

Deceptive Logic Locking for Hardware Integrity Protection against Machine Learning Attacks

Dominik Sisejkovic, Farhad Merchant, Lennart M. Reimann, and Rainer Leupers

Abstract—Logic locking has emerged as a prominent key-driven technique to protect the integrity of integrated circuits. However, novel machine-learning-based attacks have recently been introduced to challenge the security foundations of locking schemes. These attacks are able to recover a significant percentage of the key without having access to an activated circuit. This paper address this issue through two focal points. First, we present a theoretical model to test locking schemes for key-related structural leakage that can be exploited by machine learning. Second, based on the theoretical model, we introduce D-MUX: a deceptive multiplexer-based logic-locking scheme that is resilient against structure-exploiting machine learning attacks. Through the design of D-MUX, we uncover a major fallacy in existing multiplexer-based locking schemes in the form of a structural-analysis attack. Finally, an extensive cost evaluation of D-MUX is presented. To the best of our knowledge, D-MUX is the first machine-learning-resilient locking scheme capable of protecting against all known learning-based attacks. Hereby, the presented work offers a starting point for the design and evaluation of future-generation logic locking in the era of machine learning.

Index Terms—hardware security, logic locking, attack resilience, hardware integrity, machine learning.

I. INTRODUCTION

Logic Locking (LL) is a premier technique to safeguard the integrity of hardware designs throughout the Integrated Circuit (IC) supply chain [1], [2]. This obfuscation technique performs functional and structural design alterations through the insertion of key-dependent logic. Nevertheless, the security of LL has been challenged in the past decade through different key-recovery attacks. While all attacks assume the existence of the design in locked netlist format, the differentiating factor is the availability of an *oracle*, i.e., an activated IC instance. The oracle is used to acquire golden Input/Output (I/O) patterns. Thus, all attacks can be classified into Oracle-Guided (OG) [3] and Oracle-Less (OL) attacks [4]–[8]. The OG model is manifested in high-volume commercial fabrication, where it is assumed that an attacker is able to purchase an unlocked instance of the design [9]. In comparison, the OL scenario is relevant in a low-volume setup where an attacker is not able to get a copy of the activated IC. This is, for example, the case in the development of security-critical systems with unique and highly confidential hardware requirements [10], [11].

Moreover, recent findings strongly suggest that the OL model plays a more important role in realistic attack scenarios. These include the following. (i) The secret key can be exposed

by hardware Trojans that leak the key once the IC is activated regardless of the LL scheme [12]. (ii) Various OG attacks rely on having access to the internal states of a design; otherwise, the attacks become impractical. However, this assumption is unrealistic, since legitimate IC vendors never leave a scan-chain open and typically use some form of authentication [13]. (iii) Recent works have demonstrated the extraction of keys from activated ICs even in the presence of a tamper- and read-proof memory through probing or fault-injection attacks [14]–[16]. These observations exclude the necessity for some of the most efficient OG attacks as the key is retrievable from the *activated* IC regardless of the LL scheme.

Furthermore, it is important to understand the exact security implications of LL. First, it can render a fabricated IC inoperable, as the correct operation is only ensured if the activation key is provided. However, as discussed in [17], achieving inoperability can be done without LL. Hence, the primary objective of LL is to ensure concealing the design’s functionality. This security property is known as *functional secrecy* [17]. In this context, even though we measure the success of key-recovery attacks in terms of key accuracy, acquiring the key itself is not the main threat model. Since the design concept of locking ensures that the key impacts the functional and structural alteration of the design, the key-accuracy metric acts as a *proxy* to measuring functional secrecy. Hereby, security can be defined in the form of Exact Functional Secrecy (EFS), Approximate Functional Secrecy (AFS), and Best-Possible Approximate Functional Secrecy (BPAFS); for both the OG and OL model [17]. EFS embodies the security against a *perfect* reconstruction of the design under attack. In that regard, AFS offers a stricter criterion; it requires *approximation* resiliency, i.e., security against reconstructing a design that disagrees with the original up to a selected ϵ fraction of its input space. BPAFS represents a relaxed form of AFS, where the attacker has a certain *a priori* knowledge about the circuit’s size and depth. Evidently, AFS-OG implies AFS-OL and EFS-OG implies EFS-OL, respectively. Furthermore, AFS implies EFS; and BPAFS-OL is equivalent to AFS-OL.

Note that functional secrecy can be provably achieved. For example, an adaptation of the SFLL scheme is provably secure under EFS-OG/OL. However, BPAFS-OG/OL has only been achieved through universal circuits [17]–[19], which suffer from impractical overheads. Therefore, achieving BPAFS-OL with traditional, low-cost LL is still an open problem. Hence, this goal remains in the main focus of this work.

Moreover, with the proliferation of Machine Learning (ML) across various domains, ML is slowly being introduced into hardware security as well. A few recent works have evaluated

ML models for attacking LL [8], [20]–[24]. Nevertheless, so far, there has been a gap in designing LL that is resilient against ML-based attacks as well as in the theoretical means of uncovering the source of ML-exploitable leakage.

Contributions: To fill the theoretical and practical gap on ML resilience in an OL setting, in this work, we design a learning-resilient locking scheme—from theory to practice. Hereby, we use the term *learning resilience* to evaluate schemes against learning-based attacks that capture *locking-induced structural residue* that can lead to *key-related information leakage*. The contributions are as follows:

- The introduction of the first theoretical concept to test learning resilience in logic locking.
- The Deceptive Multiplexer (D-MUX) logic locking; a novel scheme based on MUX insertion that is resilient against existing ML-based attacks.
- We introduce a novel oracle-less Structural Analysis Attack on Multiplexer-based locking (SAAM) based on a major pitfall in existing MUX-based logic locking.
- We empirically evaluate the resilience of D-MUX against the oracle-less attacks SAAM, SWEEP, and SnapShot.
- We evaluate the cost of D-MUX in terms of area, power, and delay through a theoretical and empirical assessment.
- Finally, based on the lessons learned, we analyze the structural aspects of related work. Hereby, we identify novel structural leakage points in existing schemes.

To the best of our knowledge, this work is the first to shed light on the theoretical concepts of learning resilience in logic locking as well as propose the first empirically evaluated LL scheme that is resilient against learning-based attacks.

This work is organized as follows. Section II introduces the background on LL. Learning resilience is discussed in Section III. The structural-analysis attack is described in Section IV. Section V introduces D-MUX. The resilience and cost evaluation of D-MUX is presented in Section VI and Section VII, respectively. The related work is summarized in Section VIII. Limitations and opportunities are discussed in Section IX. Finally, Section X concludes the paper.

II. PRELIMINARIES

Logic locking can be categorized into two groups: combinational and sequential [25]. Combinational LL focuses on manipulations in the combinational path of circuits, whereas sequential LL obfuscates the state transition graph of a design. In the following, we only focus on combinational LL.

A. Logic Locking

To showcase the working principles of LL, let us consider EPIC [26]—one of the first LL schemes. EPIC is based on the insertion of key-controlled XOR and XNOR (XOR + INV) gates (known as *key gates*) at randomly selected locations in the gate-level design. An EPIC-locked circuit is shown in Fig. 1. In this example, the original circuit (Fig. 1 (a)) is locked through the insertion of an XOR+INV gate bound to the key input k_1 (Fig. 1 (b)). Only if $k_1 = 1$, the key gates act as a *buffer*, thereby restoring the original functionality of the circuit. Otherwise, for a $k_1 = 0$, the key gates disrupt

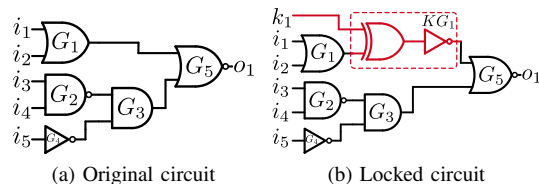


Fig. 1: Example: logic locking using EPIC.

the intended functionality by *inverting* the output of gate G_1 , leading to faulty output values. In terms of EPIC, it is assumed that a simple removal of the key gates is prevented since an adversary must guess if the inverter is part of the key gate or the original functionality. In the past years, a wide range of combinational LL schemes has been introduced based on XOR/XNOR gates [26], [27], multiplexers [28]–[31], AND/OR gates [32], and others [19], [33]–[35].

Logic Locking in the IC Supply Chain: The IP owner is the trusted entity that is introducing a legitimate product on the semiconductor market. The design is either done in-house or parts of the design services, e.g., placement and routing, are subcontracted to an external design house. In both cases, the final layout is sent to the foundry. Both the external design house and the foundry are considered untrusted entities.

The IP owner can utilize LL to conceal the design after the first logic synthesis round. Thus, LL is applied to the synthesized *gate-level netlist*. After LL, the netlist is typically *resynthesized* to integrate the induced changes. Hence, we can differentiate the *pre-resynthesis* and *post-resynthesis* netlist. After fabrication, the IC is returned to the IP owner for activation. The key is configured through a non-volatile memory. Since the key is only known to the IP owner, LL protects the design while in untrusted hands. Henceforth, the term *netlist* refers to a gate-level netlist.

Attack Model: The attack model includes the following assumptions. (i) The adversary has only access to the locked netlist (OL model). (ii) The adversary is aware of the algorithmic details of the LL scheme. (iii) The location of the key inputs in the netlist is known. In the rest of this work, we refer to the netlist under attack as the *target* netlist.

B. Oracle-less Attacks

The following provides an overview of existing OL attacks. Note that even though these attacks can achieve high accuracy, they do not defeat locking policies that are provably secure in the OG model [17].

The desynthesis attack [6] tries to extract the secret key through a Hill Climbing search of a series of new synthesis rounds using a randomly selected key. Hereby, the search is guided by the similarity between the locked and resynthesized netlist. However, this attack is not scalable to large keys and assumes a robust understanding of the synthesis tool.

The redundancy attack [7] attempts to recover the correct key by pruning out incorrect key values when these introduce a significant level of functional redundancy in the netlist. Similarly, the Topology-Guided Attack (TGA) exploits the fact that basic functions in a logic cone are often repeated multiple times in a netlist [5]. Thereby, TGA can recover the correct key

of a particular key gate by looking at equivalent functions that are constructed using all possible hypothesis key values. Even though both are OL, these attacks are designed and evaluated to exploit XOR/XNOR-based LL schemes. Moreover, resilient LL against both attacks has been proposed in [5], [36].

SAIL [23] deploys ML algorithms to retrieve the key-driven local structures of the pre-resynthesis netlist based on the post-resynthesis target. Once the reconstruction is done, SAIL extracts the correct key solely based on the properties of XOR/XNOR-based LL. Thus, it is not relevant for evaluation in this work. Moreover, recently, the UNSAIL locking technique has been proposed to thwart the SAIL attack [37].

In the context of MUX-based locking and ML-based attacks, two attacks play a major role: **SWEEP** and **SnapShot**. Both are described in more detail in the following.

1) **SWEEP**: This constant propagation attack exploits the differentiating circuit characteristics induced by hard-coding a single key-bit value during resynthesis [4]. An overview of the subsequent steps of the attack is presented in the following.

Training Set Generation: The training set is assembled by either using a set of available locked benchmarks with known key values or by relocking the target netlist with new keys.

Feature Extraction: The training data is resynthesized by hard-coding *each key input* to the correct and incorrect value (logic 0 or 1). This procedure results in two synthesis reports for each key value. The relevant features are extracted from the reports and consolidated in the form of a feature matrix.

Feature Weighting: In this stage, the attack generates a set of optimum weights that evaluate the correlation between the entries in the feature matrix and the known correct keys.

Deployment: Finally, the design features are extracted from the target by repeating the previous steps for an unknown key. To extract the key value, SWEEP utilizes the generated weights from the previous step and the target features to deduce a value for each key input. The value can be 0, 1, or X. If X is given, the attack was not able to make a safe decision about the correct value. Moreover, SWEEP allows an adjustable margin m to be selected by the attacker. m controls the freedom of SWEEP to make "wild" guesses for values that are similar (close to each other). By default, $m = 0$.

Evaluation Metrics: SWEEP relies on two metrics: *accuracy* and *precision*. Accuracy is defined as the percentage of correctly extracted key-bit values out of the entire key, *regardless of potential X values*, i.e., $(N_{correct}/N_{total}) \cdot 100\%$. Precision is defined as the percentage of correct keys, where every potential X value is regarded as a correct guess, i.e., $((N_{correct} + N_X)/N_{total}) \cdot 100\%$.

Note that SWEEP is relevant for evaluation in this work as it is an OL attack and well suited for MUX-based locking.

2) **SnapShot**: The SnapShot attack utilizes an ML model to predict key values based on structural features extracted from the target netlist [8]. Compared to the mentioned attacks, SnapShot has the ability to make a *direct* key guess by "just" looking at the post-resynthesis netlist, without having to extract synthesis features or reconstruct the pre-resynthesis netlist. The attack flow is described in the following.

Training Set Generation: Similar to SWEEP, the training set can be generated in two ways, resulting in the Generalized

Set Scenario (GSS) and the Self-Referencing Scenario (SRS). In GSS, the ML model learns to predict key values based on a generalized set of locked circuits (not including the target). In SRS, a novel set is generated by copying the target multiple times and relocking each copy with an additional key.

Extraction: In this step, SnapShot extracts the *key-affected* netlist subgraphs that serve as training samples for the ML model. For each key input, the extraction procedure follows the key wire until a gate is encountered. This gate is considered as the central key gate. Starting from this gate, the netlist is traversed in a Breadth-First Search (BFS) fashion towards the primary I/Os. While traversing, each gate is mapped to an integer value based on an encoding table. The extracted representation is known as Locality Vector (LV).

ML Model Design: The model is trained based on the labeled LVs to predict key values for unseen, target LVs. In principle, any ML model can be used in SnapShot. However, due to the BFS-nature of the extraction, the LVs embody the structural (image) representation of the subgraphs; thus making them suitable for processing with Convolutional Neural Networks (CNNs). However, since this prediction problem is of a novel nature, the original work deploys neuroevolution to automatically design suitable CNN architectures. In addition, a shallow artificial neural network is evaluated for comparison.

Deployment: Once the model is trained, SnapShot is deployed to predict the key of the locked netlist. For this task, the *unlabeled* locality vectors of the target are extracted using the same procedure as for the training set. Finally, the target localities are presented to the ML model for key prediction.

Evaluation Metric: SnapShot uses the Key Prediction Accuracy (KPA) metric for attack evaluation. *The KPA is equivalent to the accuracy definition used in SWEEP except that all bits are always guessed (no X values are allowed)*.

Note that other attacks utilize ML-models as well; however, only in the form of OG attacks [20]–[22].

III. THE CONCEPT OF LEARNING RESILIENCE

To evaluate LL for learning resilience, we analyze what structural changes are induced by a scheme, thereby considering two netlist variants. The first variant includes netlists that only consist of a *single gate type*. The second variant covers netlists that consist of a randomly selected and well distributed amount of *all gate types*. These variants represent two ends of the spectrum of possible netlist structures: regular and irregular. *Note that the regularity describes the repetition of equivalent logic structures throughout the netlist*. The rationale of looking at these variants is that a structural key-related leakage is likely to repeat for similar structures in a design. In theory, any netlist can be placed between the two variants. For example, designs that exhibit very regular and repeating structures are closer to the first variant. Examples include sbox, adder, multiplier, and decoder (tree of multiplexers) implementations. On the other hand, a very irregular structure (e.g., specific control logic) is closer to the second variant. The success of various attacks that exploit this regularity suggests that hardware designs are mostly regular and closer to the first case. Using this spectrum, we can devise two theoretical tests

for learning resilience. The first variant is represented by the AND Netlist Test (ANT) and the second variant is represented by the Random Netlist Test (RNT). Hereby, the concept is that ANT and RNT evaluate for potential structure-related leakage at the *two ends of the spectrum* (regular and irregular).

A. Learning-Resilience Test

The core idea of a learning-resilience test is to offer a means to uncover potential design flows of LL that lead to exploitable structural information leakage. Hereby, the test is designed around the two spectrum ends. Both ANT and RNT follow the *same test procedure*. The only difference is in the type of netlists used in the process. The test is manifested in the form of a learning game between two parties: Trusted (T) and Untrusted (U). The goal of U is to learn how to predict the key based on the locked netlists provided by T. *If a locking scheme fails at least one test, i.e., U is able to make an educated guess about the key (guessing accuracy is higher than 50%), the locking scheme is conclusively vulnerable, otherwise it might be learning resilient.* Therefore, these tests can only make a concrete decision about whether a scheme is *not resilient*. One iteration of the game is presented in Fig. 2, consisting of the following steps:

- 1) T randomly generates a netlist (Net) that consists of any number of gates, where each node can have any number of input or output connections.
- 2) T locks the netlist (Net_L) with a selected scheme using one randomly selected key (Key).
- 3) T sends the locked netlist to U for analysis.
- 4) U guesses the key based on the collected observations.
- 5) U sends the guessed key (Key_G) to T for comparison.
- 6) T responds whether the guess is correct or not for each key bit value individually. Note that T does not provide the reason why a bit is correctly or falsely guessed.

Hereby, we assume the following: (i) T always uses the same LL scheme and (ii) the original netlist is different in each iteration. Hereby, (i) ensures that U is able to learn about the LL scheme, otherwise T could always select different schemes. (ii) prevents learning from *different key values* for *identical locations* in a netlist. Note that in this game, U does not have to know the exact type of the used LL scheme, just that it is always the same one. The game is repeated until U terminates the game or achieves a desired success rate.

At the end of an iteration, U is able to update his knowledge with new observations taken from the response provided by T. These observations are used to make guesses about the key in new iterations. If the LL scheme leaks information, U is able to extract meaningful observations leading to the ability to make educated guesses. If the guessing accuracy of U continuously increases, then the scheme-under-attack leaks information. If the guessing accuracy remains 50%, i.e., equal to a random guess, then the collected observations exhibit no evident correlation to the correct key. However, note that this does not prove the absence of leakage, as previously discussed.

Furthermore, it is noteworthy that the only knowledge U has about the netlist is that its original form contains either AND gates exclusively or randomly selected gates. However,

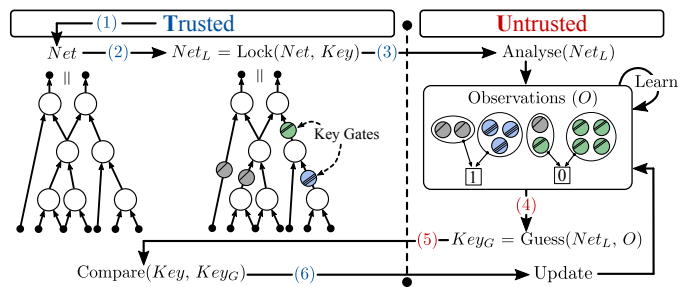


Fig. 2: Learning-resilience test.

U does not know how the gates are initially connected. This is not in line with an actual attack scenario, since an adversary should not be informed about the original gate types of a design. However, this limitation has a favorable connotation; it amplifies the changes induced by an LL scheme. *Therefore, any conclusions made about a scheme in ANT or RNT hold also in a general, for the attacker less invasive case when the adversary is not aware of the original structure of the netlist.*

1) *ANT*: In ANT, T generates netlists consisting exclusively of AND gates, representing the regular end of the spectrum. This feature maximizes the exposure of the inner workings of the LL mechanism, since predictable selections or changes can more easily be spotted in repeating structures. *Evidently, for this particular test, any type of gate would suffice.* In practice, to avoid a potential bias of an LL scheme toward specific gate types, the test should be repeated for all primary types.

2) *RNT*: In RNT, T generates a netlist consisting of a variety of randomly selected gate types. Here, a particular LL scheme might not leak, i.e., reveal information about the key due to the nature of the netlist. For example, the fact that all gate types are present in the netlist might have a favorable implicating on the security of the scheme. If a scheme fails at least one test, it is regarded as conclusively vulnerable, since its resilience depends on the structure of the netlist rather than the key. Moreover, evaluating a scheme through both tests can reveal the cause for information leakage and suggest a specific setting in which resiliency might be achieved.

3) *Termination Condition*: A concrete termination condition cannot be defined in advance. The exact number of iterations depends on the ability of U to extract usable observations, as exemplified in Section III-B and Section III-C.

4) *Resynthesis and Security*: Typically, resynthesis is performed after locking is done to further integrate any LL-induced changes. However, as discussed in [8], *an LL scheme must not depend on synthesis to obtain a sufficient security level.* Otherwise, the scheme clearly leaks information if resynthesis is not performed, as its security *depends* on specific synthesis transformations (often part of proprietary software). Moreover, this dependency has an important implication: the security of the system would not solely depend on the key. This can enable novel attacks; for example, the attacker can exploit this dependency by reverse engineering specific synthesis transformations to neutralize the security aspects of a scheme, as performed in the mentioned SAIL attack [23]. To overcome this requirement, the proposed tests *do not include resynthesis.*

5) *Circuit Family*: Both tests operate on two spectrum ends. In practice, designs might be placed somewhere between

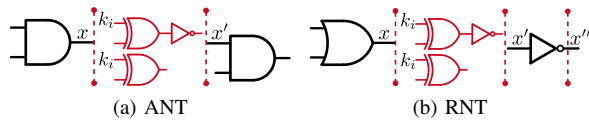


Fig. 3: ANT and RNT for XOR/XNOR-based locking.

these extremes. In this case, a locked netlist might *not* exhibit leakage under suitable structural conditions; implying that the scheme has no identifiable leakage *for this one specific case*. Thus, its security would be predicated on the structural characteristics, i.e., the circuit family, of the underlying netlist. Hence, security conclusions about the challenged scheme cannot be generalized across all netlist structures. Evidently, this must not be the objective for the design of generally secure LL. Therefore, it is beneficial to generalize the evaluation by looking at the spectrum ends. Nevertheless, designing LL that is specific to a circuit family is a valid objective. In this context, both ANT and RNT can help uncover *for which circuit family the evaluated scheme might exhibit resilience*. This allows for the design of LL schemes which are structurally secure under the assumption of a certain circuit family; which can be a reasonable use case in practice. Examples of such cases are presented in Section III-B and III-C.

B. Testing the XOR/XNOR Scheme

Let us first consider EPIC as a representative of XOR/XNOR-based locking. In EPIC, the location selection is random. However, EPIC has been successfully attacked with ML-based attacks [23]. Therefore, EPIC is leaking information through the introduced change, as argued in [8].

1) *ANT Observations*: In ANT, **T** locks netlists that consist exclusively of AND gates. After many iterations, **U** is able to learn that an XOR implies 0, while XNOR implies 1. This can be understood through the visualization in Fig. 3 (a). In ANT, **U** is aware that the original netlist exclusively contains AND gates. Thus, it is possible for **U** to isolate the *exact* location that is affected by an XOR or XNOR between two AND gates. In the given example, it is evident that the value x entering a key gate must be preserved once processed. Therefore, the adversary has to guess the value of k_i that ensures this preservation, i.e., $k_i = ? \Rightarrow x = x'$. Due to the nature of the scheme and ANT, it is always true that $k_i = 0$ for XOR and $k_i = 1$ for XNOR.

Verdict: **U** is able to easily learn about the key by isolating the limited region around a key gate. *Therefore, EPIC fails ANT and is considered vulnerable to learning-based attacks.*

2) *RNT Observations*: In RNT, **T** generates a netlist using a variety of gate types for each iteration. Thus, the possibility exists that an XOR key gate is placed before an original inverter. This leads to the occurrence of observations where an XOR+INV is associated both with the key value 0 and the value 1. In other words, sometimes it is more difficult for **U** to correctly isolate the key gates, as shown in Fig. 3 (b). Here, it is not clear where the correct cut-off line for the restoration of x is drawn, i.e., whether $x = x'$ or $x = x''$. Therefore **U** is more likely to make *false guesses* when these situations occur.

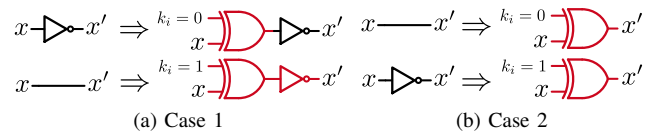


Fig. 4: XOR/XNOR learning resilience cases.

Verdict: In RNT, **U** can sometimes make a false guess. However, since **T** generates netlists with an evenly distributed number of gate types, these contradicting cases do not occur often. Consequently, it is expected that the guessing accuracy in RNT is higher than 50%. Therefore, *EPIC fails RNT*.

Discussion: Based on the tests, we can conclude that XOR/XNOR LL might exhibit learning resilience only in *very specific* cases, including: XOR key gates are only placed in front of existing inverters (Fig. 4 (a)) or XOR gates replace existing inverters for key bit 1, i.e., XNOR gates are not used (Fig. 4 (b)). This concept has been implemented in Truly Random LL (TRLL) [38]. Even though both cases lead to contradicting observations, the learning resilience is only enabled by the *inherent structural features* of the design. Therefore, TRLL fails ANT since no inverters are available to be replaced or coupled with an XOR gate; disintegrating TRLL into traditional XOR/XNOR LL. However, TRLL passes RNT since a sufficient amount of gate types is available to support the scheme. These examples also suggest how to achieve learning resilience; *equivalent local regions are labeled with the same key-bit*. Here, the attacker has a 50% chance to guess correctly. Nevertheless, it still remains a challenge to support this concept and satisfy both tests using XOR/XNOR gates.

C. Testing the Twin-Gate Scheme

To present other aspects of the tests, we propose a new LL scheme dubbed *Twin-Gate* (Algorithm 1). Twin-Gate receives the correct key and a netlist as input. First, Twin-Gate prepares the set of all gate types with more than one input (T_{multi} , line 1) and the set of all single-input gate types (T_{single} , line 1). Afterwards, for each key bit (line 2), Twin-Gate performs three steps: the random selection of a *true* node (line 3), the creation of a *false* node (lines 5-11), and the assembly of a replacement MUX (lines 13-14). The true node represents an original node from the netlist. The false node is the additional *twin* that is added as a pair to the true node. Depending on the type of the selected true node, the algorithm selects a suitable false node. The suitability is defined by ensuring that both nodes have the same number of inputs and a different type. Therefore, if the true node is of type INV or BUF, its twin can only be BUF or INV (line 6). In case the true node is a multi-input node, the algorithm randomly selects a type from the available gate types (line 8). Once the false node is created (line 11), a MUX is assembled by coupling its inputs to the true and false node outputs (line 13). The current key input ($K[i]$) acts as selection bit and determines which gate output is forwarded. A locked example is shown in Fig. 5 (a) based on the original netlist shown in Fig. 1 (a). Here, the true node G_3 is replaced by the pair (KG_1, KG_2) . The correct gate output is selected by the multiplexer through the key bit k_1 .

Algorithm 1 Twin-Gate locking scheme

```

Input: Activation key  $K$ , netlist  $Net$ 
Output: Locked netlist
1:  $T_{multi} \leftarrow \{\text{AND, OR, } \dots, \text{XOR}\}$ ;  $T_{single} \leftarrow \{\text{NOT, BUF}\}$ 
2: for  $i = 0$  to  $|K|$  do
3:    $n_{true} \leftarrow \text{RandomSelection}(Net)$ 
4:   /* Create false node */
5:   if  $\text{TypeOf}(n_{sel}) \in \{\text{NOT, BUF}\}$  then
6:      $F_{types} \leftarrow \{t \in T_{single} \mid t \neq \text{TypeOf}(n_{true})\}$ 
7:   else
8:      $F_{types} \leftarrow \{t \in T_{multi} \mid t \neq \text{TypeOf}(n_{true})\}$ 
9:   end if
10:   $f_{type} \leftarrow \text{RandomSelection}(F_{types})$ 
11:   $n_{false} \leftarrow \text{CreateNodeOfType}(f_{type})$ 
12:  /* Assemble MUX node */
13:   $m_{enc} \leftarrow \text{CreateMUXFor}(n_{true}, n_{false}, K[i])$ 
14:  ReplaceNodes $(n_{true}, m_{enc})$ 
15: end for
16: return  $Net$ 

```

1) *ANT Observations:* In ANT, Twin-Gate resolves the *direct leakage* problem of XOR/XNOR-based locking, as now \mathbf{U} is not able to guess anything about the nature of the key based on the induced gate types if only considering the local change. Let us consider the example in Fig. 5 (b). Here, Twin-Gate replaces an original AND gate with the pair (OR, AND). Due to the nature of Twin-Gate, it is possible to isolate which gates are part of the LL scheme. In this case, identifying the correct key k_i is equivalent to correctly guessing whether $x' = x_1$ or $x' = x_2$. Here the adversary is not able to guess whether AND or OR is the true node by analyzing the samples individually, since both are valid options. However, it turns out that \mathbf{U} is able to learn about the key based on the type and distribution of *all* key gates in the netlist. For example, after many iterations of the game, \mathbf{U} can learn that AND is *never* a false gate, as otherwise two ANDs are fed to a MUX; implying that any key bit value is correct. Since AND is never a false gate, \mathbf{U} can easily learn which key bit is correct for a given pair. This effect is amplified by the nature of AND-based netlists that exclude AND gates as possible false gates.

Verdict: Twin-Gate fails ANT as the selection of false nodes directly depends on the original netlist structure. However, RNT uncovers specific cases where the scheme might be learning resilient, as discussed in the following.

2) *RNT Observations:* Since each iteration of RNT is based on a netlist with randomly selected gate types, \mathbf{U} has a disadvantage in guessing the correct key value. Now, every gate type is equally likely to be selected as true or false node, making it difficult to extract meaningful observations.

Verdict: In theory, Twin-Gate passes RNT. However, in practice, typical circuits do not have a perfectly balanced amount of each gate type, making it possible to identify potential false nodes. Moreover, the presented discussion showcases how even a global leakage can be identified in LL schemes using the proposed tests. Hereby, global refers to the leakage being uncovered due to a frequency analysis over all gates.

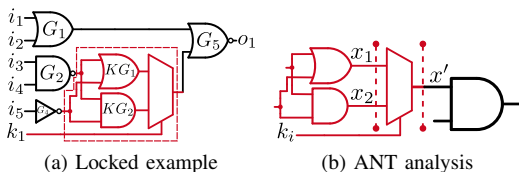


Fig. 5: Twin-Gate locking scheme.

D. Lessons Learned

Based on the analysis, we can conclude that a learning-resilient scheme must have the following properties: (i) its security must not depend on resynthesis, (ii) the induced change must not depend on the inherent structural features of the original netlist, and (iii) the induced change and its location should not depend on the value of the connected key bit. Hereby, (iii) is manifested in the following; the observations made by \mathbf{U} must be either random (noisy data) or two identical observations must suggest two different key values (e.g., two identical "images" are labeled differently). For example, XOR/XNOR-based locking is vulnerable because, in most cases, it cannot fulfill the latter.

Testing whether a human adversary (or, ultimately, an ML system) is able to challenge a scheme through *learning-based attacks* is a fundamentally important step in designing LL, as it tests the basic security foundation; *can the adversary guess the key based on the induced change?* Therefore, the introduced theoretical tests offer a simple yet powerful approach to uncover elemental security flaws in LL, *even before an implementation or empirical results are available.*

IV. SAAM: STRUCTURAL ANALYSIS ATTACK ON MUX-BASED LOGIC LOCKING

Through ANT and RNT, we identified an important vulnerability of MUX-based LL that can determine the correct key bit value for a single MUX. The MUX takes two inputs: the true (T) and false (F) wire. Both values of the key bit result in a functionally valid netlist. However, by analyzing both gates that drive T and F, we can easily identify the correct true wire as follows. Let us consider the example shown in Fig. 6 (a). Here, the original netlist from Fig. 1 is locked with a MUX. In this setting, the adversary can create two cases: for each value of the selecting key bit, the MUX is removed and the selected wire is forwarded. If the MUX selects T as the correct wire (Fig. 6 (b)), the resulting netlist displays no structural faults. However, if the MUX selects F as the correct wire (Fig. 6 (c)), the wire T (output of gate G2) remains an *unconnected* (dangling) wire. Moreover, this is *never* the case for the wire F, as it is randomly selected out of the other wires in the netlist. Thereby, the gate output that drives F (gate G1) is *always* connected to some other gate as well (gate G5). *Therefore, the MUX input that remains unconnected when not selected is the true wire.* To the best of our knowledge, this major design fallacy in MUX-based locking has not been documented properly yet; therefore, we summarize it in the form of a simple attack named SAAM; **Structural Analysis Attack on MUX-based locking**. SAAM is presented in Algorithm 2. SAAM iterates through each key input and checks which input of the connected MUX has one

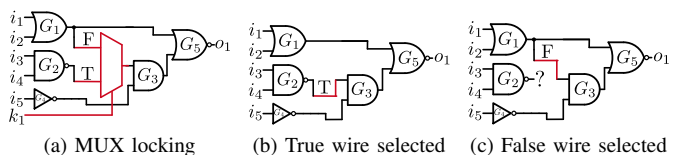


Fig. 6: Example: MUX-based locking and SAAM.

TABLE II: Cost summary for D-MUX locking strategies.

Locking Strategy	Number of Gates	Min. Key Bits per Iteration
S_1	$l \cdot MUX $	2
S_2, S_3	$l \cdot MUX $	1
S_4	$l \cdot (2 \cdot MUX)$	1

outputs, the case when one of the input nodes is not forwarded to an output via the MUX remains valid.

3) S_3 - *OneMultiOutOneBitOneMux*: In this strategy, only one multi-output node f_i is selected (OneMultiOut) to be locked using a single key bit k_i (OneBit) driving one MUX (OneMux), as presented in Fig. 7 (c). S_3 enables two configurations as shown in Table I (c). Even though these configurations are, in principle, equivalent to the previous strategy, S_3 differs from S_2 only in the fact that the output node must be selected from the multi-output input node. Hereby, regardless of which key bit is chosen, both input nodes remain connected.

4) S_4 - *AnyOutOneBitTwoMux*: This strategy selects two input nodes $\{f_i, f_j\}$ from the set of all available nodes (AnyOut). Afterwards, a single key bit k_i is selected (OneBit) to drive two MUXs simultaneously (TwoMux). For each input node, one output node is selected. The MUXs are configured to always forward the results of both input gates (otherwise SAAM is applicable in case a single-output gate is selected as input), as presented in Fig. 7 (d). S_4 allows two configurations as shown in Table I (d). Regardless of the key-bit value, all nodes remain connected; thereby forming a valid path.

Interestingly, driving each MUX in S_4 with an individual key is also a viable option. However, this configuration is susceptible to a simple reduction attack in case the input nodes are of the single-output type, thereby enabling SAAM. In this specific case, for $\{k_i, k_j\} = \{0, 1\}$ the node f_j remains unconnected, whereas for $\{k_i, k_j\} = \{1, 0\}$ the node f_i remains unconnected. Therefore, the adversary just has to guess whether $\{k_i, k_j\}$ is $\{0, 0\}$ or $\{1, 1\}$. This is equivalent to using S_4 with only one key. Note that S_4 is the only strategy which is always applicable.

B. Cost Model

The cost of a particular locking strategy can be expressed in terms of the number of gates inserted per key K of length l , where $l = |K|$. The cost summary for all strategies is presented in Table II, assuming the following:

- A single key bit input $k_i \in K$ is used only once.
- Every MUX is implemented with the same number of gates denoted as $|MUX|$. Typically, a multiplexer is implemented by using one inverter, two AND gates, and one OR gate, i.e., $|MUX| = 4$.

This model can be used to easily steer the final cost in the D-MUX locking procedure depending on the available strategies.

C. The D-MUX Locking Scheme

In this section, we discuss the construction of the D-MUX LL scheme. The scheme is presented in Algorithm 3. The scheme takes the following inputs: a set of available strategies L_s , the correct key vector K , the original netlist Net , the maximum input node iteration variable I_{max} , and the maximum output node iteration variable O_{max} . Hereby,

Algorithm 3 D-MUX locking scheme.

Input: Available strategies L_s , key K , netlist Net , max input node iterations I_{max} , max output node iterations O_{max}

Output: Locked netlist

```

1:  $F_{single} \leftarrow \text{ExtractSingleOutputNodes}(Net)$ 
2:  $F_{multi} \leftarrow \text{ExtractMultiOutputNodes}(Net)$ 
3:  $K_{list} \leftarrow \text{ToList}(K)$  // Convert key to list
4: while  $|K_{list}| > 0$  do
5:   /* Select a candidate strategy */
6:    $\text{RandomShuffle}(L_s)$ 
7:    $fallback \leftarrow \text{TRUE}; S_{sel} \leftarrow \emptyset$ 
8:   for  $S_i$  in  $L_s$  do
9:     /* Enough keys available? */
10:    if  $S_i == S_1$  &&  $(|K_{list}| < 2)$  then
11:      continue
12:    end if
13:    /* Enough nodes available? */
14:    if  $S_i \in \{S_1, S_2, S_3\}$  &&  $(|F_{multi}| < 2)$  then
15:      continue
16:    end if
17:     $S_{sel} \leftarrow S_i; fallback \leftarrow \text{FALSE}$ 
18:    break
19:  end for
20:  if  $fallback$  then
21:     $S_{sel} \leftarrow S_4$ 
22:  end if
23:  /* Search for valid input/output nodes */
24:   $\{\{f_i, f_j\}, \{g^i, g^j\}, done\} \leftarrow \text{FindPairs}(S_{sel}, F_{single}, F_{multi}, I_{max}, O_{max})$ 
25:  if  $!done$  then // Check if the search was successful
26:    continue
27:  end if
28:  /* Apply selected strategy */
29:   $K_{i,j} \leftarrow \text{GetAndRemoveFrom}(K_{list}, S_{sel})$ 
30:   $\{M_{i,j}\} \leftarrow \text{CoupleToMUXs}(f_i, f_j, g^i, g^j, K_{i,j})$ 
31:   $\text{RegisterToNetlist}(Net, \{M_{i,j}\})$ 
32: end while
33: return  $Net$ 

```

I_{max} defines the maximum number of reselections of input node pairs and O_{max} the maximum number of reselections of output node pairs for already selected input nodes. Note that S_4 is always available to the locking algorithm, therefore it is not included in L_s . The final output is the locked netlist.

Before the main locking loop, the scheme prepares two separate sets: all single-output (F_{single}) and multi-output (F_{multi}) nodes (line 1 and 2). Furthermore, the key K is represented as a list K_{list} (line 3) to easily track the amount of used keys. The main locking loop repeats until no keys are left (line 4). Each iteration starts by randomly shuffling the available strategies (line 6). The shuffling ensures a random pick in the next step; all strategies in L_s are iterated and multiple checks are performed to ensure the existence of enough key bits as well input nodes of the desired type (line 8 to 19).

Once a candidate S_i is selected, the scheme proceeds with finding a valid pair of nodes (line 24). If successful, based on the nature of S_i , the scheme performs the following steps: retrieve the necessary key bits (set $K_{i,j}$, line 29) and couple all nodes through one or multiple MUXs (set $\{M_{i,j}\}$, line 30). Finally, all changes are registered in the netlist Net (line 31).

Scheme Variants: L_s enables the creation of a generalized (gD-MUX) and enhanced (eD-MUX) scheme variant. In eD-MUX, $L_s = \{S_1, S_2, S_3\}$. Thus, the scheme first tries to use the less costly strategies if possible. In case none of these strategies are viable, the scheme *falls back* to S_4 (line 21). In gD-MUX, $L_s = \emptyset$. This forces the scheme to exclusively use S_4 . The difference between gD-MUX and eD-MUX mirrors the cost discrepancy between the cases when the target netlist only supports S_4 (worst case) and when it supports all other strategies (best case).

Algorithm 4 FindPairs Function

Input: Strategy S_i , Set of single-output nodes F_{single} , set of multi-output nodes F_{multi} , max input node iterations I_{max} , max output node iterations O_{max}

Output: Valid input nodes $\{f_j, f_i\}$, valid output nodes $\{g^j, g^i\}$, success indicator $done$

```

1:  $\{F_1, F_2\} \leftarrow \{\emptyset\}$ 
2: /* Select correct nodes */
3: if  $S_i \in \{S_1, S_2\}$  then
4:    $\{F_1, F_2\} \leftarrow \{F_{multi}, F_{multi}\}$ 
5: else if  $S_i \in \{S_3\}$  then
6:    $\{F_1, F_2\} \leftarrow \{F_{multi}, F_{single}\}$ 
7: else
8:    $\{F_1, F_2\} \leftarrow \{F_{single} \cup F_{multi}, F_{single} \cup F_{multi}\}$ 
9: end if
10: /* Prepare input and output node references */
11:  $\{f_i, f_j, g^i, g^j\} \leftarrow \emptyset$ ;  $done \leftarrow FALSE$ 
12: for  $iter_{in} = 0$  to  $I_{max}$  do
13:   /* Select first and second input node */
14:    $\{f_i, f_j\} \leftarrow \{RndSel(F_1), RndSel(F_2)\}$ 
15:   /* Select output nodes */
16:   for  $iter_{out} = 0$  to  $O_{max}$  do
17:      $\{g^i, g^j\} \leftarrow \{RndSel(OutsOf(f_i)), RndSel(OutsOf(f_j))\}$ 
18:      $\{R_1, R_2\} \leftarrow \{IsInOutCone(f_j, g^i), IsInOutCone(f_i, g^j)\}$ 
19:     if  $(g^i \neq g^j) \ \&\& \ !R_1 \ \&\& \ !R_2$  then
20:        $done \leftarrow TRUE$ 
21:       break
22:     end if
23:   end for
24:   if  $done$  then
25:     break
26:   end if
27: end for
28: return  $\{\{f_i, f_j\}, \{g^i, g^j\}, done\}$ 

```

Node Selection: In all strategies, the input and output nodes must be carefully selected to avoid cycles. Therefore, D-MUX selects the nodes (line 24) according to the function $FindPairs(S_i, F_{single}, F_{multi}, I_{max}, O_{max})$ described in Algorithm 4. The function takes five inputs: the locking strategy S_i , the set of single-output nodes F_{single} , the set of multi-output nodes F_{multi} , and the maximum iteration variables I_{max} and O_{max} . The output consists of the valid input ($\{f_i, f_j\}$) and output nodes ($\{g^i, g^j\}$) as well as an indicator ($done$) if the search has been successful. The function works as follows. Two node sets are assigned (F_1 and F_2) depending on the requirements of S_i (line 1 to 9). These sets are used for the input node selection in the rest of the function. The selection loop (line 12) is repeated until two valid input and output nodes are found or the maximum iteration is reached. In each iteration, two candidate input nodes (f_i and f_j) are randomly selected from F_1 and F_2 , respectively. Next, the function tries to find two valid output nodes. The validity is defined by three requirements (line 19): (i) the nodes are different, (ii) f_j is not in the output cone of g^i , and (iii) f_i is not in the output cone of g^j . (ii) and (iii) prevent the creation of cycles. The node selection is limited to a given amount of iterations, since a valid selection is not always possible. If a selection is not made, the function returns an indication that the search has failed. In this case, the algorithm repeats the search.

The selection returns two output nodes; even though some strategies require only one node. However, since every node has at least one output, this does not impact the final result.

VI. RESILIENCE EVALUATION

This section evaluates D-MUX against all relevant attacks, as discussed in Section II-B. For the evaluation, we used the ISCAS'85 [39], ITC'99 [40], and RISC-V Ariane core [41] benchmarks listed in Table III. All evaluations have been

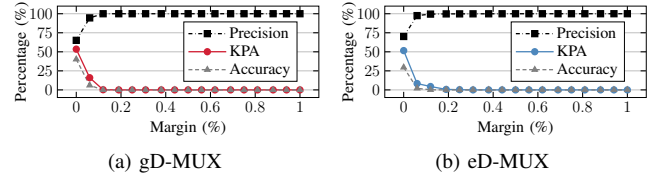


Fig. 8: SWEEP attack evaluation on D-MUX.

performed on an AMD Ryzen 9 3900X processor with 64GB of RAM and an Nvidia GeForce RTX 2080 Ti graphic card.

A. SWEEP

Attack Setup: To evaluate D-MUX using SWEEP, we applied the SWEEP tool provided by [4]. The setup is as follows. First, all benchmarks from Table III (a) have been copied 100 times and locked with random 64-bit keys; generating a data set of 700 locked benchmarks. Second, for each benchmark, the data set has been divided into a test set consisting of the 100 locked targets and a training set consisting of the remaining 600 other benchmarks. The attack is repeated for a range of margins, i.e., $m \in [0.0, 0.01]$ with a step of size $6.25 \cdot 10^{-3}$. This is repeated for both D-MUX variants.

Metrics: The attack was evaluated through accuracy, precision, and KPA (see Section II-B1). The accuracy in SWEEP is calculated using the whole key length, regardless of X values. The KPA only considers the bits that are different from X.

Results: The average attack results across all benchmarks and margins are presented in Fig. 8. At $m = 0$, the average KPA is approx. 50% as expected. However, the accuracy is lower than 50% as its formulation takes X values into account. As soon as the margin value is slightly increased ($m > 0$), both the accuracy and the KPA fall to 0% while the precision rises to 100%. *This is due to the fact that with an increased margin, SWEEP is not able to report any 'guessed' key bits (all bits are marked as X).* Based on the evaluation, we can conclude that SWEEP is not able to extract a meaningful correlation from the data regardless of the acceptable margin.

B. Learning Resilience

1) **ANT and RNT:** In ANT, U is able to identify the inserted MUX gates by following each key input. Thereby, the only observation that can be made is that both inputs to the MUX are AND gates. However, due to the nature of the locking strategies deployed by D-MUX, it is not possible to distinguish the true and the false gate driving the MUX. Even after many iterations of the game, the guessing accuracy will remain 50%. Therefore, *D-MUX passes ANT*, thereby being evaluated as *possibly learning resilient*. The same conclusion

TABLE III: Benchmark circuits used for evaluation.

(a) ISCAS'85		(b) ITC'99		(c) Ariane RISC-V	
IC	#Gates	IC	#Gates	IC	#Gates
c1355	546	b15	6931	iscan	240
c1908	880	b21	7931	commit	1584
c2670	1193	b22	12128	brunit	1655
c3540	1669	b17	21191	decoder	2169
e5315	2307	b18	49293	peselect	3333
c6288	2416	b19	98726	brpredict	4669
c7552	3512			alu	7412

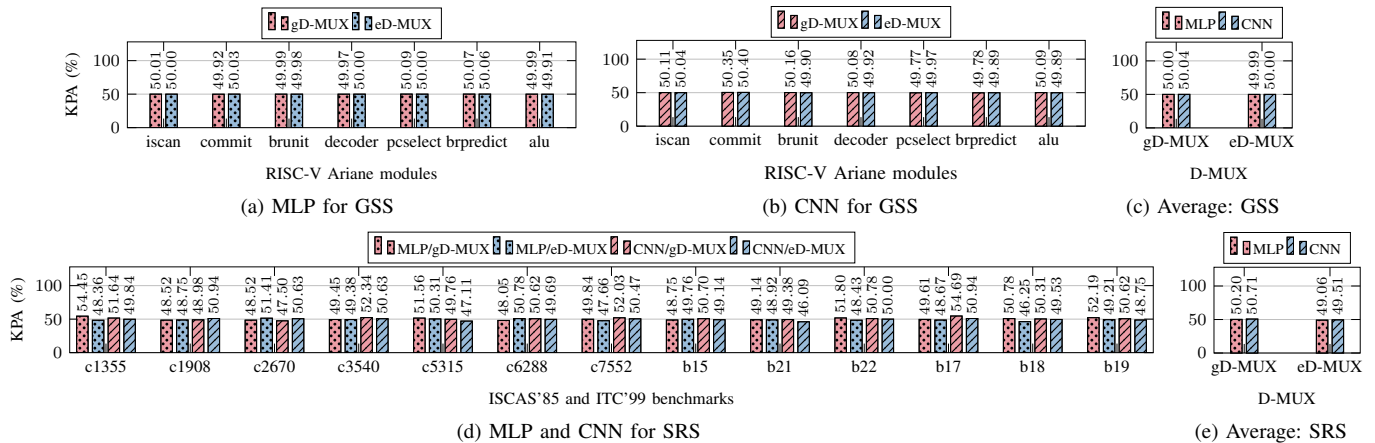


Fig. 9: SnapShot attack evaluation on D-MUX for the generalized set and self-referencing scenario.

is drawn in RNT as well. The reason why D-MUX passes both tests lies within the way the key value is mapped into the circuits; instead of adding key-bounded functionality, D-MUX reconfigures the existing one.

2) *SnapShot*: So far, SnapShot represents the only oracle-less ML-based attack on LL that is applicable to MUX-based locking. SnapShot predicts key bits based on labeled LVs that structurally represent a subcircuit associated to a key-bit value. To adapt SnapShot to D-MUX, we have to introduce a new LV extraction procedure that includes all structurally relevant components of the netlist for each key input.

D-MUX Locality Vector: The LV must include all relevant data, which could potentially leak information about the key. For D-MUX, we can observe two cases: one or two MUX gates are inserted per key bit. Both cases decompose a netlist into multiple components as shown in Fig. 10:

- k_i : the key-bit value (in case the vector is labeled).
- l_b : the backward path; the input cones of f_i and f_j .
- l_{kg} : the key gate in the form of one MUX gate (regardless of whether one or two MUX gates are inserted).
- l_{sb}^i : the backward path specific to g_i^i .
- l_{sb}^j : the backward path specific to g_j^j (if not existent, l_{sb}^j is filled with zeros).
- l_f : the forward path; the output cones of g_i and g_j .

The extraction traverses each backward and forward path until a selected depth is reached, thereby applying the fan-in and fan-out values as suggested in [8]. By applying the BFS, one locality vector can be extracted for each key input. One example is given in Fig. 11. Here, the maximum forward and

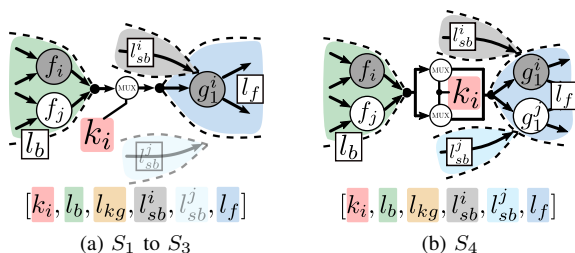


Fig. 10: Locality vectors for D-MUX.

backward depth is set to 5 ($D_b = 5$ and $D_f = 5$). Each section of the vector represents one of the components listed above.

Attack Scenarios: We evaluated both D-MUX variants using both attack scenarios as discussed in Section II-B2. In GSS, the training set is assembled using all benchmarks from Table III (a) and (b). Each benchmark is copied and locked 1,000 times using 64-bit keys, resulting in 13,000 locked netlists and 832,000 labeled localities. The test set consists of all benchmarks from Table III (c), each locked 1,000 times, yielding 448,000 unlabeled localities. SRS generates a training set by relocking the each target 1,000 times to predict the key of the initial target. This is repeated 20 times to get an average value for all benchmarks in Table III (a) and (b).

ML Model: We adapted the Mutli-layer Perceptron (MLP) and the CNN model from [8] to accommodate the higher complexity of D-MUX as follows; (i) in MLP, the input layer has a higher number of nodes (same as locality vector size), (ii) we increased the number of available internal layers in the CNN evolution to 14, and (iii) the number of epochs is set to 100 in the exploratory phase of the CNN.

Results: All evaluation results are presented in Fig. 9. For all scenarios (MLP/CNN and GSS/SRS) the key prediction accuracy is consistently around 50%. Thus, we can conclude that D-MUX is efficient in protecting against SnapShot, thereby forcing the attack to perform random guesses about the key.

Evolved CNNs: The evolved, best-performing networks tend to have few or no convolutional layers, thus degrading the CNN to a basic MLP. This suggests that no feature extraction is possible and a learning-by-heart approach yields best results.

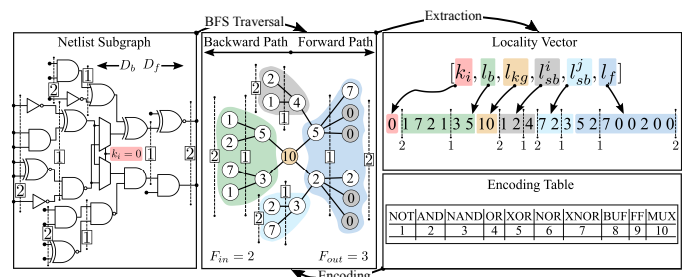


Fig. 11: Example: locality vector extraction for D-MUX.

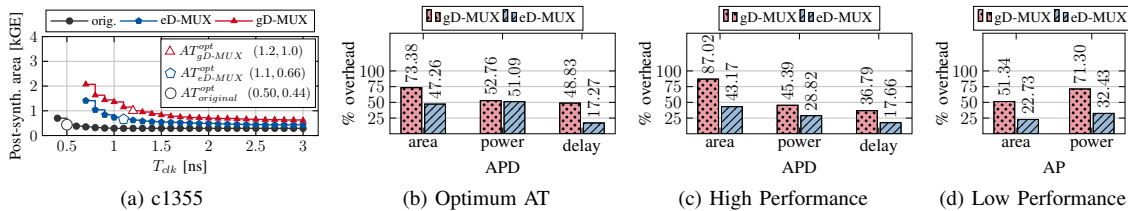


Fig. 12: D-MUX cost evaluation.

VII. COST EVALUATION

The following evaluation gathers a realistic cost analysis using the UMC 90 nm CMOS process, based on the ISCAS'85 and ITC'99 benchmarks from Table III. The Synopsys Design Compiler (DC) was used for logic synthesis.

All benchmarks were locked with 64-bit keys for gD-MUX and eD-MUX separately. Each benchmark is synthesized for a range of clock periods (T_{clk}) to generate the Area-Timing (AT) graph. As an example, we only present the AT graph for the benchmark c1355 in Fig. 12 (a). The graph includes the optimum AT point (AT^{opt}) for the original netlist and the locked variants. AT^{opt} is the Pareto-optimal point with the most *efficient* implementation (minimum $Area \cdot T_{clk}$). The lowest T_{clk} is determined by the critical path, while the highest T_{clk} is selected to exhibit a reasonable AT curve. To include locking and synthesis variations, the average value across 20 synthesis rounds for 20 locked netlist instances is taken.

The average cost across all benchmarks are presented in Fig. 12 (b)-(d). The power overhead is approximated using DC, thus it correlates with the area overhead increase. We present the evaluations at three comparisons points as follows.

Optimum AT: The average Area-Power-Delay (APD) overhead for optimum AT points is presented in Fig. 12 (b). This evaluation shows the average difference between *the most efficient implementations* per benchmark. Hereby, gD-MUX and eD-MUX induce an average area overhead of 73.38% and 47.26% respectively. In terms of delay, eD-MUX exhibits only 17.27% overhead compared to the 48.83% of gD-MUX. This is due to the fact that the AT optima of eD-MUX are typically closer in terms of T_{clk} to the original AT optima.

High performance: This evaluation shows the cost for *the lowest possible* T_{clk} (Fig. 12 (c)). The results are similar to the optimum AT case. Most importantly, the delay overhead is limited to 36.79% for gD-MUX and 17.66% for eD-MUX. This suggests that in practice the locking mechanisms have a relatively low impact on the critical path of the designs.

Low performance: This evaluation looks at the case *when no excessive synthesis-optimizations have been applied yet* since enough margin in terms of T_{clk} is still available, yielding a fair comparison (Fig. 12 (d)). The longest shown clock period has been selected for comparison for each benchmarks (e.g., $T_{clk} = 3$ ns in Fig. 12 (a)). As before, gD-MUX exhibits a higher area cost (51.34%) compared to eD-MUX (22.73%).

In conclusion, the evaluation shows that gD-MUX is costlier in terms of APD compared to eD-MUX for all benchmarks; hence offering a more cost-effective option.

VIII. ANALYSIS OF RELATED WORK

The following overview of related work focuses on four categories: MUX-based, routing-based, lookup table (LUT)-based, and ML-resilient LL. Hereby, we offer an analysis of potential structural leakage points in these schemes and indicate new perspectives on improving LL in the ML era.

MUX-Based LL: As mentioned in Section IV, existing MUX-based LL techniques are susceptible to SAAM [28]–[30]. However, a seemingly structurally secure scheme is known as cyclic obfuscation [31]. This scheme has been designed to thwart SAT-based attacks by creating combinational cycles (later successfully challenged by CycSAT [42]). To showcase its mode of operation, let us consider the locked example in Fig. 13 (exact replica from [31]). Here, all false wires are marked red. The scheme receives two inputs: the number of cycles and the cycle length. To insert a cycle, the scheme first searches for a path of a selected length between two gates; for example, the path (w_0, w_3) . First, a feedback wire is added between these endpoints (MUX M_0). Second, all edges between (w_0, w_3) must be made *removable* to disable the identification of the inserted cycle. This procedure is done by traversing all gates G_i alongside the path, thereby two MUXs are added if the fan-out(G_i) = 1 (M_2 and M_3) and, otherwise, only one MUX is inserted (M_1). Note that it is suggested that the same key input (k_2) drives both MUXs in the former case. With this configuration, all edges of the created cycle are removable. Even though this scheme checks for dangling wires, it still suffers from major leakage points. First, cyclic locking avoids dangling wires by checking the fan-out size of the gates driving the added MUXs. However, this check is not performed for the MUXs that are inserted *outside* the path, i.e., for false wires. The presented example showcases this fault; the wire r_2 is only driving M_3 , thus it has to be the true wire. Moreover, this error can also occur for the insertion of the MUX that creates the initial cycle. In the example, the fan-out of the gate that outputs w_0 is not checked. Evidently, if its fan-out is 1, w_0 must be the *true* wire. These critical cases have not been taken into account in the proposal [31]. Furthermore, we have verified this leakage in the implementations as well; a large number of

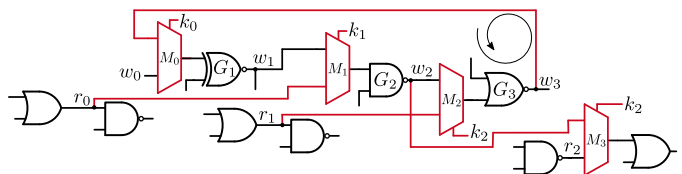


Fig. 13: Cyclic logic locking; false wires are marked red.

MUXs in virtually all cyclic-locked benchmarks that have been provided by the original authors online exhibit this fault [43]. The same is true for the cyclic locking that is available in the same author’s tool NEOS [44]. This fault enables the application of SAAM. This problem is somewhat resolvable if dummy gates are deployed; however, this requires (trusted) layout manipulations to be cost-effective [31]. The second structural fault hides in the insertion of *two* MUXs when the fan-out of a gate equals 1. Since w_2 is the only wire that simultaneously drives *both* MUXs, according to the scheme, it must be the *true* wire of the MUX placed in the cycle. Hence, both MUXs can be removed. This issue persists even if both MUXs have separate key bits. Third, since a specific cycle length is induced by the scheme, in case a cycle is identified in the netlist, there exists a high likelihood that the longest backward edge is the *false* one. Here, the length can be measured in terms of how many gates the edge “jumps over” from its *first* entry point to the *last* exit point in a *topological* order. For example, topologically (when ignoring backward edges) the first entry point in Fig. 13 is M_0 and the last exit point is G_3 . Hence, the backward edge from w_3 to w_0 is clearly the longest one and its length is *defined* as an input parameter. Finally, another issue occurs when not enough false wires are available (verified with NEOS). Here, the algorithm inserts a MUX that is driven by two inputs that originate from the *same* MUX. Evidently, this structural composition immediately reveals the correct key-bit value. Consequently, due to the presented faults, cyclic obfuscation fails both ANT and RNT. To the best of our knowledge, the introduced faults have not been discussed before. Hence, we urge the authors to address these issues in the implementation and the theoretical work. Some faults can be patched by utilizing the D-MUX locking strategies. Moreover, based on the presented concepts of cyclic locking, multiple next-generation cyclic LL schemes have been proposed [45]–[47], some of which try to mitigate the shortcomings of the initial proposal. However, using the presented evaluation concepts, a detailed analysis might reveal new structural faults in these schemes as well. Hence, we leave this task for future work.

Routing-Based LL: A promising approach to thwarting structural attacks lies in a new class of LL policies known as routing-based obfuscation. Even though these techniques have been designed to tackle SAT-based attacks, their structural composition might offer a fruitful ground to protect against structural ML-based attacks. The main building blocks of routing-based LL are *key-programmable routing blocks* (keyRBs), which can implement a variety of topologies. For example, InterLock [48] utilizes keyRBs that are composed of Switch-Boxes (SwBs) to embed selected timing paths of a predefined length. For example, the marked path in Fig. 14 (a) can be embedded into the keyRB by connecting multiple SwBs as shown in Fig. 14 (b). Each SwB hosts two 2-input gates (f_1 and f_2). These represent original gates from the netlist. For a correct key ($\{k_0, k_2, k_3, k_4\} = \{0, 0, 0, 0\}$), the embedded gates remain properly connected. Similar to D-MUX, this mechanism can create paths that are equally likely for each key value. However, some structural cases remain unresolved. For example, for $\{k_1, k_2\} = \{0, 1\}$, both functions become

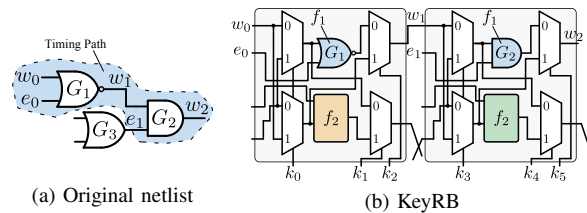


Fig. 14: InterLock: embedding a timing path into a KeyRB.

dangling gates, i.e., e_0 might lead to an *unconnected* gate; thus suggesting that these key combinations are not valid. This problem is further exacerbated in case any wire in the selected path drives multiple gates. For example, if w_1 is connected to another gate (besides G_2), it is clear that $\{k_1, k_2\} = \{0, 1\}$ may not appear. These cases have not been clarified in the proposed work [48]. Therefore, there is still the possibility to resolve such structural cases within InterLock. Moreover, we see a good potential to reuse the locking strategies presented in D-MUX for furthering the concept of InterLock. Thus, we leave this to future work.

A similar scheme is known as Banyan locking [49]. This scheme expands the SwBs by adding decoy logic. For example, the SwB can be implemented as a 2-input-2-output function, where each output can be driven by the original gate, decoy gate, buffered input wire, or a constant. This configuration leads to potential structural concerns. For example, in ANT, all SwBs would contain only AND gates, thus reducing the key space as multiple keys would result in the same output. Thus, Banyan LL might result in structural leakage depending on the structural characteristics of the netlist. Moreover, as Banyan LL implements a set of gates in addition to multiple N-to-1 MUXs for a single SwB, its overhead is likely to be a major concern for wide adoption. Even its authors left the cost evaluation to future work [49]. Nevertheless, similar to InterLock, there is a potential to adapt Banyan LL based on the introduced D-MUX locking strategies.

LUT-Based LL: A promising concept of LL is manifest in LUT-based locking [50], [51]. This LL type replaces parts of the target design with generic lookup tables that are configured *after* fabrication. This approach has one major advantage compared to traditional LL policies; the LUTs exhibit a generic memory layout topology, regardless of their functionality. Thus, if LUT-based LL is deployed, part of the design is not available to the attacker (e.g., the untrusted foundry); resulting in a structurally secure design. However, this feature comes with a few disadvantages. First, utilizing LUTs for LL can lead to considerable overheads, especially in the traditional CMOS technology. Second, LUT-based LL requires the inclusion of a trusted testing facility, which might be a concern for some IP owners. Third, to lower the cost aspects, the integration of novel technologies is required (such as the non-volatile spin transfer torque magnetic technology). Hence, another burden is put on the design flow. Nevertheless, as LUTs are able to absorb parts of the design *without leaving LL-induced structural traces*, they provide a promising direction to design resilient schemes. Interestingly, ANT provides an important observation; if the structure of the netlist is very regular (e.g.,

AND tree), the LUTs must absorb *more than one gate* to not leak information; otherwise it is clear what functionality is implemented. This can be simulated by generalizing the Twin-Gate scheme to an N-Gate scheme, where N is the number of all possible two-input gates. This suggests that the circuit family might play a role in LUT-based LL.

ML-Resilient LL: UNSAIL inserts key gates with the goal to create confusing training data for ML-based attacks (such as SAIL) [37]. For a selected LL policy, UNSAIL locks the netlist using half of the available key. Next, the design is synthesized. At this point, a dictionary of observations is generated by extracting the key-gate subgraphs from both locked designs *before* and *after* synthesis. Next, UNSAIL performs a dictionary-guided key-gates insertion, thereby using the remaining half of the key. Thus, UNSAIL-locked netlists ensure that equivalent subgraphs are linked to different key values; resulting in confusing training data. Even though deceptive observations are crucial to thwarting ML-based attacks, we can identify a few disadvantages of UNSAIL. First, the major issue of UNSAIL is its deep reliance on the usage of a synthesis tool. As discussed in Sections III-A4 and III-D, in the context of learning resilience, the security of LL should not depend on the specifics of synthesis transformations, especially since these are typically induced by closed-source, third-party software. However, note that often the synthesis tool *is not* part of the attack model, i.e., it is regarded as trustworthy. Hence, the dependency on the tool must not be seen as a disadvantage in these cases. Second, in the optimal case, an UNSAIL-locked netlist should not be changed through additional synthesis. Nevertheless, resynthesis and layout generation might lead to new leakage points that reverse the effects of UNSAIL. Moreover, tampering with the synthesis can result in sub-optimal designs. Third, the resiliency of UNSAIL has only been demonstrated for schemes that buffer or invert a single wire, such as XOR/XNOR-locking and single-wire MUX gates (MUXs that are driven by the same wire twice, except that one input is inverted). These schemes fail the theoretical tests. Finally, UNSAIL only lowers the accuracy of the attacks for a certain amount of percentage points. In comparison, D-MUX is independent of synthesis transformations and exhibits maximal resilience against the introduced attacks.

A potentially ML-resilient scheme is known as Scalable Attack-Resistant Obfuscation (SARO) [52]; operating in two steps. First, it splits the design into smaller partitions with the aim to maximize structural netlist changes. Second, SARO deploys a Truth Table Transformation (T3) to lock the created partitions. T3 is designed to induce randomized design alterations, thus increasing the complexity of deploying pattern-recognition attacks (such as SAIL). Even though the effectiveness of SARO has been evaluated against a variety of attacks, its resilience against ML-based attacks is only assumed based on a proposed metric that tries to capture the level of functional and structural changes in the design. Thus, its effectiveness against ML-based attacks remains to be evaluated.

IX. LIMITATION AND OPPORTUNITIES

Even though we could not identify any leakage in D-MUX, this does not conclusively prove its absence. Thus, D-MUX

has to be further evaluated against upcoming attacks [53].

Another interesting problem lies in the evaluation of structural leakage that is beyond the capability of the human observer. For example, LL might leak key-related information only on a *topological* level. Evidently, comparing the topology of netlists is not easily done for humans. Moreover, this problem has not been addressed yet in detail as LL mostly induces local changes. Nevertheless, the proposed tests offer the capability to pinpoint fundamental security vulnerabilities that are easily exploitable. Whether a global leakage exists will require the inclusion of suitable topology-matching procedures, hence opening up challenging questions in LL.

We have seen that the security of LL is often predicated by the structural characteristics of the target netlists, i.e., the circuit family. However, a comprehensive analysis of the existence of a structural bias across a wide range of designs has yet to be performed. A first indication is provided by the success of SnapShot [8]. Here, a range of benchmarks was successfully attacked, suggesting that there exists a bias towards repeating structures in hardware designs. The results of such a study could provide some insights on what LL should focus on—in case a general solution is difficult to provide.

Finally, we have identified the potential to transfer the D-MUX strategies to OG-resilient LL, such as routing-based LL.

X. CONCLUSION

This work addresses the challenges of logic locking in the context of learning-based attacks. Hereby, the first theoretical concept for evaluating learning resilience in logic locking has been introduced; exposing critical information leaks even before concrete locking implementations are in place. Based on the theoretical insights, we developed D-MUX; the first deceptive logic-locking scheme, thereby investigating and correcting a major fallacy in existing MUX-based locking. The resiliency of D-MUX was evaluated through the extensive application of the latest machine-learning-based and MUX-targeting attacks. The cost of D-MUX has been modeled theoretically and evaluated using a concrete technology node. Finally, we have discussed related locking concepts, thereby providing novel insights into potential information leakage as well as promising research directions. With the presented work, we establish the cornerstones for the design of next-generation logic locking in the era of machine learning.

REFERENCES

- [1] M. Yasin *et al.*, “Evolution of logic locking,” in *2017 IFIP/IEEE VLSI-SoC*, Oct 2017, pp. 1–6.
- [2] M. Yasin *et al.*, *Trustworthy Hardware Design: Combinational Logic Locking Techniques*. Springer International Publishing, 2020.
- [3] K. Zamiri Azar *et al.*, “Threats on logic locking: A decade later,” in *2019 GLSVLSI*. New York, NY, USA: ACM, 2019, p. 471–476.
- [4] A. Alaql *et al.*, “Sweep to the secret: A constant propagation attack on logic locking,” in *2019 AsianHOST*, 2019, pp. 1–6.
- [5] Y. Zhang *et al.*, “A novel topology-guided attack and its countermeasure towards secure logic locking,” *JCEN*, Oct 2020.
- [6] M. E. Massad *et al.*, “Logic locking for secure outsourced chip fabrication: A new attack and provably secure defense mechanism,” *CoRR*, vol. abs/1703.10187, 2017.
- [7] L. Li *et al.*, “Piercing logic locking keys through redundancy identification,” in *2019 DATE*, March 2019, pp. 540–545.

- [8] D. Sisejkovic *et al.*, “Challenging the security of logic locking schemes in the era of deep learning: A neuroevolutionary approach,” vol. 17, no. 3, 2021. [Online]. Available: <https://doi.org/10.1145/3431389>
- [9] C. Pilato *et al.*, “ASSURE: RTL locking against an untrusted foundry,” *IEEE TVLSI*, vol. 29, no. 7, pp. 1306–1318, 2021.
- [10] Defense Science Board Task Force, “High performance microchip supply,” *Annual Report. Defense Technical Information Center (DTIC), USA*, 2005. [Online]. Available: <https://apps.dtic.mil/sti/citations/ADA435563>
- [11] D. Šišković *et al.*, “A secure hardware-software solution based on RISC-V, logic locking and microkernel,” in *SCOPES*, ser. SCOPES '20. New York, NY, USA: ACM, 2020, p. 62–65.
- [12] A. Jain *et al.*, “TAAL: Tampering attack on any key-based logic locked circuits,” 2019. [Online]. Available: <https://arxiv.org/abs/1909.07426>
- [13] J. Lee *et al.*, “A low-cost solution for protecting IPs against scan-based side-channel attacks,” in *Proceedings of the 24th IEEE VTS*, ser. VTS '06. USA: IEEE Computer Society, 2006, p. 94–99.
- [14] M. T. Rahman *et al.*, “The key is left under the mat: On the inappropriate security assumption of logic locking schemes,” Cryptology ePrint Archive, Report 2019/719, 2019. <https://eprint.iacr.org/2019/719>.
- [15] S. Engels *et al.*, “The end of logic locking? A critical view on the security of logic locking,” *Cryptology ePrint*, vol. 2019, p. 796, 2019.
- [16] A. Jain *et al.*, “ATPG-guided fault injection attacks on logic locking,” 2020. [Online]. Available: <https://arxiv.org/abs/2007.10512>
- [17] K. Shamsi *et al.*, “On the impossibility of approximation-resilient circuit locking,” in *HOST*, 2019, pp. 161–170.
- [18] H. Zhou *et al.*, “Resolving the trilemma in logic encryption,” in *ICCAD*, 2019, pp. 1–8.
- [19] M. Yasin *et al.*, “Provably-secure logic locking: From theory to practice,” ser. CCS '17. New York, NY, USA: ACM, 2017, p. 1601–1618.
- [20] H. Chen *et al.*, “GenUnlock: An automated genetic algorithm framework for unlocking logic encryption,” in *IEEE/ACM ICCAD*, 2019, pp. 1–8.
- [21] K. Z. Azar *et al.*, “NNgSAT: Neural network guided SAT attack on logic locked complex structures,” 2020. [Online]. Available: <https://arxiv.org/abs/2009.02208>
- [22] R. Karmakar *et al.*, “A particle swarm optimization guided approximate key search attack on logic locking in the absence of scan access,” in *2020 DATE*, 2020, pp. 448–453.
- [23] P. Chakraborty *et al.*, “SAIL: Machine learning guided structural analysis attack on hardware obfuscation,” in *AsianHOST*, Dec 2018, pp. 56–61.
- [24] F. Tehranipoor *et al.*, “Deep RNN-oriented paradigm shift through BOCANet: Broken obfuscated circuit attack,” in *2019 GLSVLSI*, ser. GLSVLSI '19. New York, NY, USA: ACM, 2019, p. 335–338.
- [25] K. Shamsi *et al.*, “IP protection and supply chain security through logic obfuscation: A systematic overview,” *ACM Trans. Des. Autom. Electron. Syst.*, vol. 24, no. 6, Sep. 2019.
- [26] J. A. Roy *et al.*, “EPIC: Ending piracy of integrated circuits,” in *2008 DATE*, March 2008, pp. 1069–1074.
- [27] M. Yasin *et al.*, “On improving the security of logic locking,” *IEEE TCAD*, vol. 35, no. 9, pp. 1411–1424, Sept 2016.
- [28] S. M. Plaza *et al.*, “Solving the third-shift problem in IC piracy with test-aware logic locking,” vol. 34, no. 6, pp. 961–971, 2015.
- [29] Y. Lee *et al.*, “Improving logic obfuscation via logic cone analysis,” in *2015 16th LATS*, March 2015, pp. 1–6.
- [30] J. Rajendran *et al.*, “Fault analysis-based logic encryption,” *IEEE Trans. on Computers*, vol. 64, no. 2, pp. 410–424, 2015.
- [31] K. Shamsi *et al.*, “Cyclic obfuscation for creating SAT-unresolvable circuits,” in *GLSVLSI '17*. ACM, 2017, p. 173–178.
- [32] S. Dupuis *et al.*, “A novel hardware logic encryption technique for thwarting illegal overproduction and hardware Trojans,” in *2014 IEEE 20th IOLTS*, July 2014, pp. 49–54.
- [33] M. Yasin *et al.*, “SARLock: SAT attack resistant logic locking,” in *IEEE HOST 2016*, May 2016, pp. 236–241.
- [34] Y. Liu *et al.*, “Strong Anti-SAT: Secure and effective logic locking,” in *2020 21st ISQED*, 2020, pp. 199–205.
- [35] B. Shakya *et al.*, “CAS-Lock: A security-corrupibility trade-off resilient logic locking scheme,” *TCHES*, no. 1, pp. 175–202, Nov. 2019.
- [36] L. Li *et al.*, “Shielding logic locking from redundancy attacks,” in *2019 IEEE 37th VLSI Test Symposium (VTS)*, April 2019, pp. 1–6.
- [37] L. Alrahis *et al.*, “UNSAIL: Thwarting oracle-less machine learning attacks on logic locking,” *IEEE TIFS*, vol. 16, pp. 2508–2523, 2021.
- [38] N. Limaye *et al.*, “Thwarting all logic locking attacks: Dishonest oracle with truly random logic locking,” *IEEE TCAD*, pp. 1–1, 2020.
- [39] F. Brglez *et al.*, “Combinational profiles of sequential benchmark circuits,” in *IEEE ISCAS*, May 1989, pp. 1929–1934 vol.3.
- [40] F. Corno *et al.*, “RT-level ITC'99 benchmarks and first ATPG results,” *IEEE Design Test of Computers*, vol. 17, no. 3, pp. 44–53, 2000.
- [41] F. Zaruba *et al.*, “The cost of application-class processing: Energy and performance analysis of a Linux-ready 1.7-GHz 64-bit RISC-V core in 22-nm FDSOI technology,” *IEEE TVLSI*, vol. 27, no. 11, pp. 2629–2640, Nov 2019.
- [42] H. Zhou *et al.*, “CycSAT: SAT-based attack on cyclic logic encryptions,” in *2017 IEEE/ACM ICCAD*, 2017, pp. 49–56.
- [43] D. Forte *et al.*, “Obfuscation benchmarks,” Accessed May, 2021. [Online]. Available: <https://trust-hub.org/#/benchmarks/obfuscation>
- [44] K. Shamsi *et al.*, “NEOS: Netlist encryption and obfuscation suite,” Accessed May, 2021. [Online]. Available: <https://bitbucket.org/kavehsh/neos/src/master/>
- [45] S. Roshanifefat *et al.*, “SRCLock: SAT-resistant cyclic logic locking for protecting the hardware,” in *GLSVLSI*. ACM, 2018, p. 153–158.
- [46] A. Rezaei *et al.*, “Cyclic locking and memristor-based obfuscation against CycSAT and inside foundry attacks,” in *DATE*, 2018, pp. 85–90.
- [47] X.-M. Yang *et al.*, “LOOPLock 2.0: An enhanced cyclic logic locking approach,” *IEEE TCAD*, pp. 1–1, 2021.
- [48] H. M. Kamali *et al.*, “InterLock: An intercorrelated logic and routing locking,” in *ICCAD '20*. New York, NY, USA: ACM, 2020.
- [49] J. Sweeney *et al.*, “Modeling techniques for logic locking,” in *2020 IEEE/ACM ICCAD*, 2020, pp. 1–9.
- [50] A. Baumgarten *et al.*, “Preventing IC piracy using reconfigurable logic barriers,” *IEEE D&T*, vol. 27, no. 1, pp. 66–75, 2010.
- [51] H. Mardani Kamali *et al.*, “LUT-Lock: A novel LUT-based logic obfuscation for FPGA-bitstream and ASIC-hardware protection,” in *IEEE ISVLSI*, 2018, pp. 405–410.
- [52] A. Alaql *et al.*, “Scalable attack-resistant obfuscation of logic circuits,” 2020. [Online]. Available: <https://arxiv.org/abs/2010.15329>
- [53] L. Alrahis *et al.*, “GNNUnlock: Graph neural networks-based oracle-less unlocking scheme for provably secure logic locking,” in *IEEE/ACM DATE*, 2021, pp. 780–785.



Dominik Sisejkovic received the B.Sc. and M.Sc. degree in computing from the University of Zagreb, Croatia, in 2014 and 2016 respectively. In 2016, he started working as a Ph.D. student and research assistant at RWTH Aachen University. His research interest includes hardware security and machine learning for security. He was directly involved in the design and implementation of the logic-locking framework that was applied for the production of the first logic locked RISC-V processor on the market.



Farhad Merchant received his Ph.D. from the Indian Institute of Science, Bangalore (India), in 2016. His Ph.D. thesis title was “Algorithm-Architecture Co-design for Dense Linear Algebra Computations”. He worked as a postdoctoral research fellow at NTU Singapore, from March 2016 to December 2016. He joined Institute for Communication Technologies and Embedded Systems, RWTH Aachen University, in December 2017 as a postdoctoral research fellow in the Chair for Software for Systems on Silicon.



Lennart Reimann received his Bachelor’s and Master’s degree in Electrical Engineering, Information Technology, and Computer Engineering from RWTH Aachen University in 2016 and 2019, respectively. After his Master’s, he started working toward his Ph.D. as a research assistant at ICE under the supervision of Prof. Leupers. His research interest includes Hardware Security, Secure ASIC design, Cryptographic accelerator design, etc.



Rainer Leupers received the M.Sc. (Dipl.-Inform.) and Ph.D. (Dr. rer. nat.) degrees in Computer Science with honors from TU Dortmund in 1992 and 1997. From 1997–2001 he was the Chief Engineer at the Embedded Systems Chair at TU Dortmund. In 2002, he joined RWTH Aachen University as a professor for Software for Systems on Silicon. His research comprises embedded software development tools, multicore processor architectures, hardware security, and system-level electronic design automation.