

W. Liu is with College of Electronic and Information Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing, People's Republic of China

C. Gu, S. Yu and M. O'Neill are with Center for Secure Information Technologies, Queen's University Belfast, Belfast, United Kingdom.

In addition, this paper will present several new insights and assumptions for the first time. For example, in the literature, it is usually assumed that the HT problem was discussed from the perspective of the copyright owner or the designer, *i.e.*, the copyright owner is assumed to be trustworthy and the HT was inserted in the untrustworthy design and fabrication processes. However, when returning to the essential definition of HTs, HT is a Trojan/backdoor in the form of hardware, which is not limited to the copyright owner's perspective. Moreover, the number of users is much larger than the number of the copyright owner, and the users are usually in weaker positions compared to the copyright owners. Therefore, it is also necessary to consider the Trojans from the user's perspective. For the users, the hardware Trojan/backdoor implanted by the copyright owners or designers in the device can also be regarded as a HT, such as some super privileged structures which can be remotely controlled by the copyright owner, the hardware that transmits user's data back to the copyright owner, or the hardware that can record the user's keystrokes, and so on.

There have already been some reported accidents of potential HT attacks to gain control of devices, steal secret information, or even destroy a system. In Sept. 2007, Israel launched a successful airstrike on a nuclear reactor in Syria, while Syria's advanced air defense system did not respond throughout the operation (8). In 2008, Adee (8) speculated that the Syria's air defense system had been deactivated by a built-in kill switch, which could be accessed and activated remotely. Since the practical HTs used in industrial fields and military are often highly confidential, researchers cannot accurately determine whether they are HTs and their implementation details. However, it still reveals the concerns of various communities about the destructive power of HTs. In 2016, Yang *et al.* (9) proposed a small malicious HT, named A2, where they implemented a privilege escalation attack in the OR1200 processor by running a set of seemingly harmless commands. Such lightweight analog malicious backdoors are extremely difficult to detect. In Jan. 2018, the Free Software Foundation revealed that Intel computers have a built-in subsystem, called the Intel Management Engine (ME), which can take full control over the computer, and even has access to the main memory (10; 11). The ME structure can be a serious threat to the users' privacy and security. However, users do not have the ability to audit, examine or disable it (10; 11). From the user's point of view, this could also be considered as a HT.

To date, several review and survey papers on HT detection techniques or HT taxonomy have been published. Bhunia *et al.* (1) analyzed the threats of HT attacks, Trojan models and classifications, and protection approaches. They mainly focused on various defense techniques against HTs, including HT detection techniques, design-for-security (DFS) approaches, and runtime monitoring techniques. Tehranipoor and Koushanfar (2) presented a classification of HTs and a survey of Trojan detection techniques. In particular, they presented existing detection mechanisms and DFS methodologies. Chakraborty *et al.* (3) presented a Trojan taxonomy and a review of state-of-the-art HT detection techniques. Rostami *et al.* (12) systematized various hardware security related attacks,

including HTs, reverse engineering (RE), IC overbuilding and intellectual property (IP) piracy, side-channel analysis (SCA), *etc.* Jacob *et al.* (13) reviewed HT vulnerabilities in the IC's life cycle and HT detection techniques. Karri *et al.* (14) proposed a Trojan taxonomy in terms of activation mechanism, effects, abstraction level, insertion phase, and location. The above surveys all focus on HT detection or HT taxonomy, and are published before 2014, while a large number of HT works that have emerged in the past six years are not included.

Different from all existing surveys, this paper presents a survey of HT design and implementation based on practical attack scenarios from an attacker's perspective. The differences between this survey and these existing review/survey papers are summarized as follows:

- (1) **The HT design and implementation methods are systematically studied and analyzed, focusing on the attacker's insertion methods, capabilities, evading detection techniques, and challenges when the attacker designs and implements a HT.** To the best of the authors' knowledge, this is the first survey of HT design and implementation methods from an attacker's perspective, instead of HT detection techniques from the defender's perspective.
- (2) **We present adversarial models that show adversary's methods, adversary's capabilities and adversary's challenges to insert HT into a chip in seven practical HT implementation scenarios:** in-house design team attacks, third-party intellectual property (3PIP) vendor attacks, computer-aided design (CAD) tools attacks, fabrication stage attacks, testing stage attacks, distribution stage attacks, and field programmable gate array (FPGA) Trojan attacks. Note that, the contribution of this paper is not to provide a new HT taxonomy. Trojan taxonomies have been widely mentioned in existing review literatures. The goal of this paper is to analyze the attacker's considerations during Trojan insertion in various practical scenarios, including the technical options, advantages and disadvantages, trade-offs, anti-detection capabilities, and so on.
- (3) **HT design and implementation methods under each adversarial model are reviewed in terms of seven aspects/metrics: HT attack scenario, motivation, feasibility, detectability (anti-detection capability), protection and prevention suggestions, overhead, and case studies.** The *feasibility* and *detectability* are two main concerns from the attackers' perspective. The *protection and prevention suggestions* and *overhead* are two metrics from the defenders' perspective. Note that, the existing HT literatures, including survey works, mostly focus on HT detection and defenses. Therefore, in this paper, we do not discuss the HT detection and defense techniques in details, but only give brief suggestions for the Trojan detection. Instead, we will discuss the feasibility and the anti-detection capability in details when inserting Trojans from the attacker's perspective.
- (4) **Future potential directions on HT designs and defenses are also discussed,** including HT benchmarks

and evaluation methods, machine learning-based Trojan detection methods and HTs targeting machine learning models, attacks and defenses from chips to complex systems, universal Trojan and automatic Trojan insertion VS automatic Trojan (IC vulnerability) analysis tools, multi-stage HT attacks and defenses, split manufacturing, low overhead runtime HT monitoring techniques, logic obfuscation for HT prevention, FPGA Trojan attacks and defenses.

- (5) **This paper presents several new insights and assumptions for the first time.** On the one hand, researchers should not only consider the HTs implanted during the untrustworthy design and fabrication stages from the copyright owner’s perspective, but should also consider the Trojans inserted by the copyright owner from the user’s perspective. On the other hand, existing works usually only consider Trojans to be inserted in the design stage, CAD tools, and fabrication stage. In this paper, for the first time, we also systematically discuss the HT attacks in the testing stage and in the distribution stage. Moreover, the emerging FPGA Trojan attacks are also systematically discussed.

The rest of this paper is organized as follows. The attack models are described in Section II. In-house design team attacks are analyzed in Section III. 3PIP vendor attacks are presented in Section IV. CAD tools attacks are described in Section V. Fabrication stage attacks are presented in Section VI. Testing stage attacks are described in Section VII. Distribution stage attacks are analyzed in Section VIII. FPGA Trojans are presented in Section IX. Future directions are discussed in Section X. Finally, conclusions are provided in Section XI.

II. ATTACK MODELS: ADVERSARY’S METHODS, CAPABILITIES, AND CHALLENGES

In this section, we will present the attack models in terms of the adversary’s methods, adversary’s capabilities, and adversary’s challenges in seven practical HT implementation scenarios.

A malicious attacker in any stage of the IC supply chain can insert HTs. The most common concern is that HTs can be inserted during fabrication by untrusted foundries. A malicious designer in the IC design team could also manipulate the design and have flexibility to implement various HTs. Similarly, 3PIP core is another possible source of HTs (15). Other entities, *e.g.* CAD tool vendors, IC vendors, and users, although have less chance to insert a HT, but are still feasible to implement HT attacks. HT design and implementation methods are diverse, *e.g.* an attacker can design a HT based on the desired attack function, triggering mechanism, insertion stage, *etc.*

As the HT design and implementation methods significantly depend on practical application scenarios and the attackers’ intentions, in this paper, we present adversarial models in terms of adversary’s methods, adversary’s capabilities, and adversary’s challenges in different HT implementation scenarios, as shown in Figure 1. Related works usually consider the testing phase to be trusted, while in this paper we consider that

the testing phase may also be untrustworthy. Strictly speaking, in the testing phase, untrusted testing organizations are not able to insert Trojans, but can only collude with the malicious factory or designers to hide the inserted HTs, *i.e.*, making the HTs evade detection. Similarly, it is generally considered that the distributor cannot insert a HT because the distributor is usually unaware of the design. However, a distribution stage attacker can RE a chip so as to pirate the chip. The attacker can also directly replace the circuit with a Trojan-inserted circuit during transportation. Therefore, in this paper, we also discuss the attackers from the testing stage and the distribution stage. Specifically, we divide the practical HT implementation scenarios into seven different adversarial models: 1) in-house design team attacks; 2) 3PIP vendor attacks; 3) CAD tools attacks; 4) fabrication stage attacks; 5) testing stage attacks; 6) distribution stage attacks; 7) FPGA Trojan attacks. Moreover, we analyze the HT implementation methods under each adversarial model in terms of the following seven aspects/metrics:

- (1) **HT attack scenario:** a description of the HT attack scenario, including the HT types, trigger mechanisms, payloads, *etc.*
- (2) **Motivation:** the motivation of an attacker, including the malicious functions that an attacker wants to achieve.
- (3) **Feasibility:** the practicality of the attacks, including the resources available for an attacker, the HT design methods that an attacker can adopt. Similar to cryptography and cryptanalysis, the attacker is assumed to have significant resources, but they are restricted by the rule that the benefit from the Trojan attack should exceed the resources expended (16). The HTs should also be practical and effective under practical scenarios and be easy to control so that an attacker can employ them to perform attacks easily.
- (4) **Detectability:** anti-detection capability of the Trojan, *i.e.*, how to evade the state-of-the-art defenses from the attacker’s perspective. In other words, the detection mechanisms available for the described HT and how likely the HT will be detected.
- (5) **Protection and prevention suggestions:** guidelines for designers about protection and prevention, including challenges and opportunities from the designer’s perspective, suggestions that would help designers to protect the circuits better against Trojan insertions, and how the attack models will affect the future secure hardware design. As mentioned in Section I, since most of the existing works have discussed HT detection techniques, in this paper, we do not discuss the HT detection techniques in details, but only give brief suggestions for the Trojan detection (referred to as *protection and prevention suggestions*). Instead, we will discuss in details the anti-detection capability of a HT and the attacker’s considerations of evading detection from the attacker’s perspective (referred to as *detectability*).
- (6) **Overhead:** cost for Trojan detection from the defender’s perspective, in terms of power, area, and performance.
- (7) **Case studies:** examples of HT design and implementations.

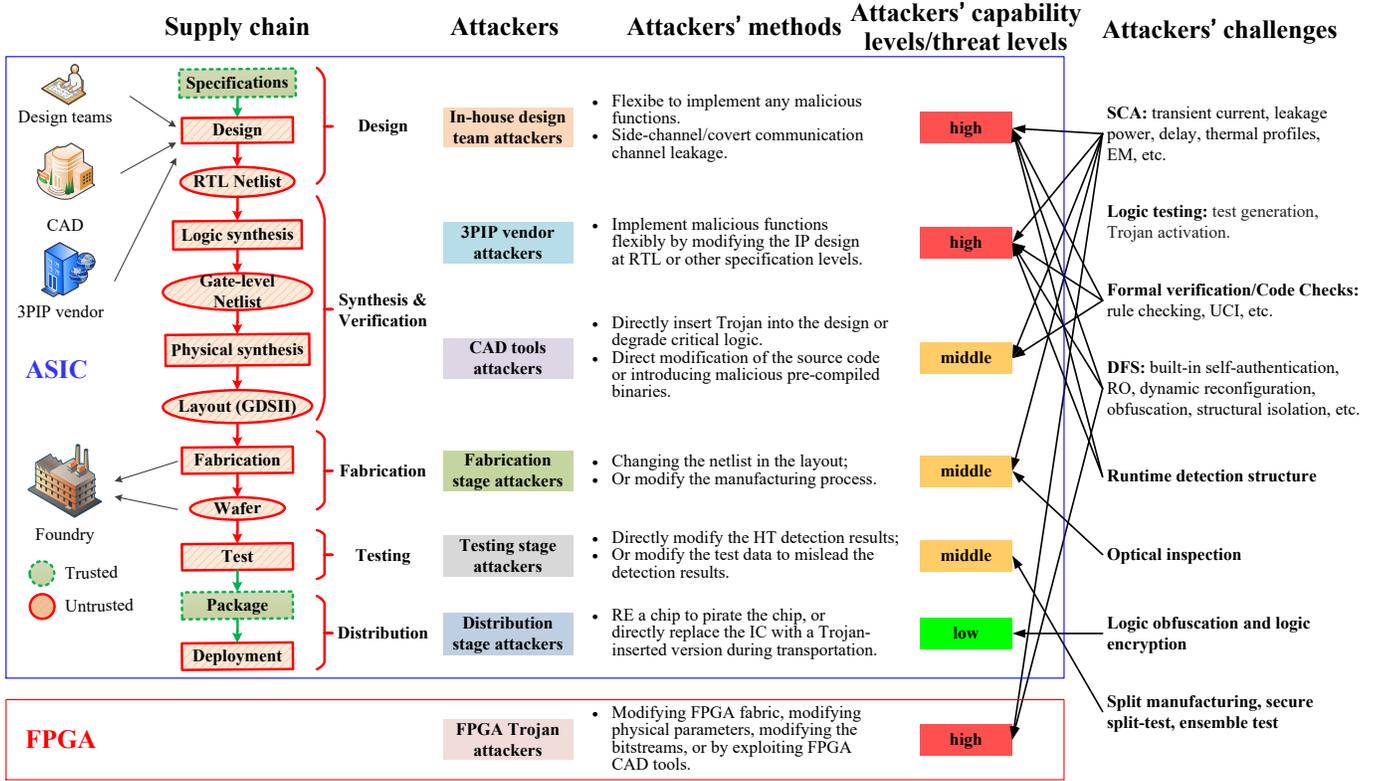


Fig. 1. HT attack models in terms of the adversary's methods, adversary's capabilities, and adversary's challenges.

In the following sections, HT design and implementation methods are reviewed and analyzed under the above seven practical HT implementation scenarios, respectively, in terms of the above seven aspects. Particularly, Figure 2 presents the HT attack scenario, motivation, available resources, feasibility, detectability (anti-detection capability), case studies from the attacker's perspective and the protection and prevention suggestions, overhead from the defender's perspective under different attack models, which will be discussed in the following sections.

III. IN-HOUSE DESIGN TEAM ATTACKS

HT attack scenario: This attack model is the one that most commonly referred to in the literature. Rogue designers in an outsourced or in-house design team can easily implement stealthy malicious modifications in the RTL design since the attackers can get the source files and codes, as shown in Figure 2. Trojans inserted by the malicious designer can implement any possible payloads with various trigger methods.

Motivation: The attacker in the design stage who insert a HT into the IC may want to steal confidential information from the deployed ICs, or cause malfunction of the ICs.

Feasibility: The attackers can manipulate the circuit with high flexibility to implement any malicious functions. The trigger is expected to be undetectable by functional tests. A feasible approach is to use a specifically designed input sequence, e.g. an abnormal condition, or a rare event. However, a trigger that relies on physical access may be restricted in practical applications. Therefore, some internal signals, e.g. a

counter, a specified temperature or voltage, can be used as an activation mechanism for the Trojan, such as the RS232-T200 HT (17). Another type of trigger, which is more aggressive, configures the Trojan as always-on. In this case, the payload of the Trojan is required to be hidden, e.g. sending secret information undetected by functional tests, such as the Advanced Encryption Standard (AES)-T200 HT (17). However, the always-on Trojans may introduce high power consumption, which could be easily to be discovered by SCA methods (18). The reasons are as follows. In general, HT is triggered by rare events so as to evade the detection of defense techniques. As a result, the HT is latent during most of the time. It does not affect the logic values of the circuit nodes, and rarely generates transition activities, so the power consumption introduced by the Trojan is very low. However, for always-on Trojan, on the one hand, the payload usually does not directly affect the digital value of the circuit node, so as not to be easily detected by logic tests. On the other hand, it has no triggering conditions and is always on, so its circuit transition activities will be relatively high. Thus it will be easily detected by SCA method based on dynamic current or power consumption. Note that, as a special case, parametric Trojans can also be considered as always-on, which does not necessarily lead to high power consumption. In conclusion, the design stage attack has a high feasibility, good practicality, and is easy to implement.

Detectability: First, we discuss the available detection methods from the designer's point of view. The insertion of HTs in the RTL can be revealed by formal code checks,

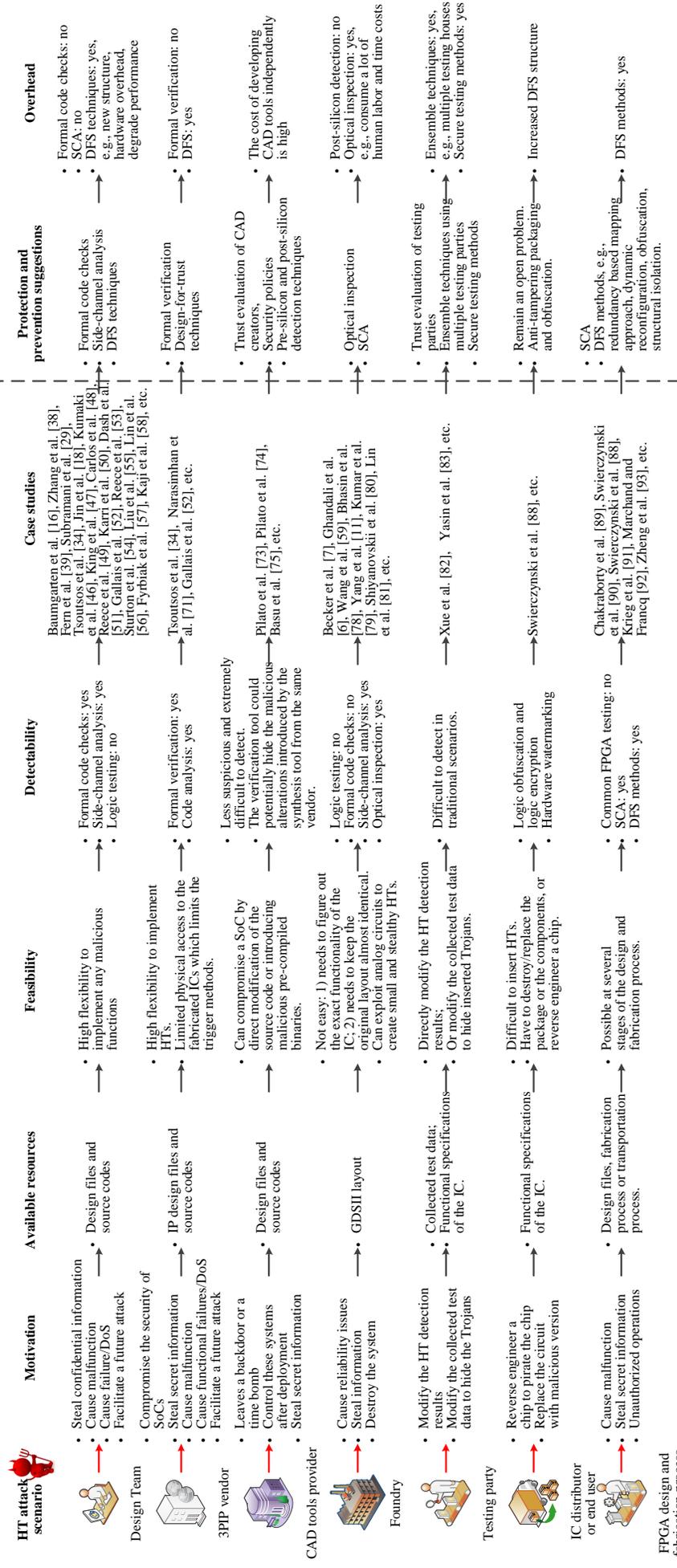


Fig. 2. HT attack scenario, motivation, available resources, feasibility, detectability (anti-detection capability), case studies under different attack models from the attacker's perspective, and protection and prevention suggestions, overhead from the defender's perspective.

which requires a comprehensive security policy to counter all possible threats. Using SCA, the impacts of Trojans on circuits' delay, leakage power, transient current, thermal profiles, electromagnetic emanation (EM), *etc.*, can be characterized for HT detection (19; 20; 21; 22; 23; 24; 25; 26). In addition to these traditional side-channel signals (power, delay, EM, *etc.*), emerging side-channel signals also include: impedance (27), backscattering (28), channel noise in wireless channels (29), *etc.* Traditional functional test fails to detect HTs, therefore, a few HT detection works have also been proposed to generate test patterns that target rarely activated nodes or events in a circuit (30; 31; 32; 33). However, the huge number of gates and states in modern ICs restrict the accuracy and the scalability of these methods.

Second, we discuss how to evade detection from an attacker's perspective. A malicious designer would normally have expertise in IC design. Thus they could insert optimized HT designs that balance the area and power overhead. Since RTL modification could affect all taped-out ICs, a golden model may not exist for use in timing and power analysis based detection methods. Moreover, process variation can also help hide HTs without hardware overhead. Functional analysis may be useful for HT detection at the RTL level. However, a stealthy HT can still evade detection using rare trigger conditions (34). Overall, HTs can be carefully inserted with design optimizations by exploiting rarely activated/observed conditions and introducing ultra low power/delay overhead to evade both post-silicon detection (30) and DFS techniques (35; 36).

Last, the emergence of new types of Trojan detection techniques will also pose challenges for Trojan design methods. Normally, HTs are stealthy with rare trigger events. As a result, HTs are usually not sensitized with test patterns during functional tests (37). Therefore, researchers can focus on such rare events and hidden corners for Trojan detection. Hicks *et al.* (37) proposed such a method, known as Unused Circuit Identification (UCI), which searches for unused components of an IC during design-time testing and marks them as potentially malicious. The UCI algorithm can detect many of the existing HTs reported in the literature, including most of the benchmarks in Trust-Hub (17), which poses a new challenge for HT implementation.

Protection and prevention suggestions: The security challenges faced by the designers are obvious. Attackers at this stage have very high flexibility to implement any malicious function (38). Furthermore, the earlier the Trojans are inserted, *e.g.* in the higher-level specifications or the RTL code, the harder for Trojan detection at later stages since it is impossible to obtain a golden RTL model (39).

From the designer's perspective, considering the motivations (goals) and capabilities of the design-stage attackers, the following countermeasures can be used against RTL Trojans: formal code checks (pre-silicon detection), SCA (post-silicon detection), and DFS techniques. Nahiyani *et al.* (40) propose a technique to analyze and quantify the vulnerabilities in a finite state machine (FSM). The state transition table of a FSM, including don't-care transitions and states, is extracted from a gate-level netlist, and then used for vulnerability analysis.

Xiao *et al.* (41) propose a built-in self-authentication (BISA) method to prevent HT insertion during physical design, which uses functional cells to fill all the unused spaces in the circuit layout. As the unused spaces in a circuit are the most likely insertion area for Trojans, any component changes in the BISA structure could be detected. The BISA structure is vulnerable to several attacks, therefore, Shi *et al.* (42) proposed obfuscated BISA structure to enhance its security. Other DFS methods include the ring oscillator (RO) based technique (35; 36), split manufacturing (43; 44), and so on. As a supplementary methods, runtime HT detection approaches, *e.g.* chaos theory based runtime power consumption monitoring (45), runtime data anomaly detection based on change-point (46), have also been proposed.

Overhead: Formal code checks and SCA do not introduce overhead to the circuit. However, the DFS techniques usually add new structure to the design thus will bring hardware overhead to the circuit. The power and area of the circuit will increase, and the performance of the circuit may be degraded due to the operations of the DFS structure.

Case studies: A summary of HT designs at RTL level proposed in the literature is shown in Table I.

(1) Cyber Security Awareness Week (CSAW). Since 2008, the annual embedded system challenge (ESC) competition, which is held as a part of CSAW, is well-known for its HT competition. The competition targets HT design and insertion techniques, Trojan detection approaches and design hardening mechanisms. Many researchers have reported their Trojan designs implemented for the CSAW ESC competition (16; 18; 49; 50; 60).

Several Trojan design and implementations have been presented by Baumgarten *et al.* (16) at CSAW ESC, including the following: information leakage through RS232 end sequence; RS232 multiple transmission rates; DoS; thermal leakage; information leakage through amplitude modulation (AM) transmission; 50MHz transmission; LED transmission. Jin *et al.* (18) presented eight RTL HTs to compromise the security of an Alpha encryption module. Santos and Fei (49) presented a backdoor Trojan, a bomb counter Trojan and a power sink Trojan to weaken a 8051 processor performing RC-5 encryption. Reece *et al.* (50) presented DoS and data leakage Trojans to attack the Intel 8051 micro-controller unit (MCU) which would probably run a data-sensitive encryption algorithm. Karri *et al.* (51) presented several case studies of HTs from ESC, including the following categories: key/information leakage through VGA display, the RS232 protocol or temperature; synthesis tool based Trojan; DoS Trojan.

In summary, these Trojans' payloads can be classified into three types as follows:

- leak sensitive data/internal signals;
- change the function of the design, or cause DoS;
- destroy the chip.

and the triggers of these Trojans can be divided into three types:

- input pattern triggered if the attacker can physically access the device;

TABLE I
AN OVERVIEW OF HARDWARE TROJAN DESIGNS TARGETED AT THE RTL LEVEL

Works	Benchmark	Trigger	Payload	Overhead	Detectability
Baumgarten <i>et al.</i> (16)	Alpha device (ESC2008)	Always-on; Internal	Information leakage; denial-of-service (DoS)	Power: 0.192%~1.026%	Functional: unlikely; SCA: likely
Zhang <i>et al.</i> (38)	OpenRisc	External; Internal	N/A	N/A	Functional: unlikely; SCA: likely; UCI: unlikely
Fern <i>et al.</i> (39)	bus protocols, ARM processor	Always-on	Leakage; Unprivileged access	Area: AXI4: 0.5%~2.1%(FF), 0.4%~3%(LUT); SoC: 0.9%(FF), 1.2%(LUT)	Functional: unlikely; SCA: unlikely; Formal verification: likely
Subramani <i>et al.</i> (29)	802.11a/g transceiver	Always-on	Information leakage	0.5dB~0.75dB extra power	Functional: unlikely; SCA: likely; Formal verification: likely
Tsoutsos <i>et al.</i> (34)	Data Encryption Standard (DES), XTEA, PRNG	External	Change functionality	N/A	Functional: unlikely; SCA: likely; Static analysis: unlikely
Jin <i>et al.</i> (18)	Alpha encryption	External; Internal	Leakage; Compromise functions; Destroy the chip	Area: -9.4%~3.3%(FF); 0.024%~6.8%(LUT)	Functional: unlikely; SCA: likely
Kumaki <i>et al.</i> (47)	AES	Predefined rule/input keyword	Leakage	Area: 0.37%; Power: 0.13%	Functional: unlikely; SCA: likely
King <i>et al.</i> (48)	Leon3 processor	A sequence of bytes; Predetermined bootstrap	Privilege escalation; Login backdoor; Stealing passwords	Area: 0.075%	Functional: unlikely; SCA: likely
Carlos <i>et al.</i> (49)	8051	External; Internal; Always-on	Leakage; Disables/enables functions	Area: 0.2%	Functional: unlikely; SCA: likely
Reece <i>et al.</i> (50)	8051	External; Internal	DoS; Leakage	Area: 0.15%~0.4%; Leakage power: 0.146%~0.399%; Dynamic power: -0.433%~0.93%	Functional: unlikely; SCA: likely
Karri <i>et al.</i> (51)	crypto-core <i>etc.</i> (ESC works)	External; Internal; Always-on	Leak secret key/info. through the RS232 protocol, through temperature using SCA, or through VGA display; DoS	N/A	Functional: unlikely; SCA: likely; UCI: likely
Dash <i>et al.</i> (52)	modern computers	A certain temperature	N/A	Area: high; Power: high	Functional: unlikely; Path delay-based: unlikely; Power-based: likely
Gallais <i>et al.</i> (53)	Rivest Shamir Adleman (RSA), AES	Specific instructions; Particular input	Leak info./secret key	N/A	Functional: unlikely; SCA: likely
Reece <i>et al.</i> (54)	AES	External; Internal; Always-on	Leakage; Drains the battery	Area: 90nm 0.16%, 45nm 0.78%; Leakage power: 90nm 0.53%, 45nm 0.46%; Dynamic power: 90nm 2.59%, 45nm 0.49%	Functional: unlikely; SCA: likely
Sturton <i>et al.</i> (55)	Leon3 processor	External	Change functionality	N/A	Functional: unlikely; Formal analysis: likely; UCI: unlikely
Liu <i>et al.</i> (56)	wireless cryptographic IC	Always-on	Leakage	Area: 0.005%, 0.025%; Power: 0.4%, 0.1%	Functional: unlikely; SCA: likely
Lin <i>et al.</i> (57)	crypto-processor	Always-on	Convey secret information off-chip	N/A	Functional: unlikely; SCA: likely
Fyrbiak <i>et al.</i> (58)	AES	Always-on; Conditionally	Cancel self-tests; key leakage	N/A	Functional: unlikely; SCA: likely; Formal verification: likely
Kaji <i>et al.</i> (59)	UART	Always-on	Facilitate data injection attack	N/A	Functional: unlikely; SCA: likely

- triggered internally by an internal event or sequence;
- always-on.

(2) **Crypto-cores.** HTs can be carefully designed to compromise the security of widely used crypto-cores, which may be of particular interest to an attacker. In (47), a malicious circuit was developed to connect the encryption module and the decryption module in an AES core. The Trojan is triggered when a predetermined condition is satisfied and then half-encoded data is sent from the encryption module to the decryption module by a specific Trojan path (47). Therefore, plain text is directly sent to the output. Moreover, when a predefined keyword is input to the AES core, which is transferred to a controller through the Trojan path, the secret key is outputted directly (47).

In (56), key leakage HTs are demonstrated in a wireless cryptographic IC that containing an AES module and an ultra-wideband (UWB) transmitter (TX) module. The impact of malicious components is carefully hidden below the side-channel margins. The key is leaked through parameter modulation, *e.g.* frequency or amplitude of the wireless transmission (56). The adversary is able to retrieve the 128-bit key leaked by a 128-bit ciphertext block through a transmission power waveform sent by the UWB TX.

In (53), HTs have been proposed to introduce or amplify

side-channel leakage of a cryptographic software. Particularly, they implement several alterations to cause information leakage through faulty computations or the variations in the power consumptions and latency of some instructions (53). Software-based Trojan activation mechanisms are proposed and the side-channel leakage of Rivest Shamir Adleman (RSA) and AES implementation were illustrated (53). Lin *et al.* (57) propose a HT, which conveys secret information off-chip by employing power side-channels. By using a spread-spectrum technique, the information is leaked below the noise level of the AES circuit so as to evade detection. Each key bit is modulated by a simple XOR operation with a pseudo-random number (PN) sequence (57). Fyrbiak *et al.* (58) propose a framework to RE the gate-level netlists and insert HTs to weaken cryptographic circuits.

(3) **Exploiting unspecified specifications, or creating covert channels.** Attackers can also carefully design HTs to hide in unspecified specifications or behaviors. Fern *et al.* (39) highlighted that current SoC bus implementations are vulnerable to HTs which can hide in the partially specified specifications or behaviors. They present a Trojan which introduces a covert channel by modifying bus signals of unspecified behaviors. The Trojan communication channel is demonstrated on a SoC design which runs a multi-user Linux OS to allow

an attacker get root user's data without permissions (39). It is shown that there are some redundant bus signals which will be ignored by the verification test, thus can be exploited for HT implementation (39). Subramani *et al.* (29) demonstrate a HT in wireless network by exploiting the forward error correction block to create a covert channel. Similarly, Kaji *et al.* (59) propose a data injection attack by exploiting HT to create specific electromagnetic waves as a covert channel.

(4) Remote activation. Since an attacker may have limited physical access to deployed devices, triggering HTs remotely is an ideal choice. Dash *et al.* (52) proposed a method to remotely activate HTs through a stealthy temperature channel. An analog HT trigger is implemented on modern computers which can be remotely activated when the infected circuit reaches a predefined temperature (52). The temperature of the target computer can be raised remotely by sending a large amount of network requests to the computer.

(5) Evading UCI detection. To evade UCI detection, Sturton *et al.* (55) constructed malicious circuits that have hidden behaviors. Particularly, this class of malicious circuits satisfies the following property: for all signal pairs (s, t) , there is at least one input that could make $s \neq t$ and would not trigger the hidden HT (55). This property ensures that the UCI technique will not mark the circuitry between s and t as a potential HT. Exhaustive enumerations of all circuits satisfying that class are performed, and the search results are used to construct an attack on a processor, *i.e.* the Leon3 processor. This HT allows user-level programs enter into supervisor mode to take control of the system by using a secret knock (55).

Zhang and Xu (38) proposed a HT design methodology from three aspects in order to bypass existing defenses, especially the UCI technique. First, to evade functional tests, carefully selected rare trigger conditions are used to make the HTs remain dormant during testing. Second, to evade UCI detection, they combine the trigger input selection and the code writing style to mask HTs as useful circuits. Third, they introduce a metric, namely un-controllability, to represent the difficulty level of setting the value of a signal (38).

IV. 3PIP VENDOR ATTACKS

HT attack scenario: In this adversarial model, the 3PIP used by a design house or a system-on-chip (SoC) developer may contain HTs, as shown in Figure 2. This is a general threat since SoCs are usually integrated with many 3PIPs with the purpose of reducing the cost and accelerating the time-to-market (1; 12; 2; 3). HTs could be inserted at each type of the IP, *e.g.* soft for RTL-level, firm for netlist-level, hard for Graphic Database System II (GDSII) cores (15). The SoC developer, who integrates design blocks and modules, often treats the 3PIP as black boxes. These unknown IP cores are usually unmodified and integrated into the final design, which can lead to an effective attack to compromise the SoC.

Motivation: Inserting a Trojan in 3PIPs is an effective and stealthy way for an attacker to compromise the security of SoCs. The attacker from a 3PIP vendor may want to insert a HT in the IP which serves as a backdoor to steal secret information from the integrated SoC, or cause functional

failures of the SoC. The attacker can also insert a HT to facilitate a future attack. For example, implant a hardware backdoor to support the software or system attacks.

Feasibility: Untrusted 3PIP vendors can easily introduce stealthy malicious modifications to a design through insertion, deletion or modification of original circuits or functions in a stealthy manner. This type of attacker can get the IP design files and the source codes. Therefore, an attacker can flexibly implement malicious functions by modifying the IP design at RTL or other specification levels.

However, the 3PIP attackers also face some obstacles. First, it is difficult for the 3PIP attacker to physically access to the fabricated ICs to trigger the HT. Therefore, the HTs are normally triggered internally. The HTs can also be designed to be always-on. Second, as the attackers insert HT in the 3PIP without knowing the overall design of the IC, it is not easy to carry out an attack successfully. In conclusion, the 3PIP Trojan has a high feasibility, good practicality, and is easy to implement. The only limitation is the method of triggering.

Detectability: We now discuss the available detection techniques from the defender's perspective, and the anti-detection capability from the attacker's perspective. Pre-silicon detection methods, *e.g.* formal verification, code analysis, are usually utilized to detect HTs in 3PIP cores (61). Previous researches (62; 32; 63) have proposed hardware description language (HDL) code analysis, or structural analysis techniques for soft IP cores. A SoC integrator can analyze the IP source codes and find potentially suspicious components by analyzing the reachability and controllability (64). However, the complexity of such analysis method is extremely high, which increases with the circuit size exponentially (65). In formal verification methods, IPs are verified by proof-checking tools to avoid including unintended functionalities (66; 67; 68; 69). Design-for-trust techniques have also been proposed. For example, Liu *et al.* (70) detected malicious HTs by applying security constraints to the task scheduling step of the SoC design process. Rajendran *et al.* (61) involved design constraints by using high-level synthesis to detect Trojans and then isolated the Trojan-infected 3PIPs. It is shown that using a variety of vendors can prevent collusion of multiple IPs from one vendor.

Fortunately for the attacker, verifying the trust of IP cores obtained from untrusted third-party entities is challenging due to several issues. For the case of 3PIPs, the methods depending on a golden chip/model are not suitable anymore. Moreover, the complete implementation of a 3PIP is invisible. It is difficult to provide sufficient coverage by general functional simulations due to incomplete functional specifications. It is also difficult to predefine comprehensive security rules to cover all the possible risks. There can always exist HTs which can satisfy proof-checking constraints thus evading detection. Lastly, the HDL and Coq (71) (an interactive theorem prover/proof-assistant) representations of a circuit may not be completely equivalent. Even if the Coq representations of the circuit are verified to be trustworthy, it cannot guarantee that the corresponding HDL code is trustworthy (61). A smart attacker can ensure that the functional specification of the design is unchanged or the modifications are undetectable.

Protection and prevention suggestions: Pre-silicon detec-

tion methods, *e.g.* formal verification methods, can be used to detect HTs that are inserted in 3PIP cores, while post-silicon detection methods usually cannot detect 3PIP Trojans. The DFS techniques, such as using a variety of vendors, can also be used to prevent 3PIP Trojans.

Overhead: The formal verification methods will not introduce overhead to the circuit, while the DFS techniques usually bring some overhead to the design.

Case studies: Tsoutsos *et al.* (34) presented HTs which do not violate the functional specifications. Multiple levels of malicious nested finite state machines (FSMs) are introduced to the design. The threat scenario is that the SoC integrator which receives the malicious IP only applies static analysis methods on the HDL code of the IP, without actually simulating or implementing the design (dynamic analysis) (34). Such modifications are hard to detect without exhaustive testing of all system states.

While HTs are generally considered to be malicious, they may also go in the opposite direction, *e.g.* be exploited in a constructive way. In (72), a hardware IP protection technique which embeds a HT as a FSM was proposed. By using a sequential Trojan that acts like a time-bomb, illegal SoCs containing pirated evaluation copies of the IP could stop performing the specified functionality. Specifically, on occurrence of a rare sequence, the Trojan stops the normal usage of the IP (72). Therefore, an expiry date on the usage of the IP can be set up based on the Trojan.

So far, most research works consider side-channels as undesired signals such that people need to protect devices from sophisticated SCA attacks. However, Gallais *et al.* (53) used side-channel leakage introduced by a HT as a watermark for IP protection, which can be detected by SCA. A unique signal is embedded into the side-channel signal of a circuit which acts as a watermark. This enables circuit designers to detect unauthorized use of their circuits. They illustrate this by designing an integer-based multiplier (53). When a specific pair of operands arrives, the pipeline will be stalled for several clock cycles. Since this pair of inputs are hard to guess, it allows the designer to verify his own IP by analyzing the power profile (53).

V. CAD TOOLS ATTACKS

HT attack scenario: Untrusted commercial CAD tools supplied by different vendors can also introduce malicious circuits to a design, which reflects the synthesis & verification stage attacks, as shown in Figure 2. CAD tools attackers can directly insert Trojan circuits into the design or degrade critical logic, *e.g.* the random number generator (RNG) used in a cryptographic circuit.

Motivation: The attacker from the CAD tools providers may want to insert HTs in the design files, which leaves an undetectable backdoor or a time bomb in these designs. The attacker can also control these systems after deployment or steal secret information from those systems.

Feasibility: Although this attack model is less possible compared with design attacks and fabrication attacks, it is still feasible. A CAD tools attack is more powerful and stealthy

than design attacks and fabrication attacks. A SoC designer has to design chips by relying on CAD vendors. By compromising the CAD tools or the running scripts, the attacker can introduce malicious modification to IPs from the HDL codes to the generated netlist (16). In conclusion, the CAD tool attacks are feasible, having good practicality and stealthiness, but are not easy to implement.

Detectability: Since this attack happens during the synthesis stage on generally trusted tool suits, it is not suspicious and extremely hard to detect (16). On the other hand, Trojans inserted by CAD tools are difficult to detect or remove since they are coupled with other design units (73). Furthermore, a SoC designer often uses a suite of CAD tools supplied from the same vendor, which means the verification tool could potentially hide the malicious alterations introduced by the synthesis tool from the same vendor (1; 16).

Protection and prevention suggestions: A trust evaluation of CAD creators, and security policies, which are currently lacking, need to be established to defeat the malicious tampering by CAD tools (16). It is suggested to use reliable CAD tools or use self-developed CAD tools. Pre-silicon and post-silicon Trojan detection techniques are also needed to be applied.

Overhead: The cost of developing CAD tools independently is high. However, there will be security threats when using third-party CAD tools.

Case studies: Pilato *et al.* (74) demonstrate the CAD threats by compromising a high-level synthesis tool to insert three HTs. The payloads of these Trojans are adding latency, compromising the security of crypto-cores, and draining energy, respectively. Similarly, Pilato *et al.* (73) use high-level synthesis to inject a benign HT, which serves as a IP watermark to prevent piracy and counterfeiting. Basu *et al.* (75) investigate the CAD attacks from all the CAD tools (from synthesis, design, verification to test), and show that all these CAD tools can launch potential attacks. They demonstrate the CAD-attacks on a ARM Cortex processor.

VI. FABRICATION STAGE ATTACKS

HT attack scenario: This attack model represents the threat of untrusted foundries. Nowadays, most modern ICs are manufactured worldwide in untrusted foundries due to budget considerations. The foundry receives the complete design (physical layout geometry file) and its specifications. However, the IC designer has little or no control over the foundries. A fabrication attacker could modify the manufacturing process by directly inserting a HT into the chip or changing the manufacturing process steps to cause reliability issues in the SoCs.

Motivation: A fabrication stage attacker may want to cause reliability issues in the SoCs, steal information from the ICs, or even directly destroy the system.

Feasibility: The foundries have complete access to the layout of the design, which provides them with opportunities to flexibly add or remove components of ICs by modifying the layout. Since the foundry has no access to the RTL code, the modifications can only be achieved in the layout by changing

the netlist, or modify the manufacturing process by changing design masks in order to not affect the functions of the design (7), as shown in Figure 2.

Generally, it is not easy for an attacker to insert HTs during fabrication. First, the attacker has to figure out the exact functionality of the circuit (in the form of the GDSII file). The attacker also needs to find the necessary space to add extra gates and connections. Second, the attacker needs to keep the layout (place & route) almost identical, to avoid being detected by optical inspection. In conclusion, the fabrication stage attacks are feasible, having good practicality and stealthiness, but are not easy to implement.

Detectability: A Trojan inserted during fabrication is difficult to discover by functional tests and verification performed on the HDL. When the layout of the circuit remains unchanged during Trojan insertion, it is almost impossible to detect these Trojans by using optical inspection. An attacker can insert HTs based on the modification of the electrical characteristics while the metal, active area and polysilicon layer remain unchanged (7).

Fabricating a golden chip in a trusted factory for HT detection is difficult in practice. Thus, the detection technique can only compare the golden simulated model and the fabricated chip under test. However, Yang *et al.* (9) show that an analog HT can be much smaller and more stealthy than a digital HT. The trigger is implemented by diverting charge from unlikely signal transitions, which makes the Trojan invisible to side-channel detection.

Protection and prevention suggestions: Optical inspection is considered as a reliable way to detect layout-level HTs, while SCA is also a general method to detect this kind of HT. Since the RE based optical inspection needs a lot of human efforts, it is also helpful to use machine learning based image analysis method for automatic analysis (76). Besides, a golden simulated model, if exists, will be helpful for post-silicon detection techniques (77), *e.g.* SCA.

Overhead: Post-silicon detection does not introduce overhead to the circuit, while the optical inspection will consume a lot of human labor and time costs.

Case studies: As most of the HTs reported to date in the literature are inserted at RTL level, constructing practical HTs at the layout level is still an open problem. The summary of HT design and implementation works targeting at the layout level is shown in Table II, which will be discussed in the following paragraphs.

(1) Exploit analog circuits. It has been shown that an attacker during fabrication can exploit analog circuits to create small and stealthy HTs (9). Yang *et al.* (9) leveraged analog circuits to perform a hardware attack, named A2. A circuit is constructed using capacitors to siphon charge from nearby wires in the spare spaces of a design after place & route. A victim flip-flop is changed to a desired value when the capacitors are fully charged. This attack has been implemented in an OR1200 processor by applying it to privilege escalation which can be controlled remotely (9).

(2) Parametric Trojans. Becker *et al.* (7) proposed layout-level Trojans by slightly altering the manufacturing process conditions, *i.e.* the dopant polarities of a transistor. The Trojan

can accelerate the wear-out mechanisms so as to affect the reliability of ICs. The Trojans have been inserted into two designs, *i.e.* an Intel secure RNG in an Ivy Bridge processor, and a side-channel attack resistant substitution box (S-Box) implementation (7).

Ghandali *et al.* (6) presented a parametric Trojan, which is designed through modifying the parameters of transistors, and does not require extra logic. It is triggered under rare conditions that are determined by the delays of some combinational logic paths. This design has been applied in a multiplier circuit to create a Trojan multiplier. If specific patterns are input, this Trojan multiplier will compute faulty outputs (6). This Trojan multiplier is further applied to attack a key agreement protocol, the Elliptic Curve Diffie-Hellman (ECDH). The bug attack works as follows (6). First, the first several bits of the key are guessed. Then, a point Q which can lead to a failure of the scalar multiplication is searched. After that, the attacker sends Q to the server to make a handshake which performs the ECDH protocol. If the handshake fails, it indicates that the Trojan multiplier outputs the expected multiplication error. Hence, the current guessed key is correct. More bits will be cracked successively in this way to recover the key (6).

Kumar *et al.* (79) proposed parametric manufacturing process HTs to facilitate fault-injection attacks. The Trojans are inserted by altering the doping concentration and the dopant area of predetermined transistors in a target circuit. The trigger condition of the HT is a slightly reduced supply voltage with very low probability (79). The Trojans have been utilized to inject faults into the lightweight cipher PRINCE. It is shown that they can reconstruct the secret key after around 5 fault-injections by differential cryptanalysis (79).

Shiyanovskii *et al.* (80) proposed a HT based on process reliability. The Trojan reduces the reliability of ICs by altering the conditions of the manufacturing process, to wear out CMOS transistors. Such minor changes in the manufacturing process are hard to detect.

(3) Unchanged place & route. Bhasin *et al.* (78) analyzed how to introduce a HT while the place & route remain unchanged. It is shown that when the placement density is over 80%, it is difficult to insert Trojans. They also inserted a Trojan to aid the differential fault analysis (DFA) attack. The payload of the HT is an XOR gate that alters one bit of the AES to be faulty in the 8th round. As a result, the attacker can retrieve the whole key by activating the HT for two encryption processes (78).

Wang *et al.* (60) considered new placement techniques and delay-aware Trojan insertion. A hard macro is used to prevent delay variations in FPGAs. For the Application Specific Integrated Circuit (ASIC) scenario, where the Trojan is inserted at post-layout, the placement and route of the original design is also preserved by making it a hard macro (60). It is shown that such Trojans only have small impact on path delay, which can evade on-chip monitor based DFS approaches.

(4) Trojan side-channels. Lin *et al.* (81) used side-channel leakage for HT implementations, called Trojan Side-Channels (TSCs). A hidden backdoor can be inserted at the foundry for unauthorized leakage of secret information. Power side-channels are demonstrated to leak information that can be hid-

TABLE II
SUMMARY OF HT DESIGN AND IMPLEMENTATION WORKS AT THE LAYOUT LEVEL

Paper	Benchmark	Trigger	Payload	Overhead	Detectability
Becker <i>et al.</i> (7)	RNG, AES	Always-on	Change functionality; Degrade performance; Leakage	Area: 0	Functional: unlikely; Optical inspection: unlikely
Ghandali <i>et al.</i> (6)	32-bit multiplier, ECDH key agreement protocols	Violating the delays of rare combinational logic paths	Trojan multiplier computes faulty outputs	N/A	Functional: unlikely; Visual inspection: difficult; SCA: difficult
Wang <i>et al.</i> (60)	ESC2010	Rare events	Change function	Area: 0.6%, Power: 0.4%	Functional: unlikely; SCA: unlikely
Bhasin <i>et al.</i> (78)	Cryptographic IP	External; Internal	Facilitate DFA; Leakage	Area: 0.5% (LUT)	Optical imaging: likely
Yang <i>et al.</i> (9)	OR1200	Internal	Privilege escalation; Change functionality	Area: 0.08%; Delay: 0.33%	Functional: unlikely; SCA: unlikely
Kumar <i>et al.</i> (79)	PRINCE	Slightly reduced supply voltage	Facilitate attacks	Area: 0	Functional: unlikely; SCA: likely; Optical inspection: unlikely
Shiyanoskii <i>et al.</i> (80)	SRAM	Always-on	Reduce the reliability by acceleration of the wearing out mechanisms	N/A	Functional: unlikely; Delay monitoring: difficult; RO: difficult; Wafer and package level reliability monitoring: likely
Lin <i>et al.</i> (81)	Crypto core	Always-on	Convey secret information	Area: 14 LUTs	Functional: unlikely; SCA: likely

den in the noise. Two Trojan side-channels are implemented, *i.e.* TSC based on spread-spectrum theory and TSC using specific input values (81). Moreover, the TSCs have physical encryption property, so that it can keep the information secure even if the introduced side-channel is successfully detected.

Table III presents a comparison between Trojan insertion at RTL level and layout level from the attacker’s perspective. The advantages of Trojan insertion at RTL level are having full access to the source code and high flexibility to implement any malicious function. Moreover, as an RTL modification will affect all fabricated ICs, a golden model may not exist for SCA based detection methods. The disadvantages of RTL Trojans are that they can be exposed by complete code reviews, adequate security policy checks, or SCA. On the other hand, the advantages of Trojan insertion at the layout level are that it can evade detection by functional testing and be invisible to side-channel defenses. It can also leverage analog circuits or parameter changes to introduce small and stealthy HTs. The disadvantages of Trojan insertion at the layout level are that it is not easy for the attacker to make modifications to the layout mask or change the manufacturing process. They must also keep the original place & route mostly unchanged, to avoid being detected by optical inspection.

VII. TESTING STAGE ATTACKS

HT attack scenario: In general, the manufacturing test is done by a credible test party, *e.g.* reputable semiconductor company or government agency, which could be considered as trusted. As a special case, Xue *et al.* (82) formulate untrustworthy testing parties into two attack models and illustrate that the test parties may be untrustworthy. In (83), Yasin *et al.* extract secret information from test data. These attacks indicate that the testing phase may also be insecure. The testing party is important in the IC supply chain. However, nowadays, there is usually only one test party to test the fabricated ICs. If the only testing house is not credible, or colludes with attackers from other stages (84), the testing results will no longer be trustworthy.

Motivation: An attacker during the test stage may want to modify the HT detection results or modify the test data to hide the HTs.

Feasibility: Generally, the testing party collects test data of fabricated ICs and then performs the HT detection procedure.

In this scenario, the testing agency can directly modify the Trojan detection results. In a special case, the designer is involved in the testing process. In this scenario, the test agency needs to modify the test data to mislead the final Trojan detection result. An adversarial test data generation algorithm was proposed in (82) for the untrustworthy testing houses, which can use the minimum test data modifications to cause the maximum detection errors of ICs. In conclusion, the testing stage attacks are feasible, having good practicality, but are not easy to implement, as shown in Figure 2.

Detectability: Little research has been done on testing stage defenses. Xue *et al.* (82) proposed a HT detection method based on hybrid clustering ensemble to resist untrustworthy testing houses. Three testing houses are used in the scheme, and each testing house carries out the HT detection process. Then, the three detection results is consolidated by using the hybrid clustering ensemble method to obtain the final test result. The technique can resist malicious modifications by untrustworthy testing houses, and can achieve higher detection accuracy than each of the three testing houses regardless of whether the testing house has maliciously modified the test data or not (82).

Protection and prevention suggestions: Since the motivations of the testing stage attackers are to modify the HT detection results or modify the test data to hide HTs, two methods can be applied to resist such attacks. One is trust evaluation of testing parties, and the other is the ensemble technique using multiple testing parties (82).

On the other hand, there have already been a few secure testing methods against IC piracy, which may provide a reference for secure testing of HT detection. For example, Contreras *et al.* (85) present a Secure Split-Test (SST) technique to prevent counterfeiting. The method allows the IP owner to meter the IPs by holding a lock key. During the test phase, a key is required to unlock the IP’s functionality, so that the IP owner can verify the testing results. Later, Rahman *et al.* (86) improve the above SST technique against IP piracy by simplifying the communication between the IP owner and the foundry, named CSST. In the CSST method, the IP owner controls the testing by locking the IC and the scan chains (86). Only the IP owner can understand the testing results under locking conditions, and can unlock the IC. Zhang *et al.* (87) propose a hybrid

TABLE III
COMPARISON OF TROJAN INSERTION AT RTL LEVEL AND LAYOUT LEVEL FROM THE ATTACKER'S PERSPECTIVE

	RTL level	Layout level
Pros	1) Have full access to the source code; 2) High flexibility to implement any malicious function; 3) Since an RTL modification will affect all the fabricated ICs, a golden model may not exist for SCA;	1) Can evade detection by functional testing and verification; 2) May be invisible to SCA; 3) Can leverage analog circuits or parameter changes to introduce small and stealthy Trojans;
Cons	1) Can be revealed by complete code reviews and adequate security policy checks; 2) May be exposed by SCA;	1) Have to make modifications to the layout mask or at process level which is neither easy nor flexible; 2) Must keep the original layout mostly unchanged, to avoid being detected by optical inspection;

approach that combines a dynamically obfuscated wrapper technique (referred to as DOST) and SST to protect IP rights, which allows the IC designer to control the fabrication and the testing processes. In the locked model, structural tests are performed, while in the unlock model, the functional tests can be performed (87).

Overhead: Since multiple testing houses are used, the cost will be of particular concern. It is shown in (82) that the time overhead of ensemble technique is small and acceptable, while the computational overhead is large. However, the computational overhead is distributed across multiple testing parties, which means that the ensemble technique does not increase the computational and storage overhead of each test party.

Case studies: To date, little research (82)(83) has been done on testing stage HT attacks, as described above.

VIII. DISTRIBUTION STAGE ATTACKS

As described in Section II, since the distributor is usually unaware of the IC design, it is generally considered that the distributor cannot insert a HT. However, a distribution stage attacker can RE a chip to pirate the chip, or directly replace the IC with a Trojan-inserted version during transportation. Therefore, we also describe the attackers from the distribution stage.

HT attack scenario: After IC fabrication and packaging, a distribution attacker may appear in the IC supply chain. Such a distribution attacker, which may be either an IC distributor or an user, is restricted in inserting Trojans. Instead of being able to modify logic gates, they have to destroy the package or the components, or manipulate the transport process (16), as shown in Figure 2.

Motivation: A distribution stage attacker may want to RE a chip so as to pirate the chip. The attacker may also want to directly replace the circuit with a Trojan-inserted circuit during transportation.

Feasibility: An attacker has limited flexibility at this stage and it is difficult to implement such hardware attacks. Such attackers cannot obtain the HDL code and the layout level geometry. The attacker also does not have the input/output test patterns. However, they usually have a set of specifications about the functions of the ICs. They may obtain the netlist of the design by RE, which is a difficult but feasible task. In conclusion, the distribution stage attacks have limited flexibility, some practicality, and are difficult to implement.

Detectability: Some defense techniques have been proposed to address this type of vulnerability, including anti-tampering packaging and obfuscation against SCA (16). HT attacks and defense techniques at this stage remain open problems.

Protection and prevention suggestions: Since the motivation of an distribution stage attacker is to RE or replace the circuit, logic obfuscation and logic encryption can be used against RE attacks. Some fragile hardware watermarking structures can also be used. Once the integrity of the hardware is compromised, the watermark will be broken.

Overhead: The DFS techniques will add hardware overhead to the circuit.

Case studies: Swierczynski *et al.* (88) described a HT attack on a USB flash drive. The USB flash drive is intercepted and attacked during transportation. The FPGA bitstream is manipulated such that the S-Box of the 256-bit AES design is changed to a linear function, and thus can be easily broken (88). If the attacked USB flash drive is used by a victim, the user's data can be revealed from the ciphertexts.

IX. FPGA TROJANS

HT attack scenario: With the extensive use of FPGAs in critical applications, the security of FPGA designs has become a major concern. In the past, the researches have focused on IP protection in FPGA, *i.e.* protecting the IP mapped on an FPGA from being stolen. However, little research has been conducted on security and protection of the FPGA device itself. Recently, a few FPGA HT detection techniques have been proposed, while the FPGA-oriented HT design and implementation works are relatively less.

Motivation: The attacker may want to cause malfunction of the FPGA system, steal secret information, or lead to other unauthorized operations.

Feasibility: Similar to the ASIC scenario, malicious modifications of the FPGAs are possible at several stages of the design and fabrication process. An attacker can create an FPGA Trojan by directly modifying the HDL, modifying FPGA fabric, modifying physical parameters, modifying the bitstreams, or by exploiting FPGA CAD tools (95). For example, a malicious circuit can be inserted by an adversary to monitor the logic values of internal nodes, logic modules, and the look-up tables (LUTs) (96). Once the Trojan is triggered, the FPGA can malfunction in different ways, *e.g.* the LUT values can be changed, configuration cells can be altered to perform incorrect routing, or incorrect values can be written into block-RAMs (BRAM) (96). In conclusion, the FPGA Trojan attacks are feasible, having good practicality, but are not easy to implement, as shown in Figure 2.

Detectability: These Trojans can escape common FPGA testing that cannot cover all possible triggering conditions. Existing FPGA Trojan defense techniques fall into two categories, SCA and DFS techniques. The power consumption

TABLE IV
SUMMARY OF FPGA-ORIENTED HT DESIGN AND IMPLEMENTATION WORKS

Works	Benchmark	HT type	Insertion mechanism	Trigger	payload
Chakraborty <i>et al.</i> (89)	128-bit AES on Xilinx Virtex-II	Bitstream Trojan	Bitstream modification to implement many ROs as the HT	Always-on	Temperature increases thus accelerating aging
Swierczynski <i>et al.</i> (90)	AES and 3-DES	Bitstream Trojan	Detect S-boxes in bitstreams, then modify the bitstream of S-boxes	Always-on	Weaken the cryptographic algorithm
Swierczynski <i>et al.</i> (88)	XTS-AES on Kingston DataTraveler 5000	Bitstream Trojan	Bitstream modification replacing AES S-boxes during interdiction	Always-on	Recovering plaintext
Krieg <i>et al.</i> (91)	iCE40 design flow running a instruction decoder of a CPU	CAD tool Trojan	Malicious insertion during synthesis, then activate malicious part during bitstream generation	Output of the malicious LUT	Privilege escalation
Marchand and Francq (92)	128-bit AES on SASEBO-GII Board (Xilinx Virtex-5)	Functional Trojan	Design, place and route the 12 HTs by hand	Time based, user, internal state	DoS, changing specifications, information leakage
Zheng <i>et al.</i> (93)	OpenRISC ORI200 on Xilinx Spartan-6	Functional Trojan	Implementing delay-logic arbiters as HTs	Digital value	Disable general purpose registers
Krieg <i>et al.</i> (94)	Xilinx 7	Functional Trojan	Exploiting the X-Optimism operations in an FPGA simulation model	Always-on	The signal which is '0' during simulation becomes '1' in hardware

based (92) and electromagnetic emanation (EM) based (97) SCA methods are proposed to detect FPGA HTs. Chen *et al.* (98) measure the EM of FPGA clock tree, and use principal component analysis (PCA) for signal processing. Then, back propagation (BP) neural network is used to automatically detect the FPGA HTs. Similar to fingerprint-based HT detection methods in ASIC scenarios, FPGA detection methods based on anomaly features have also been proposed. Pino *et al.* (99) propose a process variation based anomaly detection method for FPGAs which can isolate suspicious Trojan areas with inconsistent characteristics. In their later work (100), after isolating these suspicious areas, the remaining trustworthy areas, named FPGA Trust Zone, are selected to run the designs securely.

On the other hand, some DFS techniques are also proposed against FPGA Trojans. Mal-Sarkar *et al.* (96) propose a redundancy-based approach using Trojan tolerance which modifies the application mapping process to provide defenses against HTs. Swierczynski *et al.* (101) use dynamic obfuscation of cryptographic primitives to prevent the bitstream reverse engineering and modification based FPGA HTs. Bloom *et al.* (102) propose a FPGA HT defense technique, named MORPH, which uses onion-encryption for encrypted execution, and use a specific hardware abstraction layer to isolate the hardware and software. Zhang *et al.* (103) use the moving target defense principle to prevent FPGA CAD tools based Trojan insertion, in which three kinds of unpredictability are introduced into the FPGA designs.

Protection and prevention suggestions: Considering the diverse motivations (goals) and strong capabilities of the FPGA Trojan attackers, the defense against FPGA Trojans is still an open problem. SCA method can be used with the help of a golden model or built-in consistency verification structures. DFS methods, *e.g.* redundancy based mapping approach, dynamic reconfiguration, obfuscation, structural isolation, are also promising protection methods against FPGA Trojans.

Overhead: The DFS techniques will bring some hardware overhead in terms of logic resources (area), power, and performance.

Case studies: The summary of FPGA-oriented HT design and implementation works is shown in Table IV. The HT Type is based on the FPGA HT taxonomy proposed in (95). Note that, the FPGA Trojans implemented by using direct HDL

modification are not included in this table, because those HTs are not specifically for the FPGA, but just using the FPGA device as a code verification platform.

Chakraborty *et al.* (89) insert HTs into FPGA by directly modifying the unencrypted bitstream file. They implement a number of ROs as the HT in a 128-bit AES circuit, which can cause the temperature increases thus lead to accelerated aging. Since this Trojan is inserted during the bitstream configuration, it does not leave traces in the logic and place & route phases (89). Swierczynski *et al.* (90) propose an FPGA bitstream Trojan implementation scheme, which detects the S-boxes in bitstreams, then modifies the bitstream of S-boxes to weaken the AES and 3-DES algorithms. As mentioned in Section VIII, Swierczynski *et al.* (88) propose a interdiction based FPGA Trojan insertion, which modifies the bitstream to replace AES S-boxes. They demonstrate their work on XTS-AES on Kingston DataTraveler 5000 to recover plaintext. Krieg *et al.* (91) propose an FPGA CAD tool Trojan, including malicious insertion during synthesis, and malicious part activation during bitstream generation. They evaluate the scheme using iCE40 design flow running a instruction decoder of a CPU to launch a privilege escalation attack. Marchand and Francq (92) design, place and route 12 functional FPGA Trojans by hand on 128-bit AES on SASEBO-GII Board (Xilinx Virtex-5), which can lead to DoS, changing specifications, or information leakage. Krieg *et al.* (94) implemented a Trojan trigger by exploiting the X-Optimism operations (on unknown 'X') in an FPGA simulation model. They generated a trigger signal which is '0' during simulation phase and '1' in implemented hardware. FPGA HTs can also be used with a benign purpose. Zheng *et al.* (93) propose an functional Trojan to disable particular general purpose registers, which works as a security mechanism for FPGA systems. They implement delay-logic arbiters as HTs and evaluate on OpenRISC ORI200 on Xilinx Spartan-6. These efforts demonstrate the flexibility of FPGA HTs.

X. FUTURE DIRECTIONS

In this section, we will discuss the potential future HT implementation and detection techniques.

A. HT benchmarks and evaluation methods

A common concern is that whether a real HT has been found in industry. Due to the sensitive nature of industry IP,

it is unlikely that such HTs will be reported publicly. As such, standard benchmarks to evaluate HT implementations and defenses are highly needed. The Trust-HUB benchmark (17) developed by Tehranipour *et al.* is well-known for its hardware security related benchmarks. Trust-HUB (17) currently provides the largest database of HT benchmarks and has been widely used in the literature. For example, Reece and Robinson (54) evaluated 18 AES HTs supplied from the Trust-HUB database, in terms of power and area. It was shown that when spending enough effort on optimizing the HT, the introduced overhead could be very small.

Furthermore, in order to create dynamic Trojan benchmarks, Cruz *et al.* (104) proposed an automatic HT insertion framework, which can insert HTs with validated trigger conditions and payloads in gate-level designs. It allows configurations, *e.g.* the type of the Trojan, Trojan trigger probability, and choices of payload (104). Although this powerful Trojan automation design and implantation tool has been emerged, Trojan design and detection is still a game process. Once the defenders know how these automatic tools generate HTs, Trojans inserted by these tools may also be easily detected. However, defenses always lag behind attacks. On the other hand, various new HT detection techniques have also been proposed. When attackers are aware of these detection methods, more powerful Trojan design methods will also appear.

To date, most of the Trojan implementation methods or Trojan detection techniques are verified under specific experimental conditions, specific stages and specific scenarios, and targeting specific circuits or Trojans. This non-uniform paradigm raises a question: which attacks (defenses) are more effective (universal)? To this end, a uniform evaluation method with comprehensive evaluation metrics is required to evaluate and analyze various HT implementation methods and defense techniques. Such a uniform evaluation method can ensure researchers and IC designers to: 1) evaluate the effectiveness of different HT attack and defense methods; 2) assess ICs' vulnerabilities; 3) carry out complete and quantitative comparative researches on HT implementations and detection methods.

B. Machine learning-based Trojan detection methods, and HTs targeting machine learning models

Recent research in this field has explored machine learning methods for HT detection (105; 106; 76; 107; 108; 77; 109; 110; 84; 82; 111). Generally, machine learning methods can be utilized for HT detection in the following aspects: providing automatic layout identification in RE-based methods (105; 106; 76), providing run-time HT detection architectures which are trained by HT attack behaviors (107; 108), providing automatic feature analysis (112), and providing golden chip-free HT detection techniques based on classification or clustering (77; 109; 110; 84; 82; 111). In particular, the machine learning method has its own specialties in feature extraction and image recognition, which makes it possible to reveal unknown HTs by monitoring suspicious behaviors and features. It can also improve detection capabilities through self-learning. Elnaggar and Chakrabarty (113) reviewed the works applying machine learning methods for hardware security. Specifically,

they summarized that the machine learning methods can be used to classify or cluster the IC's parameters, gate-level nets, or traffic data in multi-core systems, for HT detection (113).

To defeat machine learning-based Trojan detection methods, attackers may introduce adversarial HT designs which can make the detection methods produce incorrect decisions. In machine learning systems, adversarial input perturbations carefully crafted at test stage can subvert the model's predictions on the instances. Attackers can investigate the vulnerability of machine learning models to such adversarial inputs (also known as adversarial examples) to mislead the HT detection. Xue *et al.* (82) propose a data modification algorithm for untrusted testers to slightly modify the collected test data, so as to mislead the HT detection results. Such an example is illustrated in Figure 3, in which the original power trace x of an IC is detected as Trojan-infected by machine learning based HT detection method (82). After introducing an imperceptible adversarial perturbation δ to the test data, the power trace $x + \delta$ will be recognized as Trojan-free by the machine learning model (82).

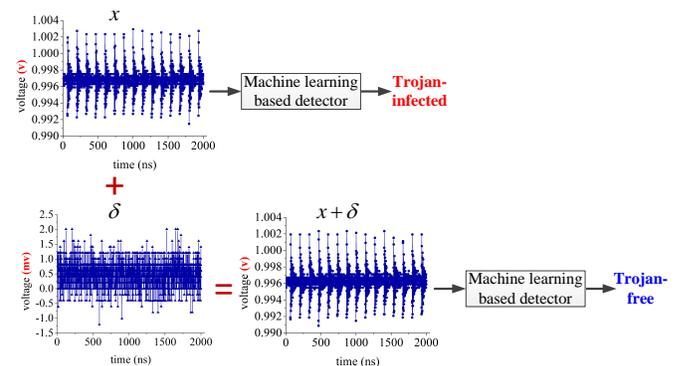


Fig. 3. Illustration of adversarial HT designs against machine learning-based Trojan detection method: the original power trace x is detected as Trojan-infected. After introducing an imperceptible perturbation δ to x , the power trace $x + \delta$ will be detected as Trojan-free.

On the contrary, there are also few works to study HT attacks targeting machine learning models and artificial intelligence (AI) chips. Clements and Lao (114) propose to insert HTs in the functional block of the neural network implementations. As a result, a desired misclassification can be achieved when a specific input trigger arrives. Ye *et al.* (115) insert a HT into the FPGA Convolutional Neural Network (CNN) accelerator to launch an attack on a CNN based image classification task. The HT can control the classification result once triggered. Odetola *et al.* (116) propose a HT attack on deep learning models without modifying the parameters or functions within the layer. They exploit statistical properties of each layer's output to trigger the HT, which makes the HT extremely stealthy. Li *et al.* (117) propose a more flexible attack framework on neural network which combines the hardware and software. Particularly, in addition to the hardware HT circuit, Trojan weights are embedded in neural networks. The Trojan is only inserted in a part of the network, and does not affect the overall accuracy, thus can ensure stealthy (117). In the above attacks, the attacker needs to have the knowledge of

the model. Hu *et al.* (118) propose memory Trojan on Deep Neural Networks (DNN), in which the Trojan logic is only inserted into the memory controller without the knowledge of the model. Targeted attacks or untargeted attacks can be achieved when the trigger image arrives.

C. Attacks and defenses from chips to complex systems

Most of the existing Trojan attacks and detections aimed at the chip level. A more practical scenario in industry is that how to implant and detect hardware Trojans on a complex SoC, or larger systems. Such a system contains many components and connections. It also contains hardware, firmware and software. This makes HT attacks more diverse, such as hardware-promoted software attacks, or software-promoted hardware attacks, or covert-channel attacks, and so on. It is important but challenging to detect HTs in such a complex system.

D. Universal Trojan and automatic Trojan insertion VS automatic Trojan (IC vulnerability) analysis tools

Most HTs reported to date are manually inserted into a specific target circuit (119). However, a more ideal situation is that arbitrary Trojan circuits with arbitrary components could be inserted into any circuits, as shown in Figure 4. There are two requirements for such practical attacks: 1) designing a universal Trojan independent of the host circuit, which is applicable for any given circuit; 2) developing automatic Trojan design and insertion tools. In order to meet these requirements, the automatic trigger and payload identification of a design at different levels are required.

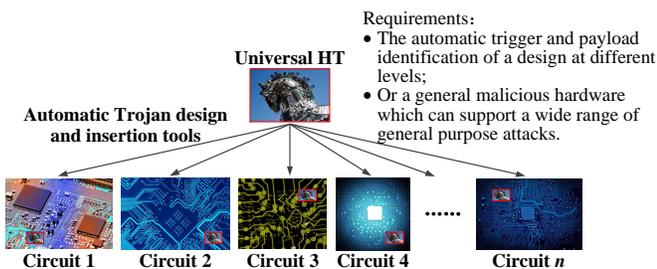


Fig. 4. Universal HT and automatic Trojan insertion.

Another type of universal Trojan is a general malicious hardware which can support a wide range of general purpose attacks. This is a more aggressive attempt. King *et al.* (48) presented two such hardware designs in Illinois Malicious Processors, and exhibited three attacks using this hardware. Through the memory access mechanism, a privilege escalation attack was implemented, which gives the attacker root access without identification or creating system logs. Under a shadow mode, a login backdoor is implemented, giving an attacker authority to log in as a root user with no password needed (48). Another function which steals passwords is also implemented.

In contrast, automatic Trojan analysis tools and automatic IC vulnerability analysis tools (120) are required. Most of the existing Trojan detection techniques are manually customized detection methods/scripts applied for specific scenarios and

specific stages. The detection process requires manual participation, and the universality is limited. It is necessary to develop mature universal tools, including automatic Trojan detection tools and automatic circuit vulnerability analysis tools, to promote DFS and Trojan detection works.

E. Multi-stage HT attacks and defenses

The majority of previously reported HTs in the literature are inserted at a single stage in the IC's life cycle. However, malicious conspiracy between multiple entities at different stages in the IC supply chain could make HT attacks more powerful. Ali *et al.* (121) described such an attack on an AES implementation. They show that such a multi-stage attack is significantly stronger than a HT attack by a single entity, both in the life cycle of ASICs and FPGAs. As a result, it would be very difficult for current defense approaches targeting individual stages to detect such a distributed attack (121). Detecting such a multi-stage Trojan is still an open problem. It is necessary to study universal Trojan detection methods independent of the stages.

F. Split manufacturing

Split manufacturing is a promising hardware security solution in the manufacturing stage where the untrusted foundries only knows part of the design information thus makes it difficult for them to insert HTs. In recent works, different split manufacturing methods are proposed, *e.g.*(44), or combined with other hardware security techniques, *e.g.*layout camouflaging (43).

G. Low overhead runtime HT monitoring techniques

A large number of Trojan detection techniques have been proposed, but it is still possible for Trojans to escape detection and activate when the chip is used in the field. Runtime HT monitoring technique (which is relatively less in existing works) is a necessary supplement to Trojan detection, and is also the last line of defense. However, existing runtime HT detection techniques suffer from high additional hardware overhead or high computational complexity. Low overhead runtime HT monitoring technique is a promising research direction, *e.g.*(45; 46).

H. Logic obfuscation for HT prevention

Logic obfuscation (logic encryption, lock locking) is a widely studied hardware security technique, which is usually used to prevent IC piracy, and IC overbuilding. It can also be used as a DFS method to prevent HT insertion. Chakraborty and Bhunia (122) propose a key based obfuscation scheme to prevent HT attacks, in which two functional modes are introduced, *i.e.*, obfuscated and normal modes. A large number of states have also been added to the obfuscated mode for obfuscation. This method prevents the attacker from finding the real rare states in the circuit (122). Dupuis *et al.* (123) propose a logic encryption approach to prevent HT insertion by minimizing the number of rare events in a circuit. Similarly, Rathor *et al.* (124) propose a logic encryption method using

key-gate topologies to remove rare-triggered nets to thwart HT. Frey and Yu (125) propose an approach using state obfuscation for HT detection. Illegal states caused by wrong keys are examined to detect HTs. They indicate that an attacker without the correct key cannot successfully modify the design without being noticed (125). Yu *et al.* (126) review the works on logic obfuscation for HT prevention and detection. They indicate that logic obfuscation can make it difficult for attackers to understand or reverse engineer the design thus can hinder the implantation of Trojans, or can facilitate the HT detection after manufacturing. Similar to ASICs, obfuscation can also be used to protect FPGA designs. Hoque *et al.* (127) propose an obfuscation based approach against bitstream modification attacks on FPGAs. Particularly, the critical functions in an FPGA design are identified and masked (obfuscated). Besides, they use a redundancy technique for obfuscation to thwart tampering (127). Potential future directions on logic encryption for HT prevention include expanding logic obfuscation from chip level to system level, and the key management in key-based obfuscation schemes (126).

Although logic obfuscation is usually used as a DFS method to prevent Trojans, the opposite application is also possible (128; 129). Vijayakumar *et al.* (128) indicated that physical design obfuscation can also be used to insert parametric Trojans. Such an example is demonstrated by Becker *et al.* (7), where the dopant polarities of transistors are changed to insert HT while making the HT difficult to be detected.

I. FPGA Trojan attacks and defenses

Compared with ASIC HTs, the works on FPGA Trojans are relatively less, both on attacks and defenses. The research of FPGA Trojan is not systematic and comprehensive at present. With the widespread use of FPGAs, the FPGA Trojan research is a valuable research direction, *e.g.*(90; 91; 103; 98).

XI. CONCLUSION

HT is an emerging threat to hardware security and information security. In the last decade, a large amount of HT detection techniques have been proposed. However, much less researches have been conducted into the design and implementation of HTs. In this paper, we provide a review of the development of HT implementations in the last decade and also make an outlook. Unlike all previous surveys or most HT works that focus on Trojan detection from the defender's perspective, for the first time, we study the Trojans from an attacker's perspective, focusing on the attacker's methods, capabilities, evading detection techniques, and challenges. We conclude that HT implantation or HT-related attacks can be launched at any stages of the IC supply chain, including the testing stage and the distribution stage that were rarely discussed in previous works. There are significant differences in the capabilities of attackers at each stage, which can be roughly divided into three levels: level 1, *i.e.*, in-house design team attackers and 3PIP vendor attackers; level 2, *i.e.*, CAD tools attackers, fabrication stage attackers, and testing stage attackers; level 3, *i.e.*, distribution stage attackers. Similar to the ASIC scenario, FPGA Trojan attacks are also feasible

at all stages of the FPGA supply chain. Some potential future directions on HT implementation and defense have emerged, which are tit-for-tat endless games. This paper would hopefully help defenders better understand the Trojan insertion so as to design reliable defense techniques, and better protect the circuits against HT attacks.

XII. ACKNOWLEDGMENTS

This work is supported by the National Natural Science Foundation of China (No. 61602241), and the Engineering and Physical Sciences Research Council (EPSRC) (EP/N508664/-CSIT2).

REFERENCES

- [1] S. Bhunia, M. S. Hsiao, M. Banga, and S. Narasimhan, "Hardware Trojan attacks: Threat analysis and countermeasures," *Proc. IEEE*, vol. 102, no. 8, pp. 1229–1247, 2014.
- [2] M. Tehranipoor and F. Koushanfar, "A survey of hardware Trojan taxonomy and detection," *IEEE Des. Test Comput.*, vol. 27, no. 1, pp. 10–25, Jan. 2010.
- [3] R. S. Chakraborty, S. Narasimhan, and S. Bhunia, "Hardware Trojan: Threats and emerging solutions," in *Proc. IEEE Int. High Level Design Validation and Test Workshop, San Francisco, USA*. IEEE, November 2009, pp. 166–171.
- [4] T. F. Wu, K. Ganesan, Y. A. Hu, H. P. Wong, S. Wong, and S. Mitra, "TPAD: Hardware Trojan prevention and detection for trusted integrated circuits," *IEEE Trans. Comput-Aided Des. Integr. Circuits Syst.*, vol. 35, no. 4, pp. 521–534, 2016.
- [5] D. Agrawal, S. Baktir, D. Karakoyunlu, P. Rohatgi, and B. Sunar, "Trojan detection using IC fingerprinting," in *Proc. IEEE Symp. on Security and Privacy, Oakland, USA*, May 2007, pp. 296–310.
- [6] S. Ghandali, G. T. Becker, D. Holcomb, and C. Paar, "A design methodology for stealthy parametric Trojans and its application to bug attacks," in *Int. Conf. on Cryptographic Hardware and Embedded Systems, Santa Barbara, USA*. Springer, August 2016, pp. 625–647.
- [7] G. T. Becker, F. Regazzoni, C. Paar, and W. P. Burleson, "Stealthy dopant-level hardware Trojans," in *Int. Workshop on Cryptographic Hardware and Embedded Systems, Santa Barbara, USA*. Springer, August 2013, pp. 197–214.
- [8] S. Adee, "The hunt for the kill switch," *IEEE Spectr.*, vol. 45, no. 5, pp. 34–39, 2008.
- [9] K. Yang, M. Hicks, Q. Dong, T. M. Austin, and D. Sylvester, "A2: Analog malicious hardware," in *Proc. IEEE Symp. on Security and Privacy, San Jose, USA*, May 2016, pp. 18–37.
- [10] "The intel management engine: An attack on computer users' freedom," <https://www.fsf.org/blogs/sysadmin/the-management-engine-an-attack-on-computer-users-freedom>, 2018.
- [11] "Intel x86s hide another CPU that can take over your machine (you can't audit it)," <https://boingboing.net/2016/06/15/intel-x86-processors-ship-with.html>, 2016.

- [12] M. Rostami, F. Koushanfar, and R. Karri, "A primer on hardware security: Models, methods, and metrics," *Proc. IEEE*, vol. 102, no. 8, pp. 1283–1295, 2014.
- [13] N. Jacob, D. Merli, J. Heyszl, and G. Sigl, "Hardware Trojans: Current challenges and approaches," *IET Comput. Digit. Tech.*, vol. 8, no. 6, pp. 264–273, 2014.
- [14] R. Karri, J. Rajendran, K. Rosenfeld, and M. Tehranipoor, "Trustworthy hardware: Identifying and classifying hardware Trojans," *IEEE Computer*, vol. 43, no. 10, pp. 39–46, 2010.
- [15] B. Shakya, T. He, H. Salmani, D. Forte, S. Bhunia, and M. Tehranipoor, "Benchmarking of hardware Trojans and maliciously affected circuits," *J. Hardw. Syst. Secur.*, vol. 1, no. 1, pp. 85–102, 2017.
- [16] A. Baumgarten, M. Steffen, M. Clausman, and J. Zambreno, "A case study in hardware Trojan design and implementation," *Int. J. Inf. Secur.*, vol. 10, no. 1, pp. 1–14, 2011.
- [17] "Trust-hub," <http://www.trust-hub.org/>, 2019.
- [18] Y. Jin, N. Kupp, and Y. Makris, "Experiences in hardware Trojan design and implementation," in *Proc. IEEE Int. Workshop on Hardware-Oriented Security and Trust, San Francisco, USA*, July 2009, pp. 50–57.
- [19] X. Wang, H. Salmani, M. Tehranipoor, and J. F. Plusquellic, "Hardware Trojan detection and isolation using current integration and localized current analysis," in *Proc. IEEE Int. Symp. on Defect and Fault Tolerance of VLSI Systems, Boston, USA*. IEEE, October 2008, pp. 87–95.
- [20] J. Li and J. Lach, "At-speed delay characterization for IC authentication and Trojan horse detection," in *Proc. IEEE Int. Workshop on Hardware-Oriented Security and Trust, Anaheim, USA*, June 2008, pp. 8–14.
- [21] K. Hu, A. N. Nowroz, S. Reda, and F. Koushanfar, "High-sensitivity hardware Trojan detection using multimodal characterization," in *Proc. Conf. on Design, Automation and Test in Europe, Grenoble, France*, March 2013, pp. 1271–1276.
- [22] K. Xiao, X. Zhang, and M. Tehranipoor, "A clock sweeping technique for detecting hardware Trojans impacting circuits delay," *IEEE Des. Test*, vol. 30, no. 2, pp. 26–34, 2013.
- [23] S. Narasimhan, D. Du, R. S. Chakraborty, S. Paul, F. G. Wolff, C. A. Papachristou, K. Roy, and S. Bhunia, "Hardware Trojan detection by multiple-parameter side-channel analysis," *IEEE Trans. Comput.*, vol. 62, no. 11, pp. 2183–2195, 2013.
- [24] A. N. Nowroz, K. Hu, F. Koushanfar, and S. Reda, "Novel techniques for high-sensitivity hardware Trojan detection using thermal and power maps," *IEEE Trans. Comput-Aided Des. Integr. Circuits Syst.*, vol. 33, no. 12, pp. 1792–1805, 2014.
- [25] M. Xue, W. Liu, A. Hu, and Y. Wang, "Detecting hardware Trojan through time domain constrained estimator based unified subspace technique," *IEICE Trans. Inf. Syst.*, vol. 97-D, no. 3, pp. 606–609, 2014.
- [26] M. Xue, A. Hu, and G. Li, "Detecting hardware Trojan through heuristic partition and activity driven test pattern generation," in *Proc. Communications Security Conf., Beijing, China*, May 2014, pp. 1–6.
- [27] D. Fujimoto, S. Nin, Y.-I. Hayashi, N. Miura, M. Nagata, and T. Matsumoto, "A demonstration of a HT-detection method based on impedance measurements of the wiring around ICs," *IEEE Trans. Circuits Syst. II-Express Briefs*, vol. 65, no. 10, pp. 1320–1324, 2018.
- [28] L. N. Nguyen, C.-L. Cheng, M. Prvulovic, and A. Zajić, "Creating a backscattering side channel to enable detection of dormant hardware Trojans," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 27, no. 7, pp. 1561–1574, 2019.
- [29] K. S. Subraman, A. Antonopoulos, A. A. Abotabl, A. Nosratinia, and Y. Makris, "Demonstrating and mitigating the risk of an FEC-based hardware Trojan in wireless networks," *IEEE Trans. Inf. Forensic Secur.*, vol. 14, no. 10, pp. 2720–2734, 2019.
- [30] R. S. Chakraborty, F. Wolff, S. Paul, C. Papachristou, and S. Bhunia, "MERO: A statistical approach for hardware Trojan detection," in *Int. Workshop on Cryptographic Hardware and Embedded Systems, Lausanne, Switzerland*. Springer Berlin Heidelberg, September 2009, pp. 396–410.
- [31] Y. Huang, S. Bhunia, and P. Mishra, "Scalable test generation for Trojan detection using side channel analysis," *IEEE Trans. Inf. Forensic Secur.*, vol. 13, no. 11, pp. 2746–2760, 2018.
- [32] M. Banga and M. S. Hsiao, "Trusted RTL: Trojan detection methodology in pre-silicon designs," in *Proc. IEEE Int. Symp. on Hardware-Oriented Security and Trust, Anaheim, USA*, June 2010, pp. 56–59.
- [33] A. Waksman, M. Suozzo, and S. Sethumadhavan, "FANCI: Identification of stealthy malicious logic using boolean functional analysis," in *Proc. ACM SIGSAC Conf. on Computer and Communications Security, Berlin, Germany*, November 2013, pp. 697–708.
- [34] N. G. Tsoutsos, C. Konstantinou, and M. Maniatakos, "Advanced techniques for designing stealthy hardware Trojans," in *Proc. ACM Annual Design Automation Conf., San Francisco, USA*, June 2014, pp. 1–4.
- [35] X. Zhang and M. Tehranipoor, "RON: An on-chip ring oscillator network for hardware Trojan detection," in *Proc. Conf. on Design, Automation and Test in Europe, Grenoble, France*, March 2011, pp. 1–6.
- [36] J. Rajendran, V. Jyothi, O. Sinanoglu, and R. Karri, "Design and analysis of ring oscillator based design-for-trust technique," in *Proc. IEEE VLSI Test Symp., Dana Point, USA*, May 2011, pp. 105–110.
- [37] M. Hicks, M. Finnicum, S. T. King, M. M. K. Martin, and J. M. Smith, "Overcoming an untrusted computing base: Detecting and removing malicious hardware automatically," in *Proc. IEEE Symp. on Security and Privacy, Oakland, USA*, May 2010, pp. 159–172.
- [38] J. Zhang and Q. Xu, "On hardware Trojan design and implementation at register-transfer level," in *Proc. IEEE Int. Symp. on Hardware-Oriented Security and Trust, Austin, USA*, June 2013, pp. 107–112.
- [39] N. Fern, I. San, C. K. Koç, and K.-T. T. Cheng, "Hiding

- hardware Trojan communication channels in partially specified SoC bus functionality,” *IEEE Trans. Comput-Aided Des. Integr. Circuits Syst.*, vol. 36, no. 9, pp. 1435–1444, 2017.
- [40] A. Nahiyan, K. Xiao, K. Yang, Y. Jin, D. Forte, and M. Tehranipoor, “AVFSM: A framework for identifying and mitigating vulnerabilities in FSMs,” in *Proc. 53rd Annu. Des. Autom. Conf. (DAC), Austin, USA*, June 2016, pp. 1–6.
- [41] K. Xiao, D. Forte, and M. Tehranipoor, “A novel built-in self-authentication technique to prevent inserting hardware Trojans,” *IEEE Trans. Comput-Aided Des. Integr. Circuits Syst.*, vol. 33, no. 12, pp. 1778–1791, 2014.
- [42] Q. Shi, M. M. Tehranipoor, and D. Forte, “Obfuscated built-in self-authentication with secure and efficient wire-lifting,” *IEEE Trans. Comput-Aided Des. Integr. Circuits Syst.*, vol. 38, no. 11, pp. 1981–1994, 2018.
- [43] S. Patnaik, M. Ashraf, O. Sinanoglu, and J. Knechtel, “A modern approach to IP protection and Trojan prevention: Split manufacturing for 3D ICs and obfuscation of vertical interconnects,” *IEEE Trans. Emerg. Top. Comput.*, pp. 1–18, 2019.
- [44] M. Li, B. Yu, Y. Lin, X. Xu, W. Li, and D. Z. Pan, “A practical split manufacturing framework for Trojan prevention via simultaneous wire lifting and cell insertion,” *IEEE Trans. Comput-Aided Des. Integr. Circuits Syst.*, vol. 38, no. 9, pp. 1585–1598, 2018.
- [45] H. Zhao, L. Kwiat, K. A. Kwiat, C. A. Kamhoua, and L. Njilla, “Applying chaos theory for runtime hardware Trojan monitoring and detection,” *IEEE Trans. Dependable Secur. Comput.*, pp. 1–14, 2018.
- [46] R. Elnaggar, K. Chakrabarty, and M. B. Tahoori, “Hardware Trojan detection using changepoint-based anomaly detection techniques,” *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 27, no. 12, pp. 2706–2719, 2019.
- [47] T. Kumaki, M. Yoshikawa, and T. Fujino, “Cipher-destroying and secret-key-emitting hardware Trojan against AES core,” in *Proc. IEEE Int. Midwest Symp. on Circuits and Systems, Columbus, USA*. IEEE, August 2013, pp. 408–411.
- [48] S. T. King, J. Tucek, A. Cozzie, C. Grier, W. Jiang, and Y. Zhou, “Designing and implementing malicious hardware,” in *Proc. USENIX Workshop on Large-Scale Exploits and Emergent Threats, San Francisco, USA*. USENIX Association, April 2008, pp. 1–8.
- [49] J. C. M. Santos and Y. Fei, “Designing and implementing a malicious 8051 processor,” in *Proc. IEEE Int. Symp. on Defect and Fault Tolerance in VLSI and Nanotechnology Systems, Austin, USA*, October 2012, pp. 63–66.
- [50] T. Reece, D. B. Limbrick, X. Wang, B. T. Kiddie, and W. H. Robinson, “Stealth assessment of hardware Trojans in a microcontroller,” in *Proc. IEEE Int. Conf. on Computer Design, Montreal, Canada*, September 2012, pp. 139–142.
- [51] R. Karri, J. Rajendran, and K. Rosenfeld, “Trojan taxonomy,” in *Introduction to Hardware Security and Trust*, M. Tehranipoor and C. Wang, Eds. Springer, USA, 2012, pp. 325–338.
- [52] P. Dash, C. Perkins, and R. M. Gerdes, “Remote activation of hardware Trojans via a covert temperature channel,” in *Int. Conf. on Security and Privacy in Communication Systems, Dallas, USA*. Springer, October 2015, pp. 294–310.
- [53] J.-F. Gallais, J. Großschädl, N. Hanley, M. Kasper, M. Medwed, and F. Regazzoni, “Hardware Trojans for inducing or amplifying side-channel leakage of cryptographic software,” in *Int. Conf. on Trusted Systems, Beijing, China*. Springer, December 2011, pp. 253–270.
- [54] T. Reece and W. H. Robinson, “Analysis of data-leak hardware Trojans in AES cryptographic circuits,” in *Proc. IEEE Int. Conf. on Technologies for Homeland Security, Boston, USA*, November 2013, pp. 467–472.
- [55] C. Sturton, M. Hicks, D. Wagner, and S. T. King, “Defeating UCI: Building stealthy and malicious hardware,” in *Proc. IEEE Symp. on Security and Privacy, Berkeley, USA*, May 2011, pp. 64–77.
- [56] Y. Liu, Y. Jin, A. Nosratinia, and Y. Makris, “Silicon demonstration of hardware Trojan design and detection in wireless cryptographic ICs,” *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 25, no. 4, pp. 1506–1519, 2017.
- [57] L. Lin, W. Burleson, and C. Paar, “MOLES: Malicious off-chip leakage enabled by side-channels,” in *Proc. Int. Conf. Computer-Aided Design, San Jose, USA*, November 2009, pp. 117–122.
- [58] M. Fyrbiak, S. Wallat, P. Swierczynski, M. Hoffmann, S. Hoppach, M. Wilhelm, T. Weidlich, R. Tessier, and C. Paar, “HAL-The missing piece of the puzzle for hardware reverse engineering, Trojan detection and insertion,” *IEEE Trans. Dependable Secur. Comput.*, vol. 16, no. 3, pp. 498–510, 2018.
- [59] S. Kaji, M. Kinugawa, D. Fujimoto, and Y.-i. Hayashi, “Data injection attack against electronic devices with locally weakened immunity using a hardware Trojan,” *IEEE Trans. on Electromagn. Compat.*, vol. 61, no. 4, pp. 1115–1121, 2018.
- [60] X. Wang, S. Narasimhan, A. Krishna, T. Mal-Sarkar, and S. Bhunia, “Sequential hardware Trojan: Side-channel aware design and placement,” in *Proc. IEEE Int. Conf. on Computer Design, Amherst, USA*, October 2011, pp. 297–300.
- [61] J. J. V. Rajendran, O. Sinanoglu, and R. Karri, “Building trustworthy systems using untrusted components: A high-level synthesis approach,” *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 24, no. 9, pp. 2946–2959, 2016.
- [62] X. Zhang and M. Tehranipoor, “Case study: Detecting hardware Trojans in third-party digital IP cores,” in *Proc. IEEE Int. Symp. on Hardware-Oriented Security and Trust, San Diego, USA*, June 2011, pp. 67–70.
- [63] J.-Y. Jou and C.-N. J. Liu, “Coverage analysis techniques for HDL design validation,” *Proc. Asia Pac. Chip Des. Lang.*, pp. 48–55, 1999.

- [64] H. Salmani and M. Tehranipoor, "Analyzing circuit vulnerability to hardware Trojan insertion at the behavioral level," in *Proc. IEEE Int. Symp. on Defect and Fault Tolerance in VLSI and Nanotechnology Systems, New York, USA*, October 2013, pp. 190–195.
- [65] J. Zhang, F. Yuan, and Q. Xu, "DeTrust: Defeating hardware trust verification with stealthy implicitly-triggered hardware Trojans," in *Proc. ACM SIGSAC Conf. on Computer and Communications Security, Scottsdale, USA*, November 2014, pp. 153–166.
- [66] E. Love, Y. Jin, and Y. Makris, "Proof-carrying hardware intellectual property: A pathway to trusted module acquisition," *IEEE Trans. Inf. Forensic Secur.*, vol. 7, no. 1, pp. 25–40, 2012.
- [67] Y. Jin and Y. Makris, "A proof-carrying based framework for trusted microprocessor IP," in *Proc. Int. Conf. on Computer-Aided Design, San Jose, USA*, November 2013, pp. 824–829.
- [68] N. Veeranna and B. C. Schäfer, "Hardware Trojan detection in behavioral intellectual properties (IP's) using property checking techniques," *IEEE Trans. Emerg. Top. Comput.*, vol. 5, no. 4, pp. 576–585, 2017.
- [69] J. Rajendran, A. M. Dhandayuthapany, V. Vedula, and R. Karri, "Formal security verification of third party intellectual property cores for information leakage," in *Proc. Int. Conf. on VLSI Design, Kolkata, India*, January 2016, pp. 547–552.
- [70] C. Liu, J. Rajendran, C. Yang, and R. Karri, "Shielding heterogeneous MPSoCs from untrustworthy 3PIPs through security-driven task scheduling," *IEEE Trans. Emerg. Top. Comput.*, vol. 2, no. 4, pp. 461–472, 2014.
- [71] "The coq proof assistant," <https://coq.inria.fr/>, 2019.
- [72] S. Narasimhan, R. S. Chakraborty, and S. Chakraborty, "Hardware IP protection during evaluation using embedded sequential Trojan," *IEEE Des. Test Comput.*, vol. 29, no. 3, pp. 70–79, 2012.
- [73] C. Pilato, K. Basu, M. Shayan, F. Regazzoni, and R. Karri, "High-level synthesis of benevolent Trojans," in *Proc. Conf. on Design, Automation and Test in Europe Conf. and Exhibition, Florence, Italy*, March 2019, pp. 1124–1129.
- [74] C. Pilato, K. Basu, F. Regazzoni, and R. Karri, "Black-hat high-level synthesis: Myth or reality?" *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 27, no. 4, pp. 913–926, 2019.
- [75] K. Basu, S. M. Saeed, C. Pilato, M. Ashraf, K. Chakraborty, and R. Karri, "CAD-base: An attack vector into the electronics supply chain," *ACM Transact. Des. Automat. Electron. Syst.*, vol. 24, no. 4, pp. 38:1–38:30, 2019.
- [76] C. Bao, D. Forte, and A. Srivastava, "On reverse engineering-based hardware Trojan detection," *IEEE Trans. Comput-Aided Des. Integr. Circuits Syst.*, vol. 35, no. 1, pp. 49–57, 2016.
- [77] M. Xue, J. Wang, and A. Hu, "An enhanced classification-based golden chips-free hardware Trojan detection technique," in *Proc. IEEE Asian Hardware-Oriented Security and Trust, Yilan, Taiwan*, December 2016, pp. 1–6.
- [78] S. Bhasin, J.-L. Danger, S. Guilley, X. T. Ngo, and L. Sauvage, "Hardware Trojan horses in cryptographic IP cores," in *Workshop on Fault Diagnosis and Tolerance in Cryptography, Alamitos, USA*. IEEE, August 2013, pp. 15–29.
- [79] R. Kumar, P. Jovanovic, W. Burleson, and I. Polian, "Parametric Trojans for fault-injection attacks on cryptographic hardware," in *Workshop on Fault Diagnosis and Tolerance in Cryptography, Busan, South Korea*, September 2014, pp. 18–28.
- [80] Y. Shiyanovskii, F. Wolff, A. Rajendran, C. A. Papachristou, D. J. Weyer, and W. Clay, "Process reliability based Trojans through NBTI and HCI effects," in *Proc. Conf. on Adaptive Hardware and Systems, Anaheim, California*, June 2010, pp. 215–222.
- [81] L. Lin, M. Kasper, T. Güneysu, C. Paar, and W. Burleson, "Trojan side-channels: Lightweight hardware Trojans through side-channel engineering," in *Int. Workshop on Cryptographic Hardware and Embedded Systems, Lausanne, Switzerland*. Springer, September 2009, pp. 382–395.
- [82] M. Xue, R. Bian, W. Liu, and J. Wang, "Defeating untrustworthy testing parties: A novel hybrid clustering ensemble based golden models-free hardware Trojan detection method," *IEEE Access*, vol. 7, pp. 5124–5140, 2019.
- [83] M. Yasin, O. Sinanoglu, and J. Rajendran, "Testing the trustworthiness of IC testing: An oracle-less attack on IC camouflaging," *IEEE Trans. Inf. Forensic Secur.*, vol. 12, no. 11, pp. 2668–2682, 2017.
- [84] R. Bian, M. Xue, and J. Wang, "Building trusted golden models-free hardware Trojan detection framework against untrustworthy testing parties using a novel clustering ensemble technique," in *Proc. IEEE Int. Conf. on Trust, Security and Privacy in Computing and Communications, New York, USA*, July 2018, pp. 1458–1463.
- [85] G. K. Contreras, M. T. Rahman, and M. Tehranipoor, "Secure split-test for preventing IC piracy by untrusted foundry and assembly," in *IEEE Int. Symp. on Defect and Fault Tolerance in VLSI and Nanotechnology Systems, New York, USA*. IEEE, October 2013, pp. 196–203.
- [86] M. T. Rahman, D. Forte, Q. Shi, G. K. Contreras, and M. M. Tehranipoor, "CSST: Preventing distribution of unlicensed and rejected ICs by untrusted foundry and assembly," in *IEEE Int. Symp. on Defect and Fault Tolerance in VLSI Nanotechnology Systems, Amsterdam, The Netherlands*. IEEE, October 2014, pp. 46–51.
- [87] D. Zhang, X. Wang, M. T. Rahman, and M. Tehranipoor, "An on-chip dynamically obfuscated wrapper for protecting supply chain against IP and IC piracies," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 26, no. 11, pp. 2456–2469, 2018.
- [88] P. Swierczynski, M. Fyrbiak, P. Koppe, A. Moradi, and C. Paar, "Interdiction in practice-hardware Trojan against a high-security USB flash drive," *J. Cryptogr.*

- Eng.*, vol. 7, no. 3, pp. 199–211, 2017.
- [89] R. S. Chakraborty, I. Saha, A. Palchoudhuri, and G. K. Naik, “Hardware Trojan insertion by direct modification of FPGA configuration bitstream,” *IEEE Des. Test*, vol. 30, no. 2, pp. 45–54, 2013.
- [90] P. Swierczynski, M. Fyrbiak, P. Koppe, and C. Paar, “FPGA Trojans through detecting and weakening of cryptographic primitives,” *IEEE Trans. Comput-Aided Des. Integr. Circuits Syst.*, vol. 34, no. 8, pp. 1236–1249, 2015.
- [91] C. Krieg, C. Wolf, and A. Jantsch, “Malicious LUT: A stealthy FPGA Trojan injected and triggered by the design flow,” in *Proc. 35th Int. Conf. on Computer-Aided Design, Austin, USA*, November 2016, pp. 1–8.
- [92] C. Marchand and J. Francq, “Low-level implementation and side-channel detection of stealthy hardware Trojans on field programmable gate arrays,” *IET Comput. Digit. Tech.*, vol. 8, no. 6, pp. 246–255, 2014.
- [93] J. X. Zheng, E. Chen, and M. Potkonjak, “A benign hardware Trojan on FPGA-based embedded systems,” in *22nd Int. Conf. on Field Programmable Logic and Applications, Oslo, Norway*, August 2012, pp. 464–470.
- [94] C. Krieg, C. Wolf, A. Jantsch, and T. Zseby, “Toggle MUX: How X-optimism can lead to malicious hardware,” in *Proc. 54th Annu. Des. Autom. Conf. (DAC), Austin, USA*, June 2017, pp. 1–6.
- [95] V. Jyothi and J. J. V. Rajendran, “Hardware Trojan attacks in FPGA and protection approaches,” in *The Hardware Trojan War*, S. Bhunia and M. Tehranipoor, Eds. Springer, Switzerland, 2018, pp. 345–368.
- [96] S. Mal-Sarkar, A. Krishna, A. Ghosh, and S. Bhunia, “Hardware Trojan attacks in FPGA devices: Threat analysis and effective countermeasures,” in *Proc. Great Lakes Symp. on VLSI, Houston, USA*, May 2014, pp. 287–292.
- [97] O. Söll, T. Korak, M. Muehlberghuber, and M. Hutter, “EM-based detection of hardware Trojans on FPGAs,” in *Proc. IEEE Int. Symp. on Hardware-Oriented Security and Trust, Arlington, USA*, May 2014, pp. 84–87.
- [98] Z. Chen, S. Guo, J. Wang, Y. Li, and Z. Lu, “Toward FPGA security in IoT: A new detection technique for hardware Trojans,” *IEEE Internet Things J.*, vol. 6, no. 4, pp. 7061–7068, 2019.
- [99] Y. Pino, V. Jyothi, and M. French, “Intra-die process variation aware anomaly detection in FPGAs,” in *IEEE Int. Test Conf., Seattle, USA*. IEEE, October 2014, pp. 1–6.
- [100] V. Jyothi, M. Thoonoli, R. Stern, and R. Karri, “FPGA trust zone: Incorporating trust and reliability into FPGA designs,” in *Proc. IEEE Int. Conf. on Computer Design, Scottsdale, USA*. IEEE, October 2016, pp. 600–605.
- [101] P. Swierczynski, M. Fyrbiak, C. Paar, C. Huriaux, and R. Tessier, “Protecting against cryptographic Trojans in FPGAs,” in *IEEE Annual Int. Symp. on Field-Programmable Custom Computing Machines, Vancouver, Canada*, May 2015, pp. 151–154.
- [102] G. Bloom, B. Narahari, R. Simha, A. Namazi, and R. Levy, “FPGA SoC architecture and runtime to prevent hardware Trojans from leaking secrets,” in *Proc. IEEE Int. Symp. on Hardware-Oriented Security and Trust, Washington, USA*, May 2015, pp. 48–51.
- [103] Z. Zhang, L. Njilla, C. A. Kamhoua, and Q. Yu, “Thwarting security threats from malicious FPGA tools with novel FPGA-oriented moving target defense,” *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 27, no. 3, pp. 665–678, 2018.
- [104] J. Cruz, Y. Huang, P. Mishra, and S. Bhunia, “An automated configurable Trojan insertion framework for dynamic trust benchmarks,” in *Proc. Conf. on Design, Automation and Test in Europe Conf. and Exhibition, Dresden, Germany*, March 2018, pp. 1598–1603.
- [105] A. A. Nasr and M. Z. Abdulmageed, “Automatic feature selection of hardware layout: A step toward robust hardware Trojan detection,” *J. Electron. Test.*, vol. 32, no. 3, pp. 357–367, 2016.
- [106] C. Bao, D. Forte, and A. Srivastava, “On application of one-class SVM to reverse engineering-based hardware Trojan detection,” in *Int. Symp. on Quality Electronic Design, Santa Clara, USA*, March 2014, pp. 47–54.
- [107] A. Kulkarni, Y. Pino, and T. Mohsenin, “SVM-based real-time hardware Trojan detection for many-core platform,” in *Int. Symp. on Quality Electronic Design, Santa Clara, USA*, March 2016, pp. 362–367.
- [108] —, “Adaptive real-time Trojan detection framework through machine learning,” in *Proc. IEEE Int. Symp. on Hardware-Oriented Security and Trust*, May 2016, pp. 120–123.
- [109] R. Bian, M. Xue, and J. Wang, “A novel golden models-free hardware Trojan detection technique using unsupervised clustering analysis,” in *Proc. Int. Conf. on Cloud Computing and Security, Haikou, China*, June 2018, pp. 634–646.
- [110] M. Xue, R. Bian, J. Wang, and W. Liu, “A co-training based hardware Trojan detection technique by exploiting unlabeled ICs and inaccurate simulation models,” in *Proc. IEEE Int. Conf. on Trust, Security and Privacy in Computing and Communications, New York, USA*, August 2018, pp. 1452–1457.
- [111] —, “Building an accurate hardware Trojan detection technique from inaccurate simulation models and unlabeled ICs,” *IET Comput. Digit. Tech.*, vol. 13, no. 4, pp. 348–359, 2019.
- [112] K. Hasegawa, M. Oya, M. Yanagisawa, and N. Togawa, “Hardware Trojans classification for gate-level netlists based on machine learning,” in *Proc. IEEE Int. Symp. on On-Line Testing and Robust System Design, Sant Feliu de Guixols, Spain*, July 2016, pp. 203–206.
- [113] R. Elnaggar and K. Chakraborty, “Machine learning for hardware security: Opportunities and risks,” *J. Electron. Test.*, vol. 34, no. 2, pp. 183–201, 2018.
- [114] J. Clements and Y. Lao, “Hardware Trojan design on neural networks,” in *IEEE Int. Symp. Circuits Syst., Sapporo, Japan*, May 2019, pp. 1–5.
- [115] J. Ye, Y. Hu, and X. Li, “Hardware Trojan in FPGA CNN accelerator,” in *IEEE 27th Asian Test Symp.*, Hefei, China, October 2018, pp. 68–73.

- [116] T. A. Odetola, H. R. Mohammed, and S. R. Hasan. (arXiv:1911.00783, 2019.) A stealthy hardware Trojan exploiting the architectural vulnerability of deep learning architectures: Input interception attack (IIA).
- [117] W. Li, J. Yu, X. Ning, P. Wang, Q. Wei, Y. Wang, and H. Yang, “Hu-Fu: Hardware and software collaborative attack framework against neural networks,” in *IEEE Comput. Soc. Annu. Symp. Very Large Scale Integr.*, Hong Kong, China, July 2018, pp. 482–487.
- [118] X. Hu, Y. Zhao, L. Deng, L. Liang, P. Zuo, J. Ye, Y. Lin, and Y. Xie, “Practical attacks on deep neural networks by memory Trojanning,” *IEEE Trans. Comput-Aided Des. Integr. Circuits Syst., Early Access*, pp. 1–14, 2020.
- [119] G. T. Becker, M. Kasper, A. Moradi, and C. Paar, “Side-channel based watermarks for integrated circuits,” in *Proc. IEEE Int. Symp. on Hardware-Oriented Security and Trust, Anaheim, USA*, June 2010, pp. 30–35.
- [120] S. Bhunia and M. Tehranipoor, *The hardware Trojan war: Attacks, myths, and defenses*. Springer, Switzerland, 2017.
- [121] S. S. Ali, R. S. Chakraborty, D. Mukhopadhyay, and S. Bhunia, “Multi-level attacks: An emerging security concern for cryptographic hardware,” in *Proc. Conf. on Design, Automation and Test in Europe, Grenoble, France*, March 2011, pp. 1–4.
- [122] R. S. Chakraborty and S. Bhunia, “Security against hardware Trojan attacks using key-based design obfuscation,” *J. Electron. Test.*, vol. 27, no. 6, pp. 767–785, 2011.
- [123] S. Dupuis, P. Ba, G. D. Natale, M. Flottes, and B. Rouzeyre, “A novel hardware logic encryption technique for thwarting illegal overproduction and hardware Trojans,” in *Proc. IEEE 20th Int. On-Line Test. Symp.*, Girona, Spain, July 2014, pp. 49–54.
- [124] V. S. Rathor, B. Garg, and G. K. Sharma, “A novel low complexity logic encryption technique for design-for-trust,” *IEEE Trans. Emerg. Top. Comput., Early Access*, pp. 1–12, 2018.
- [125] J. Frey and Q. Yu, “Exploiting state obfuscation to detect hardware Trojans in NoC network interfaces,” in *Proc. IEEE 58th Int. Midwest Symp. Circuits Syst.*, Fort Collins, USA, August 2015, pp. 1–4.
- [126] Q. Yu, J. Dofe, and Z. Zhang, “Exploiting hardware obfuscation methods to prevent and detect hardware Trojans,” in *Proc. IEEE 60th Int. Midwest Symp. Circuits Syst.*, Boston, USA, August 2017, pp. 819–822.
- [127] T. Hoque, K. Yang, R. Karam, S. Tajik, D. Forte, M. Tehranipoor, and S. Bhunia, “Hidden in plaintext: An obfuscation-based countermeasure against FPGA bitstream tampering attacks,” *ACM Trans. Design Autom. Electr. Syst.*, vol. 25, no. 1, pp. 1–32, 2020.
- [128] A. Vijayakumar, V. C. Patil, D. E. Holcomb, C. Paar, and S. Kundu, “Physical design obfuscation of hardware: A comprehensive investigation of device and logic-level techniques,” *IEEE Trans. Inf. Forensics Secur.*, vol. 12, no. 1, pp. 64–77, 2017.
- [129] G. T. Becker, M. Fyrbiak, and C. Kison, “Hardware obfuscation: Techniques and open challenges,” in *Foundations of Hardware IP Protection*, L. Bossuet and L. Torres, Eds. Springer, Cham, 2017, pp. 105–123.