

# Robust and Attack Resilient Logic Locking with a High Application-Level Impact

YUNTAO LIU, MICHAEL ZUZAK, YANG XIE, ABHISHEK CHAKRABORTY, ANKUR SRIVASTAVA,

University of Maryland, College Park

Logic locking is a hardware security technique aimed at protecting intellectual property (IP) against security threats in the IC supply chain, especially those posed by untrusted fabrication facilities. Such techniques incorporate additional locking circuitry within an IC that induces incorrect digital functionality when an incorrect verification key is provided by a user. The amount of error induced by an incorrect key is known as the **effectiveness** of the locking technique. A family of attacks known as "SAT attacks" provide a strong mathematical formulation to find the correct key of locked circuits. In order to achieve high **SAT resilience** (*i.e.* complexity of SAT attacks), many conventional logic locking schemes fail to inject sufficient error into the circuit when the key is incorrect. For example, in the case of SARLock and Anti-SAT, there are usually very few (or only one) input minterms that cause any error at the circuit output. The state-of-the-art stripped functionality logic locking (SFL) technique provides a wide spectrum of configurations which introduced a trade-off between **SAT resilience** and **effectiveness**. In this work, we prove that such a trade-off is universal among all logic locking techniques. In order to attain high effectiveness of locking without compromising SAT resilience, we propose a novel logic locking scheme, called Strong Anti-SAT (SAS). In addition to SAT attacks, removal-based attacks are another popular kind of attack formulation against logic locking where the attacker tries to identify and remove the locking structure. Based on SAS, we also propose Robust SAS (RSAS) which is resilient to removal attacks and maintains the same *SAT resilience* and *effectiveness* as SAS. SAS and RSAS have the following significant improvements over existing techniques. (1) We prove that the *SAT resilience* of SAS and RSAS against SAT attack is not compromised by increase in *effectiveness*. (2) In contrast to prior work which focused solely on the circuit-level locking impact, we integrate SAS-locked modules into an 80386 processor and show that SAS has a high application-level impact. (3) Our experiments show that SAS and RSAS exhibit better SAT resilience than SFL and their effectiveness is similar to SFL.

CCS Concepts: • **Security and privacy** → **Security in hardware**.

Additional Key Words and Phrases: logic locking, SAT attack, machine learning

## ACM Reference Format:

Yuntao Liu, Michael Zuzak, Yang Xie, Abhishek Chakraborty, Ankur Srivastava. 2020. Robust and Attack Resilient Logic Locking with a High Application-Level Impact. *ACM J. Emerg. Technol. Comput. Syst.* 1, 1, Article 1 (January 2020), 22 pages. <https://doi.org/10.1145/3446215>

## 1 INTRODUCTION

Due to the increasing cost of maintaining IC foundries with advanced technology nodes, many chip designers have become fabless and outsource their fabrication to off-shore foundries. However, such foundries are not under the designer's control which puts the security of the IC supply chain

---

Author's address: Yuntao Liu, Michael Zuzak, Yang Xie, Abhishek Chakraborty, Ankur Srivastava, University of Maryland, College Park.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM.

1550-4832/2020/1-ART1 \$15.00  
<https://doi.org/10.1145/3446215>

at risk. Untrusted foundries are capable of various malicious activities, among which “overbuilding” is the most concerning since the foundry can simply produce more copies of the IC and sell the extra copies for profit. Many design-for-trust techniques have been studied as countermeasures among which logic locking has been the most widely studied [4]. A logic locked circuit requires a secret key input and the correct key is kept by the designer and not known to the foundry. The functionality of the circuit is correct only if the key is correct. After the foundry manufactures the locked circuit and returns it to the designer, the correct key is applied to the circuit by connecting a tamper-proof memory containing the key to the key inputs. This process is called *activation*. Over the years, different types of logic locking mechanisms have been suggested. Initially, locking involved inserting XOR/XNOR gates in a synthesized design netlist [22]. Later, techniques based on VLSI testing principles have been outlined to improve logic locking schemes by manifesting high corruption at the output bits when an incorrect key is applied [20, 21]. It is noteworthy that logic locking has evolved from an academic proposal to a practical solution. Sisejkovic *et al.* proposed Inter-Lock, a cross-core logic locking technique [27], based on which a scalable logic locking solution for multi-module hardware designs was implemented it on a RISC-V processor [28, 29], which has already been made commercially available.

The Boolean satisfiability-based attack, a.k.a. SAT attack [30] was a game changer in the logic locking field. SAT provides a strong mathematical formulation to find the correct locking key of a logic locked IC which prunes out wrong keys in an iterative manner. In each iteration, an input (called the Distinguishing Input, or DI) is chosen by the SAT solver and all the wrong keys that corrupt the output of this DI are pruned out. All wrong keys are pruned out when no more DI can be found. Point function (PF)-based logic locking, including SARLock [33] and Anti-SAT [32], force the number of SAT iterations to be exponential in the key size by pruning out only a very small number of wrong keys in each iteration. However, PF-based locking schemes have a drawback that there are very few (or only one) input minterms whose output is incorrect for each wrong key. Hence the overall error rate of the locked circuit with a wrong key is very small. This disadvantage is captured by approximate SAT attacks such as AppSAT [24] and Double-DIP [25]. These attack schemes are able to find an *approximate key* (*approx-key*) which makes the locked circuit behave correctly for most (but not all) of the input values. Another kind of popular attack against logic locking schemes is removal attacks [35]. In a removal attack, the attacker tries to find the logic locking module, remove it, and replace its output with a constant 0 or 1. The key step in this attack is to identify the output wire of the locking module. This can be achieved by structural analysis assisted by calculating the signal probability skew (SPS) of each wire [35]. Locking techniques such as Anti-SAT [32] is most vulnerable to this type of attack since the correct functionality of the original circuit can be obtained by removing the Anti-SAT module and replacing its output with 0. Other types of attacks on logic locking have also been proposed, such as Hamming distance guided hill climbing attack [18], reduced-order binary decision diagram (ROBDD)-based attack [14], machine learning based structural attacks “SAIL” [6] and “SWEEP” [1], and combined structural and functional attack “SURF” [7]. These attacks are effective in recovering *most* of the key bits correctly. However, these attacks mainly targeted pre-SAT logic locking schemes and do not guarantee to find a correct key. Hence, in this paper, we focus on SAT and removal based attacks.

More recently, Yasin *et al.* proposed *stripped functionality logic locking* (SFL) which allows the designer to select a set of protected input patterns that are affected by a large percentage of wrong keys while other input patterns are affected by very few wrong keys [37]. SFL is not vulnerable to removal attack since the functionality of the original circuit for the protected input patterns has been modified in SFL. However, when the number of protected patterns increases, SAT attacks need significantly fewer iterations to find the correct key. More details of SFL and attack methods targeting SFL are introduced in Section 2.3. Essentially, SFL creates a fundamental trade-off

between **SAT resilience** (*i.e.* SAT attack complexity) and **effectiveness** (*i.e.* the amount of error injected by a wrong key). This trade-off is problematic. On the one hand, if only very few input patterns are protected, a wrong key may not inject enough error into the circuit and useful work may still be done using the chip, rendering locking **ineffective**. On the other hand, having more protected input patterns will compromise the circuit's **SAT resilience**. Moreover, as we move into the machine learning (ML) era, error-resilient applications are becoming increasingly relevant since most ML-based applications usually embody substantial amount of error resilience. Hence small amount of error in the hardware (introduced by incorrect keys and/or hardware simplification) may not necessarily impact the overall application accuracy. With SFL, if we want to ensure a very high corruption at the hardware level (for wrong keys), the resiliency to SAT would inevitably reduce. Addressing this dilemma is the main theme of our paper.

We propose *Strong Anti-SAT (SAS)* to address the challenges in achieving high effectiveness without sacrificing SAT resilience. On one hand, SAS ensures that, given any wrong (including approximate) key, the error injected by locking circuitry will have significant application-level impact. On the other hand, SAS is provably resilient to SAT attacks (*i.e.* requiring exponential time). Based on SAS, we also propose Robust SAS (RSAS), a variant of SAS that is not vulnerable to removal attacks and has the same *SAT resilience* and *effectiveness* as SAS. This makes RSAS a substantial improvement over the limitations posed by SAS. The contribution of this work is as follows.

- (1) We prove the fundamental trade-off between *SAT resilience* and *effectiveness* which is applicable to any logic locking scheme.
- (2) We demonstrate the inability of existing locking techniques to secure hardware running real-world workloads due to such a trade-off. We show that, when the longest combinational path (*i.e.* the multiplier) in a 32-bit 80386 processor is locked using SFL, the processor fails to simultaneously have high SAT complexity and high application-level impact on both PARSEC [2] and ML-based application benchmarks.
- (3) We propose *Strong Anti-SAT (SAS)* to address this challenge. In SAS, a set of input minterms that have higher impact on the applications are identified as *critical minterms*. We design the locking infrastructure of SAS such that given a wrong key, the critical minterms are more likely to introduce error in the circuit and hence result in an application-level error. We also prove that the SAT complexity is exponential in the number of key bits and does not deteriorate when the number of critical minterms increases. This is a substantial improvement over SFL.
- (4) We also propose a removal attack resistant variant of SAS, called Robust SAS (RSAS). RSAS is designed such that it achieves the same *SAT resilience* and *effectiveness* levels as SAS and if the locking module of RSAS is removed, the remaining circuit will exhibit incorrect functionality for critical minterms.
- (5) Experiment results show that, when locked using the same number of critical minterms, SAS and RSAS have higher *SAT resilience* than SFL and they have about the same level of effectiveness. In terms of area, power, and delay overhead, RSAS and SFL have similar overheads in general.

The rest of the paper is organized as follows. Sec. 2 introduces the background on SAT attack and existing logic locking schemes. We show that SFL's trade-off makes it incapable to secure real-world applications in Section 3. We then mathematically prove that the trade-off applies to all logic locking schemes in Section 4. In Section 5, SAS's hardware structure is presented and its exponential SAT attack complexity is proved in theory. The removal attack resistant variant of SAS, *i.e.* RSAS, is introduced in Section 6. Section 7 describes the methodology to choose critical

minterms. Section 8 shows the experimental results which demonstrate that when the same set of critical minterms are selected by SAS, RSAS, and SFL, SAS and RSAS achieve higher *security* than SFL while maintaining similar application-level effectiveness. Section 9 concludes the paper.

## 2 BACKGROUND

### 2.1 Attack Model

Fig. 1 illustrates the threat model we consider which is consistent with the latest papers in the logic locking field [13, 23, 24, 32]. The attacker can be either an untrusted foundry or an untrusted user who has the ability to reverse engineer the fabricated chip and obtain the locked gate-level netlist. The attacker is considered to have the following resources:

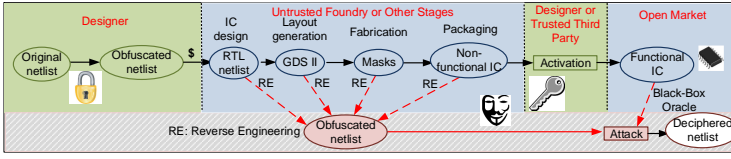


Fig. 1. The targeted attack model of logic locking

- (1) *The locked gate-level netlist of the circuit under attack.* This can be obtained by reverse engineering the GDS-II file (which the foundry has) or a fabricated chip (which can be done by a capable end user).
- (2) *An activated chip.* The attacker is considered to own an activated chip (*i.e.* the one loaded with the correct key) since such a chip can be purchased from the open market.

In general, logic locking research does not assume that the attacker is able to insert probes into the activated circuit, *i.e.* to observe the intermediate values. This is because protection schemes (*e.g.* analog shield [16]) can counter probing attacks.

### 2.2 SAT Attack

For any combinational digital circuit, the functionality can be expressed using a Boolean function  $F : \vec{X} \rightarrow \vec{Y}$  where  $\vec{X} \in \{0, 1\}^n$  and  $\vec{Y} \in \{0, 1\}^{n_o}$  are the input and output of the circuit, respectively. The logic locked circuit  $F_L$  takes one more input, the key input  $\vec{K} \in \{0, 1\}^k$ , in addition to the primary input  $\vec{X}$ , *i.e.*  $F_L : \vec{X}, \vec{K} \rightarrow \vec{Y}$ . If  $\vec{K}$  is correct, then  $\forall \vec{X}, F(\vec{X}) = F_L(\vec{X}, \vec{K})$ .  $F(\vec{X})$  may not be equal to  $F_L(\vec{X}, \vec{K})$  if  $\vec{K}$  is incorrect. As stated earlier, the key is stored tamper-proof memory and is not accessible to the attacker.

The Boolean satisfiability-based attack, a.k.a. *SAT attack* is a strong theoretical formulation to find the correct key of a locked circuit. In the context of the SAT attack, we use the *Conjunctive Normal Form (CNF)*:  $C(\vec{X}, \vec{K}, \vec{Y})$  to characterize Boolean satisfiability:  $C(\vec{X}, \vec{K}, \vec{Y}) = \text{TRUE}$  if  $\vec{X}$ ,  $\vec{K}$ , and  $\vec{Y}$  satisfy  $\vec{Y} = F_L(\vec{X}, \vec{K})$ , where  $F_L$  stands for the Boolean functionality of the locked circuit.  $C(\vec{X}, \vec{K}, \vec{Y}) = \text{FALSE}$  otherwise. SAT attacks run iteratively and prune out incorrect keys in every iteration. The attack consists of the following steps:

- (1) In the initial iteration, the attacker looks for a primary input,  $\vec{X}_1$ , and two keys,  $\vec{K}_\alpha$  and  $\vec{K}_\beta$ , such that the locked circuit produces two different outputs  $\vec{Y}_\alpha$  and  $\vec{Y}_\beta$ :

$$C(\vec{X}_1, \vec{K}_\alpha, \vec{Y}_\alpha) \wedge C(\vec{X}_1, \vec{K}_\beta, \vec{Y}_\beta) \wedge (\vec{Y}_\alpha \neq \vec{Y}_\beta) \quad (1)$$

$\vec{X}_1$  is called the *Distinguishing Input (DI)*.

- (2) The DI,  $\vec{X}_1$ , is applied to the activated circuit (the oracle) and the output  $\vec{Y}_1$  is recorded. Note that  $\vec{K}_\alpha$ ,  $\vec{Y}_\alpha$ , and  $\vec{K}_\beta$ ,  $\vec{Y}_\beta$  are not recorded. Only the DI and its correct output are carried over to the following iterations.
- (3) In the  $i^{\text{th}}$  iteration, a new DI and a pair of keys,  $\vec{K}_\alpha$  and  $\vec{K}_\beta$ , are found. The newly found  $\vec{K}_\alpha$  and  $\vec{K}_\beta$  should produce correct outputs for all the DIs found in previous iterations. To this end, we append a clause to Eq. (1):

$$C(\vec{X}_i, \vec{K}_\alpha, \vec{Y}_\alpha) \wedge C(\vec{X}_i, \vec{K}_\beta, \vec{Y}_\beta) \wedge (\vec{Y}_\alpha \neq \vec{Y}_\beta) \wedge \bigwedge_{j=1}^{i-1} (C(\vec{X}_j, \vec{K}_\alpha, \vec{Y}_j) \wedge C(\vec{X}_j, \vec{K}_\beta, \vec{Y}_j)) \quad (2)$$

In this way, all the wrong keys that corrupt the output of previously found DIs (*i.e.* the output is different from that of the activated chip) are pruned out from the search space.

- (4) SAT solves Eq. (2) repeatedly until no more DI can be found, *i.e.* Eq. (2) is not satisfiable any more.
- (5) In this case, there is no more DI. The output of the SAT attack is a key  $\vec{K}$  that produces the same output as the activated circuit to all the DIs, which can be expressed using the following CNF:

$$\bigwedge_{i=1}^{\lambda} C(\vec{X}_i, \vec{K}, \vec{Y}_i) \quad (3)$$

where  $\lambda$  is the total number of SAT iterations.

**THEOREM 2.1.** *SAT is guaranteed to find a correct key  $\vec{K}_c$  to the locked circuit.*

The proof is given in Appendix A. Note that there can be multiple correct keys: some keys can be different from but functionally equivalent to the actual key in the activated chip.

### 2.3 Existing Logic Locking Schemes

Multiple logic locking schemes have been proposed to thwart the SAT attack [33, 36, 37]. There are two ways to mitigate the SAT attack: one is to increase the time for each SAT iteration and the other is to increase the number of SAT iterations. The former requires adding SAT-hard circuitry such as AES blocks [36] or permutation blocks [9]. These techniques usually incur huge area overhead which is impractical for most circuits. The other approach is to exponentially increase the number of SAT iterations. This approach is also not perfect because a locking scheme must be rather ineffective to improve security. This is the case for Anti-SAT [32], SARLock [33], and TTLock [34]. All these techniques are vulnerable to the approximate SAT attacks (such as AppSAT [24] and Double-DIP [25]).

The state-of-the-art *stripped functionality logic locking (SFL)* [37] explores the trade-off between SAT resilience and effectiveness. SFL comprises of two parts: a functionality stripped circuit (FSC) and a restore unit (RU). The FSC is the original circuit with the functionality modified for a set of *protected input cubes*. This modification makes SFL resistant to removal attack. If the RU is removed, the FSC's functionality of the protected input cubes is different from the original circuit, thus making the attack unsuccessful. The RU stores the key, compares the circuit's input with the key, and outputs a *restore vector* which is XOR'ed with the FSC output. If the key is correct, the restore vector will fix the FSC's output and the circuit will have correct output. There are two types of SFL: SFL-HD and SFL-flex. SFL-HD leaves very specific structural traces in the FSC which have been successfully captured by a functional analysis based attack [26]. SFL-flex, on the other hand, can leave almost no structural traces in the FSC if a fault-injection-based approach

is taken to strip the functionality [23]. However, very recently, a sensitivity-based approach to identify the protected cubes of SFL-flex is proposed [31]. The sensitivity of a Boolean function on a specific input value measures how likely the output bit is to change when an input bit is flipped. The authors of [31] found that in the ISCAS85 benchmark suite, circuits with more input bits tend to have lower sensitivity on average, based on which they proposed a SAT formulation to find input minterms with high sensitivities in the stripped-functionality circuit. In many benchmarks, the stripped input minterms can be found in this way. In terms of a countermeasure, the authors of [31] proposed to find and strip the functionality of critical minterms which, after inverting their functionality, will have sensitivity values close to the average sensitivity of all the input minterms. Doing so will hide the stripped minterms’ sensitivities among other minterms and hence their sensitivities will no longer be ‘outliers’ and reduce the sensitivity-based attack to a brute-force attack. As SFL-flex remains relatively secure, provides higher flexibility in selecting protected cubes, and is more relevant to SAS, *we focus on SFL-flex* in this paper.

An SFL-flex configuration can be described using the number of protected cubes,  $c$ , and the number of specified bits of each cube,  $k$ , denoted as SFL-flex $^{c \times k}$ . The authors of [37] derived the following characteristics of a circuit locked with SFL-flex $^{c \times k}$ : (1) the fraction of input minterms whose output will be corrupted by a wrong key (*i.e.* the “error rate” of a wrong key) is  $c \cdot 2^{-k}$ ; and (2) the probability that a SAT attack finds the correct key within  $q$  iterations is  $q \cdot 2^{\lceil \log_2 c \rceil - k}$ . We illustrate this relationship in Fig. 2. As a higher SAT success probability indicates weaker SAT resilience, SFL inherently suffers from a trade-off between SAT resilience and effectiveness.

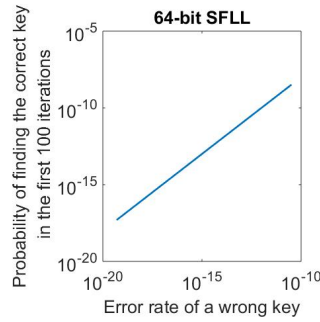


Fig. 2. The positive correlation between the error rate of wrong keys and the probability that SAT finds the correct key in certain iterations

### 3 INSUFFICIENCY OF SFL FOR REAL-WORLD APPLICATIONS

In this section, we investigate the application-level effectiveness of SFL [37]. Specifically, we lock the multiplier within a 32-bit 80386 processor since it is the largest combinational component. The application-level effects are evaluated using both generic and machine learning (ML) benchmarks. *We emphasize ML-based applications because they are inherently error-resilient and hence are more difficult to protect using logic locking.* Details of the benchmarks are listed in Table 1.

Table 1. Application benchmark details

Benchmark Type	Quantity	Content
Generic Applications	9	The PARSEC Benchmark Suite [2]
Machine Learning	5	MNIST [12], SVHN [15], CIFAR10 [11], ILSVRC-2012 [8], Oxford102 [17]

In order to evaluate the application-level impact of a logic locking scheme, we modify the GEM5 [3] simulator to reflect the effects of logic locking. Specifically, before simulation, based on

the key value, the modified GEM5 simulator calculates the set of input minterms whose output will be corrupted according to the logic locking configuration. During the simulation, when the locked module gets any of these input minterms, the simulator randomizes the locked module’s output to emulate locking-induced corruption. Under all other circumstances, the simulator works in the same way as GEM5. In this way, the circuit-level error induced by an incorrect (including approximate) key can be evaluated at the application level. This framework is illustrated in Fig. 3 which is similar to the strategy used in [5, 39].

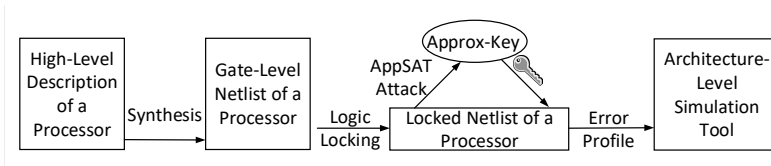


Fig. 3. Our Experimental Framework

SFLL allows the designer to explore the trade-off between effectiveness and SAT resilience. We show that a “sweet spot” does not exist. In our experiments, we lock the multiplier with various SFLL configurations, each having a different level of SAT resilience, quantified by the average SAT iterations to unlock (as the X axis in Fig. 4). The effectiveness of each locking scheme is evaluated by running the PARSEC and ML benchmarks on the locked processors loaded with approximate keys. The percentage of PARSEC benchmark runs with incorrect outcome and the accuracy loss of ML models are used as the criteria to evaluate the effectiveness of each locking configuration. The trade-off is illustrated in Fig. 4.

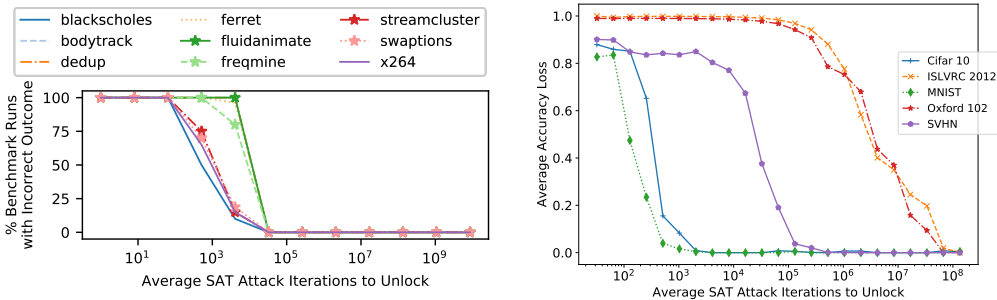


Fig. 4. SAT resiliency vs. locking effectiveness trade-off. Left: PARSEC benchmarks. Right: ML benchmarks.

From Fig. 4, we observe that the wrong keys’ impact decreases with the increase in SAT resiliency. In order to have a visible accuracy drop for the most error-resilient benchmarks, the SFLL locked processor cannot endure more than roughly 1000 SAT iterations. Such a locking scheme is extremely vulnerable since 1000 SAT iterations can be fulfilled within minutes. Therefore, a logic locking scheme that ensures high application-level impact without sacrificing SAT resiliency is needed.

#### 4 FUNDAMENTAL TRADE-OFF FOR ALL LOGIC LOCKING SCHEMES

This section generalizes the trade-off of SFLL to all logic locking schemes. We start with definitions of concepts and then prove the relationship between SAT resilience and effectiveness. Recall that we use  $F(\vec{X})$  to denote the Boolean functionality of the original circuit and  $F_L(\vec{X}, \vec{K})$  to denote the Boolean functionality of the locked circuit, where  $\vec{X}$  is the input minterm and  $\vec{K}$  is the key.

**Definition 4.1.** We say that a key  $\vec{K}$  **corrupts** an input minterm  $\vec{X}$  if and only if the locked circuit's output is different from the original circuit's output when  $\vec{X}$  is the input, i.e.  $F_L(\vec{X}, \vec{K}) \neq F(\vec{X})$ .

**Definition 4.2.** The **key error rate (KER)**  $e_{\vec{K}}$  of a key  $\vec{K}$  is defined as the number of input minterms corrupted by the key  $\vec{K}$  divided by the total number of input minterms.

Let  $\mathcal{X}_{\vec{K}}$  be the set of input minterms corrupted by  $\vec{K}$ . Then,

$$e_{\vec{K}} = \frac{|\mathcal{X}_{\vec{K}}|}{2^n} \quad (4)$$

where  $n$  is the number of bits in the input.

We use  $e_w$  to denote the average KER across all the *wrong keys*.

$$e_w = \frac{1}{|\mathcal{K}^W|} \sum_{\forall \vec{K} \in \mathcal{K}^W} e_{\vec{K}} \quad (5)$$

**Definition 4.3.** The **input error rate (IER)**  $\gamma_{\vec{X}}$  of an input minterm  $\vec{X}$  is the number of wrong keys that corrupt this minterm divided by the total number of wrong keys.

Let  $\mathcal{K}_{\vec{X}}$  be the set of wrong keys that corrupt the input minterm  $\vec{X}$  and  $\mathcal{K}^W$  be the set of all wrong keys. Then,

$$\gamma_{\vec{X}} = \frac{|\mathcal{K}_{\vec{X}}|}{|\mathcal{K}^W|} \quad (6)$$

Let  $\gamma$  denote the average IER over all the input minterms, i.e.

$$\gamma = \frac{1}{2^n} \sum_{\forall \vec{X} \in \{0,1\}^n} \gamma_{\vec{X}} \quad (7)$$

Let us illustrate the above concepts with the following example. We consider a circuit with two input bits ( $x_0, x_1$ ) and locked with a two-bit key ( $k_0, k_1$ ), as shown in Fig. 5. Table 2 is the truth table for each possible input minterm and key. If a key corrupts an input minterm, the corresponding cell is marked with ( $\mathbf{X}$ ).

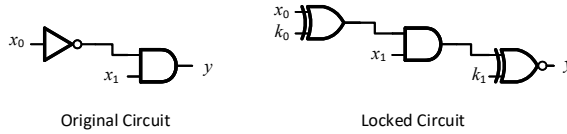


Fig. 5. An example of logic locking, with the original circuit on the left and the locked circuit on the right.

Table 2. Truth table of the locked circuit in Fig. 5

	$\vec{K} = (0, 0)$	$\vec{K} = (0, 1)$	$\vec{K} = (1, 1)$	$\vec{K} = (1, 0)$	Correct $y$	$\gamma_{\vec{X}}$	$\gamma$
$\vec{X} = (0, 0)$	1( $\mathbf{X}$ )	0	0	1( $\mathbf{X}$ )	0	$\frac{2}{3}$	$\frac{2}{3}$
$\vec{X} = (0, 1)$	1	0( $\mathbf{X}$ )	1	0( $\mathbf{X}$ )	1	$\frac{2}{3}$	
$\vec{X} = (1, 1)$	0	1( $\mathbf{X}$ )	0	1( $\mathbf{X}$ )	0	$\frac{2}{3}$	
$\vec{X} = (1, 0)$	1( $\mathbf{X}$ )	0	0	1( $\mathbf{X}$ )	0	$\frac{2}{3}$	
$e_{\vec{K}}$	$\frac{1}{2}$	$\frac{1}{2}$	0	1			
$e_w$	$\frac{2}{3}$						



In Table 2, we also calculate the KER of each key, the IER of each input minterm, and their averages. We can observe that both the average KER and average IER equal  $\frac{2}{3}$ . It turns out that this equality is universal in logic locking:

**THEOREM 4.4.** *The average KER of all wrong keys equals the average IER of all input minterms, i. e.  $e_w = \gamma$ .*

The proof is in Appendix B.

Let  $\lambda$  be the number of SAT iterations that are needed to find the correct key.

**THEOREM 4.5.** *The expected number of SAT iterations  $E[\lambda]$  is lower bounded by  $\frac{1}{\gamma}$ .*

**PROOF.** Recall that  $\gamma$  is the average IER of all input minterms, defined in Equation (7). In each SAT iteration, a distinguishing input (DI) is found and *all* the wrong keys that corrupt this DI will be pruned out from the key search space. The average number of such wrong keys for each DI is  $\gamma|\mathcal{K}^W|$ . Because some of the wrong keys may have already been pruned out in previous iterations, the average number of wrong keys pruned out in each SAT iteration is at most  $\gamma|\mathcal{K}^W|$ . SAT attack finishes when every wrong key is pruned out, hence the average number of SAT iterations can be lower bounded by:

$$E[\lambda] \geq \frac{|\mathcal{K}^W|}{\gamma|\mathcal{K}^W|} = \frac{1}{\gamma} \quad (8)$$

Hence proved.  $\square$

Theorems 4.4 and 4.5 explicitly point out that there exists an inverse relationship between  $e_w$  and the lower bound of  $E[\lambda]$ . This quantifies the trade-off between them. This trade-off applies to any logic locking scheme. Note that different input minterms may inject a different amount of error at the application level. By assigning higher IER to a few minterms with high application-level impact, we can achieve high effectiveness while maintaining high SAT resilience by keeping  $\gamma$  low and  $E[\lambda]$  high. This is the main intuition behind SAS.

## 5 THE ARCHITECTURE AND PROPERTIES OF SAS

Theorems 4.4 and 4.5 have expressed a relationship between input & key error rates and the number of SAT iterations. These quantities are related to two objectives of logic locking:

- (1) **Effectiveness:** Any incorrect key should have a high application-level error impact.
- (2) **SAT resilience:** The complexity of determining the correct key via SAT attacks should be very high.

Assuming that each SAT iteration takes constant time, these two objectives may compete with each other. In this section, we introduce *Strong Anti-SAT (SAS)* logic locking scheme which aims to achieve both objectives simultaneously. SAS guarantees that the expected complexity of SAT attack increases exponentially in the size of the locking key while having a significant impact on the accuracy of real-world applications. In SAS, instead of uniformly distributing the error across all possible inputs, we identify certain input patterns which potentially have a higher impact on the overall application-level error. We call these inputs **critical minterms**. SAS is configured in such a way that any incorrect key corrupts at least 1 critical minterm, resulting in high input error rate (IER) for critical minterms. For the other minterms, the IER is low.

### 5.1 The SAS Block

Let  $\mathcal{M}$  be the set of critical minterms and  $m = |\mathcal{M}|$  be the number of critical minterms. For the ease of implementation, we always choose  $m$  to be a power of 2. The basic locking infrastructure is the SAS block which is illustrated in Fig. 6. When the input minterm  $\vec{X}$  has  $n$  bits, the key  $\vec{K}$  has

$2n$  bits. Let  $\vec{K}_1$  and  $\vec{K}_2$  be the first  $n$  bits and the second  $n$  bits of  $\vec{K}$ , respectively, since they play different roles in SAS. In order to describe the mechanism of the SAS locking scheme clearly, we use a reverse order and start our illustration from the output side.

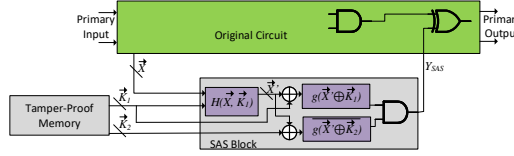


Fig. 6. The Architecture of SAS Configuration 1 with the Details of the SAS Block

$Y_{SAS}$  is the output of the SAS block. If  $Y_{SAS} = 1$ , a fault will be injected into the original circuit.  $Y_{SAS}$  is the output of an AND gate, whose inputs are two function blocks with opposite functionalities, namely  $g$  and  $\bar{g}$ .  $g$  is a function with an on-set-size of 1, *i.e.* only one input minterm will have output 1 and all others will have output 0. We call this very input value  $\vec{X}_g$ , which can be *any*  $n$ -bit value determined by the designer.  $\bar{g}$  has the opposite functionality of  $g$ , *i.e.* only when  $\vec{X} = \vec{X}_g$  will  $\bar{g}(\vec{X}) = 0$ . Due to the AND gate which produces  $Y_{SAS}$ , in order to corrupt an input minterm, the outputs of both  $g$  and  $\bar{g}$  need to be 1.

Each input of  $g$  and  $\bar{g}$  is the output of a XOR or XNOR gate and the designer can choose either one for each gate. We use “ $\oplus$ ” to denote this bit-wise operation. It can be observed that the output AND gate, the  $g$  and  $\bar{g}$  functions, and the  $\oplus$  gates constitute an Anti-SAT block. Notice that, the  $\oplus$  operations before  $g$  and  $\bar{g}$  need not be the same. However, without losing generality, we assume that they are the same for the ease of the following discussions.

A function block  $\vec{X}' = H(\vec{X}, \vec{K}_1)$  is inserted before  $g$  and  $\bar{g}$ . This block will determine the IER of  $\vec{X}$  and it works as follows. If  $\vec{X}$  is not a critical minterm, *i.e.*  $\vec{X} \notin \mathcal{M}$ , then  $\vec{X}' = \vec{X}$  *i.e.* the input minterm simply passes through  $H$ . In this case, a wrong key must satisfy two conditions to corrupt input minterm  $\vec{X}$ :  $\vec{K}_1 = \vec{X} \oplus \vec{X}_g$  and  $\vec{K}_2 \neq \vec{X} \oplus \vec{X}_g$ . Thus there are  $2^n - 1$  wrong keys (1 possible  $\vec{K}_1$  combined with  $2^n - 1$  possible  $\vec{K}_2$ ) that corrupt  $\vec{X}$ , out of the total number of  $2^n(2^n - 1)$  wrong keys. Therefore, the IER of a non-critical input minterm  $\vec{X}$ ,  $\gamma_{\vec{X}}$ , is very low:  $\gamma_{\vec{X}} = \frac{2^n - 1}{2^n(2^n - 1)} = 2^{-n}$ .

If  $\vec{X}$  is a critical minterm, *i.e.*  $\vec{X} \in \mathcal{M}$ , its IER will be  $\gamma_{\vec{X}} = \gamma_c = \frac{1}{m}$ , where  $\gamma_c$  denotes the IER of each critical minterm and  $m$  is the total number of critical minterms. This is achieved as follows. For  $\frac{2^n}{m}$  values of  $\vec{K}_1$ ,  $H(\vec{X}, \vec{K}_1) = \vec{X}_g$ . We use  $\mathcal{K}_{\vec{X}}^1$  to denote this set of  $\vec{K}_1$  values. For other  $\vec{K}_1$  values, the input minterm will still pass through  $H$ , *i.e.*  $H(\vec{X}, \vec{K}_1) = \vec{X}$ . Moreover, the  $\mathcal{K}_{\vec{X}}^1$  sets for each critical minterm are mutually exclusive and evenly partition  $\{0, 1\}^n$ , *i.e.*

$$\forall \vec{X}_1, \vec{X}_2 \in \mathcal{M}, |\mathcal{K}_{\vec{X}_1}^1| = |\mathcal{K}_{\vec{X}_2}^1|, \mathcal{K}_{\vec{X}_1}^1 \cap \mathcal{K}_{\vec{X}_2}^1 = \emptyset, \text{ and } \bigcup_{\forall \vec{X} \in \mathcal{M}} \mathcal{K}_{\vec{X}}^1 = \{0, 1\}^n \quad (9)$$

where  $n$  is the number of bits in  $\vec{X}$ ,  $\vec{K}_1$ , and  $\vec{K}_2$ . Notice that, there are many possible implementations of the partition of  $\vec{K}_1$  space that satisfies Equation (9). For example, if there are two critical minterms  $\vec{X}_1$  and  $\vec{X}_2$ , the designer can let  $\mathcal{K}_{\vec{X}_1}^1$  be the set of  $\vec{K}_1$  values whose most significant bit (MSB) is 1 and  $\mathcal{K}_{\vec{X}_2}^1$  be the set of  $\vec{K}_1$  values whose MSB is 0. In Table 3, we demonstrate how the space of  $\vec{K}_1$  is partitioned to corrupt each critical and non-critical input minterm. In the first row of Table 3, we use the indices to indicate *how many*  $K_1$  values there are that can lead a wrong key to corrupt each input minterm. The indices can accord to any ordering of the elements in  $\{0, 1\}^n$ .

The 2 configurations of SAS will be introduced in the rest of this section.

Table 3. Illustration of how  $m$  critical minterms partition the set of wrong keys

$\vec{k}_1$ of wrong keys		$\vec{k}_1$	$\dots$	$\vec{k}_{\frac{2n}{m}}$	$\vec{k}_{\frac{2n}{m}+1}$	$\dots$	$\vec{k}_{2\frac{2n}{m}}$	$\dots$	$\vec{k}_{2n}$
critical minterms	$\vec{X}_1$	•	•	•					
	$\vec{X}_2$				•	•	•		
	$\dots$						$\dots$		
	$\vec{X}_m$								• •
non- critical minterms	$\vec{X}_{m+1}$	•							
	$\vec{X}_{m+2}$		•						
	$\dots$						$\dots$		
	$\vec{X}_{2n}$								•

## 5.2 Configuration 1: SAS with One SAS Block

This configuration is illustrated in Fig. 6. In this configuration, there is one SAS block. As the critical minterms evenly partition the set of wrong keys, the IER of each critical minterm is  $\gamma_c = \frac{1}{m}$ . Below we derive the SAT resilience of this configuration assuming that the SAT solver chooses a DI uniformly at random in each iteration. This is a common assumption [23, 34, 37]. The SAT resilience is quantified using the expected number of SAT iterations  $E[\lambda]$ . To start with, we give 2 useful lemmas.

LEMMA 5.1. *Let  $\mathcal{D}^i$  be the set of DIs that have been chosen in the first  $i$  iterations and  $\vec{X}$  be a primary input minterm. If all the wrong keys that corrupt  $\vec{X}$  have already been pruned out in the previous SAT iterations, i.e.  $\mathcal{K}_{\vec{X}} \subset \bigcup_{\forall \vec{X}' \in \mathcal{D}^i} \mathcal{K}_{\vec{X}'}$ , then  $\vec{X}$  cannot be the DI of any SAT iteration beyond  $i$ .*

The proof is given in Appendix C.

LEMMA 5.2. *For SAS Configuration 1, any critical minterm must exist in the set of DIs when SAT finishes:  $\vec{X} \in \mathcal{D}^\lambda \forall \vec{X} \in \mathcal{M}$ , where  $\lambda$  is the total number of SAT iterations and  $\mathcal{D}^\lambda$  is the set of all DIs.*

The proof is given in Appendix D.

THEOREM 5.3. *The expected number of SAT iterations of SAS Configuration 1 is*

$$E[\lambda] = \frac{2^n + m}{2} \quad (10)$$

PROOF. The total number of SAT iterations equals the total number of DIs. DIs consist of critical minterms and non-critical minterms. By Lemma 5.2, all the critical minterms must be in the set of DIs for SAT to terminate. Therefore, we only need to find the expected number of *non-critical minterms* that are chosen as DIs. As illustrated in Table 3,  $\forall \vec{X}' \notin \mathcal{M}, \exists$  exactly one  $\vec{X} \in \mathcal{M}$  such that  $\mathcal{K}_{\vec{X}'} \subset \mathcal{K}_{\vec{X}}$ . By Lemma 5.1, if this  $\vec{X}$  is chosen as DI before  $\vec{X}'$ , then  $\vec{X}'$  cannot be chosen in further iterations any more. In other words,  $\vec{X}'$  will count towards the total number of iterations only when it is chosen before the critical minterm  $\vec{X}$ . By our assumption that the DI is chosen uniformly at random in each iteration,  $\vec{X}'$  has a probability of  $\frac{1}{2}$  to be chosen as DI before  $\vec{X}$  is chosen. As this is true for any non-critical minterm, the expected number of SAT iterations is  $E[\lambda] = \frac{1}{2}(2^n - m) + m = \frac{2^n + m}{2}$ .  $\square$

## 5.3 Configuration 2: SAS with Multiple Blocks

In this configuration, we have  $l$  SAS blocks as illustrated in Fig. 7. The  $n$ -bit primary input  $\vec{X}$  is shared among all the SAS blocks and there is a  $2n$ -bit key input for each SAS block. The output of each SAS block is XOR'ed with a wire in the original circuit. Therefore, a fault is injected into the original circuit if any SAS block has output 1. Let  $\mathcal{M}^j$  be the set of critical minterms for the  $j^{\text{th}}$  SAS block,  $j = 1, 2, \dots, l$ . For ease of implementation, we choose  $l$  also to be a power of 2 and  $l \leq m$ .

The relationship between  $\mathcal{M}^j$  and the total set of critical minterms  $\mathcal{M}$  is that  $\mathcal{M}^1, \mathcal{M}^2, \dots, \mathcal{M}^l$  have the same cardinality, are mutually exclusive, and evenly partition  $\mathcal{M}$ , i.e.

$$|\mathcal{M}^1| = |\mathcal{M}^2| = \dots = |\mathcal{M}^l|, \mathcal{M}^i \cap \mathcal{M}^j = \emptyset \forall i \neq j, \text{ and } \bigcup_{k=1}^l \mathcal{M}^k = \mathcal{M} \quad (11)$$

In this way, each SAS block has  $\frac{m}{l}$  critical minterms. As each critical minterm receives high IER from only one SAS block, the IER of any critical minterm is  $\gamma_c = \frac{l}{m}$ .

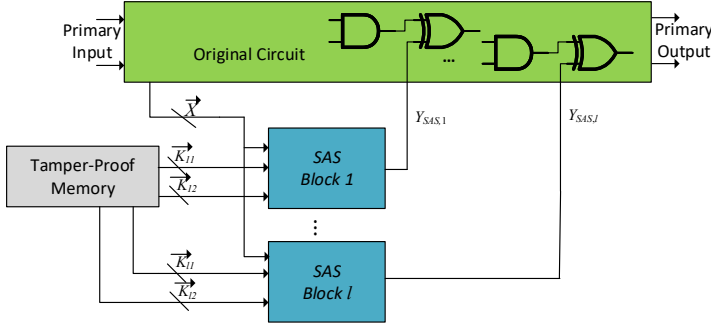


Fig. 7. Configuration 2 with  $l$  SAS blocks

LEMMA 5.4. For SAS Configuration 2, any critical minterm must exist in the set of DIs when SAT finishes:  $\vec{X} \in \mathcal{D}^\lambda \forall \vec{X} \in \mathcal{M}$ , where  $\lambda$  is the total number of SAT iterations and  $\mathcal{D}^\lambda$  is the set of all DIs.

The proof is given in Appendix E. Below, we will analyze the SAT resilience of this configuration by deriving the expected number of SAT iterations.

THEOREM 5.5. The expected number of SAT iterations of SAS Configuration 2 with  $l$  SAS blocks and  $m$  critical minterms is

$$E[\lambda] = \frac{l \cdot 2^n + m}{l + 1} \quad (12)$$

PROOF. By Lemma 5.4, every critical minterm must count toward the total number of SAT iterations. Therefore, we only need to derive the expected number of non-critical minterms that are chosen as DIs.

For any non-critical minterm  $\vec{X}' \notin \mathcal{M}$ , in the  $i^{\text{th}}$  SAS block, there exists exactly one critical minterm  $\vec{X}_i$  such that the set of wrong keys that corrupt  $\vec{X}'$  in this SAS block,  $\mathcal{K}_{i, \vec{X}'}$ , is a subset of the set of wrong keys that corrupt  $\vec{X}_i$ ,  $\mathcal{K}_{i, \vec{X}_i}$ . i.e.  $\mathcal{K}_{i, \vec{X}'} \subset \mathcal{K}_{i, \vec{X}_i}$ . As the construction of the SAS block makes this true for any individual SAS block and the critical minterms for each SAS block are mutually exclusive, there are a total of  $l$  such critical minterms. When all of these  $l$  critical minterms are chosen as DI, they will cover the entire set of wrong keys that corrupt  $\vec{X}'$ . Therefore, by Lemma 5.1, in order to include  $\vec{X}'$  in the set of DIs, it must be selected before all  $l$  critical minterms are selected. This holds for any non-critical minterm.

By our assumption that the DIs are chosen uniformly at random in each SAT iteration, the probability that each non-critical minterm will be chosen as DI is  $\frac{l}{l+1}$ . Therefore, the expected number of SAT iterations is  $E[\lambda] = \frac{l}{l+1}(2^n - m) + m = \frac{l \cdot 2^n + m}{l+1}$ .  $\square$

The properties of both configurations of SAS are summarized in Table 4.

Table 4. Properties of the 2 Configurations of SAS

Configuration	$l$	$\gamma_c$	$E[\lambda]$
1	1	$\frac{1}{m}$	$\frac{2^n+m}{2}$
2	$1 \leq l \leq m$	$\frac{l}{m}$	$\frac{l2^n+m}{l+1}$

## 6 ROBUST SAS: A REMOVAL-RESILIENT SAS VARIANT

Although SAS achieves desirable SAT resilience and high IER on critical minterms, it is still vulnerable to removal attack. In such an attack, the attacker can identify and remove each SAS block and replace their output wires with constant 0. In this way, the remaining part of the locked circuit will have correct functionality. In order to address this drawback, we introduce Robust SAS (RSAS), a variant of SAS that is resilient to removal attacks. In addition to adding an RSAS function block, RSAS modifies the functionality of the original circuit. Therefore, unlike SAS, one cannot obtain the correct functionality of the circuit by identifying and removing the RSAS block. We will introduce the architecture of RSAS and show how any SAS configuration can be converted to a functionally equivalent RSAS configuration. Due to the equivalence in functionality, an RSAS configuration will have the same **SAT resilience** and **effectiveness** as its SAS counterpart.

### 6.1 RSAS Architecture and Relationship with SAS

A circuit locked by RSAS consists of an altered original circuit and one or more RSAS block(s). Fig. 8 illustrates the RSAS configuration with one RSAS block. Given the same set of critical minterms and the same number of locking function blocks, locking a circuit with RSAS and SAS will yield the same functionality. An RSAS-locked circuit can be obtained by converting a functionally equivalent SAS-locked circuit in the following way.

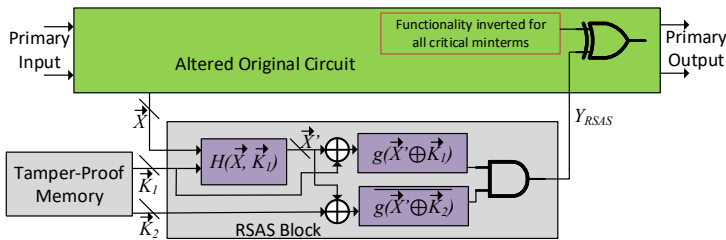


Fig. 8. A circuit locked with one RSAS block, equivalent to SAS Configuration 1

**6.1.1 Altering the original circuit.** Recall that  $l$  is the number of SAS blocks in a SAS configuration. For the  $j^{\text{th}}$  SAS block,  $j = 1, 2, \dots, l$ , the set of critical minterms it contains is denoted by  $\mathcal{M}^j$  and  $|\mathcal{M}^j| = \frac{m}{l}$ , where  $m$  is the total number of critical minterms. In order to implement RSAS, we need to modify the original circuit's functionality. Notice that, for each SAS block, there is a wire in the original circuit that is XOR'ed with the SAS block's output. For the  $j^{\text{th}}$  SAS block, we locate this wire. For each critical minterm in  $\mathcal{M}^j$ , we invert the functionality of critical minterms at this wire. This needs to be done for each  $j$  in  $j = 1, 2, \dots, l$ . This is illustrated in Fig. 9.

**6.1.2 Converting the SAS block into the RSAS block.** The RSAS block is very similar to the SAS block and there is only one difference between them. For the  $j^{\text{th}}$  SAS block,  $j = 1, 2, \dots, l$ , if the primary input is a critical minterm in  $\mathcal{M}^j$ , the output of RSAS block,  $Y_{RSAS,j}$ , is the inversion of the output of SAS block,  $Y_{SAS,j}$ . Recall that, for a SAS configuration with  $m$  critical minterms and  $l$  SAS

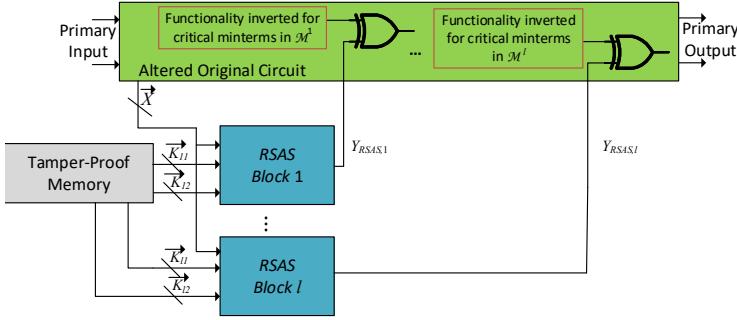


Fig. 9. A circuit locked with multiple RSAS blocks, equivalent to SAS Configuration 2

blocks, each critical minterm's IER is  $\gamma_c = \frac{l}{m}$ . Hence for a portion of  $\frac{l}{m}$  wrong keys,  $Y_{SAS,j}$  is 1. This is achieved by the  $\vec{X}' = H(\vec{X}, \vec{K}_1)$  function: if  $\vec{X}$  is a critical minterm, then the  $H(\vec{X}, \vec{K}_1)$  function makes sure that for  $\gamma_c$  portion of wrong keys, we will have  $g(\vec{X}' \oplus \vec{K}_1) = 1$ . For RSAS, the functionality for critical minterms is inverted, the portion of wrong keys that makes  $Y_{RSAS,j}$  be 1 is  $1 - \gamma_c = \frac{m-l}{m}$ . This means the functionality of  $H(\vec{X}, \vec{K}_1)$  needs to be modified in the following way: if  $\vec{X}$  is a critical minterm, then for  $1 - \gamma_c$  portion of wrong keys,  $g(\vec{X}' \oplus \vec{K}_1)$  will output 1. For non-critical input minterms,  $Y_{RSAS}$  behaves in the same as  $Y_{SAS}$ . This is illustrated in Table 5.

Table 5. Illustration of RSAS block's functionality. A '•' stands for  $Y_{RSAS} = 1$ .

$\vec{K}_1$ of wrong keys	$\vec{k}_1$	...	$\vec{k}_{\frac{2n}{m}}$	$\vec{k}_{\frac{2n}{m}+1}$	...	$\vec{k}_{2\frac{2n}{m}}$	...	$\vec{k}_{2n}$
critical minterms	$\vec{X}_1$			•	•	•		•
	$\vec{X}_2$	•	•	•				•
	...	•	•	•	•	•		•
	$\vec{X}_m$	•	•	•	•	•		•
non-critical minterms	$\vec{X}_{m+1}$	•						
	$\vec{X}_{m+2}$		•					
	...				...			
	$\vec{X}_{2n}$							•

## 6.2 SAT Resilience and Effectiveness of RSAS

In Sec. 6.1, we introduced how to convert a SAS-locked circuit into an equivalent RSAS-locked circuit. These steps essentially invert the functionality of each critical minterm at two places: the first at the wire in the original circuit where RSAS is integrated, and the other at the RSAS block's output. Since these two wires are XOR'ed, the two inversions will cancel out which makes the RSAS-locked circuit functionally equivalent to the SAS-locked circuit. Due to the equivalence in functionality, the derivations of SAS's *SAT resilience* and *effectiveness* will also hold for RSAS. Therefore, Table 4 is also the summary of these properties of RSAS.

## 7 CHOOSING CRITICAL MINTERMS

The critical minterms for injecting large errors should be selected judiciously. A careful analysis of the workload would help identify these typical minterms. Generally these minterms would be very few as compared to the overall input space of the functional modules. Here we describe how to select the critical minterms. As mentioned in Sec. 3, we use PARSEC (generic) and ML models as application benchmarks. For the PARSEC benchmarks, we arbitrarily choose critical minterms from the input minterms that exist in all the application benchmarks. We take a similar approach

for ML benchmarks. A significant part of an ML-based application is the parameters of the ML model and it turns out that the parameter values of most ML models follow a similar distribution. For example, Figure 10 shows the distribution of parameters of the LeNet (MNIST dataset) and CaffeNet (ISLVRC-2012 dataset) models. These two are the smallest benchmark and the largest benchmark, respectively. The parameter distributions are similar across ML benchmarks and many other ML models. This kind of similarity can be also found among generic applications.

We select a subset of parameter values to be critical minterms based on their application-level impact. The selected critical minterms should cause significant application-level error. Fig. 10 also shows the accuracy loss of the ML model in the following experiment: for each input minterm, we measure the accuracy loss of the ML model when every computation involving this very minterm is corrupted while no other minterm is corrupted. As the input minterm distributions are similar

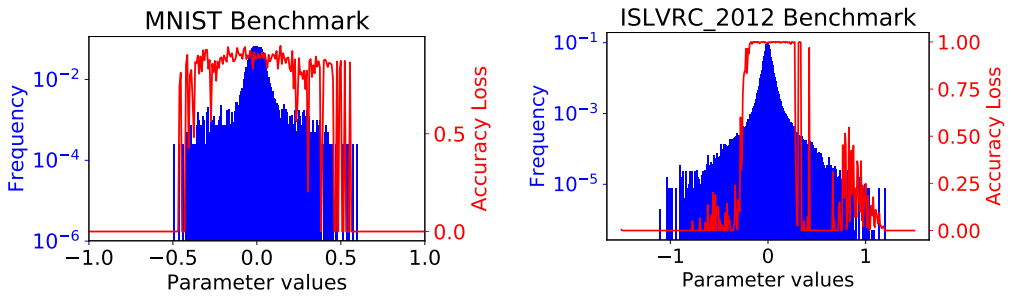


Fig. 10. Weight Distribution (Blue Histogram, Left Y Axis) and Application-Level Accuracy Loss (Red Line, Right Y Axis) of LeNet and CaffeNet when the corresponding inputs were locked

among the same type of applications, the flexibility of SAS/RSAS allows the designer to choose a configuration and a combination of critical minterms that work well in securing the intended applications without compromising SAT resilience.

It should be assumed that the application programs (written in binary code) are public knowledge according to the Kerckhoff's principle. An attacker may be tempted to take advantage of the above-mentioned strategy to identify the critical minterms. Since the binary code does not contain the locked module's input values directly, the attacker needs to observe the inputs of the locked module from the scan chain. There are scan chain obfuscation techniques which will corrupt the scan output unless the correct scan chain key is provided [10, 19, 38]. Alternatively, the designer can develop software-based defenses to detect malicious scan chain access or simply burn the scan chain before deploying the activated chips to the open market. Therefore, we do not consider this kind of attack as an immediate security threat to our technique.

## 8 EXPERIMENTS & COMPARISON WITH SFLL

This section shows the experimental results of SAS and RSAS as well as the comparison with SFLL. Recall that, as illustrated in Fig. 3, we obtain the gate-level netlists of the multiplier within a 32-bit 80386 processor by synthesizing the high-level description using Cadence RTL Compiler. Then we lock the netlist using various SAS and RSAS configurations and SFLL-flex with the same set of critical minterms. Note that the critical minterms are selected according to the method described in Sec. 7. The architecture-level simulation is conducted by a modified GEM5 [3] simulator where error is injected into the locked processor module according to the hardware error profile due to the wrong key. We conduct the following experiment to verify the SAT resilience and effectiveness of SAS and RSAS and compare them with SFLL.

## 8.1 SAT Resilience

We first verify whether the SAT resilience of SAS/RSAS (*i.e.* the actual number of SAT iterations) matches what we have derived in Sec. 5. The SAT resilience of SAS/RSAS and SFL is also compared. We lock the multiplier in a 32-bit 80386 processor with SAS and RSAS as well as SFL. Fig. 11 shows the actual and expected number of SAT iterations of multipliers locked with SAS and RSAS. These numbers are compared to the actual number of iterations of SFL. In these locking configurations, we use 14 bits of primary input for locking purposes ( $n = 14$ ) and experiment with each feasible configuration with up to 4 critical minterms. We can observe that SAS and RSAS have similar numbers of actual SAT iterations and they are both close to the expected value. When there is more than one critical minterms, SAS and RSAS exhibit higher SAT resilience than SFL. This is because the corruptibility of each critical minterm in SFL is almost 1 no matter how many critical minterms there are. This compromises its SAT resilience.

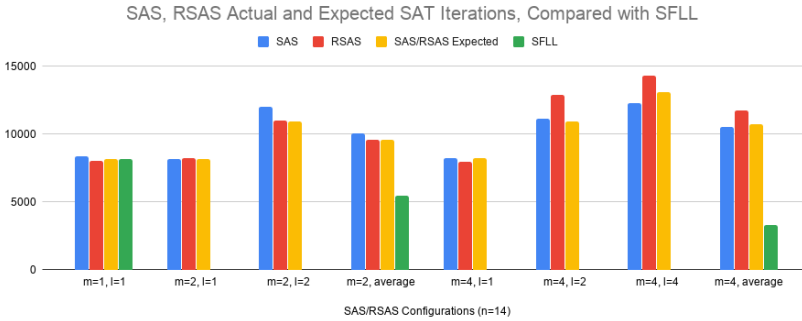


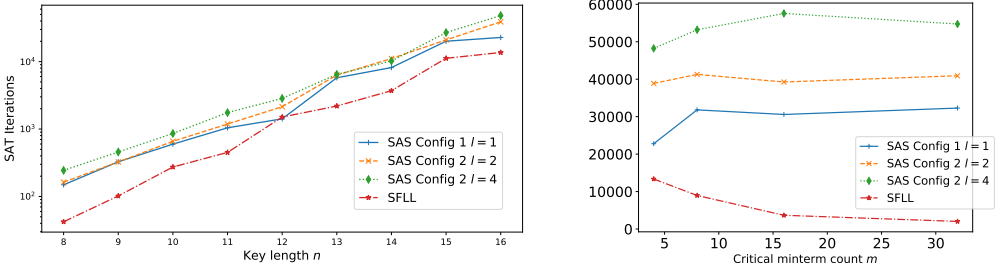
Fig. 11. Actual and expected number of SAT iterations of SAS and RSAS, compared with SFL.

Fig. 12 compares the actual SAT iterations of SAS and SFL. In Fig. 12a, it can be observed that SAS's SAT complexity is higher than that of SFL by a roughly constant factor when  $m$  is fixed at 4. Note that the same set of four critical minterms are used for each locking scheme. Among various SAS configurations, a larger  $l$  comes with higher SAT resilience as expected. In Fig. 12b, we vary the critical minterm count ( $m$ ) from 4 to 32 and demonstrate its impact on the SAT resilience of SAS and SFL. While SAS configurations become stronger with more critical minterms, SFL becomes weaker. Therefore, SAS is more SAT resilient and gives designers more flexibility when more critical minterms are needed.

## 8.2 Effectiveness

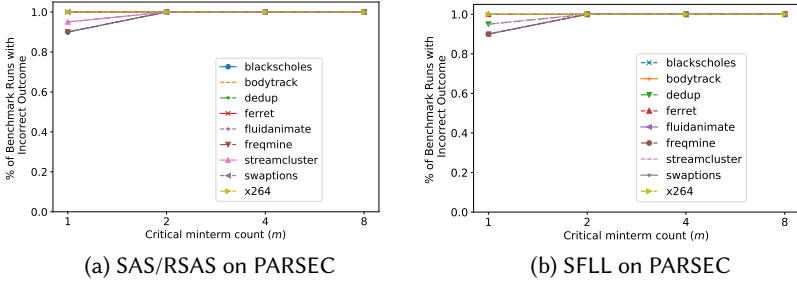
We evaluate the effectiveness of SAS/RSAS and SFL at the application level using PARSEC [2] and ML benchmarks as listed in Table 1. Due to the functional equivalence of SAS and RSAS, they will have the same architecture-level effects and we use the same functional model to perform architecture-level simulation of SAS and RSAS. In our experiments, various numbers of critical minterms are locked. The same set of critical minterms are used for SAS/RSAS and SFL in each experiment. The critical minterms are chosen according to the methods described in Sec. 7. For SAS, we choose  $l = 1$  when  $m = 1$  and  $l = 2$  when  $m \geq 2$ . Figs. 13 and 14 show that both SAS/RSAS and SFL are effective at the application level for both generic and ML-based applications. SAS/RSAS achieves high application-level effectiveness and exponential SAT resiliency at the same time. Considering that SAS/RSAS's SAT resilience is not compromised with the increase in  $m$  as opposed to SFL (as shown in Figs. 11 and 12b), SAS/RSAS is a significant improvement over SFL.





(a) Varying key length ( $n$ ), fixing # critical minterms  $m = 4$  (b) Varying # critical minterms ( $m$ ), fixing key length  $n = 16$

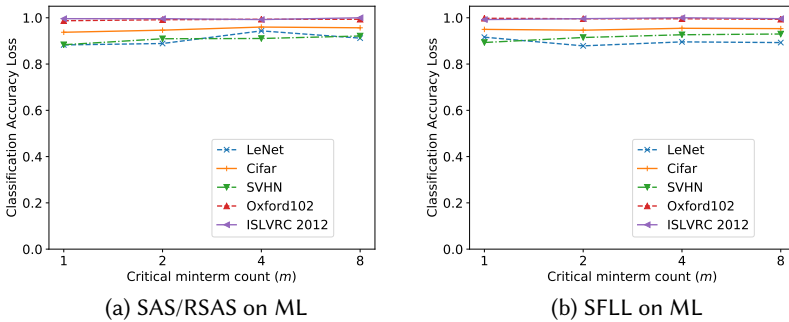
Fig. 12. The observed SAT iterations of SAS and SFL by varying key length and critical minterm count.



(a) SAS/RSAS on PARSEC

(b) SFL on PARSEC

Fig. 13. The application-level effectiveness of SAS/RSAS and SFL on PARSEC and ML benchmarks



(a) SAS/RSAS on ML

(b) SFL on ML

Fig. 14. The application-level effectiveness of SAS/RSAS and SFL on PARSEC and ML benchmarks

8.3 Area, Power, and Delay Overhead of SAS, RSAS, and SFL

Now that we have demonstrated the SAT resilience of SAS and RSAS and their application-level effectiveness, we evaluate their area, power, and delay overhead. The overhead is also compared with SFL. In our evaluation, we use 32 bits from the primary input for locking ( $n = 32$ ) and lock up to 4 critical minterms ( $m = 1, 2, 4$ ). We synthesize the original and locked circuits using Cadence RTL Compiler using SAED 90nm process. In order to account for the area of tamper-proof memory

that store keys of SAS and look-up tables of SFL, we add  $2.80\mu\text{m}^2$  for each memory bit to the area of locked designs. This is scaled from [23, 37] which used a 65nm library and reported  $1.46\mu\text{m}^2$  per bit of tamper-proof memory for the square of the feature size, *i.e.*  $1.46\mu\text{m}^2 \times (\frac{90\text{nm}}{65\text{nm}})^2 = 2.80\mu\text{m}^2$ . Figs. 15, 16, and 17 show the area, power, and delay overhead values, respectively. Compared with SFL, on average, SAS and RSAS have 1.90% and 0.73% more area overhead, 0.43% more and 0.04% less power overhead, 0.93% and 0.71% more delay overhead, respectively. These are not significant increases in overhead and should be worth the gain in SAT resilience.

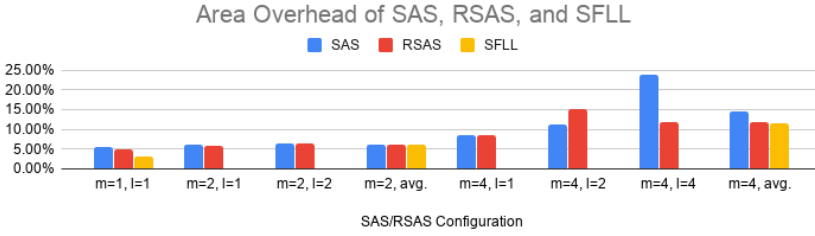


Fig. 15. Area overhead of SAS and RSAS compared with SFL

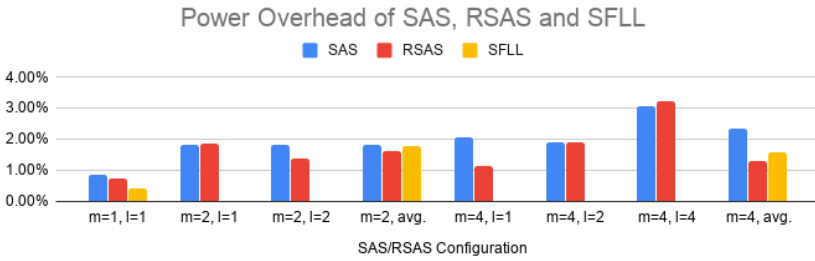


Fig. 16. Power overhead of SAS and RSAS compared with SFL

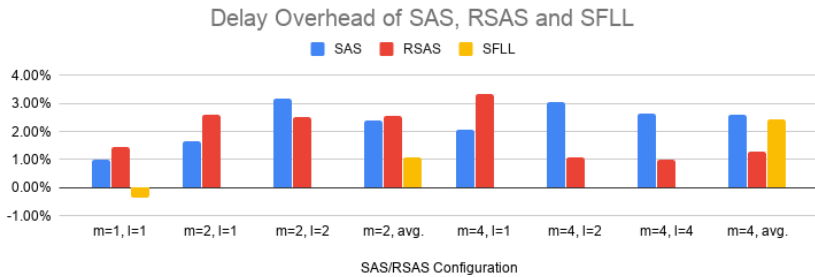


Fig. 17. Delay overhead of SAS and RSAS compared with SFL

## 9 CONCLUSION

In this work, we investigate logic locking techniques to secure both generic and error-resilient workloads running on locked processors. We motivate our work by demonstrating the insufficiency of the state-of-the-art logic locking scheme in securing such applications. We point out that this is

due to the fundamental trade-off between *SAT resilience* (SAT attack complexity) and *effectiveness* (error rate of wrong keys) of logic locking. We formally prove this trade-off. In order to address this dilemma, we propose Strong Anti-SAT (SAS) where a set of critical minterms are assigned higher corruptibility in order to ensure high application-level impact. Based on SAS, we also propose Robust SAS (RSAS) to thwart removal attacks on logic locking. RSAS is functionally equivalent to SAS and has the same SAT resilience and effectiveness. Experimental results show that SAS and RSAS secure processors against SAT attack by ensuring exponential SAT attack complexity and high application-level impact simultaneously given any wrong key. We also evaluate the area, power, and delay overhead of SAS and RSAS and compare it with SFL. It is shown that SAS and RSAS have modest increase in overhead. In summary, RSAS exhibits a higher SAT resilience than SFL when multiple critical minterms are secured, while also maintaining equivalent effectiveness and removal attack resilience. Therefore, RSAS constitutes a significant improvement over SFL-based locking.

## ACKNOWLEDGMENTS

This work is supported by AFOSR MURI under Grant FA9550-14-1-0351 and Northrop Grumman Corporation and University of Maryland Seedling Grant.

## REFERENCES

- [1] Abdulrahman Alaql, Domenic Forte, and Swarup Bhunia. 2019. Sweep to the Secret: A Constant Propagation Attack on Logic Locking. In *2019 Asian Hardware Oriented Security and Trust Symposium (AsianHOST)*. IEEE, 1–6.
- [2] Christian Bienia, Sanjeev Kumar, Jaswinder Pal Singh, and Kai Li. 2008. The PARSEC benchmark suite: Characterization and architectural implications. In *Proceedings of the 17th international conference on Parallel architectures and compilation techniques*. ACM, 72–81.
- [3] Nathan Binkert, Bradford Beckmann, Gabriel Black, Steven K Reinhardt, Ali Saidi, Arkaprava Basu, Joel Hestness, Derek R Hower, Tushar Krishna, Somayeh Sardashti, et al. 2011. The gem5 simulator. *ACM SIGARCH Computer Architecture News* 39, 2 (2011), 1–7.
- [4] Abhishek Chakraborty, Nithyashankari Gummidipoondi Jayasankaran, Yuntao Liu, Jeyavijayan Rajendran, Ozgur Sinanoglu, Ankur Srivastava, Yang Xie, Muhammad Yasin, and Michael Zuzak. 2020. Keynote: A Disquisition on Logic Locking. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 39 (October 2020), 1952–1972. Issue 10.
- [5] Abhishek Chakraborty, Yang Xie, and Ankur Srivastava. 2018. GPU obfuscation: attack and defense strategies. In *Proceedings of the 55th Annual Design Automation Conference*. ACM, 122.
- [6] Prabhuddha Chakraborty, Jonathan Cruz, and Swarup Bhunia. 2018. SAIL: Machine learning guided structural analysis attack on hardware obfuscation. In *2018 Asian Hardware Oriented Security and Trust Symposium (AsianHOST)*. IEEE, 56–61.
- [7] Prabhuddha Chakraborty, Jonathan Cruz, and Swarup Bhunia. 2019. SURF: Joint structural functional attack on logic locking. In *2019 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE Computer Society, 181–190.
- [8] J Deng, A Berg, S Satheesh, H Su, A Khosla, and L Fei-Fei. 2012. ILSVRC-2012, 2012. URL [http://www.image-net.org/challenges/LSVRC\(2012\)](http://www.image-net.org/challenges/LSVRC(2012)).
- [9] Hadi Mardani Kamali, Kimia Zamiri Azar, Houman Homayoun, and Avesta Sasan. 2019. Full-Lock: Hard Distributions of SAT Instances for Obfuscating Circuits using Fully Configurable Logic and Routing Blocks. In *Proceedings of the 56th Annual Design Automation Conference 2019*. ACM, 89.
- [10] Rajit Karmakar and Santanu Chattopadhyay. 2020. On Securing Scan Obfuscation Strategies Against ScanSAT Attack. In *2020 21st International Symposium on Quality Electronic Design (ISQED)*. IEEE, 213–218.
- [11] Alex Krizhevsky and Geoffrey Hinton. 2009. Learning multiple layers of features from tiny images. (2009).
- [12] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.
- [13] Yuntao Liu, Michael Zuzak, Yang Xie, Abhishek Chakraborty, and Ankur Srivastava. 2020. Strong Anti-SAT: Secure and Effective Logic Locking. In *Twenty-first International Symposium on Quality Electronic Design*. IEEE, 199–205.
- [14] Mohamed El Massad, Jun Zhang, Siddharth Garg, and Mahesh V Tripunitara. 2017. Logic locking for secure outsourced chip fabrication: A new attack and provably secure defense mechanism. *arXiv preprint arXiv:1703.10187* (2017).

- [15] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. 2011. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, Vol. 2011. 5.
- [16] Xuan Thuy Ngo, Jean-Luc Danger, Sylvain Guilley, Tarik Graba, Yves Mathieu, Zakaria Najm, and Shivam Bhasin. 2017. Cryptographically Secure Shield for Security IPs Protection. *IEEE Trans. Comput.* 66, 2 (2017), 354–360.
- [17] M-E. Nilsback and A. Zisserman. 2008. Automated Flower Classification over a Large Number of Classes. In *Proceedings of the Indian Conference on Computer Vision, Graphics and Image Processing*.
- [18] Stephen M Plaza and Igor L Markov. 2015. Solving the Third-Shift Problem in IC Piracy with Test-aware Logic Locking. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on* 34, 6 (2015), 961–971.
- [19] M Sazadur Rahman, Adib Nahiyani, Sarah Amir, Fahim Rahman, Farimah Farahmandi, Domenic Forte, and Mark Tehranipoor. 2019. Dynamically Obfuscated Scan Chain To Resist Oracle-Guided Attacks On Logic Locked Design. *LACR Cryptol. ePrint Arch.* 2019 (2019), 946.
- [20] Jeyavijayan Rajendran, Youngok Pino, Ozgur Sinanoglu, and Ramesh Karri. 2012. Security analysis of logic obfuscation. In *Proceedings of the 49th Annual Design Automation Conference*. ACM, 83–89.
- [21] Jeyavijayan Rajendran, Huan Zhang, Chi Zhang, Garrett S Rose, Youngok Pino, Ozgur Sinanoglu, and Ramesh Karri. 2015. Fault Analysis-Based Logic Encryption. *Computers, IEEE Transactions on* 64, 2 (2015), 410–424.
- [22] Jarrod A Roy, Farinaz Koushanfar, and Igor L Markov. 2008. EPIC: Ending piracy of integrated circuits. In *Proceedings of the conference on Design, Automation and Test in Europe*. ACM, 1069–1074.
- [23] Abhrajit Sengupta, Mohammed Nabeel, Nimisha Limaye, Mohammed Ashraf, and Ozgur Sinanoglu. 2020. Truly Stripping Functionality for Logic Locking: A Fault-based Perspective. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* (2020).
- [24] Kaveh Shamsi, Meng Li, Travis Meade, Zheng Zhao, David Z Pan, and Yier Jin. 2017. AppSAT: Approximately deobfuscating integrated circuits. In *2017 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE, 95–100.
- [25] Yuanqi Shen and Hai Zhou. 2017. Double dip: Re-evaluating security of logic encryption algorithms. In *Proceedings of the on Great Lakes Symposium on VLSI 2017*. ACM, 179–184.
- [26] Deepak Sirone and Pramod Subramanyan. 2020. Functional analysis attacks on logic locking. *IEEE Transactions on Information Forensics and Security* 15 (2020), 2514–2527.
- [27] Dominik Šišejković, Farhad Merchant, Rainer Leupers, Gerd Ascheid, and Sascha Kegreiss. 2019. Inter-lock: Logic encryption for processor cores beyond module boundaries. In *2019 IEEE European Test Symposium (ETS)*. IEEE, 1–6.
- [28] Dominik Šišejković, Farhad Merchant, Lennart M Reimann, Rainer Leupers, Massimiliano Giacometti, and Sascha Kegreiß. 2020. A secure hardware-software solution based on RISC-V, logic locking and microkernel. In *Proceedings of the 23th International Workshop on Software and Compilers for Embedded Systems*. 62–65.
- [29] Dominik Šišejković, Farhad Merchant, Lennart M Reimann, Rainer Leupers, and Sascha Kegreiß. 2020. Scaling Logic Locking Schemes to Multi-module Hardware Designs. In *International Conference on Architecture of Computing Systems*. Springer, 138–152.
- [30] Pramod Subramanyan, Sayak Ray, and Sharad Malik. 2015. Evaluating the security of logic encryption algorithms. In *Hardware Oriented Security and Trust (HOST), 2015 IEEE International Symposium on*. IEEE, 137–143.
- [31] Joseph Sweeney, Marijn Heule, and Lawrence T Pileggi. 2020. Sensitivity Analysis of Locked Circuits.. In *23rd International Conference on Logic for Programming, Artificial Intelligence and Reasoning*. 483–497.
- [32] Yang Xie and Ankur Srivastava. 2018. Anti-SAT: Mitigating SAT Attack on Logic Locking. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* (2018).
- [33] Muhammad Yasin, Bodhisatwa Mazumdar, Jeyavijayan JV Rajendran, and Ozgur Sinanoglu. 2016. SARLock: SAT attack resistant logic locking. In *Hardware Oriented Security and Trust (HOST), 2016 IEEE International Symposium on*. IEEE, 236–241.
- [34] Muhammad Yasin, Bodhisatwa Mazumdar, Jeyavijayan JV Rajendran, and Ozgur Sinanoglu. 2017. TTLock: Tenacious and traceless logic locking. In *2017 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE, 166–166.
- [35] Muhammad Yasin, Bodhisatwa Mazumdar, Ozgur Sinanoglu, and Jeyavijayan Rajendran. 2017. Removal attacks on logic locking and camouflaging techniques. *IEEE Transactions on Emerging Topics in Computing* (2017).
- [36] Muhammad Yasin, Jeyavijayan JV Rajendran, Ozgur Sinanoglu, and Ramesh Karri. 2016. On improving the security of logic locking. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 35, 9 (2016), 1411–1424.
- [37] Muhammad Yasin, Abhrajit Sengupta, Mohammed Thari Nabeel, Mohammed Ashraf, Jeyavijayan JV Rajendran, and Ozgur Sinanoglu. 2017. Provably-Secure Logic Locking: From Theory To Practice. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 1601–1618.
- [38] Dongrong Zhang, Miao He, Xiaoxiao Wang, and Mark Tehranipoor. 2017. Dynamically obfuscated scan for protecting IPs against scan-based attacks throughout supply chain. In *2017 IEEE 35th VLSI Test Symposium (VTS)*. IEEE, 1–6.

[39] M. Zuzak and A. Srivastava. 2019. Memory Locking: An Automated Approach to Processor Design Obfuscation. In *2019 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*. 541–546. <https://doi.org/10.1109/ISVLSI.2019.00103>

## A PROOF OF THEOREM 2.1

PROOF. This can be proved by contradiction: suppose the key returned by the last step of SAT attack is a wrong key. This implies that there must exist an input minterm  $\vec{X}$  such that

$$C(\vec{X}, \vec{K}_c, \vec{Y}_c) \wedge C(\vec{X}, \vec{K}, \vec{Y}) \wedge (\vec{Y}_c \neq \vec{Y})$$

where  $\vec{K}$  is the actual key,  $\vec{Y}_c$  is the output with returned key  $\vec{K}_c$  and  $\vec{Y}$  is the correct output according to the actual key  $\vec{K}$ .  $\vec{X}$  cannot be a previously found DI because otherwise  $\vec{K}_c$  will not satisfy (3). We can see that  $\vec{X}$  qualifies for a DI: just assign  $\vec{K}_\alpha = \vec{K}_c$  and  $\vec{K}_\beta = \vec{K}$ . This means that (2) is still satisfiable and contradicts the criteria that no more DI can be found before the SAT attack goes to the final step.

Hence proved.  $\square$

## B PROOF OF THEOREM 4.4

PROOF. Recall that

$$e_w = \frac{1}{|\mathcal{K}^W|} \sum_{\forall \vec{K} \in \mathcal{K}^W} e_{\vec{K}} = \frac{1}{|\mathcal{K}^W|} \sum_{\forall \vec{K} \in \mathcal{K}^W} \frac{|\mathcal{X}_{\vec{K}}|}{2^n} = \frac{1}{2^n |\mathcal{K}^W|} \sum_{\forall \vec{K} \in \mathcal{K}^W} |\mathcal{X}_{\vec{K}}|$$

and

$$\gamma = \frac{1}{2^n} \sum_{\forall \vec{X} \in \{0,1\}^n} \gamma_{\vec{X}} = \frac{1}{2^n} \sum_{\forall \vec{X} \in \{0,1\}^n} \frac{|\mathcal{K}_{\vec{X}}|}{|\mathcal{K}^W|} = \frac{1}{2^n |\mathcal{K}^W|} \sum_{\forall \vec{X} \in \{0,1\}^n} |\mathcal{K}_{\vec{X}}|$$

Therefore, in order to prove  $e_w = \gamma$ , we only need to prove

$$\sum_{\forall \vec{K} \in \mathcal{K}^W} |\mathcal{X}_{\vec{K}}| = \sum_{\forall \vec{X} \in \{0,1\}^n} |\mathcal{K}_{\vec{X}}| \quad (13)$$

Let us consider the following bipartite graph  $G = (\mathcal{X}, \mathcal{K}^W, \mathcal{E})$  where  $\mathcal{X}$  is  $\{0,1\}^n$  which is the set of all the possible input minterms,  $\mathcal{K}^W$  is the set of wrong keys, and the set of edges  $\mathcal{E} = \{(\vec{X}, \vec{K}) | \vec{X} \in \mathcal{X} \text{ and } \vec{K} \in \mathcal{K}^W, \vec{K} \text{ corrupts } \vec{X}\}$ . Both sides of Eq. 13 denote the total number of edges in  $\mathcal{E}$  and hence must be equal.  $\square$

## C PROOF OF LEMMA 5.1

PROOF. Recall that Equation (2) gives the SAT formula for each SAT iteration:

$$C(\vec{X}_i, \vec{K}_\alpha, \vec{Y}_\alpha) \wedge C(\vec{X}_i, \vec{K}_\beta, \vec{Y}_\beta) \wedge (\vec{Y}_\alpha \neq \vec{Y}_\beta) \\ \bigwedge_{j=1}^{i-1} (C(\vec{X}_j, \vec{K}_\alpha, \vec{Y}_j) \wedge C(\vec{X}_j, \vec{K}_\beta, \vec{Y}_j))$$

To satisfy the first line, at least one of  $\vec{K}_\alpha$  and  $\vec{K}_\beta$  must be a wrong key that corrupts  $\vec{X}$ . However, if any wrong key that corrupts  $\vec{X}_i$  also corrupts at least 1 previously found DI, this wrong key cannot satisfy the second line. Therefore, such  $\vec{X}_i$  cannot be the DI in future iterations.  $\square$

#### D PROOF OF LEMMA 5.2

PROOF. Recall that  $g$  has on-set size 1. Let  $\vec{X}_g$  be the very input that makes  $g(\vec{X}_g) = 1$ .  $\forall \vec{X} \in \mathcal{M}$ , let  $\vec{K}_1 = \vec{X} \oplus \vec{X}_g$ . Then, any  $\vec{K} = (\vec{K}_1, \vec{K}_2) \in \mathcal{K}^W$  is a wrong key that only corrupts  $\vec{X}$ . Therefore,  $\vec{X}$  has to be chosen as a DI to prune out this wrong key.  $\square$

#### E PROOF OF LEMMA 5.4

PROOF. This is a natural extension to Lemma 5.2. Let  $\vec{X}$  be a critical minterm and  $\vec{X} \in \mathcal{M}^j$ . Recall that  $g$  has on-set size 1. Let  $\vec{X}_g$  be the very input that makes  $g(\vec{X}_g) = 1$ .  $\forall \vec{X} \in \mathcal{M}^j$ , let  $\vec{k} = \vec{X} \oplus \vec{X}_g$ . Then, let us consider the following wrong key  $\vec{K} = (\vec{K}^1, \vec{K}^2, \dots, \vec{K}^l) \in \mathcal{K}^W$  which is composed as follows:  $\vec{K}^j = (\vec{k}, \vec{K}_2^j) \in \mathcal{K}_j^W$  where  $\mathcal{K}_j^W$  is the set of wrong keys for the  $j^{\text{th}}$  SAS block. For any  $i = 1, 2, \dots, l$  that  $i \neq j$ ,  $\vec{K}^i \in \mathcal{K}_i^C$  where  $\mathcal{K}_i^C$  is the set of correct keys for the  $i^{\text{th}}$  SAS block. Such a key  $\vec{K}$  is a wrong key that only corrupts  $\vec{X}$ . Therefore,  $\vec{X}$  has to be chosen as a DI to prune out this wrong key.  $\square$