

ISR TECHNICAL REPORT 2013-03

Temperature Tracking: An Innovative Run-Time Approach for Hardware Trojan Detection

Domenic Forte
Chongxi Bao
Ankur Srivastava

The
Institute for
Systems
Research



A. JAMES CLARK
SCHOOL OF ENGINEERING

ISR develops, applies and teaches advanced methodologies of design and analysis to solve complex, hierarchical, heterogeneous and dynamic problems of engineering technology and systems for industry and government.

ISR is a permanent institute of the University of Maryland, within the A. James Clark School of Engineering. It is a graduated National Science Foundation Engineering Research Center.

www.isr.umd.edu

Temperature Tracking: An Innovative Run-Time Approach for Hardware Trojan Detection

Domenic Forte, Chongxi Bao, Ankur Srivastava
ECE Dept., University of Maryland, College Park, USA
dforte@umd.edu, chongxi.bao@gmail.com, ankurs@umd.edu

Abstract—The hardware Trojan threat has motivated development of Trojan detection schemes at all stages of the integrated circuit (IC) lifecycle. While the majority of existing schemes focus on ICs at test-time, there are many unique advantages offered by post-deployment/run-time Trojan detection. However, run-time approaches have been underutilized with prior work highlighting the challenges of implementing them with limited hardware resources. In this paper, we propose innovative low-overhead approaches for run-time Trojan detection which exploit the thermal sensors already available in many modern systems to detect deviations in power/thermal profiles caused by Trojan activation. Simulation results using state-of-the-art tools on publicly available Trojan benchmarks verify that our approaches can detect active Trojans quickly and with few false positives.

I. INTRODUCTION

The emerging trend of outsourcing integrated circuit (IC) design and fabrication has created new opportunities called hardware Trojan attacks that can seriously jeopardize the integrity of an electronic system. Hardware Trojans can change hardware functionality, reduce system reliability, and leak valuable/sensitive information. Since this represents a dangerous threat to all sectors that rely on ICs (military, financial, power, etc.), robust methods to detect Trojans are becoming imperative. Unfortunately, Trojan detection happens to be a very challenging problem. Attackers have so many Trojans of varying types and sizes at their disposal [1]. Trojans are often triggered by extremely rare events. Finally, measurement noise and fabrication variation can mask Trojan presence [2], [3].

To overcome these challenges, detection approaches have been proposed at three stages of the IC lifecycle: design-time, test-time, and run-time. The majority of existing schemes occur at test-time where the I/O and side channel behavior (delay, power, etc.) of suspect ICs are compared to “golden” models/ICs. Design-time approaches [4], [5] have been utilized (mostly) to support test-time detection. Run-time approaches have been investigated significantly less, but possess several advantages over their counterparts. First, test-time approaches can fail for Trojans with small well-placed triggers (eg. one gate) that do not get activated at test-time [6]. Run-time approaches on the other hand can be utilized for the *entire lifetime of the IC*. Hence, any Trojan missed during test-time can be found if the Trojan ever gets activated (eg. [7], [8]). Second, an IC with an inactive Trojan performs essentially the same functions as a Trojan-free IC. Effective run-time detection allows one to still deploy a Trojan-inserted IC and then either disable the IC entirely or bypass the Trojan logic [9], [10] if the Trojan ever gets activated. The main drawback to run-time approaches has been their large overheads [2]. For example, the path delay characterization approach proposed in [7] suffers from considerable area overhead for modern designs with millions of paths [1].

Main Contributions. In this paper, we propose innovative low-overhead approaches for online Trojan detection that take advantage of the relationship between power and temperature. When a Trojan gets activated at run-time, it can have a significant impact on the system’s power consumption which will be reflected in the system’s thermal profile. Furthermore, while direct monitoring of side channels (current, delay)

typically requires significant overheads, thermal sensors and infrastructure are already available in many modern systems for dynamic thermal management. A summary of our main contributions is as follows:

- We propose a novel framework for temperature-based Trojan detection which consists of design-time, test-time, and run-time phases. In the design phase, we statistically characterize an IC’s power/thermal dynamics and optimally place thermal sensors. The test-time phase is used to calibrate each IC due to fabrication variation. The run-time phase integrates the information from the previous phases with thermal sensor measurements to detect Trojan activation.
- We propose two mechanisms to detect Trojan activation during run-time. The first is a local sensor-based approach that uses information from thermal sensors, statistical information provided by the test phase, and hypothesis testing. The second is a global approach that exploits correlation between sensors and maintains track of the IC’s thermal profile using a Kalman filter (KF).
- We test our detection mechanisms on five publicly available Trojan benchmarks (from [11]) and use state-of-the-art Cadence simulation tools to compute power/thermal profiles. In all but one benchmark, the proposed approaches are capable of detecting Trojan activation quickly (less than 210ms) and with very few false-positives.

Outline. In Section II, we review Trojans and Trojan detection. In Section III, we discuss further the motivation for our temperature-based detection. Our specific problem and its challenges are clearly defined in Section IV. The proposed framework and detection mechanisms are discussed in Section V. Experimental results are discussed in Section VI. We conclude and discuss future work in the last section.

II. BACKGROUND

Trojan Anatomy. Trojans consist of two components [1]: (i) The Trojan *trigger* waits for a special event, such as a rare external input pattern or internal logic state, and then activates the Trojan’s attack. Before the Trojan is triggered, the IC containing the Trojan functions mainly as intended (excluding the trigger’s activity); (ii) After the Trojan is triggered, the Trojan *payload* is activated and changes the IC’s functionality.

In this paper, we refer to ICs with and without Trojans as **Trojan-inserted** and **Trojan-free** respectively. Trojan-inserted ICs whose payloads have been triggered and not triggered are referred to as **Trojan-active** and **Trojan-inactive** respectively.

Trojan Detection Schemes. Detecting hardware Trojans is a very challenging problem. Deterministic and exhaustive validation are infeasible due to the variety of Trojans available to an attacker [2]. Furthermore, since Trojans are often triggered by rare events, conventional post-manufacturing tests which only target common and repeatable faults cannot be relied upon. Thus, researchers have had to develop new schemes to outwit attackers. Three main approaches have been proposed in the literature which are classified based on their place in the IC’s lifecycle: test-time, run-time, and design-time.

1) *Test-time*: These approaches consist of additional tests that take place after conventional post-manufacturing testing. There are two types. *Logic-based* schemes develop directed test patterns that activate Trojan payloads in order to detect errors in the output. *Side Channel-based* approaches measure physical parameters, such as power consumption and path delay, of suspect ICs, and compares them with expected parameters of a “golden” model or Trojan-free IC. Unlike logic-based, the Trojan payload need not necessarily be activated because the trigger alone may impact IC delay, power, etc.

2) *Run-time Monitoring*: These approaches are utilized after the IC has been deployed and have several unique advantages. For example, if the Trojan remains inactive (i.e. not triggered) for the IC’s entire lifetime, the IC will always perform its intended functions. Hence, a valid option would be to deploy the Trojan-inserted IC and monitor its behavior. If deviation from correct behavior is ever detected, the monitor can either disable the IC entirely or bypass the Trojan logic to maintain correct operation (eg. [9], [10]). Resource overheads have been the main disadvantage of run-time monitoring [2].

3) *Design-for-Trust (DFT)*: These are design-time strategies that aid test-time approaches. Examples include [4] and [5] which use scan flip-flops to increase the probability of Trojan activation and enhance side channel analysis.

Integrated Solution. [2] suggests an integrated approach that utilizes all of the above approaches to provide more comprehensive coverage of Trojans. Logic and side channel-based detection provide coverage of small and large Trojans respectively while run-time monitors act as “last lines of defense” for Trojan-inserted ICs that bypass test-time schemes.

III. MOTIVATION

In this paper, we investigate run-time approaches based on temperature which exploit the change in power and temperature caused by activation of Trojan payload.

A. Power Consumption of Active/Inactive Trojan Payload

Prior work suggests that it can be difficult to detect certain Trojan triggers and payloads at test-time. The triggers can be very small (eg. one gate) and hidden in nets with low switching activity and leakage power [6]. Furthermore, by combining the Trojan trigger with power gating techniques [6], the Trojan’s impact on I/O and side-channels can be almost completely blocked even for very large payload circuitry. Such designs can easily bypass test-time approaches since their behavior is indistinguishable from Trojan-free designs while the Trojan is inactive. Run-time approaches on the other hand can monitor an IC for its entire lifetime and detect deviation of Trojan-active ICs from both Trojan-free and Trojan-inactive ICs.

As a motivating example, we look at power consumption of Trojan-inactive and Trojan-active ICs for a publicly available Trojan benchmark: RS232-T900 from trust-HUB [11]. RS232-T900 is a micro-UART core with a Trojan inserted in its transmitter. The Trojan is triggered by a sequence of four transmission messages and its payload prevents any further transmission. We determined power and layout information by simulating, synthesizing, and placing the design with Cadence SimVision, RTL Compiler, and Encounter tools. The average power across the IC in a $250\mu\text{s}$ time window for Trojan-inactive and Trojan-active ICs is shown in Figure 1. Comparing the two cases, one can see that the power differs in the upper left quadrant where the transmitter hardware is located. Transmission is blocked in the Trojan-active IC whereas the Trojan-inactive IC transmits and receives data at all times. Overall, there is a 40% decrease in total power consumption after the Trojan is triggered.

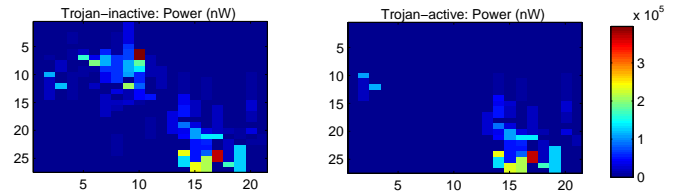


Fig. 1: Avg. power consumption (nW) in a $250\mu\text{s}$ time window across Trojan-inactive and Trojan-active ICs for the RS232-T900 benchmark.

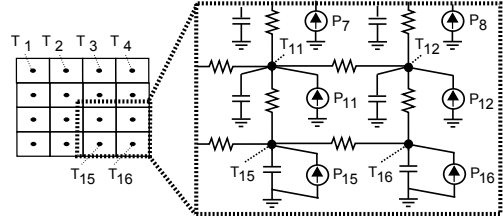


Fig. 2: IC broken into grids and RC thermal model within dotted region

B. Relationship between Power and Temperature

Temperature is a strong function of power consumption. Hence, changes in power caused by an active Trojan should also be reflected in the IC’s thermal profile. Below, we discuss the relationship between power and temperature and the advantages of temperature-based Trojan detection.

RC Thermal Model. We illustrate the link between power and temperature using the popular RC thermal model [12]. In this model, the IC is divided into grids and the temperature and power consumption of each grid at time t are represented as constants (see Figure 2). A circuit is used to estimate the IC’s thermal profile where node voltage and circuit current are analogous to temperature and power/heat flow in each grid. Voltage ground is analogous to the environment’s ambient temperature. Thermal capacitance and thermal resistance between neighboring nodes determine how heat flows between nodes/grids of the IC. Temperature for the entire IC can be determined by solving the following set of ODEs:

$$\sum_{v_j \in N_i} \frac{1}{R_{ij}} (T_i(t) - T_j(t)) + C_i \frac{dT_i(t)}{dt} - P_i(t) = 0 \quad \forall i \quad (1)$$

where $T_i(t)$ and $P_i(t)$ are the temperature and power dissipated at node i and time t ; C_i denotes thermal capacitance at node i ; R_{ij} denotes thermal resistance between nodes i and j ; and N_i is the set of all neighbors for node i . The above equations are often written in a discrete matrix form [13]:

$$\vec{T}[k] = \mathbf{A}\vec{T}[k-1] + \mathbf{B}\vec{P}[k-1] \quad (2)$$

where $\vec{T}[k]$ and $\vec{P}[k]$ are temperature and power vectors (each element corresponds to one node/grid) at discrete timestep k ; \mathbf{A} and \mathbf{B} are coefficient matrices that depend upon the RC circuit and timestep duration. The above formulation is quite flexible and has been used to simulate thermal dynamics as well as perform online temperature tracking [14], [15].

Advantages of Temperature-based Detection. Few run-time approaches take advantage of side-channel information because of the overheads involved. For example, a path delay characterization approach was proposed in [7] which, while effective, suffers from considerable area overhead for modern designs with millions of paths [1]. The current sensor approach in [8] should also come with significant area and power overheads. Temperature-based Trojan detection is an unexplored avenue which should have lower overheads compared to previous run-time approaches since many electronic systems are already equipped with thermal sensors for dynamic thermal management (DTM). For example, the AMD Operton multicore processor has 38 thermal sensors. Furthermore,

prior work has shown that the infrastructure for run-time thermal estimation has low hardware and software overheads as well. For instance, [15] used steady state Kalman filtering to track IC temperature with only five sensors. By combining the existing DTM infrastructure with new Trojan detection mechanisms, it should be possible to detect the changes in IC power/temperature caused by active Trojan payload.

IV. PROBLEM DEFINITION AND CHALLENGES

Our problem is inspired by the example discussed in Section III-A (RS232-T900). We assume that there are three possible states that an electronic system or IC can be in: Trojan-free, Trojan-inactive, and Trojan-active. Each state is defined by a set of statistical characteristics \mathbf{S}_f , \mathbf{S}_i , \mathbf{S}_a respectively. The Trojan-free and Trojan-inactive characteristics, while not necessarily identical, are close enough such that the Trojan-inserted IC can evade test-time Trojan detection methods (i.e. $\mathbf{S}_f \approx \mathbf{S}_i$). The Trojan-active characteristics \mathbf{S}_a on the other hand differ significantly from the other two. Our goal is a run-time temperature-based approach that can detect changes from \mathbf{S}_f and \mathbf{S}_i to \mathbf{S}_a after the Trojan is activated. Note, we are not concerned with Trojan-inactive ICs since, prior to Trojan activation, they essentially provide the same functionality as Trojan-free ICs. Stated formally, our problem is:

Given two hypotheses of the system's state:

$$\begin{cases} \mathcal{H}_0 & \text{The state is Trojan-free or Trojan-inactive} \\ \mathcal{H}_1 & \text{The state is Trojan-active} \end{cases}$$

Use thermal sensor observations to determine if the IC's state (characteristics) correspond to \mathcal{H}_0 (\mathbf{S}_f , \mathbf{S}_i) or \mathcal{H}_1 (\mathbf{S}_a).

The above problem has various challenges to overcome some of which are specific to temperature tracking and some of which are common to Trojan detection:

- **Golden IC/model:** Most Trojan detection approaches rely on the existence of a “golden model” to distinguish Trojan-free and Trojan-inserted ICs. In our case, we assume that the Trojan-free design is given and from it we can compute \mathbf{S}_f characteristics to function as our golden model.
- **Autonomous detection:** Since the \mathbf{S}_f characteristics are known and $\mathbf{S}_f \approx \mathbf{S}_i$, we should be able to easily track temperature for Trojan-free and Trojan-inactive designs as in prior work. The challenge is detecting active Trojan ICs because the \mathbf{S}_a characteristics are unknown. We propose two mechanisms for detecting Trojan activation at run-time.
- **Sensor Infrastructure and Noise:** Prior work has shown that sensor placement, number of sensors, sensor noise, etc. have a profound impact on temperature tracking [16]. In this paper, we vary the number of sensors to see the impact on temperature-based Trojan detection. One of our approaches uses the Kalman Filter which explicitly accounts for measurement noise.
- **Fabrication Variation (FV):** FV makes it more challenging to track temperature as well as detect Trojans. For tracking, FV results in larger uncertainty in the estimated thermal profile [15]. For Trojan detection, FV makes it difficult to distinguish between deviations in power/temperature due to manufacturing and Trojan presence [3]. In our approach, calibration is performed for each IC to ensure robustness in the face of FV.

V. TEMPERATURE-BASED DETECTION

In this section, we discuss the overall framework and algorithms for our temperature-based Trojan detection. While detection itself occurs at run-time, we also require some additional offline profiling steps to deal with the above challenges. An overview of entire approach with design, test, and run-time phases is shown in Figure 3. The details of each phase are discussed below.

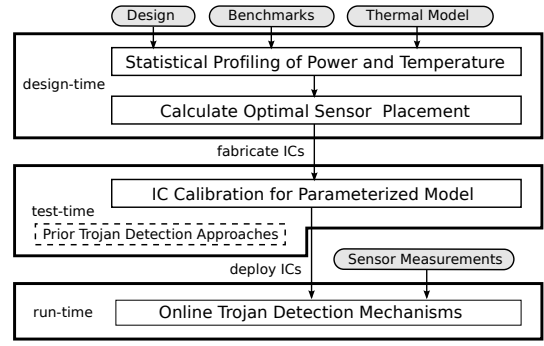


Fig. 3: Phases of the Proposed Approach

A. Design Phase

As discussed above, side channel-based Trojan detection often relies on a “golden” model. Our approach assumes that we have access to the Trojan-free design from which we can obtain statistical characteristics (\mathbf{S}_f) of switching activity, power consumption, thermal dynamics, etc. We leverage this knowledge to determine optimal sensor placement for temperature tracking/Trojan detection.

Design Profiling. We profile the Trojan-free designs using the RC thermal model (see Section III-B), benchmarks, and either state-of-the-art simulation tools or prototype ICs. The RC thermal model divides the power consumed by the design into grids and uses a vector \vec{P} to represent power consumed in the grids. The statistical approaches for temperature tracking and Trojan detection used in this paper require probability distribution functions (pdfs) to summarize the design’s expected power and/or temperature. Benchmarks that are representative of the design’s expected workload along with simulation tools are used to estimate the pdfs. Alternatively, IC prototypes that are verified as Trojan-free can be used. We approximate the pdfs as Gaussian with mean vectors $\vec{\mu}_p$ and $\vec{\mu}_T$ and covariance matrices \mathbf{Q}_p and \mathbf{Q}_T (for power and temperature respectively).

Sensor Placement. We adopt the sensor placement approach from [16] which uses the temperature covariance matrix \mathbf{Q}_T . Specifically, sensors are placed in a greedy fashion to minimize the following cost function

$$cost = \sum_{\forall grids:i} \max \left(0, 1 - \sum_{\forall sensors:j} q_{i,j} \right) \quad (3)$$

where $q_{i,j}$ denotes the element at location i and j of matrix \mathbf{Q}_T (i.e. correlation in temperature between IC grids i and j).

B. Test Phase

When the design phase is complete, we fabricate the ICs. At this point, we assume that some Trojan-inserted ICs are fabricated and inserted into the supply chain.

IC Parameter Calibration. Our test phase is mainly used to deal with fabrication variation. Fabrication variation (FV) results in ICs that have different physical, electrical, and performance parameters from the nominal design. As discussed above, FV makes it more challenging to accurately detect Trojans and track temperature. To ensure robustness, we must have accurate power/thermal statistics ($\vec{\mu}_p$, $\vec{\mu}_T$, \mathbf{Q}_p , \mathbf{Q}_T) for each IC under test (ICUT). One can accomplish this by applying test vectors to the ICUT, measuring power consumption, and estimating the pdfs after fabrication. For example, gate-level characterization has been applied in prior work [3] to successfully profile IC gate parameters. Temperature-based approaches which utilize infra-red cameras and Expectation Maximization are also applicable [17].

Test-time Detection. Prior test-time approaches could also be used at this point to remove some of the Trojan-inserted

ICs. The proposed run-time approach would then detect the remaining Trojan-inserted ICs as they are activated in the field.

C. Run-time Phase

As discussed in Section IV, our main problem is to decide the correct hypothesis (state of the system): \mathcal{H}_0 or \mathcal{H}_1 . In other words, is the IC Trojan-free/Trojan-inactive or Trojan-active? In this section, we propose two mechanisms to solve the problem. The first is a local sensor-based approach that uses an hypothesis testing (HT) framework. The second is a global approach that exploits correlation between sensors and maintains track of the IC's thermal profile with a Kalman filter.

1) *Hypothesis Testing (HT) Approach:* For simplicity, let us suppose we have one sensor measurement at timestep k denoted by $S[k]$ from which we shall decide the state (we'll consider more sensors later). In an hypothesis testing framework, one assumes that $S[k]$ can only come from one of two pdfs: S_0 or S_1 which correspond to temperature in Trojan-free/Trojan-inactive and Trojan-active ICs respectively. We shall choose the correct state as the one with the highest probability of occurrence given $S[k]$ (i.e. $\text{argmax} Pr(H_x|S[k]), x \in \{0, 1\}$). By applying Bayes rule, it can be shown the optimal decision is [13]:

$$decision = \begin{cases} \mathcal{H}_0, & Pr(S[k]|S_0)Pr(S_0) > Pr(S[k]|S_1)Pr(S_1) \\ \mathcal{H}_1, & Pr(S[k]|S_0)Pr(S_0) < Pr(S[k]|S_1)Pr(S_1) \end{cases} \quad (4)$$

where $Pr(x)$ and $Pr(x|y)$ denote the prior probability of x and probability of x given y . Basically, an IC shall be regarded as Trojan-free/inactive (Trojan-active) if the measurement comes from S_0 (S_1) with higher probability. While this methodology is theoretically sound, it is difficult to directly apply to our problem for two reasons. First, one cannot assume that all measurements come from single stationary pdf (i.e. one that is time invariant) since the IC temperature varies with time. Second, even if the pdfs were stationary, we do not have access to the Trojan-active design and therefore cannot accurately estimate S_1 , $Pr(S[k]|S_1)$, or $Pr(S_1)$. We get around these issues by making several simplifying assumptions.

Stable State Temperature. The first issue no longer presents a problem when the IC's temperature has reached a stable state. Put simply, if an IC's power consumption is similar for a long period of time, the IC's temperature will end up converging to a "stable state" [18] where measurements of its thermal state are actually samples from a stationary pdf. In this paper, we run benchmarks with random inputs on the ICUT until reaching the stable state. We then take measurements and approximate S_0 as Gaussian with mean μ_0 and variance σ_0^2 .

Trojan pdf Estimate: To overcome the second issue, we must make some assumptions about S_1 . Basically, we exploit the fact that S_0 and S_1 must be slightly different. We assume that the mean of S_1 (μ_1) differs from S_0 's known mean (μ_0) by some fixed percentage difference $S\%$ and both possess the same variance ($\sigma_0 = \sigma_1$.) If $S\% < 0$ ($S\% > 0$), Trojan activation causes the IC to lose (gain) some functionality. This is the best we can do to apply the theory because we have little if any knowledge of the actual Trojan-active design.

Single sensor Decision Rule. With the above assumptions, one can easily come up with a decision using Eqn. (4). For simplicity, we assume that $Pr(S_0) = Pr(S_1) = .5$ and compute the conditional probabilities with:

$$Pr(S[k]|S_x) = \frac{1}{\sigma_x \sqrt{2\pi}} \exp\left(-\frac{(S[k] - \mu_x)^2}{2\sigma_x^2}\right), \quad x \in \{0, 1\} \quad (5)$$

Multi-sensor Decision Rule. For multiple sensors, we can easily extend the above rule by collecting the sensor measurements as a vector $\vec{S}[k]$ and using multivariate Gaussian pdfs. For simplicity, we take a simple ad-hoc approach instead. We

evaluate Eqn. (4) for each sensor and come to decision (\mathcal{H}_0 or \mathcal{H}_1) by majority voting.

Overheads. A high-level overview of the HT approach and its overheads is shown in Figure 4. HT's offline overhead includes computing stable state pdfs (S_0 and S_1 for each sensor). At run-time, the probabilities of sensor measurements for hypotheses \mathcal{H}_0 and \mathcal{H}_1 are computed and compared resulting in a 1-bit vote for each sensor. Majority voting is used to combine the decisions of z sensors and obtain a final decision. All these operations are simple and the overall complexity depends only on the number of sensors z .

2) *Kalman Filter based Approach:* This is a global sensor-based approach that exploits correlation between sensors and uses a Kalman Filter (KF) to dynamically track the system's thermal profile at run-time. An autocorrelation based metric then decides between hypotheses \mathcal{H}_0 and \mathcal{H}_1 (Trojan-free/Trojan-inactive and Trojan-active).

Temperature Tracking Via Kalman Filter. We track IC temperature at run-time using the standard Kalman filtering (KF) approach developed in prior work [14], [15]. For simplicity, we only consider dynamic power, but leakage power can be handled as well [19]. The KF relies on a state-space equation to model the random dynamics of the state being estimated and on a measurement equation to relate measurements with the state being estimated. The state-space equation for temperature tracking is the discrete form RC thermal model equation discussed in Section III-B (copied below for convenience)

$$\vec{T}[k] = \mathbf{A}\vec{T}[k-1] + \mathbf{B}\vec{P}[k-1] \quad (6)$$

The above equation assumes that the current thermal state $\vec{T}[k]$ depends on the previous thermal state $\vec{T}[k-1]$ (Markovian assumption) and also local power dissipation $\vec{P}[k-1]$. Due to variations in the voltage supply noise, system workload, etc., the power \vec{P} is random at each timestep and $\vec{T}[k]$ cannot be precisely computed with the state-space model alone. To improve the estimate, the KF uses measurements collected by thermal sensors and the following measurement model

$$\vec{S}[k] = \mathbf{H}\vec{T}[k] + \vec{v}[k] \quad (7)$$

where $\vec{S}[k]$ is a vector of sensor measurements at timestep k ; \mathbf{H} is a transformation matrix based on the sensor placement; and $\vec{v}[k]$ is a Gaussian random vector with zero mean and known covariance \mathbf{R} [14] representing measurement noise.

The KF estimates the thermal state of a chip as follows. $\vec{P}[k]$ is modeled as a Gaussian random vector with known mean $\vec{\mu}_p$ and covariance \mathbf{Q}_p (which we determine in the design/test phases). KF estimation is then performed recursively with *predict* and *update* steps. In the *predict* step, the KF uses $\vec{\mu}_p$ and the previous temperature estimate to predict the IC's new thermal state. In the *update* step, the KF corrects this estimate based on new sensor measurements. The following equations are used at each timestep

$$\text{predict:} \quad \vec{T}[k|k-1] = \mathbf{A}\vec{T}[k-1] + \mathbf{B}\vec{\mu}_p \quad (8)$$

$$\mathbf{C}[k|k-1] = \mathbf{A}\mathbf{C}[k-1]\mathbf{A}^T + \mathbf{B}\mathbf{Q}_p\mathbf{B}^T \quad (9)$$

$$\text{update:} \quad \vec{e}[k] = \vec{S}[k] - \mathbf{H}\vec{T}[k|k-1] \quad (10)$$

$$\vec{T}[k|k] = \vec{T}[k|k-1] + \mathbf{K}[k]\vec{e}[k] \quad (11)$$

$$\mathbf{K}[k] = \mathbf{C}[k|k-1]\mathbf{H}^T(\mathbf{R} + \mathbf{H}\mathbf{C}[k|k-1]\mathbf{H}^T)^{-1} \quad (12)$$

$$\mathbf{C}[k|k] = (\mathbf{I} - \mathbf{K}[k]\mathbf{H})\mathbf{C}[k|k-1] \quad (13)$$

$\vec{T}[k|k]$ and $\vec{T}[k|k-1]$ are estimates of the temperature at time k computed with and without sensor information respectively; $\mathbf{C}[k|k-1]$ and $\mathbf{C}[k|k]$ are the error covariance matrices associated with $\vec{T}[k|k-1]$ and $\vec{T}[k|k]$; $\vec{e}[k]$ is the KF *residual* which reflects the discrepancy between the predicted and actual measurements; \mathbf{I} is the identity matrix; $\mathbf{K}[k]$ represents

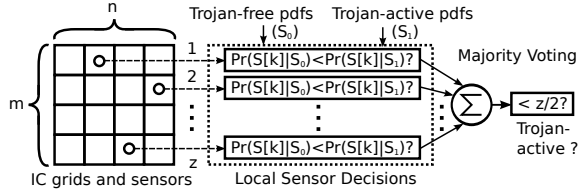


Fig. 4: Local Hypothesis Testing (HT) approach with z sensors and $m \times n$ grid

the Kalman gain at the k th step and is chosen to minimize the error in $\hat{T}[k|k]$.

Steady State Kalman Filter. When the statistical characteristics of \bar{P} and measurement noise are stationary (or do not change for relatively long time), the KF stabilizes which means $\mathbf{C}[k|k-1]$, $\mathbf{C}[k|k]$, and $\mathbf{K}[k]$ converge to static values. This is referred to as the KF *steady state*. During the steady state, even though the temperature may change with time, the error associated with the estimates remains the same. The steady state allows one to create low overhead implementations of the KF [15] by replacing $\mathbf{C}[k|k-1]$, $\mathbf{C}[k|k]$, and $\mathbf{K}[k]$ in Eqns. (8) to (13) by constants.

Autocorrelation-based Detection Rule. While the KF can be used to accurately track temperature, we also need a rule to decide on the correct state (\mathcal{H}_0 or \mathcal{H}_1). Our decision rule is based on the KF residual and uses the autocorrelation function of the residual process. In the KF, residual $\bar{e}[k]$ represents the discrepancy between the predicated temperature and thermal sensor measurements. If $\bar{e}[k]$ is small (large), the two agree (disagree) on the thermal state. Assuming the state-space model/parameters and the sensor noise covariance are reasonably accurate, the autocorrelation of the residual should be close to zero on average [20]. When a Trojan gets activated, the state-space model (which does not account for the power of an active Trojan) becomes less accurate and should cause the autocorrelation to diverge from zero.

We use the following method to detect a Trojan at timestep x . We record the residual in the N previous timesteps of the KF ($\bar{e}[x-N]$, $\bar{e}[x-N+1]$, \dots , $\bar{e}[x]$) and then compute the following cost function

$$\hat{a}[x] = \frac{1}{N} \sum_{i=x-N+1}^x (\bar{e}[i] \cdot \bar{e}[i-1]^T) \quad (14)$$

This cost is the average autocorrelation of the residual process in the last N timesteps. To decide if a Trojan is activated, we define thresholds a_T and V . If $(|\hat{a}[x]| > a_T)$ for more than V consecutive timesteps, we assume a Trojan has been activated (i.e. state \mathcal{H}_1). a_T and V are parameters that tune the aggressiveness of the decision rule.

Overheads. A high-level overview of the KF approach and its overheads is shown in Figure 5. The main offline overheads are estimating the power statistics ($\bar{\mu}_p$ and \mathbf{Q}_p) computing the steady state Kalman gain matrix \mathbf{K} . During run-time, the KF performs some matrix-vector multiplications and vector additions/subtractions (Eqns. (8), (10), (11)) at each timestep k . The size of the matrices/vectors and complexity of these operations depend on the number of grids in the RC thermal model (mn) and the number of sensors (z). Running averages of the autocorrelation \hat{a} for the last V timesteps must kept and the final state is chosen based on V thresholding operations

3) **Qualitative Comparison:** The salient differences between the above two detection mechanisms are as follows:

- **Stable State Assumption:** Hypothesis testing (HT) requires the system under test to be in a stable thermal state so that the sensor measurements can be compared with stationary pdfs. The Kalman Filter (KF) approach works with the system in any thermal state since it tracks the system's

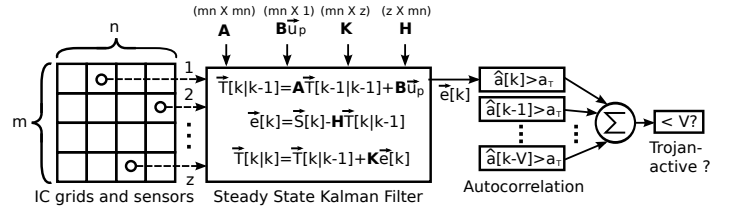


Fig. 5: Global Kalman Filter (KF) approach with z sensors and $m \times n$ grid

thermal profile at all timesteps.

- **Sensor Correlation:** HT is a local approach that compares each sensor measurement with its corresponding stable state pdf in an independent fashion. Correlation between the sensors is not exploited and the final decision is made based on a majority vote. The KF is a global approach that exploits the correlation between sensors to accurately track temperature and detect Trojans.
- **Run-time Overheads:** The KF approach clearly has larger run-time overheads than the HT approach (see Figs. 4 and 5). While the HT approach primarily works with scalar values and computes 1-bit decisions, the KF approach requires matrix-vector storage and computations which are more expensive ($O(mn)$).

VI. EXPERIMENTS AND DISCUSSION

A. Setup

Benchmarks. We tested our Trojan detection schemes on five publicly available Trojan benchmarks (from trust-HUB [11]):

- 1) **RS232-T900:** was discussed in Section III-A.
- 2) **s38417-T300:** contains a Trojan trigger with activation probability of $1.7e-44$. Once activated, the payload leaks the value of a specific net through a 29 stage ring-oscillator.
- 3) **BasicRSA-T200:** is an RSA encoder with a Trojan triggered by a specific plaintext input. The payload permanently disables encoding of the plaintext.
- 4) **MC8051-T300:** is an implementation of the 8051 micro-processor with a Trojan. The Trojan is triggered when a specific string is sent through the UART and the payload blocks new messages from the UART.
- 5) **MC8051-T600:** The Trojan is activated by an external interrupt and disables 8051 instructions containing jumps.

We determined power and layout information in the above benchmarks by simulating, synthesizing, and placing each design with Cadence SimVision, RTL Compiler, and Encounter tools for two different testbench instances: one which activates the Trojan (i.e. Trojan-active) and one which does not (i.e. Trojan-inactive). The difference in power consumption between the two is shown in Table I for all the benchmarks. In all cases but one (MC8051-T300) there is a % difference larger than 25%.

Temperature-based Trojan Detection. We divided the IC into 20 by 16 grids (320 distinct regions). In the design phase, we computed power and temperature statistics using 250ms of data generated by Cadence (from the Trojan-free designs) and the RC thermal model [12]. With the resulting statistics, we placed sensors as discussed in Section V-A. The number of sensors we tested were 4, 16, and 32. For simplicity, we ignored fabrication variation and therefore did not implement the test phase. For the run-time phase, we computed “real” dynamic thermal profiles using the RC thermal model. A steady state Kalman filter (KF) implementation was used to estimate the thermal profile for Trojan-active and Trojan-inactive cases. Sensor measurements were made by overlaying noise onto the “real” thermal profile. We assumed sensor noise variance of 0.1 which seems like a worst-case for state-of-the-art thermal sensors [21]. For KF-based Trojan detection, we

TABLE I: % difference in total power consumption between Trojan-inactive and Trojan-active in 250ms experiment

Benchmark	% difference
RS232-T900	-39.97%
s38417-T300	54.33%
BasicRSA-T200	-28.40%
MC8051-T300	-1.5%
MC8051-T600	-72.16%

TABLE II: Average true positive rate t_+ , false positive rate f_+ , and detection time t_{dec} (seconds) for 100 trials. Note detection time is given in seconds and only includes true positives. ‘-’ indicates no true positives.

	# sensors	RS232-T900			s38417-T300			BasicRSA-T200			MC8051-T300			MC8051-T600		
		t_+	f_+	t_{dec}	t_+	f_+	t_{dec}	t_+	f_+	t_{dec}	t_+	f_+	t_{dec}	t_+	f_+	t_{dec}
HT	4	100%	79%	4.9E-2	100%	45%	1.9E-3	100%	66%	1.7E-2	68%	66%	1.1E-1	100%	64%	5.2E-2
	16	100%	0%	1.7E-1	100%	0%	4.3E-3	100%	0%	3.9E-2	0%	0%	-	100%	0%	1.6E-1
	32	100%	0%	2.1E-1	100%	0%	5.1E-3	100%	0%	4.8E-2	0%	0%	-	100%	0%	2.0E-1
KF	4	100%	1%	9.2E-4	100%	2%	3.6E-2	100%	1%	8.4E-3	0%	0%	-	100%	2%	3.4E-2
	16	100%	0%	6.4E-4	100%	5%	2.6E-2	100%	1%	5.8E-3	1%	1%	1.1E-2	100%	5%	2.4E-2
	32	100%	0%	5.8E-4	100%	1%	2.2E-2	100%	0%	3.2E-3	1%	1%	1.2E-2	100%	3%	2.2E-2

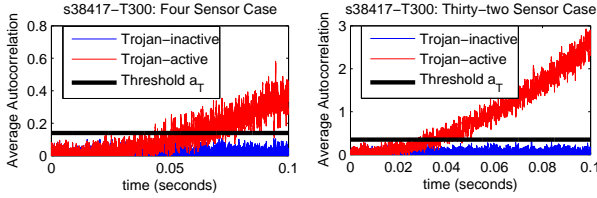


Fig. 6: Average autocorrelation (\hat{a}) over time with 4 and 32 sensors for s38417-T300

stored $N = 50$ residuals and chose autocorrelation thresholds $V = 10$ and $a_T = 0.14, 0.24, .035$ for 4, 16, and 32 sensors respectively (based on data from the Trojan-free designs). For the hypothesis testing (HT) approach, we used mean difference $S\% = \pm 2.5\%$ to estimate the Trojan-active stable state pdf S_1 . Except where specified, one can assume the experiments were conducted while the ICs were in stable thermal states.

B. Results

We conducted 100 trials with random sensor noise on both the Trojan-inactive and Trojan-active ICs. We recorded the following data: average true positive rate t_+ , average false positive rate f_+ , and average time t_{dec} to obtain a true positive. The results are shown for all 5 benchmarks and both detection mechanisms in Table II. ‘HT’ and ‘KF’ denote the hypothesis testing and Kalman Filter based approaches.

Hypothesis Testing (HT). HT was able to detect the active Trojans (true positives) in all the benchmarks with 100% accuracy except for MC8051-T300. MC8051-T300 had the smallest difference in power consumption between Trojan-active and Trojan-inactive cases (see Table I) and thus there was little deviation in the thermal profile. False positives on the inactive Trojans were only an issue in the 4 sensor case and went to 0% with additional sensors. Increasing the number of sensors also resulted in slower Trojan detection. Put simply, it took a longer time for the majority of sensors to agree on a true positive when there were more sensors. On average, 4 sensors could detect Trojans 69.5% faster than 32 (ignoring MC8051-T300), but with higher false positive rate.

Kalman Filtering (KF). The KF was also very successful with true positives in every benchmark but one. Once again, the deviation in power/temperature was too small to detect for MC8051-T300. The KF had a very low false positive rate in all instances ($\leq 5\%$). In contrast to HT, increasing the number of sensors from 4 to 32 improved detection time by 37.3% on average (ignoring MC8051-T300). Basically, the more measurements the better the resolution of thermal profile and autocorrelation \hat{a} . To illustrate, Figure 6 shows \hat{a} of Trojan-inactive (blue) and Trojan-active (red) ICs with 4 and 32 sensors for s38417-T300. The Trojan-inactive autocorrelation stays below the threshold a_T (black line) while the Trojan-active autocorrelation diverges from zero and exceeds the threshold. The Trojan-active case crosses the threshold more quickly in the 32 sensor case.

Note that the results in Table II correspond to the case where the ICs were in stable thermal states, but we also ran trials at room temperature. For the latter, the KF yielded similar results which we have not shown due to space limitations.

Comparing HT and KF. While both approaches worked well, the KF achieved better results. The KF found active Trojans 71.6% faster on average than HT and was effective even with only 4 sensors. Also, while the HT approach could only operate with the ICs in a stable thermal state, the KF worked in all scenarios. The advantage of HT is its lower overheads.

VII. CONCLUSION AND FUTURE WORK

In this paper, we proposed novel temperature-based approaches for online Trojan detection. The experimental results on real Trojan benchmarks showed that our proposed schemes were effective in all instances but one. Failure only occurred when the deviation in power/temperature was too small to detect. While we find these results promising, we shall continue making improvements to the proposed approaches. In future work, we plan on evaluating the impact of fabrication variation and further investigating the implementation overheads.

REFERENCES

- [1] M. Tehranipoor and F. Koushanfar, ‘‘A survey of hardware trojan taxonomy and detection,’’ *IEEE Des. Test*, 2010.
- [2] S. Bhunia, M. Abramovici, D. Agrawal, P. Bradley, M. Hsiao, J. Plusquellic, and M. Tehranipoor, ‘‘Protection against hardware trojan attacks: Towards a comprehensive solution,’’ *IEEE Des. Test*, 2012.
- [3] F. Koushanfar and A. Mirhoseini, ‘‘A unified framework for multimodal submodular integrated circuits trojan detection,’’ *IEEE TIFS*, 2011.
- [4] H. Salmami, M. Tehranipoor, and J. Plusquellic, ‘‘A novel technique for improving hardware trojan detection and reducing trojan activation time,’’ *IEEE TVLSI*, 2012.
- [5] —, ‘‘A layout-aware approach for improving localized switching to detect hardware trojans in integrated circuits,’’ in *Proc. WIFS*, 2010.
- [6] S. Wei, K. Li, F. Koushanfar, and M. Potkonjak, ‘‘Hardware trojan horse benchmark via optimal creation and placement of malicious circuitry,’’ in *Proc. DAC*, 2012.
- [7] J. Li and J. Lach, ‘‘At-speed delay characterization for ic authentication and trojan horse detection,’’ in *Proc. HOST*, 2008.
- [8] S. Narasimhan, W. Yueh, X. Wang, S. Mukhopadhyay, and S. Bhunia, ‘‘Improving ic security against trojan attacks through integration of security monitors,’’ *IEEE Des. Test*, 2012.
- [9] M. Abramovici and P. Bradley, ‘‘Integrated circuit security: new threats and solutions,’’ in *Proc. Annual CSIR Workshop*, 2009.
- [10] M. Hicks, M. Finnicum, S. King, M. Martin, and J. Smith, ‘‘Overcoming an untrusted computing base: Detecting and removing malicious hardware automatically,’’ in *IEEE Symp. Security and Privacy*, 2010.
- [11] trust HUB.org, <http://trust-hub.org/resources/benchmarks>.
- [12] K. Skadron, M. Stan, K. Sankaranarayanan, W. Huang, S. Velusamy, and D. Tarjan, ‘‘Temperature-aware microarchitecture: Modeling and implementation,’’ *ACM TACO*, 2004.
- [13] Y. Zhang and A. Srivastava, ‘‘Adaptive and autonomous thermal tracking for high performance computing systems,’’ in *Proc. DAC*, 2010.
- [14] S. Sharifi, C. Liu, and T. Rosing, ‘‘Accurate temperature estimation for efficient thermal management,’’ *Proc. ISQED*, 2008.
- [15] Y. Zhang, A. Srivastava, and M. Zahran, ‘‘On-chip sensor-driven efficient thermal profile estimation algorithms,’’ *ACM TODAES*, 2010.
- [16] Y. Zhang, B. Shi, and A. Srivastava, ‘‘A statistical framework for designing on-chip thermal sensing infrastructure in nano-scale systems,’’ in *Proc. ISPD*, 2010.
- [17] Y. Zhang and A. Srivastava, ‘‘Statistical characterization of chip power behavior at post-fabrication stage,’’ in *Proc. IGCC*, 2011.
- [18] D. Forte and A. Srivastava, ‘‘Energy and thermal-aware video coding via encoder/decoder workload balancing,’’ in *Proc. ISLPED*, 2010.
- [19] Y. Zhang and A. Srivastava, ‘‘Leakage-aware kalman filter for accurate temperature tracking,’’ in *Proc. IGCC*, 2011.
- [20] J. L. Crassidis and J. L. Junkins, *Optimal Estimation of Dynamic Systems*. CRC Press, 2004.
- [21] F. Sebastiano, L. Breems, K. Makinwa, S. Drago, D. Leenaerts, and B. Nauta, ‘‘A 1.2 v 10 μ w npn-based temperature sensor in 65nm cmos with an inaccuracy of $\pm 0.2^\circ$ c from -70° c to 125° c,’’ in *Proc. ISSCC*, 2010.