

Hardening Circuit-Design IP Against Reverse-Engineering Attacks

Animesh Chhotaray
University of Florida
chho58@ufl.edu

Thomas Shrimpton
University of Florida
teshrim@ufl.edu

Abstract—Design-hiding techniques are a central piece of academic and industrial efforts to protect electronic circuits from being reverse-engineered. However, these techniques have lacked a principled foundation to guide their design and security evaluation, leading to a long line of broken schemes. In this paper, we begin to lay this missing foundation.

We establish formal syntax for *design-hiding (DH) schemes*, a cryptographic primitive that encompasses all known design-stage methods to hide the circuit that is handed to a (potentially adversarial) foundry for fabrication. We give two security notions for this primitive: function recovery (FR) and key recovery (KR). The former is the ostensible goal of design-hiding methods to prevent reverse-engineering the functionality of the circuit, but most prior work has focused on the latter. We then present the first provably (FR,KR)-secure DH scheme, *OneChaff_{hd}*.

A side-benefit of our security proof is a framework for analyzing a broad class of new DH schemes. We finish by unpacking our main security result, to provide parameter-setting guidance.

Index Terms—cryptography; provable security; design hiding; hardware obfuscation; logic locking; logic encryption; IC camouflaging

I. INTRODUCTION

Modern integrated-circuits (ICs, or “chips”) are the product of a globally distributed supply-chain [1]. Much of this is driven by economics: the cost of building and operating a chip fabrication facility is exorbitant, so circuit designers are forced to send their digital intellectual property (IP) to external foundries for fabrication. In 2018, just ten foundries accounted for more than 95% of the chip-fabrication market [2].

Given the lack of choice they have concerning who will fabricate and package their IP into chips, IP authors are motivated to protect their high-value circuit designs from being reverse-engineered and stolen by untrusted foundries. Along with IP theft, inserting difficult-to-detect malicious implants (or *trojans*) into ICs and producing out-of-contract, counterfeit ICs are the three major security issues in the IC supply chain.

We focus on the threat of IP theft because (intuitively) successful trojan insertion and IC counterfeiting are more difficult tasks when the adversary is denied access to the original circuit design. What’s more, there are reports that estimate the cost of IP theft, and subsequent counterfeiting to be more than a hundred billion dollars in 2011 alone [3], [4]. Note that while there are legal restrictions for the foundries to not steal the circuit design of IP authors, the fact that they can steal and sell counterfeit chips in the black market without

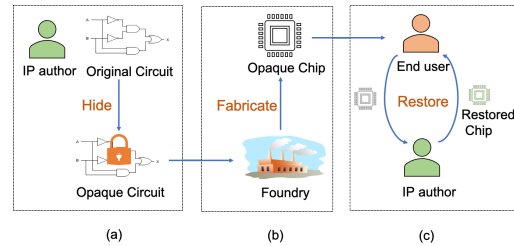


Fig. 1: Simplified view of the setting for Design-Hiding schemes. The IP author is involved in hiding (a) and restoring (c) chips. The (potentially malicious) foundry controls fabrication (b) of chips.

getting caught gives them incentive to fabricate counterfeit chips.

IP theft by the foundry is enabled when it obtains (effectively) the gate-and-wire layout for the circuit to be fabricated. Thus, the last decade has seen a surge in research on methods to “hide” the circuit IP [5]. We refer to such methods as *design-hiding (DH) schemes*. An important constraint on DH schemes is that the IP author still needs the foundry to fabricate something useful; simply applying traditional encryption to the gate-and-wire layout and handing the resulting ciphertext to the foundry does not suffice. (We discuss relation to other cryptographic primitives like program obfuscation, multi-party computation and function secret-sharing in Section B.)

Framing the problem. Fig. 1 gives a simplified picture of the setting in which DH techniques are typically deployed. An IP author attempts to hide (Fig. 1(a)) the underlying functionality F that its plaintext circuit computes by presenting the foundry with an *opaque* version of the circuit. Upon receiving the opaque circuit, the foundry is meant to fabricate it (Fig. 1(b)) into physical chips that compute whatever the opaque circuit does. These chips will likely not compute F , but this functionality can be restored (Fig. 1(c)) by the IP author. Intuitively, the IP author’s security goal is that only they can produce a chip that correctly computes F , i.e., the IP author’s restoration mechanism is required, and only they can execute it properly. Chips that have been restored may enter the market and are available for purchase by the end-user.

Characterizing prior approaches. Broadly speaking, prior

work on DH schemes falls into three categories: logic locking [6]–[8], IC camouflaging [8]–[10], and split manufacturing [11], [12]. Logic-locking¹ techniques are employed when an IP author must outsource the entire fabrication process to an untrusted, third-party foundry. Here, the foundry is provided a complete circuit description (modulo, perhaps, a small block of uninitialized write-once, tamper-resistant memory) to fabricate. This is the setting most commonly considered in the literature on circuit-design hiding. Split manufacturing primarily aims to prevent an untrusted foundry from inserting hardware-trojans into the IP author’s circuit [13]; as the name suggests, the fabrication of chips is spread across multiple parties, one of which is (typically) the IP author or a trusted proxy. When the trust model is altered, so that the IP author trusts the foundry but not the end user, IC camouflaging methods [10], [14] can be employed. Here, the foundry fabricates an opaque circuit that contains several “camouflaged” logic cells. An end user that purchases fabricated chips from the market does not know the functionality of these camouflaged cells and hence will fail to recover the hidden design.

Our focus will be on adversarial foundries. As such, IC camouflaging is out of scope. Split manufacturing targets adversarial foundries, but it requires (effectively) that the IP author has fabrication capabilities. This is most often not the case, in practice.

Thus, we drill down on the logic-locking approach to hiding IP from untrusted foundries, and building DH schemes that are provably secure in the attack model assumed by logic-locking schemes. Here, the foundry has *unrestricted access*² to the opaque circuit and full control over the entire fabrication process, making this a challenging setting. We note that the foundry may also play the role of an end-user by purchasing packaged (restored) chips from the market. By running the chips on inputs of its choosing, it may learn the value of the original circuit design on a subset of inputs. Finally, we note that although our focus is on the logic-locking setting, our syntax also covers IC camouflaging and split manufacturing techniques. Please see Section A to get a historical context of logic-locking schemes.

Advancing the state of the art. The design of logic-locking schemes (2008–2021) yielded schemes that were efficient, in terms of power, performance and area overhead of the opaque circuit, but ultimately insecure. We attribute this to a myopic view of what it means for a scheme to be secure, combined with a natural eagerness to focus upon efficiency. We are not alone in this belief: in their 2019 paper that broke the “provably secure” SFLL-HD scheme [15], Sirone and Subramanyan state

“Our results reinforce the observation that all logic locking schemes appear to be vulnerable to attack.

¹This is one of several monikers used in the hardware community, others include hardware obfuscation, logic encryption, design withholding and encryption.

²We use the racially neutral term “unrestricted access” instead of “whitebox access” to a circuit/function. Similarly, we use “oracle access” instead of “blackbox access.”

We assert this is because the logic locking community has not adopted notions of provable security from cryptography.” [16]

Indeed, the paper that presented SFLL-HD claimed provable security because it (provably) resisted *three* prominent attacks; but not Sirone and Subramanyan’s FALL attack. This provides a clear reminder that a scheme should claim security — explicitly scoped by a principled and well defined security notion— only if it provably thwarts *all* (suitably efficient) attacks that are admitted by that notion.

In response to this state of affairs, we give a principled, provable-security foundation to the area. Our main contributions are as follows:

- 1) First, we formalize DH schemes as an abstract, syntactic object, i.e., a definition of what a DH scheme *is*. A DH scheme will be defined by two component algorithms, *Hide* and *Restore*, that correspond to the hiding and restoration phases discussed, above. In addition, we elevate the fabrication step, of turning circuit descriptions (necessarily visible to the foundry) into fully packaged chips (which require considerable resources and expertise to “open” and analyze), to an explicit algorithm *Fab*. This has important implications for security that have not previously been surfaced.
- 2) We establish two security notions that capture the capabilities, and goals (KR and FR), of an adversarial foundry. These notions make formal the attack model considered in prior works, and attends to details that have sometimes been quietly elided, e.g., a priori knowledge about the hidden function F that the foundry may have. We note that the SAT attack [17], and other key-recovery attacks [16]–[26], are admitted by our formal KR-security notion. (Some of these attacks are against DH schemes [15], [27]–[32] that were designed (post 2015) to thwart SAT attack.) The KR-security notion and our formal FR-security notion share the same abstract attack model, and Theorem 1 shows that KR *insecurity* immediately implies FR *insecurity*. Thus, the existence of these attacks implies that prior logic-locking schemes cannot achieve our notion of FR security.
- 3) We give the first DH scheme, *OneChaff*_{hd}, that provably protects a broad class of functions (or their circuit representation) against reverse-engineering attacks. Along the way, we observe that certain types of “simple” functions (e.g., functions with small domain) cannot be protected by any DH scheme.

In Section II we will discuss our contributions in more detail; but, before that, let us mention some of the immediate next steps that future work might take.

Next Steps. Our work *initiates* a provable-security exploration of DH schemes, putting the goals of prior work on a solid foundation, and showing that these goals are achievable. As such, it completes an important first step towards the ultimate goal of providing IP authors with efficient schemes that provably protect broad classes of practically relevant circuit

designs. But there is much to be done between here and there, and the following list captures some of the research challenges that our results suggest.

- 1) **HIDING MORE REALISTIC CIRCUITS.** We consider DH schemes for stateless circuits, as these were the target of most prior work. However, many real-world circuits are stateful, and comprised of multiple stateless sub-circuits. There are very few DH schemes [33]–[36] for stateful circuits and all of them have been shown to be broken by a recent attack [37]. Thus, extending our formalisms to such circuits is an important next step. One may treat each stateless sub-circuit as an independent circuit, and try to use $OneChaff_{hd}$ to prevent FR-attacks on each of these. For reasons that we give in Section II, $OneChaff_{hd}$ focuses on hiding Boolean circuits, i.e., those with a single output bit. As we will explain, extending $OneChaff_{hd}$ to stateless circuits with multiple output bits requires some care. Intuitively, one can take the transitive fan-in cone³ (TFC) of any output bit and use $OneChaff_{hd}$ to hide the Boolean (sub)circuit that is the indicated TFC. But determining which TFC(s) to hide may not be straightforward, as we will see.
- 2) **MORE PRACTICAL SECURITY GOALS.** Capturing the security goals of prior work, our FR-security notion demands that the foundry recover the *full* input-output behavior of the hidden function. But in practice, this is likely too strong. In particular, it discounts attacks that recover a “good enough” approximation of the hidden function. There are many potential ways to address this. For example, one might modify our FR-security experiment to demand input-output correctness on a subset of the function’s domain; but how should this subset be determined, and how large should it be? Perhaps better, one could add an explicit *Test* algorithm the syntax of a DH scheme, and then modify the FR-security experiment to declare an attack successful only if it fools the *Test* algorithm into saying that the foundry’s dishonestly produced chip (not circuit) is functionally correct. This would allow the IP author (via the DH scheme) to specify what counts as a “win” for the adversarial foundry, rather than having the security notion assert it.
- 3) **EXPLORING THE EFFECTS OF A PRIORI KNOWLEDGE.** Earlier works implicitly assumed that the adversary has *no* a priori information about the hidden function. Our FR-notion addresses this shortcoming, by restricting the to-be-hidden function to come from a specified set (which is a parameter of the notion); intuitively, this set captures the ability of the foundry to narrow the search space. In our analysis of $OneChaff_{hd}$, we assume that the adversary knows the exact hamming weight of the hidden function, and we will justify this choice. But a thorough exploration of the effects of a priori knowledge, on $OneChaff_{hd}$ or

³The transitive fan-in cone of an output bit in a function is the smallest subgraph in the DAG (circuit) representation that connects the primary inputs to the output bit.

any other realization of a DH scheme, is likely to be illuminating.

- 4) **EMPIRICAL EVALUATIONS.** We prove that our $OneChaff_{hd}$ construction is FR-secure for a wide range of functions. Still, a “head-to-head” empirical evaluation of published attacks against $OneChaff_{hd}$ and existing schemes, for a range of parameter settings and families of circuits, would be interesting. Among other things, this would help to establish the tightness of our security results, and aid in guiding security-sensitive choices in practice. Another dimension for such an empirical evaluation is to explore the power, performance and area overhead of the opaque circuits created by $OneChaff_{hd}$ vs. those produced existing DH schemes. For reasons of scope and focus, we do not explore this evaluation in the current work.

II. OVERVIEW OF CONTRIBUTIONS

We now provide a more detailed overview of our contributions, before engaging with the technical core of the paper.

Formal foundations for DH schemes: Syntax. Notably absent from the area is a provable-security foundation for the design and analysis of DH schemes. Very few papers in the area offer anything along these lines. The works that do [5], [8], [15], [38], [39] fall short of what is needed, e.g., by giving syntactic descriptions that are imprecise or clearly mismatch existing schemes.

Such a foundation begins with a precise definition of a DH scheme as a syntactic object, i.e., what are the component algorithms that must be realized in a concrete scheme.

So, we begin by providing a formal syntax (in Section IV) for DH schemes, and our formalization captures all currently known methods of design-stage circuit hiding. Specifically, a DH scheme is a pair of algorithms (*Hide*, *Restore*) that abstract the portions (a) and (c) of Fig. 1, respectively. Loosely, the *design-hiding* algorithm *Hide* takes as input a circuit C_F (and some design parameters θ), and it returns an opaque circuit C_L , along with the associated hiding key K_O . The *design-restoring* algorithm *Restore* takes an opaque *chip* \mathbb{C}_L , a hiding key K_O and design parameters θ as inputs; it returns either a restored chip \mathbb{C}_F or an error symbol \perp , i.e. an indication that restoring has failed.

The opaque circuit is transformed into a chip by a separate *chip-fabrication* algorithm *Fab* that takes a circuit C_L and design parameters as inputs and returns a chip \mathbb{C}_L . Notice that we use the heavy typeface in \mathbb{C}_L to distinguish between unrestricted access to circuits (e.g. C_L) and oracle access to chips. This is necessary as otherwise the foundry can purchase a restored chip and use invasive attacks [40] to read the hiding key from (tamperproof) memory. In this work, we consider such attacks to be out of scope as protecting against invasive attacks will likely require design of special hardware like active shields [41] and are hence, orthogonal to the development of DH schemes. We note that no prior work considered the fabrication process, which turns circuits into *chips*, as a first-class syntactic primitive. The effect is

that fabrication-specific security issues could not be cleanly surfaced. We will see that making the fabrication process explicit uncovers an important connection between the security of DH schemes, and detecting stealthy hardware trojans (in packaged chips). More on this in a moment.

Formal foundations: Security notions. Given a precise description of what a DH scheme *is*, we next define formal notions of what it means for a DH scheme (however it is realized) to be *secure*. An intuitive definition of reverse-engineering the opaque circuit is to recover from it, the hidden IP F by any means. But literature has tended to focus on attacks (and countermeasures) that attempt to recover the secret hiding key K_O . Thus, we give two formal notions of security (in Section V): *function recovery* (FR) and *key recovery* (KR). In both notions, the adversary is provided unrestricted access to the opaque circuit, and various oracles that abstractly capture the powers of a foundry.

In the KR notion, the adversary’s goal is to return a key K that is equivalent to K_O , in the sense that when one restores the functionality of an honestly fabricated chip using either K or K_O , we get restored chips with identical functionalities.

The FR notion captures a stronger attack model. In it, the adversary’s goal is to find *any* chip that is functionally equivalent to F . As one expects, KR-insecurity implies FR-insecurity: if you can recover a key K equivalent to K_O , then you can win the FR game by returning an honestly fabricated chip restored with K . The converse is not necessarily true, i.e., reverse-engineering the hidden functionality of the opaque circuit does not necessarily require recovering something equivalent to the hiding key.

We note that certain kinds of functions cannot be protected by *any* DH scheme in the logic-locking setting, where the foundry may purchase honest chips and thereby learn input-output pairs of its choosing. For example, if the domain of the chip is small, the functionality of the chip can be recovered by querying the chip on its entire domain. In the case of Boolean functions, those whose decision-tree representations have small depth/size cannot be hidden [42], nor can those whose Fourier spectra contain relatively few significant components [43]. So, while our security notions are agnostic to structural characteristics of the function(s) one wishes to hide, our security results will surface this concern.

We also note that our security notions allow for fully malicious foundries that may fabricate arbitrary, “dishonest” chips, and submit these to be restored with the secret hiding-key K_O . The chip may have been fabricated from the opaque circuit, but (say) with an embedded hardware trojan that outputs K_O when triggered on a particular input. Unless knowledge of the secret K_O suffices to allow the *Restore* algorithm to detect such a trojan (and alert the IP author not to proceed), the restored chip can be run by the foundry (acting as user), leaking K_O and allowing it to win the FR game. Given the state of the art in trojan detection, we know of *no remotely practical* DH scheme that can be FR-secure against fully malicious foundries. Thus, we restrict our attention to

designing DH schemes that are secure against honest-but-curious foundries, i.e., ones that will try to reverse-engineer the functionality of the IP, but will only fabricate chips that adhere to the IP author’s opaque circuit. This is in keeping with all prior work on DH schemes.

A new family of DH schemes: *OneChaff*. We introduce a family of DH schemes that we call *OneChaff* (see Section VI), and analyze a particular scheme *OneChaff*_{hd} in this family. In *OneChaff*_{hd}, the *Hide* algorithm takes inspiration from SFL-flex [15] as it encodes a single n -input-bit Boolean function \underline{H} (one “chaff” function) and an uninitialized lookup table $\overline{\text{Tab}}$ in the opaque circuit. While SFL-flex allows arbitrary H , in our *OneChaff*_{hd} scheme, the chaff H matches the hidden function F , except on $\Delta \in \mathbb{N}$ uniformly chosen inputs. These are the so-called *distinguishing inputs* (DIs) for the pair (H, F) . The hiding key K_O encodes the correct input-output behaviors on the DIs. (Practical *Hide* algorithms will have $\Delta \ll 2^n$.) On input of a key K , chip \mathbb{C}_L , and design parameters θ , the *Restore* algorithm in *OneChaff*_{hd} loads the key K into the (write-once, tamper-resistant) uninitialized lookup table of the chip. Under honest operation, key K is equal to K_O .

Proving security of *OneChaff*_{hd} for Boolean functions.

After giving our foundations for DH schemes and introducing the *OneChaff* family of schemes, the remainder of this work is spent showing that *OneChaff*_{hd} provably prevents full recovery of *Boolean* functions in the presence of honest-but-curious adversarial foundries. While most real-world circuits do not compute functions returning a single bit, several prominent logic-locking schemes [15], [27]–[29] only aim to hide Boolean functions. Moreover, no provably secure scheme exists for circuits implementing functions from this “base” class. We note that for circuits with multiple output bits, one can attempt to hide Boolean sub-functions that are determined by the transitive fan-in cone (TFC) of individual output bits.

Our main security result (Theorem 3) gives an upper-bound on the probability that a computationally bounded, honest-but-curious foundry manages to win the FR game against *OneChaff*_{hd}. To the best of our knowledge, this is *the first positive provable-security result* on DH schemes.

Security holds for Boolean functions that are not “simple” in the sense we mentioned earlier (no DH scheme can hide those) under some assumptions about the *a priori* knowledge that the adversary has about F . All prior schemes assume that the adversary has *no a priori* knowledge of F . Such an assumption is unrealistic and also makes the adversary weak as the initial “guess” space of the adversary is the set of all Boolean functions. In our FR analysis of *OneChaff*_{hd}, we assume that the adversary knows a priori the hamming weight h , i.e., the number of inputs that cause F to output one. This narrows the initial “guess space” to Boolean functions that have hamming weight of h . Also, the hamming-weight parameter allows us to capture the fact that the number of functions in the guess space of the adversary increases exponentially in h . Hence, functions with hamming weight close to 2^{n-1} will be more

secure compared to functions with hamming weights close to zero or 2^n . This is also intuitive and in agreement to a result from learning theory that states that a random Boolean function (with sufficiently large domain) will not be simple, as it will lack the highly concentrated Fourier spectral structure that is typically needed for a function to be learnable. Note that a random Boolean function will have expected hamming weight close to 2^{n-1} .

Our analysis essentially bounds the number of functions that remain in the adversary’s guess space after some number of true input-output observations $(X_1, F(X_1)), \dots, (X_q, F(X_q))$ are (adaptively) obtained. Intuitively, if a large number of functions remain in the guess space, then the probability of winning the FR-game will be small, and conversely if the adversary is able to eliminate all but a few functions, the adversary’s winning probability will be close to one.

From Theorem 3, and the analysis leading to it, we can glean some useful observations. In particular, the IP author should use *OneChaff*_{hd} to protect Boolean functions (or Boolean sub-functions) that have large domains, and hamming weights not too close to 0 or 2^n . Functions with small domains cannot be hidden by any DH scheme, at least not without severely restrictive assumptions on the adversary. When the hamming weight tends towards 0 or 2^n , the function tends towards a constant function. Intuitively, as the (known) hamming weight of the hidden function moves away from 2^{n-1} towards either 0 or 2^n , the number of possible functions decreases. This makes it more likely, although not necessarily “likely”, that the hidden function can be guessed after seeing the opaque circuit and some true $(X, F(X))$ pairs. Finally, the IP author should choose to make Δ as large as is feasible. Intuitively, if Δ is small, the number of functions that remain in the adversary’s guess space after it gets access to the opaque circuit (that encodes chaff H and a lookup table $\overline{\text{Tab}}$) will also be small compared to large Δ . Note that the new guess space will contain only functions that have a hamming weight of h (due to its a priori knowledge of the hamming weight of F) and that also differ from H on Δ DIs (by construction).

III. PRELIMINARIES

Basic notation. When X, Y are strings, we establish the following notations. We write $X \parallel Y$ to denote the concatenation of X and Y ; $X[i]$ for the i -th element of X ; and $|X|$ to denote the length of X . We extend the last two notations to ordered objects (e.g., a sequence, list, table).

When T is any ordered object, we write $T[i]$ for the i -th element, and $|T|$ to denote the number of elements in T . In pseudocode, our convention will be: all such T are initialized to $T[i] = \perp$ for all values of i , where \perp is a distinguished symbol. Likewise, all sets will be initially empty. We use the notation $\langle V \rangle$ to denote the encoding of object V as a bit string. The method of encoding is left implicit, and it is silently overloaded to accommodate whatever is the type of V .

When m is an integer, we use the standard notation $[m]$ to denote the set $\{1, 2, \dots, m\}$. We write $v_1, v_2, \dots, v_r \leftarrow^s V$ to

denote sampling (with replacement) $r > 0$ uniform elements of V , where V is some non-empty set. An unembellished \leftarrow denotes deterministic assignment. This notation is also used for randomized algorithms, i.e., $x \leftarrow^s A(\dots)$ means that algorithm A runs on its indicated inputs, and halts with an output that is assigned to x . (In this case, the distribution on output values is determined by A , and is not necessarily uniform.)

We write $A^{\mathcal{O}_1, \mathcal{O}_2, \dots}$ to denote that algorithm A has oracle access to the superscripted oracles. For randomized algorithms X that may return a distinguished symbol \perp , the support of the algorithm $\text{Sup}(X)$ is defined only over non- \perp outputs. An *adversary* is a randomized algorithm.

Functions and their representations. When \mathcal{D}, \mathcal{R} are non-empty sets, we write $\text{Func}(\mathcal{D}, \mathcal{R})$ for the set of all functions $F: \mathcal{D} \rightarrow \mathcal{R}$. We write $\text{Func}(n, m)$ as shorthand for functions with $\mathcal{D} = \{0, 1\}^n$ and $\mathcal{R} = \{0, 1\}^m$. When $m = 1$, $\text{Func}(n, 1)$ is the set of all Boolean functions. We will use three representations of functions in this work: *circuits*, *chips* and *lookup tables*. While circuits and chips will be used to describe the design phase of the IC supply chain and are important to denote the type of access (oracle or unrestricted) to a circuit that an adversary gets as part of the threat model, lookup tables will be used primarily for functional analyses.

Formally, a *circuit* C_F is a directed acyclic graph that implements some mapping $F \in \text{Func}(\mathcal{D}, \mathcal{R})$. Access to circuits will always be unrestricted. This captures a reality of our setting, in which the foundry is handed a description of a circuit to be fabricated. Once a circuit C_F is fabricated, we will refer to it as a *chip* and use the heavy typeface \mathbb{C}_F to make this distinction clear. Crucially, access to \mathbb{C}_F is *not* unrestricted; rather an adversary can only use \mathbb{C}_F to make *oracle* queries. This syntactic choice is to make invasive attacks [40] on chips (to leak secrets) out of scope as they are orthogonal to the (algorithmic) development of design-hiding schemes — the central primitive in our work.

The lookup-table representation of a function F with domain \mathcal{D} and range \mathcal{R} is a table $T_F = \langle (X_1, Y_1), (X_2, Y_2), \dots \rangle$, where $X_i \in \mathcal{D}$ and $Y_i = F(X_i)$.

In all representations of functions, when the underlying mapping is implicit/understood, we will omit the subscript.

When F, G are two n -bit to m -bit functions, the *hamming distance* between F and G is the number of inputs X_i on which the value of the functions differ. We use $\mathcal{I}^\neq(F, G)$ to denote the set of such *distinguishing inputs*, i.e., $\mathcal{I}^\neq(F, G) = \mathcal{I}^\neq(G, F) = \{x \in \{0, 1\}^n \mid F(x) \neq G(x)\}$. Formally, $\text{hd}(F, G) = |\mathcal{I}^\neq(F, G)|$. We write $F \equiv G$ whenever $\text{hd}(F, G) = 0$.

We will also find it useful to define the *hamming weight* of Boolean functions. When $F \in \text{Func}(n, 1)$, the *hamming weight* of F is defined as the number of inputs that map to one, i.e., $\text{hw}(F) = |\{X \in \{0, 1\}^n \mid F(X) = 1\}|$. Also, for Boolean functions, we will use sets $\mathcal{X}_i(F) = \{X \mid F(X) = i\}$, where $i = 0$ or 1 , to denote the set of inputs for which F map to 1 or 0. Notice that these two sets fully define F as we can

construct the truth table of F using $\mathcal{X}_0(F)$ and $\mathcal{X}_1(F)$.

For any $F \in \text{Func}(n, m)$, we define the Δ -neighborhood of F as $N_\Delta(F) = \{H \in \text{Func}(n, m) \mid \text{hd}(F, H) = \Delta\}$.

IV. DH SCHEMES

We begin by defining design-hiding scheme as a syntactic object. Loosely speaking, the syntax describes the inputs and outputs of the core algorithms that any DH scheme must realize, as well as what it means for a DH scheme to operate correctly.

Definition 1. Fix integers $n, m > 0$. A design-hiding (DH) scheme $\Pi = (\text{Hide}, \text{Restore})$ for $\text{Func}(n, m)$, with key length $k_o: \{0, 1\}^* \rightarrow \mathbb{N}$ is a tuple of algorithms with the following syntax.

- *The randomized design-hiding algorithm Hide:*
Inputs: a circuit C_F implementing $F \in \text{Func}(n, m)$, and a string of design parameters θ .
Outputs: the distinguished symbol \perp (“error”), or a tuple consisting of (1) a secret key $K_O \in \{0, 1\}^{k_o(\theta)}$, (2) and an opaque circuit C_L , where $L \in \text{Func}(n + k_o(\theta), m)$.
Requirements: for every input (F, θ) , either $\Pr[\text{Hide}(C_F, \theta) = \perp] = 1$ or $\Pr[\text{Hide}(C_F, \theta) = \perp] = 0$.
- *The deterministic design-restoring algorithm Restore:*
Inputs: a key $K \in \{0, 1\}^{k_o(\theta)}$, a string of design parameters θ , and a chip \mathbb{C}_L .
Outputs: a restored chip \mathbb{C}_F or \perp . When K_O — the key that is used by *Hide* to produce C_L — is the key input to *Restore*, we refer to the resulting chip \mathbb{C}_F as an honestly-restored chip.
Requirements: for every input $(K, \theta, \mathbb{C}_L)$, it must be that $\mathbb{C}_F \leftarrow \text{Restore}(K, \theta, \mathbb{C}_L) \neq \perp$ implies that $F \in \text{Func}(n, m)$.

We assume that if an algorithm is called on a point outside of its domain, in particular if any of its inputs are \perp , then the algorithm returns \perp . \diamond

The syntax that we have just established is fashioned to capture the techniques — logic locking, IC camouflaging, and split manufacturing — that an IP author uses to (a) protect the “privacy” of its high-value circuit-design C_F from adversarial entities in the post-design phase of the IC supply chain, and (b) make functionally-correct chips available to the end user. Let us elaborate using the logic-locking setting.

The IP author often outsources the fabrication of C_F into physical chips \mathbb{C}_F to third-party foundries. (We will formalize this transformation in the next section.) The foundries are potentially malicious entities. Hence, the IP author cannot give the foundries the circuit design C_F (in plaintext). The design-hiding algorithm *Hide* abstracts the mechanism by which the IP author turns C_F into an opaque circuit C_L that “hides” the functionality F using a secret key K_O . The circuit C_L takes $k_o(\theta) + n$ bits of input, where n is the length of the input to F . Here, the additional $k_o(\theta)$ bits encode the key K_O , which will be used to restore the functionality of the chips (that the foundry produces) to the original, intended functionality F .

Specific instantiations of *Hide* include the locking process in *logic-locking* schemes like random logic-locking (RLL) [6] and its variants [44]–[47], SAT-attack-resistant schemes [15], [27]–[29], [32], etc. We insist that for a given pair (F, θ) , either *Hide* works always (returns valid K_O and L) or it always fails (returns \perp).

The meaning of the parameter string θ depends heavily on the particular instantiation of the design-hiding scheme. For example, it may encode the number of “key gates” to be inserted in an RLL scheme [6], or the target hamming distance in SFLL-HD [15]. It may also encode other design constraints, such as the maximum size and depth of a locked circuit, restrictions on gate types, and so on.

We envision that the IP author will securely store the key K_O , and send the opaque circuit C_L to the foundry, instead of the (plaintext) circuit C_F . The foundry will fabricate and package one or more opaque chips; if the foundry is *honest*, then each chip will implement L . We formalize this in a moment.

Intuitively, the design-restoring algorithm *Restore* abstracts the mechanism by which a fabricated and packaged chip \mathbb{C}_L is restored to its original, intended functionality F . Loosely, this entails fixing the $k_o(\theta)$ “key bits” in the input to \mathbb{C}_L to K_O . In logic locking, the opaque chips are restored by (at least) having a key installed in some tamper-proof, one-time writable memory unit within the chip.

Fabrication and DH-scheme correctness. In order to define the correctness for a DH scheme, we need some mechanism for turning an opaque circuit into a chip. This is exactly the role of the fabrication process that the foundry is meant to provide. Thus, let *Fab* be a randomized *chip-fabrication algorithm*: it takes as inputs a circuit C_L and a string of design parameters θ , and it outputs either a chip \mathbb{C}_L , or the error symbol \perp .

Definition 2. A DH scheme $\Pi = (\text{Hide}, \text{Restore})$ is correct with respect to chip fabrication *Fab* if, for any (F, θ) and any $(K_O, C_L) \in \text{Supp}(\text{Hide}(C_F, \theta))$, we have

$$\Pr[\mathbb{C}_G \leftarrow \text{Restore}(K_O, \theta, \text{Fab}(C_L, \theta)) : (\mathbb{C}_G \neq \perp) \Rightarrow (G \equiv F)] = 1,$$

where, the probability is over the coins of *Fab*. \diamond

In words, this requirement asserts that whenever (F, θ) is a pair that results in an opaque circuit C_L with associated key K_O , it must be the case that an honestly-restored chip computes F exactly.

We will find it useful to define an *honest* chip-fabrication algorithm. Loosely, we say that a chip-fabrication algorithm *Fab* is honest if the chip that it produces computes exactly what it is supposed to compute. Formally, this requires that for any circuit C_L , and design-parameters θ , $\Pr[\mathbb{C}_{\hat{L}} \leftarrow \text{Fab}(C_L, \theta) : (\mathbb{C}_{\hat{L}} \neq \perp) \Rightarrow (\hat{L} \equiv L)] = 1$, where the probability is over the coins of *Fab*. We note that the correctness of a DH scheme does not require honest *Fab*; indeed, $\text{Fab}(C_L, \theta)$ may produce a chip $\mathbb{C}_{\hat{L}}$ such that for

some $Y \in \text{Dom}(F)$, $\hat{L}(K, Y) \neq L(K, Y)$ for (say) a small set of keys $K \neq K_O$. For example, a hardware trojan [48] can be embedded in the description of Fab such that a fabricated opaque chip $\mathbb{C}_{\hat{L}}$ implements $L(K_O, \cdot)$ honestly; on keys $\hat{K} \neq K_O$, \hat{L} may leak information about K_O .

V. SECURITY NOTIONS

We consider two notions of security for DH schemes. Both notions deal with an adversarial foundry that attempts to recover the original functionality that the opaque circuit hides. The adversary gets unrestricted access to the opaque circuit, as this is something that a real foundry would receive in order to carry out fabrication. It is also given oracles that capture various capabilities that a foundry is likely to have.

Function recovery. We begin with a notion of function recovery for a design-hiding scheme $\Pi = (\text{Hide}, \text{Restore})$. The pseudocode for the FR experiment is given in Fig. 2. The experiment is parameterized by the DH scheme Π , some design parameters θ that Π uses, and an honest chip-fabrication algorithm Fab . It takes as inputs: an adversary’s attack algorithm A , and a set $\mathcal{F} \subseteq \text{Func}(n, m)$. Intuitively, the set \mathcal{F} captures the *a priori* uncertainty/knowledge that the adversary has about the function F that the DH scheme Π is used to protect. Clearly, there must be some uncertainty, since otherwise, the adversary already knows the functionality of the IP author’s circuit design. Note that the design of Π does not depend on \mathcal{F} .

The experiment begins by sampling F uniformly from \mathcal{F} . A circuit C_F (that implements F) and the design parameters θ are inputs to the hiding algorithm Hide , which returns the key K_O , and the opaque circuit C_L . The adversary is given C_L and θ as inputs, and it is provided oracles named FAB, RESTORE, and RUN. These oracles model processes that an adversarial foundry can employ while trying to recover F .

When the foundry is *honest-but-curious*, it will *only* produce chips that are fabricated from the opaque circuit that the IP author generates using the DH scheme Π , and the design parameters θ . In the *fully-malicious* setting, the adversary can produce any *arbitrary* chip using arbitrary circuits and arbitrary design parameters.

For our security experiments, we define an honest-but-curious adversary as one that *always* queries FAB with zero arguments, i.e., $Z = \varepsilon$. On the other hand, a fully-malicious adversary can call FAB with any circuit C_N , and any design parameters γ of its choice, i.e., $Z = (C_N, \gamma)$.

When $Z = \varepsilon$, the FAB oracle runs Fab with C_L (generated in the FR experiment) and θ as inputs to get an *honest* chip \mathbb{C}_L . Otherwise, FAB parses Z to get (C_N, γ) and runs the Fab algorithm on (C_N, γ) to get an *arbitrary* chip \mathbb{C}_N of the adversary’s choice. We use the sets \mathcal{H} and \mathcal{N} to keep track of the honest and arbitrary chips, respectively.

The RESTORE oracle models the adversary’s ability to obtain honestly-restored chips, i.e., chips that are restored with the secret key K_O . We allow the adversary to query RESTORE on any chip that it obtained from FAB, i.e., chips in set $\mathcal{H} \cup$

\mathcal{N} . We use \mathcal{U} to keep track of honestly-restored chips. The RESTORE oracle does not return the restored chips in order to prevent the adversary from reading the secret key from the description of the chip. Instead, we return the index of the restored chip in \mathcal{U} . Notice that we do not restrict the adversary to run the deterministic *Restore* algorithm *locally* on any triple $(K, \gamma, \mathbb{C}_N)$ of its choice.

The RUN oracle captures the foundry’s ability to see the output of any honestly-restored chip \mathbb{C}_P (that is stored in \mathcal{U}) on any input of its choice. The oracle takes as input the index j of \mathbb{C}_P in \mathcal{U} and $X \in \{0, 1\}^n$, and returns the value $P(X)$. (By notation, \mathbb{C}_P implements function P .)

The goal of the adversary in the FR experiment is to output a chip \mathbb{C}_G as its guess for F . The adversary is said to win the FR experiment if $F \equiv G$. Notice that we prevent trivial wins — by returning honestly-restored (honestly-fabricated) chips — by not allowing the adversary direct access to chips that are restored by the RESTORE oracle.

We define the *FR advantage* of the pair (\mathcal{F}, A) against DH scheme Π , design parameters θ , and honest-fabrication algorithm Fab to be

$$\text{Adv}_{(\Pi, \theta), Fab}^{\text{FR}}(\mathcal{F}, A) = \Pr \left[\text{Exp}_{(\Pi, \theta), Fab}^{\text{FR}}(\mathcal{F}, A) = 1 \right],$$

where the probability is over the indicated experiment. We say A is (t, q_f, q_s, q_r) -resource when its time complexity is t , and it makes q_f queries to the FAB oracle, q_s queries to the RESTORE oracle, and q_r queries to the RUN oracle. By convention, an FR adversary does not make pointless queries to any oracle, i.e. queries that cannot increase its advantage.

Key recovery. The notion of key recovery (KR) is similar to the FR notion, except the adversary’s goal is to recover a key K , and the key is then used to get an honestly restored chip \mathbb{C}_G that is obtained by running the honest Fab algorithm on (C_L, θ) . In the final step of the KR notion, we check whether the function G that the chip \mathbb{C}_G implements is functionally equivalent to F . We define the *KR advantage* of A as

$$\text{Adv}_{(\Pi, \theta), Fab}^{\text{KR}}(\mathcal{F}, A) = \Pr \left[\text{Exp}_{(\Pi, \theta), Fab}^{\text{KR}}(\mathcal{F}, A) = 1 \right],$$

where the probability is over the indicated experiment. The resources are the same as those for the FR advantage.

Note that one might think it more natural to define key recovery as determining *the* secret key K_O . We define it as we do because some of the existing DH schemes, including RLL [6] and strong logic obfuscation (SLO) [44], admit multiple keys $K \in \{0, 1\}^{k_o(\theta)}$ that map to the hidden function F , i.e., $F \equiv \text{Restore}(K_O, \theta, \mathbb{C}_L) \equiv \text{Restore}(K, \theta, \mathbb{C}_L)$. (These keys are said to constitute an equivalence class of the hidden function.) Thus, our KR notion captures (for example) the SAT attack of Subramanyan et. al [17], which recovers some key in the equivalence class of the hidden function. It also captures other key-recovery attacks [16], [18], [21] that exploit structural and functional characteristics of the opaque circuit.

It is intuitively clear that FR security implies KR security for any design-hiding scheme. We formalize this relation in

Experiment $\text{Exp}_{(\Pi, \theta), \text{Fab}}^{\text{FR}}(\mathcal{F}, A)$ $F \leftarrow \mathcal{F}; i \leftarrow 0$ $(K_O, C_L) \leftarrow \text{Hide}(C_F, \theta)$ $\mathbb{C}_G \leftarrow A^{\text{FAB, RESTORE, RUN}}(C_L, \theta)$ Ret $[G \equiv F]$	Experiment $\text{Exp}_{(\Pi, \theta), \text{Fab}}^{\text{KR}}(\mathcal{F}, A)$ $F \leftarrow \mathcal{F}; i \leftarrow 0$ $(K_O, C_L) \leftarrow \text{Hide}(C_F, \theta)$ $K \leftarrow A^{\text{FAB, RESTORE, RUN}}(C_L, \theta)$ $\mathbb{C}_G \leftarrow \text{Restore}(K, \theta, \text{Fab}(C_L, \theta))$ Ret $[G \equiv F]$	
oracle $\text{FAB}(Z)$: if $Z = \varepsilon$ then //make intended chip $\mathbb{C}_L \leftarrow \text{Fab}(C_L, \theta)$ $\mathcal{H} \leftarrow \mathcal{H} \cup \{\mathbb{C}_L\}$ Ret \mathbb{C}_L else //make arbitrary chip $\langle C_N, \gamma \rangle \leftarrow Z$ $\mathbb{C}_N \leftarrow \text{Fab}(C_N, \gamma)$ $\mathcal{N} \leftarrow \mathcal{N} \cup \{\mathbb{C}_N\}$ Ret \mathbb{C}	oracle $\text{RESTORE}(\mathbb{C}_{\hat{L}})$: if $\mathbb{C}_{\hat{L}} \notin \mathcal{H} \cup \mathcal{N}$ then Ret \perp $\mathbb{C}_P \leftarrow \text{Restore}(K_O, \theta, \mathbb{C}_{\hat{L}})$ if $\mathbb{C}_P \neq \perp$ then $i \leftarrow i + 1; \mathcal{U}[i] \leftarrow \mathbb{C}_P$ Return i	oracle $\text{RUN}(j, X)$: $\mathbb{C}_P \leftarrow \mathcal{U}[j]$ Ret $P(X)$

Fig. 2: Experiments for function-recovery (FR) and key-recovery (KR) notions for DH scheme $\Pi = (\text{Hide}, \text{Restore})$, given $\mathcal{F} \subseteq \text{Func}(n, m)$ and design parameters θ , when chips are fabricated according to Fab .

Experiment $\widetilde{\text{FR}}_{(\Pi, \theta)}(\mathcal{F}, B)$ $F \leftarrow \mathcal{F}$ $(K_O, C_L) \leftarrow \text{Hide}(C_F, \theta)$ $\mathbb{C}_G \leftarrow B^{\text{TRUE}}(C_L, \theta)$ Ret $[G \equiv F]$	oracle $\text{TRUE}(X)$: Ret $F(X)$
----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-------------------------------------------------------

Fig. 3: FR-game in the honest-but-curious setting. All three oracles in the original FR-security experiment are replaced with TRUE.

Theorem 1.

Theorem 1. Fix integers $n, m > 0$, $\mathcal{F} \subseteq \text{Func}(n, m)$, and design parameters θ for a DH scheme $\Pi = (\text{Hide}, \text{Restore})$. Let Fab be honest. For any KR-adversary A , there is an FR-adversary B such that

$$\text{Adv}_{(\Pi, \theta), \text{Fab}}^{\text{KR}}(\mathcal{F}, A) \leq \text{Adv}_{(\Pi, \theta), \text{Fab}}^{\text{FR}}(\mathcal{F}, B).$$

Here, A is (t, q_f, q_s, q_r) -resource where $t = O(q_f + q_s + q_r)$, and B is $(O(t), O(q_f), O(q_s), O(q_r))$ resource. \blacklozenge

The proof of this theorem is straightforward. Given the KR-adversary A guesses K , the FR adversary B can generate \mathbb{C}_G by fixing the $k_o(\theta)$ input bits in the opaque circuit C_L with K (using the deterministic algorithm Restore), and then using Fab to fabricate the restored circuit. With \mathbb{C}_G , B wins the FR experiment with probability no less than the probability with which A wins the KR experiment.

The converse, i.e., KR security implies FR security, is not true. Consider that the opaque circuit that the hiding algorithm outputs is an encoding of the original circuit, and the secret key is sampled uniformly at random from $\{0, 1\}^{128}$. While admittedly pathological, this example suffices to make the point.

Simplified FR-notion in the honest-but-curious setting. Re-

call that in the security experiments, the goal of the adversary is to recover the full-functionality of F . In the honest-but-curious setting, the foundry effectively will have oracle access to F and unrestricted access to L . (We discuss the fully-malicious setting in Section C.) Let's see this in the context of the FR-security definition.

Let C_L be the opaque circuit that the IP author generates using Hide . In the honest-but-curious setting, the foundry uses Fab honestly. By definition of honest Fab , any opaque chip that Fab produces will implement L . Thus, running (the deterministic) Restore algorithm on the honestly-fabricated opaque chips with K_O will result in chips with identical functionalities. Since the adversary gains no additional information about F from multiple queries to the oracles FAB and RESTORE, we can fix $q_f = 1$ and $q_s = 1$, where the only (useful) query that A makes to FAB is $\hat{Z} = \langle C_L, \theta \rangle$, and the only (useful) query that A makes to RESTORE is $\widehat{\mathbb{C}}_L = \text{Fab}(\hat{Z})$. Notice that this results in \mathcal{U} to store a single restored chip $\mathbb{C}_F = \text{Restore}(K_O, \theta, \widehat{\mathbb{C}}_L)$ at index one. Using its q_r queries $(1, X_i)$ to RUN oracle, the adversary will learn q_r pairs $(X_i, F(X_i))$. Thus, in the honest-but-curious setting, we can replace all the oracles with a single TRUE oracle, that takes X_i as input and returns $F(X_i)$. We show these changes in the security experiment $\widetilde{\text{FR}}_{(\Pi, \theta)}(\mathcal{F}, B)$; See Fig. 3.

Notice that the changes are all either syntactic (i.e., no change to the advantage) or conservative (i.e., cannot decrease advantage) with respect to the original FR experiment. This observation gives rise to the following simple, but useful lemma.

Lemma 1. Fix DH scheme Π and parameters θ . For every honest-but-curious FR-adversary A that is (t, q_f, q_s, q_r) -resource, there exists a closely related adversary B such that

$$\Pr \left[\text{Exp}_{(\Pi, \theta), \text{Fab}}^{\text{FR}}(\mathcal{F}, A) = 1 \right] \leq \Pr \left[\widetilde{\text{FR}}_{\Pi, \theta}(\mathcal{F}, B) = 1 \right]$$

where B makes at most q queries to TRUE (see Fig. 3), and

runs in time $O(t)$. \blacklozenge

As discussed in Section B, almost all existing DH schemes [6], [15], [27]–[30], [44]–[46], [49] have been shown to be KR-insecure, i.e., these schemes have large (close to one) KR advantage due to various key-recovery attacks [16]–[26]. Per our Theorem 1, if a DH scheme is KR-insecure, then it is also FR-insecure, i.e., large FR advantage. The few (recent) schemes [31], [32] that have not been shown to be broken were designed specifically to thwart SAT attack. In order to prove FR-security, we need to show that a DH scheme thwarts *all* FR-attacks — not specific KR-attacks. In the next section, we will construct an FR-secure DH scheme.

VI. A FRAMEWORK FOR DESIGNING FR-SECURE DESIGN-HIDING SCHEMES

In the FR-security experiment as well as in $\widetilde{\text{FR}}$, the adversary’s guess space is the set $\mathcal{F} \subseteq \text{Func}(n, m)$ before it gets full access to the opaque circuit. As the adversary learns new information via its unrestricted access to the opaque circuit C_L and oracle access to F , it can prune its guess space to a smaller set $\mathcal{M}_1 \subseteq \mathcal{F}$ of functions. Let \mathcal{D}_1 be the probability distribution of \mathcal{M}_1 . Then, the probability that B wins $\widetilde{\text{FR}}$ will be upperbounded by $2^{-H_\infty(\mathcal{D}_1)}$, where $H_\infty(\mathcal{D}_1)$ is the min-entropy of \mathcal{D}_1 . In this section, we will describe a framework that an IP author can use to build a DH scheme (with formal descriptions of *Hide* and *Restore*) such that it can concretely define the distribution \mathcal{D}_1 , and find concrete FR-security guarantees of its construction.

Abstractly, the opaque chip \mathcal{C}_L can be viewed as a circuit that encodes a set $\mathcal{R}_L \subseteq \text{Func}(n, m)$ of n -bit to m -bit functions (not necessarily distinct) that are selected by *Hide* (either implicitly or explicitly), and each key K in set $\mathcal{K} = \{0, 1\}^{k_o(\theta)}$ is associated with some function in \mathcal{R}_L . Minimally, correctness will require that the hiding key K_O is associated with the true function F . We define a chaff function $H \in \text{Func}(n, m)$ as one that is not functionally equivalent to F ; we refer to the inputs on which H differs (in output) from F as “distinguishing” inputs.

In our framework, we insist that the set \mathcal{R}_L be “extractable” from the opaque circuit, i.e., given C_L and the description of *Hide*, the IP author knows the full-functionality of each function in \mathcal{R}_L . This is a reasonable assumption as during the design of circuit, the functionality of a circuit is fixed first as part of the system specification and architectural design, and then the topology of the circuit is decided. In RLL [6] and its variants [31], [32], [44]–[46], the functionality of the chaff functions is *not* decided during the design of *Hide*. Rather, the chaff functions are an artifact of random structural modifications — the positions of key gates in the original circuit. So, for a non-pathological circuit, i.e., circuits with sufficiently complex functionality and reasonably large domain, it is almost impossible to extract the functionality of the chaff functions. Consequently, it is difficult to evaluate the FR security of such schemes in our framework. As mentioned earlier, RLL and its variants (designed before SAT attack) as

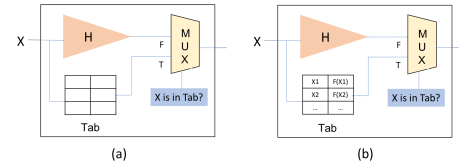


Fig. 4: a) An opaque circuit generated by any OneChaff scheme. b) An honestly-fabricated-and-restored chip. Adversary gets unrestricted access to (a) and oracle access to (b).

well as SAT-attack-resistant schemes [15], [27]–[30] (designed post SAT attack) are FR-insecure due to various key-recovery attacks [16]–[26].

A. OneChaff: a family of DH schemes

We present a family of DH schemes called *OneChaff* with the following common features:

- the design-hiding algorithm returns an opaque circuit that contains the circuit of a single *chaff* H with $\Delta \in \mathbb{N}$ distinguishing inputs and an uninitialized lookup table (see Fig. 4(a)) that is used for restoring the functionality of the hidden function;
- the hiding key encodes the set of distinguishing inputs and the correct value of the hidden function on the distinguishing inputs;
- the design-restoring algorithm, in an honest run, initializes the lookup table with the hiding key (see Fig. 4(b)).

Notice that the description of OneChaff leaves H unspecified; specific OneChaff schemes will define H explicitly.

The opaque circuit that OneChaff generates consists of: the circuit of a single Boolean function, and a lookup table that is effectively a restore unit. Structurally, this is similar to the opaque circuits that SAT-attack-resistant DH schemes [15], [27]–[29] generate. However, the functionalities of the Boolean function and the restore unit are different across schemes. In SFLL-HD [15] and TTLock [29], the encoded Boolean function H is such that its functionality differs from F only on non-key inputs X_i that are at a hamming-distance of h (in SFLL-HD, $0 \leq h \leq 2^n$; in TTLock, $h = 0$) from the correct key K_O ; on the remaining inputs, H and F have identical functionalities. In AntiSAT [27] and SARLock [28], the original function F is encoded in the opaque circuit with the assumption that the circuit-synthesis tools will “meld” the functionalities of F and the restore unit—a comparator that checks whether the key-input to the opaque circuit is identical to the hard-coded correct key (K_O)—in such a way that the adversary will not be able to locate F .

Several key-recovery attacks [16], [21], [26], [50] have exploited structural as well as functional weaknesses in the opaque circuits that the SAT-attack-resistant DH schemes [15], [27]–[29] generate to recover the correct key (and hence the entire functionality of the original circuit). In our construction *OneChaff*_{hd}, which we describe next, we chose the functionality of the Boolean function and the restore unit in such a way

that we can prove its security for a wide class of functions against all attacks admissible in the standard threat model of logic-locking.

Let $\theta = (n, m, \Delta)$, where $n, m > 0$, and $\Delta \in [2^n - 1]$.

Description of *Hide*. Given F , the *Hide* algorithm in *OneChaff_{hd}* selects F as the anchor function and samples a single chaff uniformly from its Δ -neighborhood, i.e., $H \leftarrow_s N_\Delta(F)$. Equivalently, the chaff H is initialized to F and then, on a random subset $\mathcal{I}^\neq(F, H) = \{X_1, X_2, \dots, X_\Delta\}$ of the domain of F , the value of $H(X_i) \leftarrow_s \{0, 1\}^m \setminus \{F(X_i)\}$; when $m = 1$, $H(X_i) \leftarrow \neg F(X_i)$.

Formally, the opaque circuit C_L returned by *Hide* computes the function $L: (\{0, 1\}^n \cup \{\zeta\})^\Delta \times (\{0, 1\}^m \cup \{\zeta\})^\Delta \times \{0, 1\}^n \rightarrow \{0, 1\}^m$ defined by

$$L(((U_1, V_1), (U_2, V_2), \dots, (U_\Delta, V_\Delta)), X) = \begin{cases} V_i & \text{if } \exists U_i = X \\ H(X) & \text{otherwise} \end{cases}$$

where the distinguished (non-string) symbol ζ is understood to mean “uninitialized”. We implicitly assume that all string-valued U_i are distinct so that L is well defined, and that if any $U_i = \zeta$, then all $U_i = \zeta$. Thus, L computes H and a family of functions $\mathcal{R}_L = N_\Delta(H)$, one of which is F .

In practice, we envision C_L to be a circuit that encodes H and an uninitialized table $\widetilde{\text{Tab}}$, whose gate-and-wire representation allows for a sub-circuit C_H (computing H) to be easily parsed out; overloading the $\langle \cdot \rangle$ notation, we will sometimes write $C_L = \langle H, \widetilde{\text{Tab}} \rangle$ to reflect this. This way, one can evaluate C_H given C_L and, by loading $\widetilde{\text{Tab}}$ in C_L with a list of Δ pairs in $\{0, 1\}^n \times \{0, 1\}^m$, any of the functions in $N_\Delta(H)$.

Description of *Restore*. The *Restore* algorithm parses the input (design) parameters θ to determine n, m, Δ . It then expects a key $K = \langle (U_1, V_1), \dots, (U_\Delta, V_\Delta) \rangle$ where each $(U_i, V_i) \in \{0, 1\}^n \times \{0, 1\}^m$, and a chip \mathbb{C}_L that was fabricated from a circuit $C_L = \langle H, \widetilde{\text{Tab}} \rangle$. When these expectations are met, *Restore* loads the pairs (U_i, V_i) into the uninitialized lookup table. (When not, it returns \perp .) We write $\mathbb{C}_L \uparrow K$ for this, the notation suggesting that K is “uploaded” into \mathbb{C}_L . Thus, when provided with an honestly fabricated chip from the opaque circuit C_L returned by $\text{Hide}(F, \theta)$, and the corresponding key K_O , the restored chip $\mathbb{C}_L \uparrow K_O$ properly computes the hidden function F . Note that, we abstract away the physical mechanism by which uploading the (U_i, V_i) is implemented, and assume that if this mechanism is able to detect upload failure, *Restore* returns \perp .

B. *OneChaff_{hd}* is efficient

***Hide* is efficient.** Note that *Hide* in *OneChaff_{hd}* is abstracting the process by which an IP author goes from the “idea” (i.e., informal specification) of a functionality F to an opaque circuit described in some format that is suitable for fabrication (e.g., a GDSII file). This process typically involves writing a program for the desired functionality in some high-level

language like Verilog, and then converting the program to a circuit description by an EDA tool. In *OneChaff_{hd}*, this process can be augmented to sample Δ random domain points, and incorporate these to yield a program (hence a circuit) for the chaff function H . More specifically, the IP author can follow a two-step approach in order to create H from F . Let $m = 1$. In the first step, it creates a set of Δ random n -bit inputs and finds the value of F on those inputs. Next, it sets the value of H on the Δ inputs X_i to $\neg F(X_i)$; on the remaining inputs H and F have identical values. Hence, runtime of *Hide* to create H will be $O(n\Delta) + T(H)$, where $T(H)$ is the time needed to create the circuit for H . Note that the IP author would have anyways spent $T(F)$ time to build the circuit for F in the honest setting. So, effectively, the additional runtime of *Hide* in *OneChaff_{hd}* is $O(n\Delta) + T(H) - T(F)$. (Time to create K_O will also be $O(n\Delta)$ as $|K_O| = n\Delta$.) When $\Delta \ll 2^n$, $T(H)$ will be roughly same as $T(F)$. Thus, the efficiency of *Hide* should be roughly same as that of the original process, at least when $\Delta \ll 2^n$.

Also, note that such an implementation does not leak any additional information about F apart from H when the adversary gets full view of the opaque circuit $\langle H, \widetilde{\text{Tab}} \rangle$.

***Restore* is efficient.** Since *Restore* involves uploading the $n\Delta$ -bits hiding key K_O to $\widetilde{\text{Tab}}$, it’s runtime will be $O(n\Delta)$.

C. FR security of *OneChaff_{hd}*

In $\widetilde{\text{FR}}$, given unrestricted access to $C_L = \langle H, \widetilde{\text{Tab}} \rangle$, the IP author knows that the foundry can (at best) extract the set $\mathcal{R}_L = N_\Delta(H)$ and compute $\mathcal{M}_0 = \mathcal{F} \cap N_\Delta(H)$ to reduce the size of its original guess space \mathcal{F} . Using queries to TRUE, the foundry can further reduce the guess space to $\mathcal{M}_1 \subseteq \mathcal{M}_0$ by removing functions in \mathcal{M}_0 whose distinguishing-input sets contain query inputs.

In practice, the IP author will not know the set \mathcal{F} because that captures the foundry’s initial guess at the set of possible functions F that may be what the IP author intends to hide. (Recall that the FR-notion samples $F \leftarrow_s \mathcal{F}$, suggesting that the foundry’s initial set is the correct set.) More plausibly, the author may assume that the foundry’s initial guess is based upon knowledge of likely “properties” of F . These properties may be gleaned from discussions with the IP author, statements in fabrication contracts, historical and market information, etc. The IP author may use these assumed-known properties in its description of *Hide*, *Restore*, or into the parameters θ . In particular, it can use these assume-known properties to compute its estimates of the sets \mathcal{M}_0 and \mathcal{M}_1 that can help the foundry reduce the guess space. We will use these observations in our upcoming analysis.

Scoping the set \mathcal{F} . In what follows, we will focus on the case that $m = 1$, i.e., $\mathcal{F} \subset \text{Func}(n, 1)$. We focus on Boolean functions for a few reasons. First, even if the circuit representation of the hidden function has $n' \geq n$ bits of input and $m > 1$ bits of output, we can take as a first consideration whether or not one can securely hide the transitive fan-in cone (TFC) of any particular output bit. For a collection of output

bits that have disjoint transitive fan-in cones, one can consider hiding these in parallel.

Second, focusing on Boolean functions makes the analysis less complicated. In particular, for all distinguishing inputs $U_i \in \mathcal{I}^\neq(F, H)$, $V_i = \neg F(U_i)$; otherwise $V_i \leftarrow_s U_m \setminus F(U_i)$. Note that when $m > 1$, the TFCs of different output bits may not be disjoint. In such a case, the adversary can potentially use information it learns about one TFC to learn about the functionality of a different TFC.

Third, barring SFLL-flex [15], in the remaining SAT-attack-resistant DH schemes (AntiSAT [27], SARLock [28], TTLock [29] and SFLL-HD [15]), F is Boolean.

Simple functions cannot be hidden. We borrow the definition of “simple” functions from learning theory; simple functions are those that can be efficiently learned via a reasonably small number of input-output values of F . Now, if \mathcal{F} consists predominantly of simple functions, then *no* DH scheme will be secure⁴ in hiding functions sampled from \mathcal{F} . For example, functions that have a small domain can be learned by brute force. In addition, results from computational learning theory tell us that functions whose decision-tree representation have small depth/size can be learned via the Kushilevitz-Mansour algorithm [42]. More generally, methods exist to (approximately) learn Boolean functions whose Fourier spectra are sparse (e.g., dominated by relatively few Fourier coefficients) [43].

To loosely capture a measure of the density of “simple” Boolean functions within a given set \mathcal{F} , we give the following definition.

Definition 3. [*“Simple” functions in \mathcal{F} .*] Let $\mathcal{F} \subseteq \text{Func}(n, 1)$ be a set of Boolean functions. Let $t, q \geq 0$ be integers, and let $\delta \in [0, 1]$ be a real number. Let $\mathcal{F}_{t,q,\delta} \subseteq \mathcal{F}$ be a subset such that the following holds: \exists a q -query, t -time adversary A such that, $\forall f \in \mathcal{F}_{t,q,\delta}$, when $g \leftarrow_s A^{f(\cdot)}$ we have $\Pr[x \leftarrow_s \{0, 1\}^n : g(x) = f(x)] \geq \delta$. Furthermore, define $\epsilon_{t,q,\delta} = |\mathcal{F}_{t,q,\delta}|/|\mathcal{F}|$. In particular, $\epsilon_{t,q,1}$ is the fraction of \mathcal{F} that can be exactly learned (by some A) with q input-output values and time-complexity t . \diamond

Our security bounds for $\text{OneChaff}_{\text{hd}}$ will reflect the term $\epsilon_{t,q,1}$, although we stress that corresponding bounds for any DH scheme would also have to reflect this term (perhaps not explicitly) because any “simple” function will not be hideable.

In general, specifying a set \mathcal{F} for which $\epsilon_{t,q,1}$ is small enough for practically meaningful security statements (for reasonable t, q) is challenging, as this would require results of the following kind: There exist no adversary that can learn any function in $\mathcal{F} \setminus \mathcal{F}_{t,q,1}$ with q queries in time t . We are not aware of any such results. Also, note that it is not sufficient for \mathcal{F} to be large (although it is necessary to avoid simply guessing F), as one can specify large sets of functions with sparse Fourier spectra.

⁴Shamsi et al. [5] made a similar observation and gave impossibility results on logic locking when \mathcal{F} consists of entirely simple functions. Our security experiments are more generic as they allow \mathcal{F} to contain functions of varying degrees of “simplicity”.

We conjecture that if \mathcal{F} is sufficiently “unstructured”, then $\epsilon_{t,q,1}$ will be small enough to not dominate the FR-security bounds we will prove, for practically reasonable t, q . For example, a random Boolean function lacks the highly concentrated spectral structure — the number of non-zero Fourier coefficients is $(1 - o(1))2^n$ [51] — that leads to efficient learnability from input-output pairs.

Under the above conjecture, we will focus on $\mathcal{F} = \mathcal{F}_{[h_0, h_1]} \subseteq \text{Func}(n, 1)$ that consists of all Boolean functions with hamming weight at least h_0 and at most h_1 . Recall that the set \mathcal{F} in our FR experiment is meant to reflect the adversary’s a priori “knowledge” about the function hidden in the opaque circuit, and that our FR-notion samples F uniformly from \mathcal{F} . At the extremes, setting $h_0 = 0, h_1 = 2^n$ considers the case that the adversary has no a priori knowledge; setting $h_0 = h_1$ considers an adversary that knows exactly the hamming weight of the hidden function F . In our analysis of $\text{OneChaff}_{\text{hd}}$, we will conservatively assume the latter.

A random Boolean function has an expected hamming weight close to 2^{n-1} . Hence, when \mathcal{F} comprises functions with hamming weights close to 2^{n-1} , then a randomly sampled function from \mathcal{F} is unlikely to be easily learned (i.e., “simple”) given the current state of the art, and will have very small $\epsilon_{t,q,1}$. We conjecture this based on the fact that a random Boolean function almost certainly does not have a concentrated Fourier spectrum [51] thereby making efficient learning-theoretic attacks out-of-scope.

A practical warm-up: FR security with $q = 0$. Chips that are used in critical infrastructure (e.g., military devices) will most likely require considerable effort to obtain. In this case, we can assume that the foundry cannot get access to honestly restored chips, and thereby learn the true input-output behavior of F on inputs of its choice. Thus, we begin our analysis of $\text{OneChaff}_{\text{hd}}$ for $\mathcal{F}_{[h,h]}$ in the $q = 0$ case.

Given Lemma 1, recall that we were left to find an upper bound on $\Pr[\widetilde{\text{FR}}(B) = 1]$, where B has some time-complexity t and asks q queries to its TRUE oracle. Let us fix $q = 0$, and consider an arbitrary adversary B_0 with these resource bounds.

Recalling the notation that \mathcal{R}_L is the set of functions realizable by the opaque circuit $L = \langle H, \widehat{\text{Tab}} \rangle$, we have $\mathcal{R}_L = N_\Delta(H)$; $F \in \mathcal{R}_L$. The adversary B_0 knows a priori that $F \in \mathcal{F}_{[h,h]}$. Hence, for all $\hat{F} \in \mathcal{R}_L$, $\text{hw}(\hat{F}) \in [\text{hw}(H) - \Delta, \text{hw}(H) + \Delta]$. To get an intuition of the claim on $\text{hw}(\hat{F})$, consider two cases: $\text{hw}(H) = 0$ and $\text{hw}(H) = 2^n$. When $\text{hw}(H) = 0$, $\text{hw}(\hat{F}) = \text{hw}(H) + \Delta = \Delta$; when $\text{hw}(H) = 2^n$, $\text{hw}(\hat{F}) = \text{hw}(H) - \Delta = 2^n - \Delta$.

While attacking $\text{OneChaff}_{\text{hd}}$, the optimal strategy for B_0 is to return the circuit implementation of the most-likely function in $\mathcal{F}_{[h,h]} \cap \mathcal{R}_L$. (We will use $\mathcal{F}_{[h,h]}^j, \mathcal{R}_L^j$ to denote the subset of functions in $\mathcal{F}_{[h,h]}$ and \mathcal{R}_L , respectively, that have correct values on some $0 \leq j \leq q$ inputs.) But as sampling is done in a uniform and independent fashion in $\widetilde{\text{FR}}$, all functions in $\mathcal{M}_0 = (\mathcal{F}_{[h,h]}^0 \cap \mathcal{R}_L^0)$ are equally likely.

Then, $\Pr \left[\widehat{\mathbf{FR}}(B_0) = 1 \right] = |\mathcal{M}_0|^{-1}$. We will proceed to give a lowerbound on $|\mathcal{M}_0|$ in the $q = 0$ case.

Claim 1. *Let $\text{hw}(H) = h + \alpha$, where $\alpha < (2^{n-1} - h)$, and $\Delta < h$. Then,*

$$|\mathcal{M}_0| = \left| \left(\mathcal{F}_{[h,h]}^0 \cap \mathcal{R}_L^0 \right) \right| = \binom{h + \alpha}{(\Delta + \alpha)/2} \binom{2^n - h - \alpha}{(\Delta - \alpha)/2}.$$

◆

When $\Delta < h$ and $\text{hw}(H) < 2^{n-1}$, we found (in the proof of Claim 1) that *more than half* $(\Delta/2 + \alpha/2)$ of the distinguishing inputs will belong to $\mathcal{X}_1(H)$. Hence, we can assume, without any loss in the FR advantage of the adversary, that the adversary will make all of its queries from $\mathcal{X}_1(H)$ until it finds all $(\Delta + \alpha)/2$ distinguishing inputs in $\mathcal{X}_1(H)$. (Later, we will show that when $q < \min(\text{hw}(H)/4, \Delta^2/64 \ln n)$, the adversary can find *at most* half of the distinguishing inputs in $\mathcal{X}_1(H)$, without loss.)

The lower bound on $|\mathcal{M}_0|$ in Claim 1, a function of $\text{hw}(H)$ as $\alpha = (h - \text{hw}(H))$, gives us an upper bound on the FR advantage of adversary that attacks $\text{OneChaff}_{\text{hd}}$ without making any queries to any of the three oracles. But, $\text{hw}(H)$ is not a parameter of the FR-security experiment. In the next claim, we give a lower bound on the hamming weight of H using the Hoeffding lemma; the bound is a function of h , Δ and n and all three are parameters of the FR experiment.

Claim 2. *Let $\Delta < h$ and $0 < (\Delta + h) < 2^{n-1}$. Then, $\text{hw}(H) \geq h + \Delta(1 - h/2^{n-1})$, with probability at least $(1 - 1/\Delta^2)$.*

◆

We will assume that the bound on $\text{hw}(H)$ is tight and we will reflect the uncertainty in the claim/theorem statements.

Now, we are ready to give the upper bound on the FR advantage of adversary A in the zero-query setting. We use the standard relation $\binom{p}{r} \geq (p/r)^r$ in the bound in Claim 1 to make the final bound in Theorem 2 easier to interpret. We also use the weaker bound $\text{hw}(H) > h$ (instead of the tighter bound in Claim 2) to make the final bound even more interpretable. Note that practical $\text{OneChaff}_{\text{hd}}$ schemes will have $\Delta \ll 2^{n-1}$ as the run time of *Hide* as well as *Restore* will be linear in Δ . Thus, the loss in using $\text{hw}(H) > h$ instead of $\text{hw}(H) \geq h + \Delta(1 - h/2^{(n-1)})$ will not be significantly large.

Theorem 2. *Fix $m = 1$ and integers $n, h, \Delta > 0$ such that $\Delta < h$, and $\Delta + h < 2^{(n-1)}$. Let $\epsilon_{t,0,1}$ be as defined in Definition 3. When an honest-but-curious adversary A attacks the FR security of $\Pi = \text{OneChaff}_{\text{hd}}$ without making any queries to any oracle, then, it achieves*

$$\text{Adv}_{(\Pi, \theta), \text{Fab}}^{\text{FR}}(\mathcal{F}_{[h,h]}, A) \leq \left(\frac{\Delta^2}{2^n h} \right)^{\Delta/2} + \epsilon_{t,0,1},$$

with probability at least $(1 - 1/\Delta^2)$.

◆

In the next section, we will see that when we account for the adversary's access to the oracles, more specifically, the TRUE oracle, the upper bound on the FR advantage

will be $\left(\frac{\Delta^2}{2^n \cdot (h-q)} \right)^{\Delta/4} + \epsilon_{t,q,1}$. Observe that the queries to TRUE increase the FR advantage of an adversary significantly (compared to the zero-query setting) as the exponent $\Delta/2$ decreases to $\Delta/4$; the denominator also is smaller compared to the zero-query setting. Since the first term is very similar to the bound in Theorem 3, we defer the unpacking of Theorem 2 to the next section.

FR analysis of $\text{OneChaff}_{\text{hd}}$ with $q > 0$. The ability to learn true input-output pairs, via queries to TRUE, provide a way for the adversary to verify guesses at portions of the key K_O . Recall that the key encodes $(X_1, F(X_1)), \dots, (X_\Delta, F(X_\Delta))$ for $X_i \in \mathcal{I}^\neq(F, H)$, and the opaque circuit C_L allows for local computation of $H(X)$. Thus, as a first step in analyzing the FR security of $\text{OneChaff}_{\text{hd}}$ in the $q > 0$ case, we derive a bound on the number of points in $\mathcal{I}^\neq(F, H)$ that an adversary uncovers in its q queries to the TRUE oracle.

Let $\mathcal{Q}_j = \{x_1, x_2, \dots, x_j\}$ be the the first j queries to TRUE, and let random variable $\mathcal{Q}_j^{\text{key}} = \mathcal{Q}_j \cap \mathcal{I}^\neq(F, H)$ denote the queries in \mathcal{Q}_j that uncover a portion of the distinguishing inputs. Observe that $|\mathcal{Q}_0^{\text{key}}| = 0$, and for $j > 0$ the value of $|\mathcal{Q}_j^{\text{key}}|$ depends only upon $|\mathcal{Q}_{j-1}^{\text{key}}|$ and the query x_j ; in particular that $|\mathcal{Q}_j^{\text{key}}| = |\mathcal{Q}_{j-1}^{\text{key}}| + 1$ if $x_j \in \mathcal{I}^\neq(F, H)$, and $|\mathcal{Q}_j^{\text{key}}| = |\mathcal{Q}_{j-1}^{\text{key}}|$ if not.

Let I_j be the indicator random variable indicating that the event $x_j \in \mathcal{Q}_j^{\text{key}}$ occurs. We claim that $\Pr[I_j = 1] = (\Delta - |\mathcal{Q}_{j-1}^{\text{key}}|) / (\text{hw}(H) - (j - 1))$. To see this, observe that the number of uncovered points in $\mathcal{I}^\neq(F, H)$ is precisely $(\Delta - |\mathcal{Q}_{j-1}^{\text{key}}|)$ and, given how those points were sampled, any of the remaining, unqueried points in $\mathcal{X}_1(H)$ are equally likely to be in $\mathcal{I}^\neq(F, H)$. Given this, we can prove the following lemma.

Lemma 2. *Let $\Delta < h$ and $(\Delta + h) < 2^{n-1}$. Then, we have $\mathbb{E}[|\mathcal{Q}_q^{\text{key}}|] = \frac{q\Delta}{\text{hw}(H)}$, and, $|\mathcal{Q}_q^{\text{key}}| < \mathbb{E}[|\mathcal{Q}_q^{\text{key}}|] + \sqrt{4q \ln \Delta}$ with probability at least $1 - (2/\Delta^2)$.*

◆

Notice that when $q < \Delta^2/(256 \ln \Delta)$, $\sqrt{4q \ln \Delta} < \Delta/8$; when $q < \text{hw}(H)/8$, $q\Delta/\text{hw}(H) < \Delta/8$. Thus, with probability $1 - (2/\Delta^2)$, the adversary will uncover *no more than* quarter of the distinguishing inputs in $\mathcal{I}^\neq(F, H)$ when $q < \min(\text{hw}(H)/8, \Delta^2/(256 \ln \Delta))$.

The adversary can increase its FR advantage by removing functions from $(\mathcal{F}_{[h,h]}^0 \cap \mathcal{R}_L^0)$ that are not correct on any input in $\mathcal{Q}_q^{\text{key}}$. In fact, it can remove all functions that do not agree with F on any point in the query set \mathcal{Q}_q .

Extending Claim 1 to account for $\mathcal{Q}_q^{\text{key}}$, we give a lower bound on $|\mathcal{M}_1| = \left| \left(\mathcal{F}_{[h,h]}^q \cap \mathcal{R}_L^q \right) \right|$ in Claim 3. Since one wins against $\text{OneChaff}_{\text{hd}}$ if and only if the hidden function F is completely recovered, a lowerbound on $|\mathcal{M}_1|$ will give us an upperbound on $\Pr \left[\widehat{\mathbf{FR}}(B) = 1 \right]$, and in turn (by Lemma 1) an upper bound on the FR advantage of an adversary A that attacks $\text{OneChaff}_{\text{hd}}$.

Claim 3. *Let $\text{hw}(H) = h + \delta$, where $\delta = \Delta(1 - h/2^{n-1})$,*

and $\Delta < h$. Then,

$$|\mathcal{M}_1| \geq \binom{(h + \delta - q)}{(\Delta + \delta)/2 - |\mathcal{Q}_q^{\text{key}}|} \binom{(2^n - h - \delta)}{(\Delta - \delta)/2},$$

with probability at least $(1 - 2/\Delta^2)$. \blacklozenge

Comparing the bounds in Claims 1 and 3, we can see that the adversary B can remove a large number — at least $\binom{2^n - h - \delta}{(\Delta - \delta)/2}$ — of functions from $(\mathcal{F}_{[h,h]}^0 \cap \mathcal{R}_L^0)$ using its queries to TRUE.

Upperbound on the FR security of $\text{OneChaff}_{\text{hd}}$. Recall that the lowerbound on $|\mathcal{M}_1|$ gives us an upperbound on the FR advantage of an adversary that attacks $\text{OneChaff}_{\text{hd}}$, where \mathcal{R}_L^q denotes the set of functions in the Δ neighborhood of H that are correct on all the queries that the adversary makes to TRUE. The bound in Claim 3 is a little complex and hard to interpret. In the following Corollary to Claim 3, we derive a weaker, but easier to interpret bound on $|\mathcal{M}_1|$. In the derivation, we use the standard relation $\binom{p}{r} \geq (p/r)^r$, and the assumptions: $q < \min(h/8, \Delta^2/256n)$ and $(h < h + \Delta < 2^n)$ to get the bound.

Corollary 1. *Let $0 \leq q < \min(h/8, \Delta^2/256n)$. Let $\Delta < h$ and $0 < (\Delta + h) < 2^{n-1}$. Then,*

$$|\mathcal{M}_1| = \left| \left(\mathcal{F}_{[h,h]}^q \cap \mathcal{R}_L^q \right) \right| \geq \left(\frac{2^n(h-q)}{\Delta^2} \right)^{\Delta/4},$$

with probability at least $(1 - 2/\Delta^2)$. \blacklozenge

We are now prepared to state the upper bound on the FR-advantage of an adversary attacking $\text{OneChaff}_{\text{hd}}$. Since $\text{Adv}_{(\Pi, \theta), \text{Fab}}^{\text{FR}}(\mathcal{F}_{[h,h]}, A) \leq 1/|\mathcal{M}_1| + \epsilon_{t,q,1}$, the final bound follows directly from Corollary 1. Recalling that $\epsilon_{t,q,1}$ captures the probability that the sampled hidden function F is a “simple” function (see Definition 3), we have the following theorem.

Theorem 3. *Fix $m = 1$ and integers $n, h, \Delta, q > 0$ such that $\Delta < h$, $0 < \Delta + h < 2^{n-1}$, and $0 \leq q < \min(h/8, \Delta^2/256n)$. Let an honest-but-curious adversary A attack the FR security of $\text{OneChaff}_{\text{hd}}$ using resources $(t, q_f = 1, q_s = 1, q_r = q)$. Let $\epsilon_{t,q,1}$ be as defined in Definition 3. Then,*

$$\text{Adv}_{(\text{OneChaff}_{\text{hd}}, \theta), \text{Fab}}^{\text{FR}}(\mathcal{F}_{[h,h]}, A) \leq \left(\frac{\Delta^2}{2^n(h-q)} \right)^{\Delta/4} + \epsilon_{t,q,1},$$

with probability at least $(1 - 2/\Delta^2)$. \blacklozenge

See Theorem 4 (Section D) for a much tighter, but harder to interpret upperbound on $\text{Adv}_{(\Pi, \theta), \text{Fab}}^{\text{FR}}(\mathcal{F}_{[h,h]}, A)$. Let us unpack Theorem 3 a bit. We will carve out three criteria from the theorem in terms of: $|\text{Dom}(F)|$, $\text{hw}(F)$ and Δ , that are needed for FR security of $\text{OneChaff}_{\text{hd}}$. Note that an IP author can configure Δ while running *Hide*.

The FR security of $\text{OneChaff}_{\text{hd}}$ increases polynomially in the size of the domain of the hidden function, as well as the hamming weight of the function (as the hamming weight approaches 2^{n-1}). The rate of increase depends on the value

of Δ — large values of Δ result in a faster rate of increase. All three criteria are quite intuitive. Functions with small domains will be simple because they can be learnt by brute force; functions with very small hamming weights will be “close” to being constant functions. On the contrary, when F is a uniformly sampled balanced Boolean function with large domain, the min-entropy of the (uniform) distribution on \mathcal{F} will be very large (compared to functions with very small hamming weights) as $|\mathcal{F}| = \binom{2^n}{\ell} > \binom{2^n}{\ell}$ for any $\ell \in [0, 2^n] \setminus \{2^{n-1}\}$.

While efficient $\text{OneChaff}_{\text{hd}}$ schemes will require Δ to be small (as the size of K_O is $n\Delta$), the IP author should choose as large Δ as possible to maximize FR-security guarantee. This criterion is also quite intuitive as large Δ implies large number of secret distinguishing inputs (that are encoded in K_O) on which the adversary does not know the value of F .

Even when Δ is not so large, say $\Delta = 2^{10}$, security holds with probability at least $(1 - 1/2^{20})$ which is almost equal to one. And when the domain of the hidden function as well as the hamming weight is large, i.e., $h, 2^n \gg 0$, then the first term $(\Delta^2/(2^n(h-q)))^{(\Delta/4)}$ in Theorem 3, will be extremely small. (Please note that we give concrete security bounds, not asymptotic ones, precisely because it allows this kind of analysis — one can evaluate what security is possible (in our model) for any setting of the parameters.)

None of the previous works consider these factors while discussing the security of DH schemes. For example, in the evaluation of the SAT attack [17], four base-benchmark circuits (c17, ex5, apex4, ex1010) have very small domains: in c17, $n = 5$; in ex5, $n = 8$; in apex4 and ex1010, $n = 10$. Note that the authors used each base-benchmark circuit to create 21 opaque circuits. That means roughly 19% of the benchmark circuits⁵ in the test corpora of SAT attack cannot be protected by any DH scheme. It is noteworthy to also point that the open-source test corpora of SAT attack is used by other attack algorithms as well [22], [23].

VII. CONCLUSIONS

In this paper, we gave a provable security treatment to the problem of hiding circuit-design IP from adversarial foundries. We formalized a new primitive, called a design-hiding scheme, that abstractly captures all prior work on so-called logic-locking schemes (which forms the bulk of prior work in the area), and our security notions capture the intended goals of logic-locking schemes. With this formal foundation in place, we described a DH scheme called $\text{OneChaff}_{\text{hd}}$ and proved that it meets our function-recovery security notion in the honest-but-curious foundry setting. (It is unclear, given the state of the art in hardware trojan detection, how to achieve security in the fully malicious foundry setting.) Our work surfaced some important issues that had gone unmentioned in prior work, for example the matter of “simple” circuits. We see this work as an important first step towards the ultimate goal of the

⁵Out of 441 benchmark circuits, 84 circuits have very small domains in the test corpora.

community, namely to provide efficient and provably secure DH schemes to protect high-value circuit IP from being reverse engineered by untrusted foundries.

ACKNOWLEDGEMENT

We acknowledge the hard work of the reviewers, and their helpful feedback. We also acknowledge the National Science Foundation for their support under award NSF-1564444

REFERENCES

- [1] The-Intercept, "Everybody does it: the messy truth of computer supply chains." <https://theintercept.com/2019/01/24/computer-supply-chain-attacks/>, Tech. Rep., 2019.
- [2] TrendForce, "Trendforce reports top 10 ranking of global semiconductor foundries of 2018," <https://press.trendforce.com/press/20180524-3106.html>, Tech. Rep., 2018.
- [3] "Top 5 most counterfeited parts represent a \$169 billion potential challenge for global semiconductor market," <https://technology.informa.com/405654/>, 2012, accessed: 2020-08-25.
- [4] U. Guin, K. Huang, D. DiMase, J. M. Carulli, M. Tehranipoor, and Y. Makris, "Counterfeit integrated circuits: A rising threat in the global semiconductor supply chain," *Proceedings of the IEEE*, vol. 102, no. 8, pp. 1207–1228, 2014.
- [5] K. Shamsi, D. Z. Pan, and Y. Jin, "On the impossibility of approximation-resilient circuit locking," in *2019 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE, 2019, pp. 161–170.
- [6] J. A. Roy, F. Koushanfar, and I. L. Markov, "Epic: Ending piracy of integrated circuits," in *Proceedings of the conference on Design, automation and test in Europe*. ACM, 2008, pp. 1069–1074.
- [7] K. Zamiri Azar, H. Mardani Kamali, H. Homayoun, and A. Sasan, "Threats on logic locking: A decade later," in *Proceedings of the 2019 on Great Lakes Symposium on VLSI*, ser. GLSVLSI '19. New York, NY, USA: ACM, 2019, pp. 471–476. [Online]. Available: <http://doi.acm.org/10.1145/3299874.3319495>
- [8] K. Shamsi, M. Li, K. Plaks, S. Fazzari, D. Z. Pan, and Y. Jin, "Ip protection and supply chain security through logic obfuscation: A systematic overview," *ACM Transactions on Design Automation of Electronic Systems (TODAES)*, vol. 24, no. 6, pp. 1–36, 2019.
- [9] J. Rajendran, M. Sam, O. Sinanoglu, and R. Karri, "Security analysis of integrated circuit camouflaging," in *Proceedings of the 2013 ACM SIGSAC Conference on Computer & #38; Communications Security*, ser. CCS '13. New York, NY, USA: ACM, 2013, pp. 709–720. [Online]. Available: <http://doi.acm.org/10.1145/2508859.2516656>
- [10] A. Vijayakumar, V. C. Patil, D. E. Holcomb, C. Paar, and S. Kundu, "Physical design obfuscation of hardware: A comprehensive investigation of device and logic-level techniques," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 1, pp. 64–77, 2016.
- [11] J. Rajendran, O. Sinanoglu, and R. Karri, "Is split manufacturing secure?" in *2013 Design, Automation Test in Europe Conference Exhibition (DATE)*, March 2013, pp. 1259–1264.
- [12] B. Shakya, N. Asadizanjani, D. Forte, and M. Tehranipoor, "Chip editor: leveraging circuit edit for logic obfuscation and trusted fabrication," in *Proceedings of the 35th International Conference on Computer-Aided Design*, 2016, pp. 1–8.
- [13] S. Dziembowski, S. Faust, and F.-X. Standaert, "Private circuits iii: Hardware trojan-resilience via testing amplification," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '16. New York, NY, USA: ACM, 2016, pp. 142–153. [Online]. Available: <http://doi.acm.org/10.1145/2976749.2978419>
- [14] L. W. Chow, J. P. Baukus, B. J. Wang, and R. P. Cocchi, "Camouflaging a standard cell based integrated circuit," Apr. 3 2012, uS Patent 8,151,235.
- [15] M. Yasin, A. Sengupta, M. T. Nabeel, M. Ashraf, J. J. Rajendran, and O. Sinanoglu, "Provably-secure logic locking: From theory to practice," in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '17. New York, NY, USA: ACM, 2017, pp. 1601–1618. [Online]. Available: <http://doi.acm.org/10.1145/3133956.3133985>
- [16] D. Sirone and P. Subramanyan, "Functional analysis attacks on logic locking," in *2019 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2019, pp. 936–939.
- [17] P. Subramanyan, S. Ray, and S. Malik, "Evaluating the security of logic encryption algorithms," in *2015 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*, May 2015, pp. 137–143.
- [18] X. Xu, B. Shakya, M. M. Tehranipoor, and D. Forte, "Novel bypass attack and bdd-based tradeoff analysis against all known logic locking attacks," in *International Conference on Cryptographic Hardware and Embedded Systems*. Springer, 2017, pp. 189–210.
- [19] Y. Shen and H. Zhou, "Double dip: Re-evaluating security of logic encryption algorithms," in *Proceedings of the on Great Lakes Symposium on VLSI 2017*, ser. GLSVLSI '17. New York, NY, USA: ACM, 2017, pp. 179–184. [Online]. Available: <http://doi.acm.org/10.1145/3060403.3060469>
- [20] Y. Shen, A. Rezaei, and H. Zhou, "Sat-based bit-flipping attack on logic encryptions," in *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2018, pp. 629–632.
- [21] M. Yasin, B. Mazumdar, O. Sinanoglu, and J. Rajendran, "Removal attacks on logic locking and camouflaging techniques," *IEEE Transactions on Emerging Topics in Computing*, 2017.
- [22] K. Shamsi, M. Li, T. Meade, Z. Zhao, D. Z. Pan, and Y. Jin, "Appsat: Approximately deobfuscating integrated circuits," in *2017 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*, May 2017, pp. 95–100.
- [23] Y. Shen, Y. Li, A. Rezaei, S. Kong, D. Dlott, and H. Zhou, "Besat: behavioral sat-based attack on cyclic logic encryption," in *Proceedings of the 24th Asia and South Pacific Design Automation Conference*. ACM, 2019, pp. 657–662.
- [24] H. Zhou, R. Jiang, and S. Kong, "Cycsat: Sat-based attack on cyclic logic encryptions," in *Proceedings of the 36th International Conference on Computer-Aided Design*. IEEE Press, 2017, pp. 49–56.
- [25] Y.-C. Chen, "Enhancements to sat attack: Speedup and breaking cyclic logic encryption," *ACM Trans. Des. Autom. Electron. Syst.*, vol. 23, no. 4, pp. 52:1–52:25, May 2018. [Online]. Available: <http://doi.acm.org/10.1145/3190853>
- [26] F. Yang, M. Tang, and O. Sinanoglu, "Stripped functionality logic locking with hamming distance based restore unit (sfl-hd)-unlocked," *IEEE Transactions on Information Forensics and Security*, 2019.
- [27] Y. Xie and A. Srivastava, "Mitigating sat attack on logic locking," in *International Conference on Cryptographic Hardware and Embedded Systems*. Springer, 2016, pp. 127–146.
- [28] M. Yasin, B. Mazumdar, J. J. Rajendran, and O. Sinanoglu, "Sarlock: Sat attack resistant logic locking," in *2016 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE, 2016, pp. 236–241.
- [29] M. Yasin, A. Sengupta, B. C. Schafer, Y. Makris, O. Sinanoglu, and J. J. Rajendran, "What to lock?: Functional and parametric locking," in *Proceedings of the on Great Lakes Symposium on VLSI 2017*, ser. GLSVLSI '17. New York, NY, USA: ACM, 2017, pp. 351–356. [Online]. Available: <http://doi.acm.org/10.1145/3060403.3060492>
- [30] K. Shamsi, M. Li, T. Meade, Z. Zhao, D. Z. Pan, and Y. Jin, "Cyclic obfuscation for creating sat-unresolvable circuits," in *Proceedings of the on Great Lakes Symposium on VLSI 2017*, ser. GLSVLSI '17. New York, NY, USA: ACM, 2017, pp. 173–178. [Online]. Available: <http://doi.acm.org/10.1145/3060403.3060458>
- [31] S. Roshanifefat, H. Mardani Kamali, and A. Sasan, "Srclock: Sat-resistant cyclic logic locking for protecting the hardware," in *Proceedings of the 2018 on Great Lakes Symposium on VLSI*. ACM, 2018, pp. 153–158.
- [32] H. M. Kamali, K. Z. Azar, H. Homayoun, and A. Sasan, "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*, ser. DAC '19. New York, NY, USA: ACM, 2019, pp. 89:1–89:6.
- [33] R. S. Chakraborty and S. Bhunia, "Harpoon: an obfuscation-based soc design methodology for hardware protection," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 28, no. 10, pp. 1493–1502, 2009.
- [34] Y. Alkabani, F. Koushanfar, and M. Potkonjak, "Remote activation of ics for piracy prevention and digital right management," in *Proceedings of the 2007 IEEE/ACM international conference on Computer-aided design*. IEEE Press, 2007, pp. 674–677.
- [35] A. R. Desai, M. S. Hsiao, C. Wang, L. Nazhandali, and S. Hall, "Interlocking obfuscation for anti-tamper hardware," in *Proceedings of*

the eighth annual cyber security and information intelligence research workshop. ACM, 2013, p. 8.

- [36] J. Dofe and Q. Yu, "Novel dynamic state-deflection method for gate-level design obfuscation," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 2, pp. 273–285, 2017.
- [37] M. Fyrbiak, S. Wallat, J. Déchelotte, N. Albartus, S. Böcker, R. Tessier, and C. Paar, "On the difficulty of fsm-based hardware obfuscation," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 293–330, 2018.
- [38] M. E. Massad, J. Zhang, S. Garg, and M. V. Tripunitara, "Logic locking for secure outsourced chip fabrication: A new attack and provably secure defense mechanism," *CoRR*, vol. abs/1703.10187, 2017. [Online]. Available: <http://arxiv.org/abs/1703.10187>
- [39] M. Li, K. Shamsi, T. Meade, Z. Zhao, B. Yu, Y. Jin, and D. Z. Pan, "Provably secure camouflaging strategy for ic protection," in *Proceedings of the 35th International Conference on Computer-Aided Design*, ser. ICCAD '16. New York, NY, USA: ACM, 2016, pp. 28:1–28:8. [Online]. Available: <http://doi.acm.org/10.1145/2966986.2967065>
- [40] S. Engels, M. Hoffmann, and C. Paar, "The end of logic locking? a critical view on the security of logic locking," *Cryptology ePrint Archive*, Report 2019/796, 2019, <https://eprint.iacr.org/2019/796>.
- [41] J.-M. Cioranescu, J.-L. Danger, T. Graba, S. Guilley, Y. Mathieu, D. Naccache, and X. T. Ngo, "Cryptographically secure shields," in *2014 IEEE International Symposium on Hardware-Oriented Security and Trust (HOST)*. IEEE, 2014, pp. 25–31.
- [42] E. Kushilevitz and Y. Mansour, "Learning decision trees using the fourier spectrum," *SIAM Journal on Computing*, vol. 22, no. 6, pp. 1331–1348, 1993.
- [43] P. Gopalan, R. O'Donnell, R. A. Servedio, A. Shpilka, and K. Wimmer, "Testing fourier dimensionality and sparsity," *SIAM Journal on Computing*, vol. 40, no. 4, pp. 1075–1100, 2011.
- [44] J. Rajendran, Y. Pino, O. Sinanoglu, and R. Karri, "Security analysis of logic obfuscation," in *DAC Design Automation Conference 2012*, June 2012, pp. 83–89.
- [45] —, "Logic encryption: A fault analysis perspective," in *Proceedings of the Conference on Design, Automation and Test in Europe*, ser. DATE '12. San Jose, CA, USA: EDA Consortium, 2012, pp. 953–958. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2492708.2492947>
- [46] S. Dupuis, P.-S. Ba, G. Di Natale, M.-L. Flottes, and B. Rouzeyre, "A novel hardware logic encryption technique for thwarting illegal overproduction and hardware trojans," in *On-Line Testing Symposium (IOLTS), 2014 IEEE 20th International*. IEEE, 2014, pp. 49–54.
- [47] S. M. Plaza and I. L. Markov, "Protecting integrated circuits from piracy with test-aware logic locking," in *2014 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, Nov 2014, pp. 262–269.
- [48] R. S. Chakraborty and S. Bhunia, "Security against hardware trojan attacks using key-based design obfuscation," *Journal of Electronic Testing*, vol. 27, no. 6, pp. 767–785, 2011.
- [49] A. Baumgarten, A. Tyagi, and J. Zambreno, "Preventing ic piracy using reconfigurable logic barriers," *IEEE Design & Test of Computers*, vol. 27, no. 1, 2010.
- [50] M. Yasin, B. Mazumdar, O. Sinanoglu, and J. Rajendran, "Security analysis of anti-sat," in *2017 22nd Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE, 2017, pp. 342–347.
- [51] R. De Wolf, "A brief introduction to fourier analysis on the boolean cube," *Theory of Computing*, pp. 1–20, 2008.
- [52] Z. Brakerski and G. N. Rothblum, "Virtual black-box obfuscation for all circuits via generic graded encoding," in *Theory of Cryptography Conference*. Springer, 2014, pp. 1–25.
- [53] B. Barak, O. Goldreich, R. Impagliazzo, S. Rudich, A. Sahai, S. Vadhan, and K. Yang, "On the (im) possibility of obfuscating programs," in *Annual International Cryptology Conference*. Springer, 2001, pp. 1–18.
- [54] A. C.-C. Yao, "Protocols for secure computations," in *FOCS*, vol. 82, 1982, pp. 160–164.
- [55] D. Evans, V. Kolesnikov, M. Rosulek *et al.*, "A pragmatic introduction to secure multi-party computation," *Foundations and Trends® in Privacy and Security*, vol. 2, no. 2-3, pp. 70–246, 2018.
- [56] M. Bellare, V. T. Hoang, and P. Rogaway, "Foundations of garbled circuits," in *Proceedings of the 2012 ACM conference on Computer and communications security*. ACM, 2012, pp. 784–796.
- [57] Y. Lindell, "Secure multiparty computation for privacy preserving data mining," in *Encyclopedia of Data Warehousing and Mining*. IGI Global, 2005, pp. 1005–1009.
- [58] M. Bellare, V. T. Hoang, S. Keelveedhi, and P. Rogaway, "Efficient garbling from a fixed-key blockcipher," in *2013 IEEE Symposium on Security and Privacy*. IEEE, 2013, pp. 478–492.
- [59] V. Kolesnikov, P. Mohassel, and M. Rosulek, "Flexor: Flexible garbling for xor gates that beats free-xor," in *Annual Cryptology Conference*. Springer, 2014, pp. 440–457.
- [60] S. Zahur, M. Rosulek, and D. Evans, "Two halves make a whole," in *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 2015, pp. 220–250.
- [61] J. B. Nielsen, P. S. Nordholt, C. Orlandi, and S. S. Burra, "A new approach to practical active-secure two-party computation," *Cryptology ePrint Archive*, Report 2011/091, 2011, <https://eprint.iacr.org/2011/091>.
- [62] H. Carter and P. Traynor, "Opfe: Outsourcing computation for private function evaluation," *IACR Cryptology ePrint Archive*, vol. 2016, p. 67, 2016.
- [63] S. Kamara, P. Mohassel, and B. Riva, "Salus: a system for server-aided secure function evaluation," in *Proceedings of the 2012 ACM conference on Computer and communications security*, 2012, pp. 797–808.
- [64] T. P. Jakobsen, J. B. Nielsen, and C. Orlandi, "A framework for outsourcing of secure computation," in *Proceedings of the 6th edition of the ACM Workshop on Cloud Computing Security*, 2014, pp. 81–92.
- [65] D. Malkhi, N. Nisan, B. Pinkas, Y. Sella *et al.*, "Fairplay-secure two-party computation system," in *USENIX Security Symposium*, vol. 4. San Diego, CA, USA, 2004, p. 9.
- [66] B. Kreuter, A. Shelat, B. Mood, and K. Butler, "{PCF}: A portable circuit format for scalable two-party secure computation," in *Presented as part of the 22nd {USENIX} Security Symposium ({USENIX} Security 13)*, 2013, pp. 321–336.
- [67] L. Malka, "Vmccrypt: modular software architecture for scalable secure computation," in *Proceedings of the 18th ACM conference on Computer and communications security*, 2011, pp. 715–724.
- [68] E. M. Songhori, S. U. Hussain, A.-R. Sadeghi, T. Schneider, and F. Koushanfar, "Tinygarble: Highly compressed and scalable sequential garbled circuits," in *2015 IEEE Symposium on Security and Privacy*. IEEE, 2015, pp. 411–428.
- [69] E. Boyle, N. Gilboa, and Y. Ishai, "Function secret sharing," in *Annual international conference on the theory and applications of cryptographic techniques*. Springer, 2015, pp. 337–367.

APPENDIX

A. Historical context

In 2008, the first logic-locking scheme (random-logic locking (RLL) [6]) was designed. The ostensible security goal of RLL was to construct an opaque circuit that hides the full functionality F of the IP author's circuit from an adversarial foundry who has unrestricted access to the opaque circuit and oracle access to F . However, the authors of RLL claimed that their scheme is secure because it prevents an adversary from recovering the correct key using a brute-force key-search attack.

Note that the security goal of preventing key recovery (KR) misses the original intent; namely, to prevent the adversary from learning the *functionality* of the IP author's circuit. Thus, preventing *function recovery* (FR) is arguably a more pertinent security goal.

Between 2008 and 2015, several variants [6], [44]–[46], [49] of RLL were designed. Each claimed security by (empirically) showing that their scheme is resilient against few KR attacks. In 2015, Subramanyan *et al.* [17] gave a KR attack (popularly known as the SAT attack) and empirically showed that the correct key can be efficiently recovered from circuits locked with RLL and its variants.

From 2015 to 2021 (before our work), several new logic-locking schemes [15], [27]–[32] were designed, primarily to thwart the SAT attack.

At a high level, prior works use the following steps to design a logic-locking scheme:

- 1) generate an opaque circuit using some structural modifications of the original circuit;
- 2) show that the opaque circuit has low power, performance and area overhead;
- 3) claim security by empirically showing that the scheme is resilient against few KR attacks.

Strangely, these schemes (implicitly) claim that circuits with very small domain can be hidden. This is impossible since an adversary can use its oracle access to a restored chip (purchased from the market) to build the full truth table of the original circuit. Not only that, almost all logic-locking schemes have been shown to be vulnerable to efficient KR attacks, e.g., [16]–[26] that allows the foundry to define the functionality F in full.

B. Relation to Existing Cryptographic Primitives

Our security notions capture an attack on the privacy of the IP author’s circuit design. Thus, it is natural to think of related cryptographic primitives: program obfuscation, multi-party computation, encryption and function secret sharing.

PROGRAM OBFUSCATION. Design hiding is *orthogonal* to recent work in cryptography on program/circuit obfuscation (e.g. [52], [53]). Loosely, obfuscation in the latter setting makes the syntactic requirement that the obfuscated circuit have the same input-output functionality as the original circuit. In the design-hiding setting, at an abstract level, the opaque circuit implements a *set* of functions. It is only when the chip is restored, via the secret key, that the chip must faithfully compute the intended function.

MULTIPARTY COMPUTATION. Secure multi-party computation [54]–[57] allows two or more parties to “securely” compute some arbitrary function on inputs that are secrets of the respective parties.

Though MPC has become increasingly feasible to deploy in real-world systems [58]–[68], we *cannot* use it to design secure DH schemes in the logic-locking setting because of the fundamental differences in the threat models.

Logic locking involves two parties: the IP author (whose IC design needs to be protected), and the adversarial foundry (who is responsible for fabricating multiple opaque chips from the opaque circuit). The two popular settings of two-party computation that have some similarity to the logic-locking setting are: secure-function evaluation (SFE) [55], [56], [63] and private-function evaluation (PFE) [55]–[57], [62], and the threat models of both are *different* compared to logic locking.

In SFE, the function that is evaluated is public, and the input of the other party is private. Clearly, SFE is not suitable in the logic-locking setting as it requires the IC design of the IP author to be public. In PFE, the function that is evaluated is a private input of one party, and the input to the function is private to the other party. The logic-locking setting satisfies this aspect of PFE as the IC design is the private input of the IP author. However, the adversarial foundry gets a lot

of side information about the concealed IC design due to its unrestricted access to the opaque circuit and oracle access to an honestly-restored chip. This violates the PFE security goal as the party that does not have access to the function — secret input of the other party — should learn *only* the value of the function on the secret input. Hence, PFE is also not suitable in the logic-locking setting.

ENCRYPTION. In order to fabricate chips, the foundry parses the layout (described using the GDSII format) of a circuit design to replace the API calls to its hardware library with physical circuits. Thus, one cannot simply encrypt (say) the layout file before handing it to the foundry: secure encryption schemes produce ciphertexts that are indistinguishable from random bit-string, which almost certainly will not encode a valid circuit. Even if it did, that circuit would have no relationship to one realizing the intended functionality, thereby forfeiting the economic benefit of outsourced fabrication. And clearly the foundry cannot be allowed to decrypt the ciphertext, if it is not trusted in the first place.

FUNCTION SECRET SHARING. In 2015, Boyle et al. [69] introduced a cryptographic primitive called function secret sharing (FSS) that partitions a function F into multiple shares — the shares are distributed among multiple parties — such that an adversary that does not have access to all the shares learns *nothing* about F . The threat model of FSS prevents the adversary to learn the value of F on any input. In fact, if the adversary gets oracle access to F akin to the DH setting, the security of the FSS construction [69][Section 3] for point functions in the two-party setting falls apart as the adversary can use a single oracle query to identify the point function that was protected by FSS in the security experiment.

C. Fully-malicious adversaries break KR security

Hardware trojans are malicious modifications to a target circuit that are hard to detect [48]. For example, a key-leaking hardware trojan in a cryptographic IP core (say, AES) leaks the secret key, when the IP core is run on a small and specific sequence of triggering inputs. A standard assumption in trojan-insertion attacks requires the adversary to have unrestricted access to the target circuit in order to make trojans stealthy.

In our security experiments, an adversarial foundry gets unrestricted access to the opaque circuit C_L that the experiments generate using *Hide*. When the foundry is fully malicious, it can insert a key-leaking hardware trojan (T) in C_L to leak the key K_O .

Consider the following (slightly informal) KR attack in the fully-malicious setting against *any* DH scheme Π where: *Restore* fixes the $k_o(\theta)$ key-input bits to restore the original, intended functionality. (Logic-locking schemes like RLL [6] and its variants [44]–[46] and SAT-attack-resistant schemes [15], [27]–[29], [32] fit this description perfectly.)

- Query FAB with C_L^T to get a chip \mathbb{C}_L^T , where C_L^T denotes the opaque circuit C_L with trojan T .
- Query RESTORE with \mathbb{C}_L^T to get the index p of an honestly-restored chip \mathbb{C}_P . Let \mathcal{I} be the sequence of triggering inputs

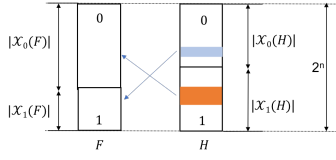


Fig. 5: Representation of F and H using $\mathcal{X}_i(\cdot)$. The orange (resp. blue) box denotes the distinguishing inputs that fall in set $\mathcal{X}_1(H)$ (resp. $\mathcal{X}_0(H)$). By construction of $OneChaff_{hd}$, the orange (resp. blue) box belongs to $\mathcal{X}_0(F)$ (resp. $\mathcal{X}_1(F)$).

for the trojan T . By design, \mathbb{C}_P contains the trojan T and leaks the key K_o when \mathbb{C}_P is run on inputs in \mathcal{I} .

- Query RUN with inputs $(p, \mathcal{I}[1]), (p, \mathcal{I}[2]), \dots, (p, \mathcal{I}[v])$ to recover K_o . Here, $v = |\mathcal{I}|$.

This attack breaks the KR security of Π . In order to thwart this attack, *Restore* has to detect arbitrary, key-leaking hardware trojans. We do not know how this can be achieved in the setting of design-hiding schemes such as logic locking — the main focus of our work — where the IP author outsources the entire fabrication process to an external foundry. *Therefore, our security theorems assume that the adversary is honest-but-curious.* Recall that none of the existing DH schemes considered *Fab* as part of the syntax. Hence, by default, these schemes are designed in the honest-but-curious setting.

Note that Dziembowski et al. [13] showed that it is possible to design “trojan-resilient” circuits using split manufacturing as a DH scheme; as mentioned earlier, split manufacturing requires the IP author to fabricate a portion of the circuit and hence is not a common setting in IC supply chain; instead we focus on the stronger logic-locking setting in this work.

D. Delayed Proofs and a Theorem

Proof of Claim 1. Any Boolean function H can be described fully using sets $\mathcal{X}_0(H)$ and $\mathcal{X}_1(H)$. Recall that $\mathcal{X}_i(H)$ denotes the set of all inputs for which H 's value is i . In Fig. 5, we show the representation of F and H using their respective \mathcal{X}_i .

For $i \in \{0, 1\}$ and for any $\hat{F} \in \mathcal{R}_L$, let $S_i(H, \hat{F}) = S_i = \mathcal{I}^\neq(H, \hat{F}) \cap \mathcal{X}_i(H)$. By construction: S_1, S_0 are disjoint, and $S_i \subset \mathcal{X}_{1-i}(\hat{F})$ as $\hat{F}(X) = \neg H(X)$ for all $X \in \mathcal{I}^\neq(H, \hat{F})$.

We can build $\mathcal{X}_1(\hat{F})$ using the sets $\mathcal{X}_1(H)$ and S_i in three steps. First, initialize $\mathcal{X}_1(\hat{F}) = \mathcal{X}_1(H)$. Then, remove all elements in set S_1 from $\mathcal{X}_1(\hat{F})$, and in the final step, add all elements of set S_0 , i.e., $\mathcal{X}_1(\hat{F}) = (\mathcal{X}_1(H) \setminus S_1) \cup S_0$. ($\mathcal{X}_0(\hat{F}) = (\mathcal{X}_0(H) \setminus S_0) \cup S_1$.)

By construction, the hamming weight of all functions in \mathcal{R}_L will fall in the window $[\text{hw}(H) - \Delta, \text{hw}(H) + \Delta]$, i.e., $\text{hw}(\hat{F}) \in [\text{hw}(H) - \Delta, \text{hw}(H) + \Delta]$ for all $\hat{F} \in \mathcal{R}_L$. Since $\mathcal{F}_{[h,h]}^0$ comprises of functions with hamming weight h only, we can construct $(\mathcal{F}_{[h,h]}^0 \cap \mathcal{R}_L^0)$ by only keeping functions $\hat{F} \in \mathcal{R}_L$ whose $\text{hw}(\hat{F}) = h$. That is, for each $\hat{F} \in (\mathcal{F}_{[h,h]}^0 \cap \mathcal{R}_L^0)$, the hamming weight of \hat{F} will satisfy: $|\mathcal{X}_1(H)| - |S_1(H, \hat{F})| + |S_0(H, \hat{F})| = h$. Also, by construction

of $OneChaff_{hd}$, the total number of distinguishing inputs will be Δ . Hence, $|S_1(H, \hat{F})| + |S_0(H, \hat{F})| = \Delta$. Solving these two equations, and using $\text{hw}(H) = h + \alpha$ gives us:

$$|S_1(H, \hat{F})| = (\Delta + \alpha)/2; |S_0(H, \hat{F})| = (\Delta - \alpha)/2.$$

These equations tell us that given H , each function \hat{F} in the set $\mathcal{F}_{[h,h]} \cap \mathcal{R}_L$ has $(\Delta + \alpha)/2$ distinguishing inputs in $\mathcal{X}_1(H)$ and $(\Delta - \alpha)/2$ distinguishing inputs in set $\mathcal{X}_0(H)$. Total number of such functions will be $\text{Bin}(|\mathcal{X}_1(H)|, (\Delta + \alpha)/2) \times \text{Bin}(|\mathcal{X}_0(H)|, (\Delta - \alpha)/2)$, where $\alpha = \text{hw}(H) - h$ and $\text{Bin}(n, m) = \binom{n}{m}$. Substituting $|\mathcal{X}_1(H)| = \text{hw}(H)$ and $|\mathcal{X}_0(H)| = 2^n - \text{hw}(H)$ in the previous equation gives us the final bound in the claim. ■

Proof of Claim 2. We will use the sets $\mathcal{X}_i(F)$ and $S_i(F, H) = S_i = \mathcal{I}^\neq(F, H) \cap \mathcal{X}_i(F)$, where $i \in \{0, 1\}$, to define $\mathcal{X}_i(H)$ similar to how we used sets $\mathcal{X}_i(H)$ and $\mathcal{I}^\neq(F, H) \cap \mathcal{X}_i(H)$ to define $\mathcal{X}_i(F)$ in the proof of Claim 1. Note that the set $\mathcal{X}_i(F)$ is used as the base function in this proof; in Claim 1, we used $\mathcal{X}_i(H)$ as the base function. Following the same three steps as in the proof of Claim 1, we will get $\mathcal{X}_1(H) = (\mathcal{X}_1(F) \setminus S_1) \cup S_0$ and $|\mathcal{X}_1(H)| = |\mathcal{X}_1(F)| - |S_1| + |S_0|$.

Next, we need bounds on S_i 's. Before that, observe that the set $\mathcal{I}^\neq(F, H) = \{X_1, X_2, \dots, X_\Delta\}$ is a random variable that follows a uniform distribution, i.e., $\mathcal{I}^\neq(F, H) \leftarrow^* (U_n)^\Delta$. Since $\text{hw}(F) = h$, $\Pr[X \in \mathcal{X}_1(F)] = h/2^n$, and $\Pr[X \in \mathcal{X}_0(F)] = 1 - h/2^n$, where $X \leftarrow^* \mathcal{I}^\neq(F, H)$.

Now, the size of S_i , i.e., $|S_i|$ is a random variable that is binomially distributed, with $\mathbb{E}[|S_1|] = \mu_1 = \Delta(h/2^n)$ and $\mathbb{E}[|S_0|] = \mu_0 = \Delta(1 - h/2^n)$. By a standard Hoeffding bound, we have, for all $\epsilon > 0$, $\Pr[|S_i| \geq \mu_i + \Delta\epsilon] \leq \exp(-2\Delta\epsilon^2)$. Setting $\epsilon = \sqrt{(\ln \Delta)/\Delta}$, we get with probability at most $1/\Delta^2$, $\Pr[|S_i| \geq \mu_i + \sqrt{\Delta \ln \Delta}]$. We can write this otherwise as: $\Pr[|S_i| < \mu_i + \sqrt{\Delta \ln \Delta}]$ with probability at least $1 - 1/\Delta^2$. We will assume that the bounds on S_i are tight; we will reflect the uncertainty in the claim/theorem statements. Using the bounds on S_i in $|\mathcal{X}_1(H)| = |\mathcal{X}_1(F)| - |S_1| + |S_0|$, we get

$$\begin{aligned} |\mathcal{X}_1(H)| &= |\mathcal{X}_1(F)| - (\mu_1 + \sqrt{\Delta \ln \Delta}) + (\mu_0 + \sqrt{\Delta \ln \Delta}) \\ &= h - \Delta(h/2^n) + \Delta(1 - h/2^n) \\ &= h + \Delta - (\Delta h)/2^{n-1}, \end{aligned}$$

with probability at least $(1 - 1/\Delta^2)$. ■

Proof of Theorem 2. Let $\delta = \Delta(1 - h/2^{n-1})$. From Lemma 1, Claims 1 and 2, we get:

$$\begin{aligned} \text{Adv}_{(\Pi, \theta), Fab}^{\text{FR}}(\mathcal{F}_{[h,h]}, A) &\leq \left(\left(\frac{h + \delta}{(\Delta + \delta)/2} \right) \left(\frac{2^n - h - \delta}{(\Delta - \delta)/2} \right) \right)^{-1} + \epsilon_{t,q,1}. \quad (1) \end{aligned}$$

Since $\delta = \Delta(1 - h/2^{n-1})$, we have: $\Delta - \delta < \Delta$, $\Delta + \delta < 2\Delta$, and $h + \delta \leq h + \Delta < 2^{n-1}$. Using the above observa-

tions, and the standard inequality $\text{Bin}(p, r) \geq (p/r)^r$, where $\text{Bin}(p, r) = \binom{p}{r}$, we get,

$$\text{Bin}(h + \delta, (\Delta + \delta)/2) \geq \left(\frac{h}{\Delta}\right)^{\Delta/2}, \text{ and}$$

$$\text{Bin}(2^n - h - \delta, (\Delta - \delta)/2) \geq \left(\frac{2^{n-1}}{(\Delta/2)}\right)^{\Delta/2} \geq \left(\frac{2^n}{\Delta}\right)^{\Delta/2}.$$

Plugging these bounds in equation 1 gives us the final bound in the claim. Note that the inequalities will shift as the equation has the inverse of each binomial term. ■

Proof of Lemma 2. Let $j \leq q$ and $Z_j = \mathcal{Q}_j^{\text{key}}$. We prove this claim by induction. We have already identified the base case of $\mathbb{E}[Z_0] = Z_0 = 0$ and $\mathbb{E}[Z_1] = \Delta/\text{hw}(H)$. In the general case, we will compute $\mathbb{E}[Z_j] = \mathbb{E}[\mathbb{E}[Z_j | Z_{j-1}]]$. First,

$$\begin{aligned} \mathbb{E}[Z_j | Z_{j-1}] &= (Z_{j-1} + 1) \frac{\Delta - Z_{j-1}}{\text{hw}(H) - (j-1)} \\ &\quad + (Z_{j-1}) \left(1 - \frac{\Delta - Z_{j-1}}{\text{hw}(H) - (j-1)}\right) \\ &= Z_{j-1} + \left(\frac{\Delta - Z_{j-1}}{\text{hw}(H) - (j-1)}\right) \\ &= \left(\frac{(\text{hw}(H) - j)Z_{j-1} + \Delta}{\text{hw}(H) - (j-1)}\right) \end{aligned}$$

Now taking the expected value of both sides, we get the expression

$$\begin{aligned} \mathbb{E}[Z_j] &= \left(\frac{\text{hw}(H) - j}{\text{hw}(H) - (j-1)}\right) \mathbb{E}[Z_{j-1}] + \frac{\Delta}{\text{hw}(H) - (j-1)} \\ &= \left(\frac{\text{hw}(H) - j}{\text{hw}(H) - (j-1)}\right) \frac{(j-1)\Delta}{\text{hw}(H)} + \frac{\Delta}{\text{hw}(H) - (j-1)} \\ &= \left(\frac{j\Delta(\text{hw}(H) - (j-1))}{\text{hw}(H)(\text{hw}(H) - (j-1))}\right) \\ &= \frac{j\Delta}{\text{hw}(H)} \end{aligned}$$

Thus, for $j = q$, $\mathbb{E}[Z_q] = \frac{q\Delta}{\text{hw}(H)}$.

We now show that Z_q is tightly concentrated around its expectation. Observe that $Z_j - Z_{j-1} \leq 1$ for all $1 \leq j \leq q$, so by the Azuma-Hoeffding inequality $\Pr[Z_q \geq (1 + \epsilon)\mathbb{E}[Z_q]] \leq 2 \exp\left(-\frac{(\epsilon\mathbb{E}[Z_q])^2}{2q}\right)$. If we set the upperbound to equal

a parameter α and solve for ϵ , we obtain $\epsilon = \frac{\sqrt{2q \ln(\frac{2}{\alpha})}}{\mathbb{E}[Z_q]}$.

Pushing a bit further, we set $\alpha = 2/\Delta^2$ and arrive at the following result: $Z_q < \mathbb{E}[Z_q] + \sqrt{4q \ln \Delta}$ with probability at least $1 - (2/\Delta^2)$. ■

Proof of Claim 3. After q queries, the adversary learns $|\mathcal{Q}_q^{\text{key}}|$ distinguishing inputs. Let $\delta = \text{hw}(H) - h$. From the proof of Claim 1, we know that given H , each function \hat{F} in the set $\mathcal{F}_{[h,h]} \cap \mathcal{R}_L$ has $(\Delta + \delta)/2$ distinguishing inputs in $\mathcal{X}_1(H)$ and $(\Delta - \delta)/2$ distinguishing inputs in set $\mathcal{X}_0(H)$. Since $\mathcal{Q}_q \subset \mathcal{X}_1(H)$ with high probability (from Lemma 2) and the adversary knows $|\mathcal{Q}_q^{\text{key}}|$ distinguishing inputs in

set $\mathcal{X}_1(H)$, the total number of functions in $\mathcal{F}_{[h,h]} \cap \mathcal{R}_L$ that have values $F(X_i)$ on all inputs $X_i \in \mathcal{Q}_q$ will be $\binom{|\mathcal{X}_1(H)| - q}{(\Delta + \delta)/2 - |\mathcal{Q}_q^{\text{key}}|} \binom{|\mathcal{X}_0(H)|}{(\Delta - \delta)/2}$. ■

Proof of Corollary 1. Let $Z_q = \mathcal{Q}_q^{\text{key}}$. From Claim 2 and Lemma 2, we have, $Z_q \leq \sqrt{4q \ln \Delta} + \frac{q\Delta}{h + \Delta(1 - h/2^{(n-1)})}$, with probability at least $(1 - 2/\Delta^2)$. Since we are finding an upperbound on Z_q , we can substitute $h + \Delta(1 - h/2^{(n-1)})$ with h to get a weaker-but-easier-to-interpret bound. As discussed earlier, for practical *OneChaff*_{hd} schemes, the loss in accuracy will not be much since $\Delta \ll 2^n$. In the first term, $\ln \Delta < \ln 2^{n-1} < (n-1) \ln 2 < n$. Thus, $\sqrt{4q \ln \Delta} < \sqrt{4qn}$. In order to prove the bound, which is the inverse of the product of the two binomial terms in Claim 3, we will show that the left binomial term satisfies

$$\left(\binom{h + \delta - q}{((\Delta + \delta)/2 - (q\Delta/h + \sqrt{4qn}))}\right) \geq \left(\frac{h - q}{\Delta}\right)^{\Delta/4},$$

and the right binomial term satisfies

$$\left(\binom{2^n - h - \delta}{(\Delta - \delta)/2}\right) \geq \left(\frac{2^n}{\Delta}\right)^{\Delta/4}.$$

Under our assumption $q < h/8$, $q\Delta/h < \Delta/8$. When $q < \Delta^2/256n$, $\sqrt{4qn} < \Delta/8$. Thus, $(q\Delta/h + \sqrt{4qn}) < \Delta/4$, and

$$\left(\binom{h + \delta - q}{((\Delta + \delta)/2 - (q\Delta/h + \sqrt{4qn}))}\right) \geq \left(\binom{h + \delta - q}{(\Delta/4 + \delta/4)}\right)$$

Using $\text{Bin}(p, r) = \binom{p}{r} \geq (p/r)^r$, we get:

$$\begin{aligned} \text{Bin}((h + \delta - q), (\Delta + \delta)/4) &\geq \left(\frac{h + \delta - q}{(\Delta + \delta)/4}\right)^{(\Delta + \delta)/4} \\ &\geq \left(\frac{h - q}{\Delta}\right)^{\Delta/4}. \end{aligned}$$

Since $\delta = \Delta(1 - h/2^{n-1})$, and $h + \Delta < 2^{n-1}$, we have: $\Delta - \delta < \Delta$, and $h + \delta \leq h + \Delta < 2^{n-1}$. Using these observations, we get

$$\begin{aligned} \text{Bin}(2^n - h - \delta, (\Delta - \delta)/2) &\geq \text{Bin}(2^{n-1}, (\Delta - \delta)/2) \\ &\geq \left(\frac{2^n}{(\Delta - \delta)/4}\right)^{(\Delta - \delta)/4} \\ &\geq \left(\frac{2^n}{\Delta}\right)^{\Delta/4} \end{aligned}$$

■

Theorem 4. Fix $m = 1$ and integers $n, h, \Delta, q > 0$ such that $\Delta < h$, $0 < \Delta + h < 2^{n-1}$. Let $\delta = \Delta(1 - h/2^{n-1})$. Let an honest-but-curious adversary A attack the FR security of *OneChaff*_{hd} using resources $(t, q_f = 1, q_s = 1, q_r = q)$. Then,

$$\begin{aligned} \text{Adv}_{(\text{OneChaff}_{\text{hd}}, \theta), \text{Fab}}^{\text{FR}}(\mathcal{F}_{[h,h]}, A) &\leq \frac{1}{\binom{h + \delta - q}{(\Delta + \delta)/2 - \sqrt{4q \ln(\Delta) - q\Delta/(h + \delta)}} \binom{2^n - h - \delta}{(\Delta - \delta)/2}} + \epsilon_{t,q,1}, \end{aligned}$$

with probability at least $(1 - 2/\Delta^2)$. ♦