

Interpretable Noninterference Measurement and Its Application to Processor Designs

ZIQIAO ZHOU*, Microsoft Research, USA
MICHAEL K. REITER*, Duke University, USA

Noninterference measurement quantifies the secret information that might leak to an adversary from what the adversary can observe and influence about the computation. Static and high-fidelity noninterference measurement has been difficult to scale to complex computations, however. This paper scales a recent framework for noninterference measurement to the open-source RISC-V BOOM core as specified in Verilog, through three key innovations: logically characterizing the core’s execution incrementally, applying specific optimizations between each cycle; permitting information to be declassified, to focus leakage measurement to only secret information that cannot be inferred from the declassified information; and interpreting leakage measurements for the analyst in terms of simple rules that characterize when leakage occurs. Case studies on cache-based side channels generally, and on specific instances including SPECTRE attacks, show that the resulting toolchain, called DINOme, effectively scales to this modern processor design.

CCS Concepts: • **Security and privacy** → **Logic and verification**; • **Hardware** → *Theorem proving and SAT solving*.

Additional Key Words and Phrases: information flow, interference, declassification, interpretability

ACM Reference Format:

Ziqiao Zhou and Michael K. Reiter. 2021. Interpretable Noninterference Measurement and Its Application to Processor Designs. *Proc. ACM Program. Lang.* 5, OOPSLA, Article 141 (October 2021), 30 pages. <https://doi.org/10.1145/3485518>

1 INTRODUCTION

Noninterference [Goguen and Meseguer 1982] is a classic information flow policy that, informally, requires that an attacker’s view be unaffected by the values that should remain secret to it. Since systems often necessarily leak some information, however, a more practical goal is to insist that the interference be “small”, which in turn requires that it be measured in some way. Various methodologies have been proposed for doing so statically (e.g., Backes et al. [2009]; Phan and Malacaria [2014]; Zhang et al. [2010]), though these techniques invariably must balance a tension between measurement fidelity and scalability to complex computations.

A recent advance in this domain was due to Zhou et al. [2018], which formulated noninterference measurement in terms of a *projected model counting* problem that, in turn, was amenable to relatively efficient, *approximate* model counting methods. Their measurement approach, however, scales to programs of only modest complexity, for two reasons. Computationally, their technique relies on symbolic execution to generate a logical postcondition for the computation for which

*Work performed in part at the University of North Carolina, Chapel Hill, NC, USA.

Authors’ addresses: Ziqiao Zhou, Microsoft Research, Redmond, WA, USA, ziqiaozhou@microsoft.com; Michael K. Reiter, Duke University, Durham, NC, USA, michael.reiter@duke.edu.



This work is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License.

© 2021 Copyright held by the owner/author(s).

2475-1421/2021/10-ART141

<https://doi.org/10.1145/3485518>

noninterference is to be measured. For example, this step alone required six hours for Smaz and eight hours for Gzip, using 16 cores, for extracting postconditions to measure the risk of CRIME attacks [Kelsey 2002] against these compression libraries. More qualitatively, while their technique provides a measurement of interference, it provides the analyst little assistance in interpreting the measurement or focusing the analysis on particular aspects of the leakage.

While noninterference measurement for arbitrary computations remains out of reach, in this paper we adapt the approach of Zhou et al. [2018] to address the previous shortcomings within a particularly important and complex domain, namely information leaks arising in hardware processors. Leakage of software secrets due to processor optimizations have attracted massive attention in recent years, especially since the discovery of vulnerabilities arising due to the footprint of speculative executions in processor caches (SPECTRE [Kocher et al. 2019], MELTDOWN [Lipp et al. 2018], and variants). Even though many defenses (e.g., Tan et al. [2020]; Wang and Lee [2007]; Werner et al. [2019]; Zhou et al. [2016]) have been proposed to interfere with cache-based side channels, we are aware of no measurement methodology to compare designs and evaluate their effectiveness, working directly from their Verilog specifications. Adapting a technique like Zhou et al. [2018] to do so, moreover, appears difficult: the sheer complexity of modern processor designs both necessitates greater support to help the analyst understand the factors contributing to the leakage and poses significant scaling challenges to such techniques.

In this paper, we present a methodology to measure and interpret leaks in processors, using three key advances:

- Our methodology enables analysts to *declassify* certain information, thereby focusing the measurement on any *other* leakage that might be occurring, i.e., leakage that cannot be inferred from the declassified information. For systems as complex as modern processors, this ability is essential to permit analysts to decompose and analyze leakage in a piecemeal fashion.
- The complexity of processor designs means that once leakage is measured, the exact conditions that cause this leakage might not immediately be evident. Our methodology therefore incorporates a method of *interpreting* the leakage, i.e., providing simple rules that indicate circumstances in which leakage will (or will not) occur. These rules facilitate analyst understanding of the root causes of leakage and can guide analysts to declassify leakage that can be ignored. Each such rule is additionally accompanied by a precision and recall, so that analysts can prioritize the rules they address. These rules are expressed in terms of conditions in which leakage occurs, enabling executions to be generated that demonstrate the leakage if desired but hiding the particulars of the executions from analysts if not.
- Since generating a logical postcondition for a processor’s execution of a program en masse is intractable, we devise a method to build the postcondition one cycle at a time. To build single-cycle formulas, we abandon symbolic execution, as we found that applying it to hardware designs induces significant path explosion for even one CPU cycle. Instead, we extract the single-cycle formulas without solving for feasible paths, and then leverage a number of aggressive optimizations when stitching single-cycle formulas together to build the postcondition for the processor’s multi-cycle execution.

Due to the focus of our methodology on support for declassification and interpretability, we call our tool that realizes it DINoME (for “Declassification and Interpretability for Noninterference Measurement”).

To evaluate DINoME, we apply it to evaluate leakage during execution on a RISC-V BOOM core [Celio et al. 2017], a state-of-the-art public domain processor design. Our improvements to generating logical postconditions for execution permit DINoME to do so for more than 100 cycles of this core. This, in turn, permits us to evaluate leakage from cache-based side channels

(PRIME+PROBE [Osvik et al. 2006] and FLUSH+RELOAD [Yarom and Falkner 2014]) in various scenarios, including cryptographic key leakage in sliding-window based modular exponentiation (e.g., Aciüzmez [2007]; Percival [2005]), leakage of secrets due to speculative execution, and how this leakage is (incompletely) mitigated by proposed improvements such as SCATTERCACHE [Werner et al. 2019] and PHANTOMCACHE [Tan et al. 2020]. In each case, we not only measure interference but also generate rules to explain why the leakage occurs, and in some cases refine our view of the leakage using declassification. Our performance evaluation of DINoME indicates that these types of analyses complete in times ranging from seconds to under 15 minutes (using horizontal scaling), after an initial phase to assemble the logical postcondition of up to (only) two hours on (only) a single core.

The rest of this paper is structured as follows. We discuss related work in Sec. 2, and provide both background on the framework on which we build [Zhou et al. 2018] and our introduction of declassification to it, our first contribution, in Sec. 3. We present our method for interpreting leakage in Sec. 4. We address implementation challenges in Sec. 5, and then evaluate DINoME through several case studies in Sec. 6. We discuss DINoME's performance in Sec. 7, its limitations in Sec. 8, and our conclusions in Sec. 9.

2 RELATED WORK

To our knowledge, DINoME is the first work to measure information leakage from an executable hardware specification instantiated with a software program, in a manner that supports declassification and interpretation of its leakage results.

Timing side-channel analysis. Constant-time verification (e.g., Almeida et al. [2016]; Barthe et al. [2014]; Blazy et al. [2019]; Gleissenthall et al. [2019]; Zhang et al. [2015]) is a commonly used technique to analyze timing side channels. Software-level verification (e.g., Almeida et al. [2016]; Blazy et al. [2019]) checks whether a software program runs in a constant time under specified hardware assumptions. For example, a software-level analysis [Almeida et al. 2016] might conclude that a variable leaks if it is used in a branch condition or as an address in memory access. In a different approach, hardware-level verifiers (e.g., Gleissenthall et al. [2019]; Zhang et al. [2015, 2018]) can formally verify the existence of timing side channels using cycle-precise logic derived from hardware specifications. These works check for timing dependencies on secret variables but do not quantify secret leakage due to timing variations in different executions.

Hardware leakage modeling. Some works use simplified hardware models instead of real designs (e.g., Chattopadhyay et al. [2017]; Doychev et al. [2013]; Malacaria et al. [2018]), which makes the computation target feasible but requires more domain knowledge and manual effort to construct the model. Black-box analysis of real systems avoids the use of domain knowledge through a data-driven method that uses sampled data in a real system for estimating the leakage (e.g., Nilizadeh et al. [2019]; Oleksii et al. [2020]; Song et al. [2001]). In contrast, DINoME measures leakage from hardware specifications written in a hardware design language.

Quantitative information flow. QIF (e.g., Gray [1991]; Smith [2009, 2011]) represents information leakage through a numeric measurement; most mainstream QIF works (e.g., Chapman and Evans [2011]; Phan and Malacaria [2014]; Zhang et al. [2010]) use entropy as their measure [Seidenfeld 1986]. The use of entropy for measuring QIF in actual systems can lead to significant costs, due to the need to compute the input preimage per output value. In addition, real implementations tend to use the most conservative min-entropy measure; e.g., QIF-Verilog [Guo et al. 2019] propagates a min-entropy label per gate and accumulates the leakage across all gates, which overestimates leakage due to its conservative leakage accumulation, especially in large, complex hardware designs

(e.g., a CPU core). Entropy also does not distinguish between leaking a few bits in many executions or leaking more bits in a few cases. Alternatives to entropy-based leakage—e.g., differential privacy [Dwork et al. 2006], noninterference measurement [Zhou et al. 2018], classifier-based measurement [Chapman and Evans 2011], and quantitative hyperproperties [Sahai et al. 2020; Yasuoka and Terauchi 2014]—measure the attacker’s ability to distinguish some secret values from others. Those metrics do not accommodate declassification or leakage interpretability, our main concerns here.

Declassification. To rule out allowed leakage and focus on targeted leakage, information flow control research supports declassification policies to specify the secret information permitted to transfer to observable variables (e.g., Banerjee et al. [2008]; Chong and Myers [2004]; Ferraiuolo et al. [2017]; Giacobazzi and Mastroeni [2018]; McCall et al. [2018]; Sabelfeld and Myers [2003]; Sabelfeld and Sands [2009]). However, while this work omits declassified information from its analysis, it does not quantitatively measure the remaining leakage in light of what the attacker can already infer from the declassified information. In contrast, our work adapts information leakage measurement to account for such inferences.

Leakage interpretability. To interpret quantitative leakage, domain-specific works (e.g., SPEECH-MINER [Xiao et al. 2020], CACHEBAR [Zhou et al. 2016]) use customized measures following a specific attack templates, forgoing general measures. Although those customized measures are more understandable when interpreting a specific attack vector, they are blind to leakage from different attacks not considered. One crucial improvement our work makes in evaluating information leakage is to generate an interpretable model to explain how leakage occurs. Already an emerging topic in machine learning (e.g., Chen et al. [2018]; Molnar [2019]), interpretability is especially important in security evaluation, since it is not easy to draw a clear threshold to indicate when a system is secure enough, even with a perfect measure. Many methods for measuring leakage in software (e.g., Chattopadhyay and Roychoudhury [2018]; Godefroid et al. [2012]; Wang et al. [2009]; Zhou et al. [2018]) generate a code path to help the analyst understand leakage. However, leakage in hardware-software joint codebases often exploits interactions between the two, which can manifest in many code-dependent paths. We are aware of no comparable work that explores an interpretable ML model to explain information-flow leakage, though the method we use to extract explanations in Sec. 4.3 builds from previous work in interpretable ML (e.g., Friedman and Popescu [2008]; Ribeiro et al. [2016, 2018]).

3 NONINTERFERENCE AND DECLASSIFICATION

We begin in Sec. 3.1 by providing background on the noninterference measurement methodology of Zhou et al. [2018]. We then discuss how we extend this methodology to support declassification, our first contribution, in Sec. 3.3.

3.1 Background on Noninterference Measure

To analyze the leakage from a procedure $proc$ ¹, the procedure is modeled as having four different sets of formal parameters: a set $Vars_{\bar{s}}$ of secret input variables; a set $Vars_{\bar{c}}$ of attacker-controlled input variables; a set $Vars_{\bar{t}}$ of other input variables; and a set $Vars_{\bar{o}}$ of attacker-observable output variables. The actual parameter values assigned to those variables in an invocation of $proc$ are given by maps $\bar{s} : Vars_{\bar{s}} \rightarrow Vals_{\bar{s}}$, $\bar{c} : Vars_{\bar{c}} \rightarrow Vals_{\bar{c}}$, and $\bar{t} : Vars_{\bar{t}} \rightarrow Vals_{\bar{t}}$, respectively; e.g., $\bar{t}(ivar) \in Vals_{\bar{t}}$ represents the value passed in variable $ivar \in Vars_{\bar{t}}$. The attacker-observable outputs of the procedure are defined by the map $\bar{o} : Vars_{\bar{o}} \rightarrow Vals_{\bar{o}}$. Accordingly, we denote the

¹Different from the definition used by Zhou et al. [2018], which is for a software procedure, our $proc(\bar{c}, \bar{t}, \bar{s})$ includes both the software and hardware logic.

procedure

$$\vec{o} \leftarrow \text{proc}(\vec{c}, \vec{i}, \vec{s})$$

We assume that *proc* is deterministic; a nondeterministic *proc* can be rendered deterministic by providing the random values as inputs, say \vec{i} ('coins'). A given *proc* can then be characterized by a logical postcondition $\Pi_{\text{proc}}(\vec{c}, \vec{o}, \vec{i}, \vec{s})$ that constrains how the values in \vec{o} relate to those in \vec{c} , \vec{i} , and \vec{s} in any execution. Without loss of generality, below we assume $\text{Vars}_{\vec{s}}$ contains a single variable *svar*, i.e., $\text{Vars}_{\vec{s}} = \{\text{svar}\}$.

The basic idea of the metric developed by Zhou et al. [2018] is to quantify the difficulty the attacker has in distinguishing between $\vec{s}(\text{svar}) \in S$ and $\vec{s}(\text{svar}) \in S'$ for random, disjoint sets S, S' , based on the $\langle \vec{c}, \vec{o} \rangle$ pairs possibly available to it in the two cases, denoted $Y_S, Y_{S'}$, i.e.,

$$\begin{aligned} X_S &= \left\{ \langle \vec{c}, \vec{o}, \vec{i} \rangle \mid \exists \vec{s} : \Pi_{\text{proc}}(\vec{c}, \vec{o}, \vec{i}, \vec{s}) \wedge \vec{s}(\text{svar}) \in S \right\} \\ Y_S &= \left\{ \langle \vec{c}, \vec{o} \rangle \mid \exists \vec{i} : \langle \vec{c}, \vec{o}, \vec{i} \rangle \in X_S \right\} \end{aligned}$$

Zhou et al. [2018] specifically explored the Jaccard distance between Y_S and $Y_{S'}$ to measure the difficulty an attacker would have in distinguishing between $\vec{s}(\text{svar}) \in S$ and $\vec{s}(\text{svar}) \in S'$. To better capture the importance of \vec{i} in the leakage, however, they further replaced $Y_S \cap Y_{S'}$ with $\hat{X}_{S,S'}$, where²

$$\begin{aligned} \check{X}_{S,S'} &= X_S \cup X_{S'} \\ \hat{X}_{S,S'} &= \left\{ \langle \vec{c}, \vec{o}, \vec{i} \rangle \mid \langle \vec{c}, \vec{o}, \vec{i} \rangle \in \check{X}_{S,S'} \wedge \langle \vec{c}, \vec{o} \rangle \in Y_S \cap Y_{S'} \right\} \end{aligned}$$

In this way, the number of values \vec{i} for $\langle \vec{c}, \vec{o} \rangle$ exposed in $\hat{X}_{S,S'}$ serves as the "weight" of that $\langle \vec{c}, \vec{o} \rangle$ pair. When $\langle \vec{c}, \vec{o}, \vec{i} \rangle$ is from

$$\check{X}_{S,S'} = \check{X}_{S,S'} \setminus \hat{X}_{S,S'} \quad (1)$$

an attacker can distinguish if $\vec{s}(\text{svar})$ is from S or S' . Zhou et al. [2018] thus suggested the measure \hat{J}_n , where

$$\hat{J}(S, S') = |\check{X}_{S,S'}| / |\hat{X}_{S,S'}| = 1 - |\hat{X}_{S,S'}| / |\check{X}_{S,S'}| \quad (2)$$

$$\hat{J}_n = \text{avg}_{\substack{S, S' : |S| = |S'| = n \\ \wedge S \cap S' = \emptyset}} \hat{J}(S, S') \quad (3)$$

As discussed by Zhou et al. [2018, Sec. III], when n is small, \hat{J}_n measures how frequently leakage occurs, whereas when n is large, it measures how much information about the secret leaks, when leakage occurs.

3.2 Motivating Examples

To see this measure applied to simple programs, consider the two programs with a secret shown in Fig. 1(a) and Fig. 1(b). The procedure in Fig. 1(a) returns a random value between $0 - 7$ or a fixed value 8 depending on whether \vec{s} ('secret') $\bmod 32 < 16$ if \vec{c} ('test') $\bmod 32 > 15$ and returns a fixed value 9 otherwise. The second procedure in Fig. 1(b) returns the five least significant bits of \vec{s} ('secret') & \vec{c} ('test'). Directly measuring the two procedures using \hat{J}_n leads to different leakage measures, as it should, as shown in Fig. 1(e).

Some sources of information leakage may be inevitable or intentional; e.g., a bank website may not mask the last four digits of a user's social security number when displaying it to her browser, and

²Our definition of $\hat{X}_{S,S'}$ differs from Zhou et al. [2018], which only requires $\langle \vec{c}, \vec{o}, \vec{i} \rangle \in X_S$. Ours has the same essential properties but is symmetric with respect to S and S' and so is easier to work with.

```

proc (c̄, ī, s̄)
  if (c̄('test') mod 32 > 15)
    if (s̄('secret') mod 32 < 16)
      ō('result') ← ī('random') mod 8
    else
      ō('result') ← 8
  else ō('result') ← 9

```

(a) Implicit flow

```

proc (c̄, ī, s̄)
  ō('result') ← c̄('test') & s̄('secret') & 0x1f

```

(b) Explicit flow

```

proc (c̄, ī, s̄)
  ō('result') ← c̄('test') & s̄('secret') & 0x2f

```

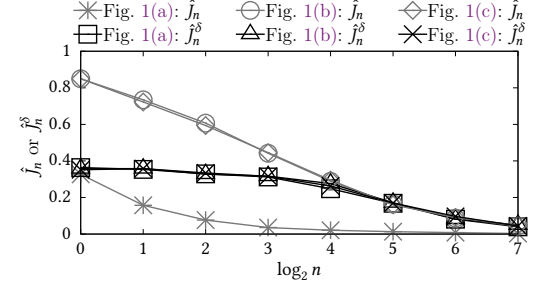
(c) Different explicit flow

```

δ(c̄, ī, s̄)
  Δ̄('info') ← s̄('secret') & 0x0f

```

(d) Declassification policy



(e) Measurement with vs. without declassification

Fig. 1. Motivating examples for declassification (Sec. 3.3) and interpretation (Sec. 4)

so the site intentionally “leaks” that portion to a malicious browser. In the context of the preceding example, now suppose the leakage of the four least significant bits of the secret is intended (similar to the SSN example). Since the \hat{J}_n curve only reflects the total interference, including the portion *intended* to leak (i.e., the four least significant bits), the \hat{J}_n curves shown in Fig. 1(e) mislead us to conclude that Fig. 1(a) is more secure than Fig. 1(b). In truth, they both additionally leak the fifth least significant bit, which is the only leakage that matters.

3.3 Declassification

To exclude such intended leakage from the analysis, it will be helpful to provide a method to exempt some identified information leakages specified by the analyst, allowing the analysis to focus on the leakage that remains. Specifically, our methodology seeks to assess the degree to which a procedure permits secrets to be distinguished by the attacker using attacker-observable and declassified information but not by the declassified information alone.

Let $\vec{\Delta} \leftarrow \delta(\vec{c}, \vec{i}, \vec{s})$ denote the allowed information exposure (e.g., for a website requiring SSN, $\vec{\Delta}$ is the last four digits), and let

$$\Pi_{proc, \delta}(\vec{c}, \vec{o}, \vec{\Delta}, \vec{i}, \vec{s}) \leftarrow \Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s}) \wedge \Pi_{\delta}(\vec{c}, \vec{\Delta}, \vec{i}, \vec{s})$$

where $\Pi_{\delta}(\vec{c}, \vec{\Delta}, \vec{i}, \vec{s})$ is a logical postcondition for δ that relates $\vec{\Delta}$ to \vec{c} , \vec{i} , and \vec{s} . Then, we can define the attacker’s accessible set Y_S^{δ} of $\langle \vec{c}, \vec{o}, \vec{\Delta} \rangle$ tuples and allowed accessible set D_S^{δ} consistent with chosen secret set S by

$$\begin{aligned}
X_S^{\delta} &= \left\{ \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \mid \exists \vec{s} : \vec{s}(svar) \in S \wedge \Pi_{proc, \delta}(\vec{c}, \vec{o}, \vec{\Delta}, \vec{i}, \vec{s}) \right\} \\
Y_S^{\delta} &= \left\{ \langle \vec{c}, \vec{o}, \vec{\Delta} \rangle \mid \exists \vec{i} : \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \in X_S^{\delta} \right\} \\
D_S^{\delta} &= \left\{ \langle \vec{c}, \vec{\Delta} \rangle \mid \exists \vec{o}, \vec{i} : \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \in X_S^{\delta} \right\}
\end{aligned}$$

Since the declassified information is allowed to leak, we are concerned only with cases where the secret is distinguishable by $\langle \vec{c}, \vec{o}, \vec{\Delta} \rangle$ but not by $\langle \vec{c}, \vec{\Delta} \rangle$. Here, we define a set $\tilde{X}_{S, S'}^{\delta}$ to include the

$proc(\vec{c}, \vec{i}, \vec{s})$
 $\vec{o}(ovar) \leftarrow \vec{s}(svar)[0 : 3]$
 (a) An artificial procedure

$\delta_{i,j}(\vec{c}, \vec{i}, \vec{s})$
 $\vec{\Delta}(dvar) \leftarrow \vec{s}(svar)[i : j]$
 (b) Declassification policy

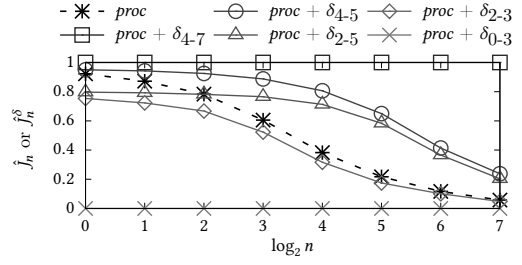


Fig. 2. Declassification example

tuples $\langle \vec{c}, \vec{o}, \vec{\Delta} \rangle$ that leak whether the secret is in S or S' , assuming $\langle \vec{c}, \vec{\Delta} \rangle$ is equivalent.

$$\check{X}_{S,S'}^\delta = \left\{ \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \left| \begin{array}{l} \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \in X_S^\delta \cup X_{S'}^\delta \\ \wedge \langle \vec{c}, \vec{\Delta} \rangle \in D_S^\delta \cap D_{S'}^\delta \end{array} \right. \right\} \quad (4)$$

$$\hat{X}_{S,S'}^\delta = \left\{ \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \left| \begin{array}{l} \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \in \check{X}_{S,S'}^\delta \\ \wedge \langle \vec{c}, \vec{o}, \vec{\Delta} \rangle \in Y_S^\delta \cap Y_{S'}^\delta \end{array} \right. \right\} \quad (5)$$

$$\tilde{X}_{S,S'}^\delta = \check{X}_{S,S'}^\delta \setminus \hat{X}_{S,S'}^\delta \quad (6)$$

Thus, we can use an alternative metric

$$\hat{J}_n^\delta(S, S') = \left| \tilde{X}_{S,S'}^\delta \right| / \left| \check{X}_{S,S'}^\delta \right| \quad (7)$$

$$\hat{J}_n^\delta = \text{avg}_{\substack{S, S' : |S| = |S'| = n \\ \wedge S \cap S' = \emptyset}} \hat{J}_n^\delta(S, S') \quad (8)$$

Returning to the examples in Fig. 1(a) and Fig. 1(b) with declassification of the four least significant bits (Fig. 1(d)), the \hat{J}_n^δ curves show the same quantitative leakage (Fig. 1(e)), as they should.

To further illustrate the impact of declassification, consider the simple procedure shown in Fig. 2(a). In this procedure, $\vec{s}(svar)$ is an 8-bit value, and $proc$ outputs the lowest 4 bits as $\vec{o}(ovar)$. The declassification policy shown in Fig. 2(b) allows the i -th to j -th bits of $\vec{s}(svar)$ to be released. We evaluate \hat{J}_n^δ with differently parameterized declassification policies in Fig. 2(c). Specifically, when the lowest 4 bits ($i = 0, j = 3$) are declassified, then the additional leakage from $proc$ is nothing, which is demonstrated by the “ $proc + \delta_{0-3}$ ” curve. When the declassification policy declassifies all but the lowest 4 bits ($i = 4, j = 7$), then the additional leakage by $proc$ is maximized, as shown by the “ $proc + \delta_{4-7}$ ” curve. Intuitively, if $\vec{o}(ovar)$ and $\vec{\Delta}(dvar)$ do not overlap (e.g., “ $proc + \delta_{4-7}$ ” and “ $proc + \delta_{4-5}$ ”), then the \hat{J}_n^δ curve should be higher than \hat{J}_n , whereas if $\vec{o}(ovar)$ includes all of $\vec{\Delta}(dvar)$ (e.g., “ $proc + \delta_{0-3}$ ” and “ $proc + \delta_{0-1}$ ”), then \hat{J}_n^δ should be lower than \hat{J}_n . A hybrid case occurs when $\vec{o}(ovar)$ includes a portion of $\vec{\Delta}(dvar)$ (e.g., “ $proc + \delta_{2-5}$ ”), where \hat{J}_n^δ is lower than \hat{J}_n when n is small but becomes larger when n is large. This is consistent with the interpretation that \hat{J}_n^δ with small n primarily reflects the number of secret values for which interference occurs [Zhou et al. 2018]; e.g., when $n = 1$, two secret values share bits 0–1 (and so cannot be distinguished by bits 0–3 after declassifying bits 2–5) in 25% of cases, but share bits 0–3 (and so cannot be distinguished using them) in only 6.25% of cases. Larger n , in contrast, better reflects the amount of leakage that occurs [Zhou et al. 2018]. For example, in a random partition of all 2^8 values into sets S and S' of equal size (i.e., $n = 2^7$), every value for bits 2–5 is represented in both S and S' with high probability.

In conjunction with the additional bits 0–1 output in \vec{o} (yielding six bits of the secret value in total), however, these bits give the attacker greater distinguishing power than do bits 0–3 alone.

4 INTERPRETING LEAKAGE

Our metric measures the additional interference of a secret with values observable by the attacker, beyond that implied by declassified information. For this to be useful to an analyst, however, we need to explain *how* this leakage occurs. Specifically, while the conditions under which leakage occurs are already present in the procedure postcondition, it is difficult to understand the formula without further help (e.g., see Sec. 6.6).

4.1 Motivating Examples for Interpretation

Consider again the motivating examples in Fig. 1(a) and Fig. 1(b). The two procedures own quite different outputs but still leak the same additional information about the secret after declassification (i.e., both leak the fifth least significant bit of the secret when \vec{c} (‘test’)’s fifth bit is 1 and nothing otherwise). To cut through the differences in code style and concrete values, DINoME derives the condition when a pair of secrets are distinguishable using paired samples of input. Thus, the interference rule for both cases becomes $|\vec{s}(\text{‘secret’})[4] - \vec{s}'(\text{‘secret’})[4]| > 0 \wedge \vec{c}(\text{‘test’})[4] = 1$. This rule shows the equivalence of these procedures’ leakages after declassification.

In addition, interpreting leakage can differentiate cases with the same *amount* of leakage but different conditions in which that leakage occurs. For example, the procedure in Fig. 1(c), which reveals the four least significant bits and the sixth bit of the secret when the sixth bit of \vec{c} (‘test’) is 1, leaks the same amount of information about a different portion of the secret under a different attack condition. A quantitative leakage measurement with the same four low-order bits declassified will not distinguish Fig. 1(c) from Fig. 1(b) (see Fig. 1(e)). Through DINoME’s interpretation, we provide a slightly different interference rule for Fig. 1(c), however: $|\vec{s}(\text{‘secret’})[5] - \vec{s}'(\text{‘secret’})[5]| > 0 \wedge \vec{c}(\text{‘test’})[5] = 1$.

Though these motivating examples seem small and readable even when using different coding styles and output values, real-world code can become difficult to understand, particularly when spanning different levels of abstraction (e.g., a processor and the code it is executing). It is here we expect our interpretation of interference to simplify investigating leakage. Learning from the previous examples, our interpretation should achieve two goals. First, the interference interpretation for the same functionality should be consistent no matter how the functionality is implemented. Second, the interference interpretation should distinguish two procedures if they leak information in different ways, even when they leak the same amount.

4.2 Noninterference and Interference Tuples

Our first step toward providing an intuitive explanation for the leakage that occurs is to train a binary classifier to classify 4-tuples $\langle \vec{c}, \vec{i}, \vec{s}, \vec{s}' \rangle$ into those that illustrate leakage occurring (i.e., that permit the attacker to distinguish \vec{s} (*svar*) and \vec{s}' (*svar*) from the resulting output \vec{o}) and those that do not. When using declassification, the interference tuples should only include those where the secrets can be distinguished using \vec{c} , \vec{o} , $\vec{\Delta}$ but not using just \vec{c} , $\vec{\Delta}$.

More specifically, we define the interference set *IS* based on (6). That is, when the attacker chooses \vec{c} , if an observable value is feasible for (\vec{i}, \vec{s}) for some \vec{i} but is never possible for (\vec{i}', \vec{s}') for any \vec{i}' that shares a declassification value with (\vec{i}, \vec{s}) , then $\langle \vec{c}, \vec{i}, \vec{s}, \vec{s}' \rangle$ is added to *IS*:

$$IS = \left\{ \langle \vec{c}, \vec{i}, \vec{s}, \vec{s}' \rangle \mid \exists \vec{o}, \vec{\Delta} : \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \in X_S^\delta \wedge \langle \vec{c}, \vec{\Delta} \rangle \in D_{S'}^\delta \cap D_S^\delta \wedge \langle \vec{c}, \vec{o}, \vec{\Delta} \rangle \in Y_S^\delta \setminus Y_{S'}^\delta \right\} \quad (9)$$

where $S = \{\vec{s}(\textit{svar})\}$ and $S' = \{\vec{s}'(\textit{svar})\}$.

The noninterference set NS should include two types of tuples. For an attacker-chosen \vec{c} , if there is a observable value \vec{o} that is feasible for an $\langle \vec{t}, \vec{s} \rangle$ pair and an $\langle \vec{t}', \vec{s}' \rangle$ pair, tuple $\langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle$ belongs to NS as it is an example where no interference occurs. In addition, for an attacker-chosen \vec{c} , if there is a declassification value $\vec{\Delta}$ that is feasible for $\langle \vec{t}, \vec{s} \rangle$ but not $\langle \vec{t}', \vec{s}' \rangle$ for any \vec{t}' , then $\langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle$ should also be added to NS , as \vec{s} and \vec{s}' can already be distinguished using the declassified value:

$$NS = \left\{ \langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle \mid \begin{array}{l} \exists \vec{o}, \vec{\Delta} : \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{t} \rangle \in X_S^\delta \wedge \langle \vec{c}, \vec{o}, \vec{\Delta} \rangle \in Y_S^\delta \cap Y_{S'}^\delta \\ \cup \left\{ \langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle \mid \exists \vec{o}, \vec{\Delta} : \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{t} \rangle \in X_S^\delta \wedge \langle \vec{c}, \vec{\Delta} \rangle \in D_S^\delta \setminus D_{S'}^\delta \right\} \end{array} \right\} \quad (10)$$

where $S = \{\vec{s}(svar)\}$ and $S' = \{\vec{s}'(svar)\}$.

Since NS and IS are large in practical scenarios, enumerating all tuples is generally infeasible. Instead, we generate samples in each set to train a machine learning model, from which explanations of the leakage will be extracted (as described below). Doing so with modern SAT solvers, however, typically results in samples that cover NS and IS unevenly, since solvers generally enumerate the next solution by simply adding a conflict constraint to block out previous solutions; as a result, the next solution found is typically close to the previous. Another drawback of using this “blocking” method to sample is that we cannot parallelize the sampling.

For this reason, we sample from NS and IS using hash-based sampling (cf., [Zhou et al. \[2018\]](#)). Specifically, we sample a limited number of solutions by adding a random universal hashing constraint to the formula given to the solver. Due to the hash function’s universality, we can run multiple samplers in parallel to generate a large number of uniformly distributed solutions. In most cases, the sizes of the sampled sets \hat{NS} and \hat{IS} differ either due to differences in the sizes of NS and IS or due to the solving difficulty of one set compared to the other. We associate a sample weight to each element so the weight of each set is equal in the training process described below.

4.3 Interpretation through a Rule-Based Method

Given \hat{NS} and \hat{IS} —i.e., $\langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle$ tuples labeled according to whether they illustrate noninterference or interference—we could train an interpretable machine-learning model and then extract rules to explain to the user what gives rise to interference. A natural such model to consider is a decision tree. In a decision tree, each decision node (i.e., interior node) is a predicate on features of a $\langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle$ tuple, and its two children correspond to a true or false evaluation of this predicate on a tuple, respectively. A $\langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle$ tuple is classified by traversing the tree from its root, following the branch from each decision node corresponding to the result of evaluating the predicate at that node on the tuple. Each leaf is labeled with an estimate of the probability that a tuple constrained by the predicates’ evaluations from the root to that leaf is in IS . We will discuss what features we include in the process of building decision trees in Sec. 4.4, but an example might be individual variables (e.g., $cvar$).

A single decision tree can easily grow to be deep and complex, and it can miss some useful combinations of predicates since each decision predicate is highly influenced by the splits above it in the tree. To make the decision tree model more powerful in finding useful predicates, we used a decision-tree ensemble called gradient boosted trees [[Friedman 2001](#)]. This process produces m trees denoted T_1, \dots, T_m , with associated weights. If we denote by $T_j(\langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle)$ the real number stored at the leaf to which $\langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle$ is assigned by T_j , then the weighted sum of $T_j(\langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle)$ for $j = 1, \dots, m$ is an estimate of the probability that $\langle \vec{c}, \vec{t}, \vec{s}, \vec{s}' \rangle \in IS$.

To interpret tree ensembles, rule-based classifiers (e.g., RuleFit [[Friedman and Popescu 2008](#)], Slipper [[Cohen and Singer 1999](#)], Pre [[Fokkema 2020](#)]) were introduced to bridge the interpretability of a decision tree with the modeling power of a tree ensemble. Our toolchain leverages SKOPE-RULES (<https://skope-rules.readthedocs.io/>) to generate logical rules from the tree ensemble. Specifically,

consider any path from the root to a leaf in a tree T_j , and let $\pi_{j,1}, \dots, \pi_{j,\ell}$ denote the predicates along that path that evaluated to true. So, for example, if the first predicate encountered in T_j , say “ $\vec{c}(cvar) = 1$ ”, evaluated to false, then $\pi_{j,1} = \vec{c}(cvar) \neq 1$. Then, SKOPE-RULES constructs a rule by conjoining $\pi_{j,1}, \dots, \pi_{j,\ell}$, with the caveat that it limits the number of predicates included in any rule by heuristically pruning them.

Each such rule has a *precision* and *recall*, which we evaluate using a validation set held out from \hat{IS} and \hat{NS} during training. That is, the *recall* of a rule is the fraction of validation samples held out from \hat{IS} for which the rule evaluates to true, and its *precision* is the fraction of validation samples (from \hat{IS} or \hat{NS}) for which the rule evaluates to true that were held out from \hat{IS} . We further prune rules by iteratively removing conjuncts from a long rule if the precision of the resulting rule is at least 95% of the original. We then rank order rules according first to precision, and then according to recall.

4.4 Feature Engineering

The utility of the rule generation described in the previous section depends critically on the features of each $\langle \vec{c}, \vec{i}, \vec{s}, \vec{s}' \rangle$ tuple exposed when training the tree ensemble, from which the predicates making up the decision nodes of each tree are formed. One factor that makes feature engineering especially critical here is that the SAT solver used to produce elements of \hat{IS} and \hat{NS} requires that the conditions defining IS and NS (i.e., conditions (9) and (10))

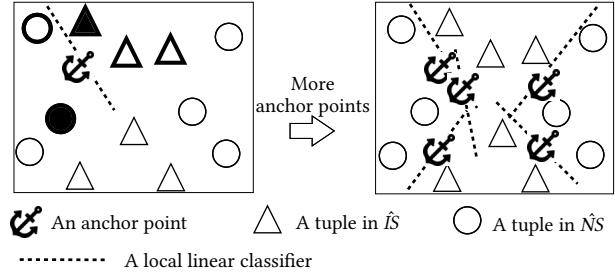


Fig. 3. Finding linear combinations of features near anchor points

be presented to the SAT solver in terms of binary variables only. As such, each solution generated by the SAT solver is expressed as an assignment to these binary variables. While for some hardware logic, a binary representation of the relevant variables is most natural, for other types of logic (e.g., on integers), it is not. For this reason, we augment each binary solution returned by the SAT solver (i.e., each $\langle \vec{c}, \vec{i}, \vec{s}, \vec{s}' \rangle$ tuple) with additional features.

- **Type-aware features:** First, we reconstruct features in a type-aware way from their binary representations. For example, if a variable was initially an integer before being reduced to a collection of binary variables in the formula presented to the SAT solver, we recover the integer value from the bit-vector solution and include it as a feature on which the tree ensemble can trained. With such type-aware features, predicates such as, e.g., $\vec{s}(svar) < 15$ can be learned in a search for simple predicates testing only a single feature, i.e., unary predicates.
- **Symmetric features:** Due to the symmetry of \vec{s} and \vec{s}' , an interference rule could be trivially transformed to another valid interference rule by exchanging \vec{s} and \vec{s}' . For example, when a rule is $\vec{s}(svar)[0] = 0 \wedge \vec{s}'(svar)[0] = 1$, there must be a rule $\vec{s}(svar)[0] = 1 \wedge \vec{s}'(svar)[0] = 0$. Thus, we create $|\vec{s}(svar)[i] - \vec{s}'(svar)[i]|$ for each bit i in $svar$.
- **Linear combinations of multiple variables:** Unary predicates will be unable to naturally capture some relationships resulting in leakage. For example, if leakage happens only when $\vec{s}(svar) > \vec{c}(cvar)$, permitting only unary predicates will result in a boundary characterized point-by-point, e.g., “ $\vec{s}(svar) \geq \theta \wedge \vec{c}(cvar) < \theta$ ” where $\theta = 1, 2, \dots$. We thus expanded our feature

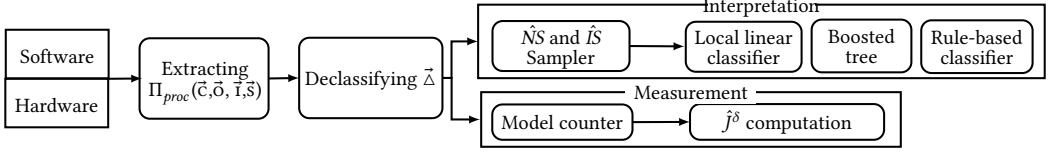


Fig. 4. DINOme workflow

set to permit linear combinations of some features (e.g., $\vec{s}(svar) - \vec{c}(cvar)$), chosen by a linear classifier.

To accommodate branching in the procedure that results in discontinuities in the boundary between sample sets \hat{IS} and \hat{NS} , we opted for a *local* linear classifier (e.g., Fan [1993]; Ribeiro et al. [2018]). That is, we pick *anchor points*, around each of which we train a local classifier that best separates the *nearby* samples in \hat{IS} and \hat{NS} . (See Fig. 3.) To select anchor points, we first find pairs of $\langle \vec{c}, \vec{i}, \vec{s}, \vec{s}' \rangle$ tuples, one from \hat{IS} and one from \hat{NS} , that are *neighbors* in one feature (i.e., after ranking all tuples by this feature, the pair are adjacent in the ranking) and then take the pair’s *midpoint* tuple as their per-feature means. We select anchors uniformly at random from these midpoints. For each anchor, we train a linear classifier using the tuples in \hat{IS} and \hat{NS} that are within a threshold Euclidean distance from the anchor. The linear combination of features used in this linear classifier is then added as another feature to each $\langle \vec{c}, \vec{i}, \vec{s}, \vec{s}' \rangle$ tuple.

5 IMPLEMENTATION

We developed DINOme³ for evaluating and interpreting leakage, described in Sec. 3–4, with an eye toward applying it to evaluate and understand leakage from hardware designs. Though our declassification and interpretation methodologies are not limited to hardware designs, we believe they will be most useful in complicated cases where developers need to understand the interactions between low-level and high-level code. To capture such cases, we define the procedure *proc* to be a hardware design, say written in Verilog, in its initial state but with a predefined program stored in its memory. DINOme enables the user to annotate the configuration by marking components of the hardware state as attacker-controlled (i.e., in $Vars_{\vec{c}}$), attacker-observable (in $Vars_{\vec{o}}$), secret (in $Vars_{\vec{s}}$), or otherwise unknown to the attacker (in $Vars_{\vec{i}}$). DINOme workflow for analyzing this “procedure” is illustrated in Fig. 4. Our system converts this “procedure,” which we continue to denote *proc*, to a cycle-accurate logical formula Π_{proc} that characterizes hardware execution of the program and that relates \vec{c} , \vec{o} , \vec{i} , and \vec{s} . The user can also declare a declassification function δ that operates on the hardware state of the system (we will give examples below), from which DINOme similarly produces a logical formula Π_{δ} that characterizes how the declassified information $\vec{\Delta}$ relates to inputs \vec{c} , \vec{i} , and \vec{s} in the execution of *proc*. From Π_{proc} and Π_{δ} DINOme generates \hat{J}_n^{δ} for varying n (see (8)) and, if requested, sample sets \hat{IS} and \hat{NS} from IS (see (9)) and NS (see (10)), respectively. These sets seed the generation of the rules for interpreting leakage, as discussed in Sec. 4.

Below we discuss particular challenges we encountered when building DINOme and how we overcame them. We focus on how to extract $\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s})$ in Sec. 5.1. In Sec. 5.2, we describe how we calculate our interference measure using projected model counting. Finally, we discuss our technique for sampling to create \hat{IS} and \hat{NS} in Sec. 5.3.

³<https://github.com/DINOme-Project/DINOme>

5.1 Extracting $\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s})$

To analyze the leakage from $proc$, we need an accurate postcondition $\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s})$ for $proc$. In practice, generating a postcondition for an arbitrary procedure is not trivial. Especially here, where our concern is detecting leakage from a processor implementation when running an application—i.e., the procedure $proc$ includes numerous cycles of a cycle-accurate implementation of the processor logic as well as the software logic—the postcondition will be quite large.

Our general strategy to construct $\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s})$ in these circumstances is to assemble it one cycle at a time. Yosys [Wolf [n.d.]] provides a framework to convert the Verilog code for a processor design to its internal register-transfer level (RTL) intermediate language, optimize or modify the design using a series of passes, and finally translate the design to targeted formula through its back-end pass. The SMT2 back-end pass defines a data structure for each hardware module representing the module’s temporary hardware state, a function to implement the module’s state transition from one cycle to the next, and an initialization function to initialize the module’s state. To incorporate the software logic of $proc$, we compile the software to its hardware-readable assembly and load the assembly into the instruction memory unit.

To mark the symbolic variables, the analyst defines a configuration file to mark as symbolic each input parameter of $proc$ (in this case, $svar$, $ivar$, and $cvar$), which can be a software variable located at a fixed location in the memory unit or a wire/register inside the hardware module. Our modified SMT2 backend pass in Yosys then tracks the constraints associated with this symbolic data throughout a cycle execution. Specifically, it outputs a logical postcondition $\tau_{proc}(\vec{h}^{t-1}, \vec{h}^t)$ that relates fully symbolized hardware state $\vec{h}^{t-1} : Vars_{\vec{h}} \rightarrow Vals_{\vec{h}}$ at the end of cycle $t - 1$ to the hardware state \vec{h}^t that results from executing cycle t . Since the hardware state includes memory units, registers, etc., $\tau_{proc}(\vec{h}^{t-1}, \vec{h}^t)$ with fully symbolized \vec{h}^{t-1} is too large to naively extend to cover multiple cycles. We also use the pass to generate initialization logic $\Psi_{proc}^0(\vec{c}, \vec{i}, \vec{s}, \vec{h}^0)$ that concretely characterizes the first-cycle starting state \vec{h}^0 (upon a reset) except for the configured symbolic inputs $svar$, $ivar$, and $cvar$.

Using the transition logic, we construct a cycle-accurate postcondition Ψ_{proc}^T representing the logic between symbolic inputs and its internal hardware state one cycle at a time, leveraging the entire hardware state as an “observable” output of the cycle.

$$\Psi_{proc}^T(\vec{c}, \vec{i}, \vec{s}, \vec{h}^T) \leftarrow \Psi_{proc}^0(\vec{c}, \vec{i}, \vec{s}, \vec{h}^0) \wedge \bigwedge_{t=1}^T \tau_{proc}(\vec{h}^{t-1}, \vec{h}^t)$$

We finally define $\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s})$ by defining \vec{o} in terms of the sequence of hardware states $\langle \vec{h}^t \rangle_{t=0}^T$ using a formula $\Gamma(\langle \vec{h}^t \rangle_{t=0}^T, \vec{o})$.

$$\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s}) \leftarrow \Psi_{proc}^T(\vec{c}, \vec{i}, \vec{s}, \vec{h}^T) \wedge \Gamma(\langle \vec{h}^t \rangle_{t=0}^T, \vec{o}) \quad (11)$$

For example, in cache-based side channels, the observable parameters are whether there is a cache hit/miss during the execution, which is constructed using the values of the $s2_hit$ register across the execution (as demonstrated in Sec. 6.3).

Applying a correct combination of techniques to simplify $\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s})$ is critical to scaling the sampling of IS and NS to create \hat{IS} and \hat{NS} and to count $\left| \tilde{X}_{S,S'}^\delta \right|$ and $\left| \tilde{X}_{S,S'}^\delta \right|$ to compute \hat{J}_n^δ . See Zhou [2020] for a discussion of these simplifications.

To correctly measure leakage, the postcondition for $proc$ must be complete and sound. Completeness means that if $\langle \vec{c}, \vec{i}, \vec{s}, \vec{o} \rangle$ is feasible for $proc$, then $\langle \vec{c}, \vec{i}, \vec{s}, \vec{o} \rangle$ satisfies $\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s})$. Soundness means that if $\langle \vec{c}, \vec{i}, \vec{s}, \vec{o} \rangle$ is infeasible for $proc$, then $\langle \vec{c}, \vec{i}, \vec{s}, \vec{o} \rangle$ does not satisfy $\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s})$. Here,

$\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s})$ is derived from the hardware transition logic τ_{proc} . Since τ_{proc} represents how the next hardware state is derived from the previous hardware state⁴ and is derived from the actual hardware design, our postcondition is consistent with the real verilog code, provided that the Yosys SMT2 backend pass is correct.

In our experiments, we selected T to ensure the termination of the execution, based on our knowledge gained by studying the CPU. A more conservative method would be to track the CPU pipeline and call the SAT solver each cycle to check whether the last instruction has certainly committed. We have confirmed that adding more cycles after the termination of the execution does not affect Π_{proc} meaningfully, as the additional cycles do not process any valid opcodes and so only trivially change the hardware state.

5.2 Measurement with Declassification using Projected Model Counting

Using CryptoMiniSAT 5.0 as the basic solver, we implemented a counter to estimate the numerator and the denominator in the measurement $\hat{J}^\delta(S, S')$ in (8).

5.2.1 Computing $\hat{J}^\delta(S, S')$. To compute $\hat{J}^\delta(S, S')$, we need to count the sizes of $\tilde{X}_{S, S'}^\delta$ and $\check{X}_{S, S'}^\delta$. Directly counting $\tilde{X}_{S, S'}^\delta$ is not easy as the set difference operation introduces a “forall” quantifier. Fortunately, since $|\tilde{X}_{S, S'}^\delta| = |\check{X}_{S, S'}^\delta| - |\hat{X}_{S, S'}^\delta|$, it suffices to count $\check{X}_{S, S'}^\delta$ and $\hat{X}_{S, S'}^\delta$ for each sample pair S, S' . Intuitively, counting $\check{X}_{S, S'}^\delta$ could be expressed as a projected model counting task [Aziz et al. 2015] over $\langle \vec{c}, \vec{o}, \vec{i}, \vec{s} \rangle$ in a quantifier-free SAT problem with two copies of Π_{proc} shown in \check{F} below. \check{F} is translated to a CNF proposition where it uses v bit variables to represent $\langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle$ and others to represent $\langle \vec{s}, \vec{s}', \vec{i}, \vec{i}' \rangle$ and auxiliary variables.

$$\begin{aligned} \check{F} \leftarrow & \left(\Pi_{proc}(\vec{c}, \vec{o}, \vec{i}, \vec{s}) \vee \Pi_{proc}(\vec{c}, \vec{o}, \vec{i}', \vec{s}') \right) \wedge \Pi_g(\vec{c}, \vec{\Delta}, \vec{i}, \vec{s}) \wedge \Pi_g(\vec{c}, \vec{\Delta}, \vec{i}', \vec{s}') \\ & \wedge \left((\vec{s}(svar) \in S \wedge \vec{s}'(svar) \in S') \vee (\vec{s}'(svar) \in S \wedge \vec{s}(svar) \in S') \right) \end{aligned} \quad (12)$$

Following Zhou et al. [2018], two random, disjoint sets S and S' of expected size n are specified with distinct strings $p, \hat{p} \in \{0, 1\}^b$ where $n = |\mathbb{S}|/2^b$ for \mathbb{S} being the domain of all possible secret values, and specifically with the constraint that for a fixed hash function, the hash of each $s \in S$ is p and the hash of each $s' \in S'$ is \hat{p} .

For $\hat{X}_{S, S'}^\delta$, we can define another projected model counting task over $\langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle$ in a quantifier-free SAT problem \hat{F} shown below. \hat{F} uses the logical postcondition Π_{proc} twice, where the first copy is for the execution with a secret $\vec{s}(svar) \in S$ and the second checks for existence of a secret $\vec{s}'(svar) \in S'$ leading to a result \vec{o} also possible with \vec{s} . \hat{F} also checks the existence of some secret (denoted by $\vec{s}''(svar)$) in the secret set S' leading to the equivalent declassification value $\vec{\Delta}$ so that we can ensure the \vec{s} and \vec{s}' cannot be distinguished by $\vec{\Delta}$.

$$\begin{aligned} \hat{F} \leftarrow & \Pi_{proc, \delta}(\vec{c}, \vec{o}, \vec{\Delta}, \vec{i}, \vec{s}) \wedge \vec{s}(svar) \in S \\ & \wedge \Pi_{proc}(\vec{c}, \vec{o}, \vec{i}', \vec{s}') \wedge \vec{s}'(svar) \in S' \\ & \wedge \Pi_g(\vec{c}, \vec{\Delta}, \vec{i}'', \vec{s}'') \wedge \vec{s}''(svar) \in S' \end{aligned} \quad (13)$$

5.2.2 Optimizations for Counting $\tilde{X}_{S, S'}^\delta$ and $\check{X}_{S, S'}^\delta$. Enumerating all solutions to (12) and (13) using a solver is intractable. To estimate the number of solutions to each instead, we used the approximate model counting technique due to Chakraborty et al. [2013], specifically the approach taken by Soos

⁴Unlike software, hardware code (e.g., verilog) does not use do-while loops within one cycle for which the number of iterations is determined dynamically. In our case studies, we found that the one-cycle logic for BOOM is correspondingly simple, enabling the completeness and soundness of τ_{proc} .

$$\begin{aligned} &\text{E-Solver with } H \text{ and } p \text{ generates } \langle \vec{c}, \vec{i}, \vec{s}, \vec{s}', \vec{o}, \vec{\Delta} \rangle \text{ satisfying} \\ &\Pi_{proc,\delta}(\vec{c}, \vec{o}, \vec{\Delta}, \vec{i}, \vec{s}) \wedge \Pi_{proc,\delta}(\vec{c}, \vec{o}', \vec{\Delta}, \vec{i}', \vec{s}') \wedge \vec{o} \neq \vec{o}' \wedge H(\langle \vec{c}, \vec{i}, \vec{s}, \vec{s}' \rangle) = p \end{aligned} \quad (16)$$

$$\begin{aligned} &\text{F-Solver cancels } \langle \vec{c}, \vec{i}, \vec{s}, \vec{s}', \vec{o}, \vec{\Delta} \rangle \text{ satisfying (16) if there is some } \vec{i}'' \text{ satisfying} \\ &\Pi_{proc,\delta}(\vec{c}, \vec{o}, \vec{\Delta}, \vec{i}'', \vec{s}') \end{aligned} \quad (17)$$

Fig. 5. Generating examples in $\hat{I}\hat{S}$ using EF-solver

and Meel [2019]. That is, by specifying a randomly selected hash function $\hat{H}^{\hat{b}} : \{0, 1\}^v \rightarrow \{0, 1\}^{\hat{b}}$ and an output $\hat{p} \in \{0, 1\}^{\hat{b}}$ as an additional constraint, we can estimate $|\hat{X}_{S,S'}^{\delta}|$ using the average value of multiple estimations of $|\hat{Z}_{S,S'}^{\hat{p}}|$ with some error ϵ and confidence δ (i.e., $|\hat{X}_{S,S'}^{\delta}| \approx |\hat{Z}_{S,S'}^{\hat{p}}| \times 2^{\hat{b}}$). Similarly, we could estimate $|\check{X}_{S,S'}^{\delta}|$ using $\check{Z}_{S,S'}^{\check{p}}$.

$$\hat{Z}_{S,S'}^{\hat{p}} = \left\{ \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \mid \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \in \hat{X}_{S,S'}^{\delta} \wedge \hat{H}^{\hat{b}}(\langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle) = \hat{p} \right\} \quad (14)$$

$$\check{Z}_{S,S'}^{\check{p}} = \left\{ \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \mid \langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle \in \check{X}_{S,S'}^{\delta} \wedge \check{H}^{\check{b}}(\langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle) = \check{p} \right\} \quad (15)$$

This optimization for model counting will limit the number of calls to the SAT solver by constraining the number of solutions available, and thus make the counting more scalable for large set size. Thus, $\hat{J}^{\delta}(S, S')$ is estimated using the average value of $1 - \frac{|\hat{Z}_{S,S'}^{\hat{p}}|}{|\check{Z}_{S,S'}^{\check{p}}|}$ for various \hat{p}, \check{p} .

Our primary departure from the implementation by Soos and Meel [2019] lies in utilizing task-specific properties in our counting tasks to reduce redundant effort in solution searching. Specifically, since $\hat{X}_{S,S'}^{\delta} \subseteq \check{X}_{S,S'}^{\delta}$, we ensure that $\hat{X}_{S,S'}^{\delta} \cap \check{Z}_{S,S'}^{\check{p}} \subseteq \hat{Z}_{S,S'}^{\hat{p}}$ in our counting by defining $\hat{H}^{\hat{b}}(\langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle)$ to be the \hat{b} -bit prefix of $\check{H}^{\check{b}}(\langle \vec{c}, \vec{o}, \vec{\Delta}, \vec{i} \rangle)$ for $\hat{b} \leq \check{b}$. Then once we have generated solutions in $\check{Z}_{S,S'}^{\check{p}}$, we speed up finding solutions in $\hat{Z}_{S,S'}^{\hat{p}}$ for $\hat{b} = \check{b}$ (and so $\hat{p} = \check{p}$) by first checking each solution in $\check{Z}_{S,S'}^{\check{p}}$ to see if it satisfies \hat{F} (i.e., is in $\hat{X}_{S,S'}^{\delta} \cap \check{Z}_{S,S'}^{\check{p}}$). Only if insufficient solutions are found with $\hat{b} = \check{b}$ is \hat{b} reduced and the solver used to generate additional solutions in $\hat{Z}_{S,S'}^{\hat{p}}$ for \hat{p} a \hat{b} -bit prefix of \check{p} .

In Sec. 6, we set the error $\epsilon = 0.4$ and confidence $\delta = 0.9$ in this method to estimate the sizes of $\check{X}_{S,S'}^{\delta}$ and $\hat{X}_{S,S'}^{\delta}$, from which $\hat{J}^{\delta}(S, S')$ is estimated using (8). For each set size n , we compute \hat{J}_n^{δ} using ≥ 100 hash functions, i.e., implicit selection of pairs S, S' of expected size n .

5.3 Sampling $\hat{N}\hat{S}$ and $\hat{I}\hat{S}$ for Interpretable Learning

Similar to the counting process, to construct $\hat{N}\hat{S}$ and $\hat{I}\hat{S}$, the sampler will select hash functions H randomly from a family and output values p randomly from its range to solve for tuples $\langle \vec{c}, \vec{i}, \vec{s}, \vec{s}' \rangle$ for which $H(\langle \vec{c}, \vec{i}, \vec{s}, \vec{s}' \rangle) = p$ (and are in $\hat{N}\hat{S}$ or $\hat{I}\hat{S}$, respectively). In the following experiments, we will generate up to 100,000 solutions for each of $\hat{N}\hat{S}$ and $\hat{I}\hat{S}$, where 70% used for training and 30% used for validation.

We cannot directly encode set difference, used in (9) and (10), using an equivalent quantifier-free formula. To implement a sampler to generate solutions in the set difference, we will use one solver (“E-Solver”) to search for candidate solutions and another (“F-Solver”) cancel candidates; this is a commonly used algorithm for an SMT solver to solve exist-forall problems (e.g., see Dutertre [2015]).

Here, we will illustrate sampling IS , while sampling NS is similar. The sampler first uses the E-Solver to generate feasible solutions $\langle \bar{c}, \bar{i}, \bar{s}, \bar{s}' \rangle$ (see (16)) that guarantee, for an attacker's chosen \bar{c} , the observable value \bar{o} derived from \bar{s} with \bar{i} could be different from an observable \bar{o}' generated by \bar{s}' with some \bar{i}' when the declassified value $\bar{\Delta}$ is the same. However, it does not guarantee the \bar{o} is never feasible for \bar{s} . To further test whether the $\langle \bar{c}, \bar{i}, \bar{s}, \bar{s}' \rangle$ is in \hat{IS} , we use the F-Solver to test whether $\langle \bar{s}', \bar{i}'' \rangle$ for some \bar{i}'' could generate \bar{o} with $\langle \bar{s}, \bar{i} \rangle$ when they share the declassification value $\bar{\Delta}$, to check whether we need to cancel the solution. That is, $\langle \bar{c}, \bar{i}, \bar{s}, \bar{s}' \rangle$ satisfying (16) but not (17) will be included in \hat{IS} .

After generating enough $\langle \bar{c}, \bar{i}, \bar{s}, \bar{s}' \rangle$ tuples in \hat{NS} and \hat{IS} , the interpretation module trains local support vector machine (SVM) classifiers [Fan et al. 2008] around each of 50 anchor points, after ruling out data whose normalized Euclidean distance (i.e., after scaling each attribute to a value between 0 and 1, use Euclidean distance divided by the number of attributes) is more than 0.2 from the anchor. Then a logistic regression model for NS and IS is learned using a gradient boosted tree implementation *xgboost* [Chen and Guestrin 2016]. To generate the interpretable models, we implemented the rule learner using SKOPE-RULES.

6 CASE STUDIES

In this section, we illustrate DINoME by describing its application to the BOOM core (<https://github.com/riscv-boom/riscv-boom>), an open-source RISC-V core that is susceptible to cache-based side channels and SPECTRE attacks. The goal of these case studies is to illustrate our methodology and to show how it can be useful to system analysts. Our method is also applicable to other side channels, not only cache-based ones. Analysts can specify the secret to protect and define their side channels using attacker-controlled and attacker-observable variables but, critically, not the specific attacker algorithm.

- We applied DINoME to evaluate cache-based side-channel leakage due to secret-dependent memory accesses. With different BOOM configurations (i.e., number of cache ways w and whether to share memory), the case studies show how \hat{J}_n^δ curves reveal the effects of the configurations on the leakage. We also implemented and evaluated two possible mitigations, namely SCATTER-CACHE [Werner et al. 2019] and PHANTOMCACHE [Tan et al. 2020], which reduce but do not eliminate the cache leakage. Our measurements using \hat{J}_n^δ illustrate which mitigation is better for a specific BOOM setting.
- We used DINoME to assess leakage via cache-based side channels from a modular exponentiation function commonly used in cryptographic algorithms. The rule-based interpretation explains how to choose attacker-controlled variables and which portion of the secret is leaked.
- We evaluated software code snippets causing speculative execution. This case study demonstrates how to use declassification to focus on leakage caused by speculative execution (i.e., by declassifying other leakage to reveal it) and how to generate an efficient interpretable rule set. We found that some software with a short speculation window is insufficient to cause memory leakage in the latest version of BOOM.

6.1 BOOM Configurations

In the following experiments, we used pocket-size hardware modules to replace the modules in the BOOM v2.2.3 configuration. A simplified diagram is shown in Fig. 6. Analyzing artificially small but otherwise faithful configurations of a system is not uncommon in model checking, for example (e.g., Ball et al. [2004]; Pnueli et al. [2002]). Specifically, we set the cache line size to $bbytes = 64B$ and the total L1 data cache size to 1KB (16 cache lines in total). We then varied the cache ways w and sets c (i.e., subject to $w \times c = 16$) in Sec. 6.3 but used a fixed setting $c = 2$, $c = 8$ for other

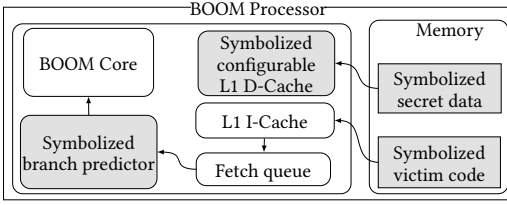


Fig. 6. BOOM configuration

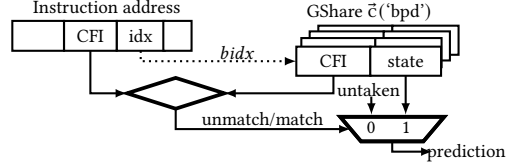


Fig. 7. GShare branch predictor's logical architecture

evaluations. BOOM only provides a configurable associative L1 cache module using a random replacement policy. To compare different cache designs, we implemented two side-channel-resistant cache modules, as described in Sec. 6.4. For the main memory, we set the memory size to 4KB and thus a memory address is only 12 bits. For evaluation purposes, we used the upper half of the memory address space as instruction memory and the lower half as data memory. To simplify the following analysis, we removed the page table walker module and assumed virtual addresses were the same as physical addresses. For the instruction fetch, we set the fetch width to 4 and configured the L1 instruction cache to a 1KB, 8-set, 2-way cache with a customized prefetching module that preloaded the software workload at the first cycle.

One feature of BOOM is that it supports speculative execution, with which we will experiment in Sec. 6.6. Speculative execution leverages a *branch predictor*, for which we used the GShare branch predictor. The logical structure of GShare is shown in Fig. 7. When a prediction request arrives for a branch instruction, the GShare predictor derives a value $bidx$ from the certain bits (denoted 'idx' in Fig. 7) in the instruction address and an instruction history register and then uses $bidx$ to index into a table to which we refer as 'bpd'. Each entry of the 'bpd' table includes a label called 'CFI' and a 2-bit 'state', of which one bit indicates whether the entry holds a strong or weak prediction and the other bit holds that prediction (i.e., whether the branch will be taken or not). If the $\text{bpd}\{bidx\}.\text{CFI}$ value matches the 'CFI' portion of the instruction address, then the predictor uses the $\text{bpd}\{bidx\}.\text{state}$ value to make a branch prediction. The GShare predictor will globally tune entries based on executions in any user's domain. Thus, an attacker can easily affect the 'bpd' table before victim's execution, and so we include 'bpd' in $\text{Vars}_{\bar{c}}$. In our evaluation, we fix the number of 'bpd' entries to 4 so that only 2 bits in the instruction address are used as 'idx' while another 2 bits ($=\log_2(\text{fetch width})$) are used as its 'CFI' label.

In the following case studies, we added the 'bpd' table in the GShare module to $\text{Vars}_{\bar{c}}$ and registers in the L1 data cache module including the cache metadata, the replacement state (i.e., the linear-feedback shift register (LFSR) for the random replacement policy), and the memory-to-cache mapping (if using a nonfixed mapping) to $\text{Vars}_{\bar{r}}$.

In cache-based side channel attacks, \bar{c} and \bar{o} are not directly represented in the hardware state or in victim's code, and so it is necessary to define them through an adversary model. We assume that the adversary has access to 16 memory blocks $block_1, block_2, \dots, block_\ell, \dots, block_{16}$ aligned to cache lines, which is sufficient to control the cache as our L1 data cache consists of only 16 cache lines in our experiments. Specifically, $\bar{c}(\text{'load'})[\ell]$ indicates whether the adversary loads (1) or flushes (0) $block_\ell$, while $\bar{o}(\text{'hit'})[\ell]$ indicates whether the adversary observes a cache hit (1) or miss (0) when accessing $block_\ell$. The following section illustrates how to automatically construct these.

6.2 Defining \bar{c} and \bar{o} for Cache-Based Side Channels

The most common cache-based side-channel attacks are PRIME+PROBE, FLUSH+RELOAD, and their variants (e.g., see Yarom and Falkner [2014]; Zhang et al. [2012]). In a PRIME+PROBE attack, the attacker loads memory blocks to fill (PRIME) cache sets, permits the victim computation to run

for a PRIME+PROBE interval, and then reads (PROBES) these same blocks to determine which were evicted by the victim computation during the PRIME+PROBE interval. In a FLUSH+RELOAD attack, the attacker FLUSHes a shared-memory block from cache and then, after a FLUSH+RELOAD interval, accesses (RELOADS) the block to determine whether the block was brought back into the cache by the victim computation.

To model side channel attacks in our framework, it is necessary to model the effects on the cache of the phases before victim execution (the PRIME and FLUSH steps) and to define $\vec{\sigma}$ to include the results of the phases after victim execution (the PROBE and RELOAD steps). To do so, we assume that the adversary has access to memory blocks $block_1, block_2, \dots, block_m$ aligned to cache lines, and we define the RISC-V assembly routine *acc* (see above) by which the adversary can access the block with index $\ell = \hat{c}$ ('blockIdx') and empty \hat{s} .

```
acc( $\hat{c}, \hat{i}, \hat{s}$ )
li s0, 0x2000000
add s1, s0,  $\ell$ 
sll s1, s1, 6
lbu a2, 0(s1)
```

Starting from hardware state $\hat{H}_\ell^0 (= \hat{i})$ that is completely symbolic, we generate the per-cycle logical postcondition $\tau_{acc}(\hat{H}_\ell^{t-1}, \hat{H}_\ell^t)$ for each $0 < t \leq \hat{T}$ as in Sec. 5.1, where we empirically choose $\hat{T} = 45$.

We use these postconditions in two ways. First, we use them to extract a constraint $\Gamma(\langle \hat{H}^t \rangle_{t=1}^T, \vec{\sigma})$ that defines the attacker's observations $\vec{\sigma}$ in terms of the hardware states $\langle \hat{H}^t \rangle_{t=1}^T$ induced by the execution (see (11)). A naive attempt to do so would be to simply include in $\vec{\sigma}$ the metadata for each cache line at every step of the execution. However, this would grant too much power to an attacker, who should not be given access to the tag values and the exact locations of blocks inside a set. Instead, we permit only a weaker attacker (cf., abstract noninterference [Giacobazzi and Mastroeni 2004]) by defining the constraint $\Gamma(\langle \hat{H}^t \rangle_{t=1}^T, \vec{\sigma})$ that represents the view of cache hits and misses immediately observable by the adversary, by:

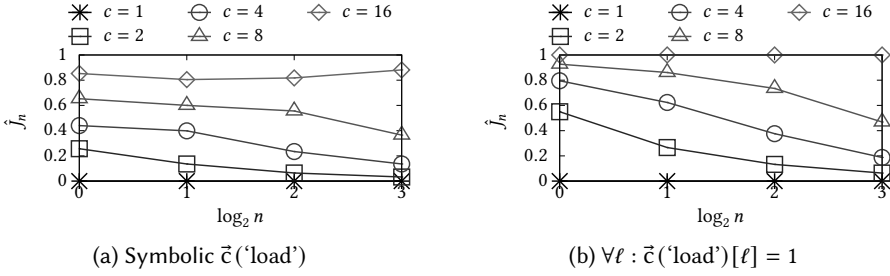
$$\vec{\sigma}(\text{'hit'})[\ell] = \left(\begin{array}{l} (\hat{H}_\ell^0 = \vec{H}^T) \wedge \left(\bigwedge_{t=1}^{\hat{T}} \tau_{acc}(\hat{H}_\ell^{t-1}, \hat{H}_\ell^t) \right) \\ \wedge \left(1 - \bigvee_{t=0}^{\hat{T}} \text{CACHEMISS}(\hat{H}_\ell^t, block_\ell) \right) \end{array} \right)$$

for $\ell = \hat{c}$ ('blockIdx'). Here, CACHEMISS is a BOOM-defined Verilog code snippet that, intuitively, checks a set of cache lines where $block_\ell$ might reside and returns 1 (in a register called `s2_hits`) if none of those cache lines has a valid tag matched with $block_\ell$ (and returns 0 otherwise). In this way, we characterize the procedure *acc* using a logical postcondition without manually modeling CACHEMISS.

Second, we permit the attacker to control which of its blocks are loaded into the cache before the victim runs. Specifically, the predicate $\Psi_{proc}^0(\vec{c}, \vec{i}, \vec{s}, \vec{H}^0)$ that controls the initial hardware state from which the victim executes is modified to constrain which of the attacker's blocks are present in cache, as communicated through a reserved variable 'load' $\in \text{Vars}_{\vec{c}}$, for which the \vec{c} ('load') is a bit vector of length m . That is, attacker block $block_\ell$ should be loaded before the victim runs if and only if \vec{c} ('load')[ℓ] = 1. To effect this in $\Psi_{proc}^0(\vec{c}, \vec{i}, \vec{s}, \vec{H}^0)$, we construct $\Psi_{proc}^0(\vec{c}, \vec{i}, \vec{s}, \vec{H}^0)$ to include

$$\vec{c}(\text{'load'})[\ell] = \left(\begin{array}{l} (\hat{H}_\ell^0 = \vec{H}^0) \wedge \left(\bigwedge_{t=1}^{\hat{T}} \tau_{acc}(\hat{H}_\ell^{t-1}, \hat{H}_\ell^t) \right) \\ \wedge \left(1 - \bigvee_{t=0}^{\hat{T}} \text{CACHEMISS}(\hat{H}_\ell^t, block_\ell) \right) \end{array} \right)$$

Of course, we rename variables to ensure no conflicts between copies of \hat{H}_ℓ^t included within the \vec{c} ('load')[ℓ] and $\vec{\sigma}$ ('hit')[ℓ] constraints.

Fig. 8. \hat{J}_n for PRIME+PROBE attacks

6.3 Cache-Based Side Channels

In this section, we evaluate cache-based side channels under different memory isolation and cache configurations.

6.3.1 Without Shared Memory. Here, we target a victim's RISC-V assembly *proc* to access a secret-indexed memory block not shared with the attacker, by setting the base address in $s0$ to a value $0x2000010$, in contrast to the one used in attacker's process *acc* (see Sec. 6.2). We experimented with different numbers of cache sets c including $c = 1$ (i.e., 16-way, 1-set, fully associative), $c = 2$ (i.e., 8-way, 2-set), $c = 4$ (i.e., 4-way, 4-set), $c = 8$ (2-way, 8-set), and $c = 16$ (i.e., 1-way, 16-set, direct-mapped). As shown in Fig. 8(a), \hat{J}_n increases when the number of sets increases. Specifically, there is no leakage ($\hat{J}_n = 0$ for all n) when $c = 1$. Using fewer cache sets, each cache set is shared by more memory blocks, and so an attacker will have more difficulty distinguishing one execution from others. When $1 < c < 16$, \hat{J}_n decreases as n grows, since the attacker can learn only $\log_2(c)$ bits about the secret and thus may be unable to distinguish secrets in large sets (i.e., large n).

```

proc ( $\vec{c}, \vec{i}, \vec{s}$ )
  li s0, 0x2000010
  add s1, s0,  $\vec{s}$ ('secret')
  sll s1, s1, 6
  lbu a2, 0(s1)

```

An example *interference* rule for *IS* generated as described in Sec. 4 with the highest precision (1.00) and a recall ≈ 0.04 in a 2-way, 8-set cache is:

$$\left\{ \begin{array}{l} \vec{s}(\text{'secret'})[2] \geq 1 \wedge \vec{s}(\text{'secret'})[1] < 1 \\ \wedge \vec{s}(\text{'secret'})[0] \geq 1 \wedge \vec{s}'(\text{'secret'})[0] < 1 \end{array} \right\} \wedge \left\{ \begin{array}{l} \vec{c}(\text{'load'})[5] \geq 1 \\ \wedge \vec{c}(\text{'load'})[13] \geq 1 \end{array} \right\} \quad (18)$$

Our approach could not directly represent $\vec{c}(\text{'load'})[\ell] \equiv \vec{s}(\text{'secret'}) \bmod c$. So, the trees in the model split the dataset based on the cache set index. In this rule, the \vec{s} and \vec{s}' conjuncts concretize the least significant 3 bits of $\vec{s}(\text{'secret'})$ (i.e., $\vec{s}(\text{'secret'}) \equiv 5 \bmod 8$) using $\vec{s}(\text{'secret'})[2] \geq 1 \wedge \vec{s}(\text{'secret'})[1] < 1 \wedge \vec{s}(\text{'secret'})[0] \geq 1$ and the lowest bit of $\vec{s}'(\text{'secret'})$ (i.e., $\vec{s}'(\text{'secret'}) \equiv 0 \bmod 2$) using $\vec{s}'(\text{'secret'})[0] < 1$. The \vec{c} conjuncts are $\vec{c}(\text{'load'})[5] \geq 1$ and $\vec{c}(\text{'load'})[13] \geq 1$; note that $13 \equiv 5 \bmod 8$. That is, an attacker could load all blocks $block_\ell$ with $\ell \equiv 5 \bmod 8$ into cache to distinguish a secret $\vec{s}(\text{'secret'}) \equiv 5 \bmod 8$ from $\vec{s}'(\text{'secret'}) \bmod 8 \in \{0, 2, 4, 6\}$.

There were many other top-ranking rules similar to (18), each focusing on one residue class of the secret value modulo c where $c = 8$ and constraining $\vec{c}(\text{'load'})[\ell] = 1$ for all ℓ with that residue class modulo c . Each such rule works for $\frac{1}{8}$ of $\vec{s}(\text{'secret'})$'s domain and $\frac{1}{2}$ of $\vec{s}'(\text{'secret'})$'s domain, thus only for $\frac{1}{8} \times \frac{1}{2} \approx 0.06$ of secret pairs. The recall rate $0.04 < 0.06$ indicates that priming the corresponding cache set ensures (i.e., precision = 1.0) the interference but is not necessary to cause it.

Analogously, we can generate rules for the *noninterference* set *NS*, as well. One example with precision 1.0 (i.e., that ensures noninterference) and recall 0.11 constrains the secret's least-significant

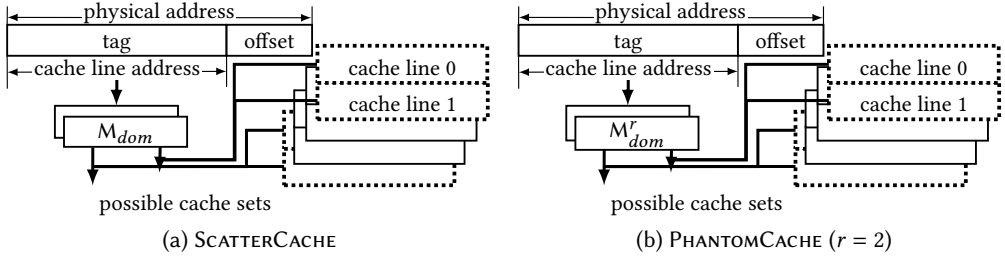


Fig. 10. Cache modules in 2-way, 4-set configure

3 bits to be the same for \vec{s} and \vec{s}' :

$$\begin{aligned} & |\vec{s}(\text{'secret'})[2] - \vec{s}'(\text{'secret'})[2]| < 1 \\ \wedge & |\vec{s}(\text{'secret'})[1] - \vec{s}'(\text{'secret'})[1]| < 1 \\ \wedge & |\vec{s}(\text{'secret'})[0] - \vec{s}'(\text{'secret'})[0]| < 1 \end{aligned} \quad (19)$$

This analysis illustrates that an attacker can easily distinguish $\vec{s}(\text{'secret'})$ and $\vec{s}'(\text{'secret'})$ when priming a cache set used by $\vec{s}(\text{'secret'})$ or $\vec{s}'(\text{'secret'})$ but not both. It is therefore safe to assume that the attacker will PRIME the cache using all its controlled memory blocks to maximize the chances for leakage. The \hat{J}_n measure under this specific attack is shown in Fig. 8(b). The worst case will leak all of the 4-bit secret when using high-granularity memory-to-cache mapping, i.e., where $c = 16$.

6.3.2 With Shared Memory. To evaluate the leakage due to shared memory (i.e., with FLUSH+RELOAD attacks), we allow the attacker to control and observe all memory blocks used by the victim by setting the base to $0x2000000$ in *proc* instead of to $0x2000010$. The \hat{J}_n curves are similar and close to 1 for all settings, indicating that the leakage does not have much correlation with w . An example rule for interference derived using the methodology of Sec. 4, having a precision of 1.0 and recall of ≈ 0.04 , is

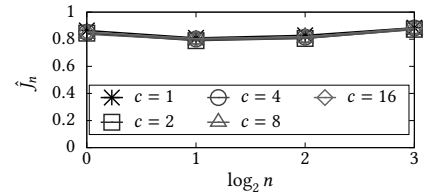
$$\vec{s}'(\text{'secret'}) < 2 \wedge \vec{s}'(\text{'secret'}) \geq 1 \wedge \vec{c}(\text{'load'})[1] < 1 \quad (20)$$

That is, if $\vec{s}'(\text{'secret'}) = 1$ then $\vec{c}(\text{'load'})[1] = 0$ results in interference. Indeed, the other top-ranked rules for this example (not shown) were roughly 32 similar rules, each one setting $\vec{c}(\text{'load'})[\ell] = 0$ for a specific secret value $\vec{s}(\text{'secret'}) = \ell$ or $\vec{s}'(\text{'secret'}) = \ell$. The intuition behind these rules is that an attacker can precisely detect if $\vec{s}(\text{'secret'}) = \ell$ by setting $\vec{c}(\text{'load'})[\ell] = 0$ (i.e., FLUSHING $block_\ell$ so he can later RELOAD it), and similarly for $\vec{s}'(\text{'secret'})$. Going further, if an attacker sets $\vec{c}(\text{'load'})[\ell] = 0$ for all ℓ , he can detect the victim's access to any $block_\ell$, where $\hat{J}_n = 1$ for all n .

6.4 Side-Channel-Resistant Cache Designs

To demonstrate the power of DINOme in comparing different implementations, we evaluate two cache designs for mitigating side channels, namely SCATTERCACHE [Werner et al. 2019] and PHANTOMCACHE [Tan et al. 2020]. Unfortunately, Verilog specifications of these are unavailable, and so we implemented two simplified cache modules (which we continue to refer to as SCATTERCACHE and PHANTOMCACHE) in BOOM following their paper designs.

SCATTERCACHE maps a memory block to a cache line using a cryptographic index derivation function computed using the block's physical address and a private key. As shown in Fig. 10(a), to

Fig. 9. \hat{J}_n for FLUSH+RELOAD attacks with symbolic $\vec{c}(\text{'load'})$

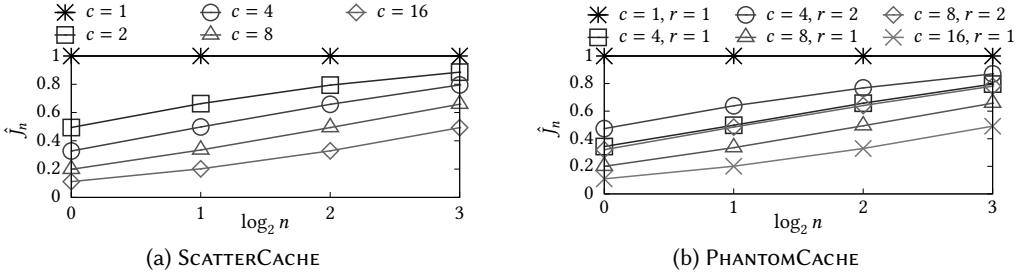


Fig. 11. Memory sharing enabled with $\forall \ell : \bar{c}(\text{'load'})[\ell] = 0$ (FLUSH+RELOAD attack)

simulate this index derivation without choosing a concrete function, we use a symbolic look-up table denoted by M_{dom} per security domain dom ($dom = 0$ denotes the victim's domain and $dom = 1$ denotes the attacker's) to store the mapping from memory address to cache line. For security domain dom , its access to memory contents at physical address $paddr$ and so with block address $baddr = \lfloor paddr/bbytes \rfloor$ is mapped to cache lines with way index k and set index $j = M_{dom}\{baddr\}\{k\}$ for $k = 0, 1, \dots, w - 1$. Similarly, for PHANTOMCACHE, we used a domain-specific memory-to-cache mapping (shown in Fig. 10(b)) represented by M'_{dom} to allow a memory block to use cache lines in up to r cache sets indexed by $M'_{dom}\{baddr\}\{k\}$ for $k = 0, 1, \dots, r$.⁵ In the following evaluation, we have $M_{dom}, M'_{dom} \in \text{Vars}_{\bar{c}}$.

6.4.1 Random Memory-to-Cache Mappings. First, we experimented without memory sharing when assuming the memory-to-cache mapping is completely unknown to the attacker. We ended up with $\hat{J}_n = 0$ for all n in both SCATTERCACHE and PHANTOMCACHE. The attacker cannot tell which memory blocks are accessed by the victim, as a memory block could be mapped to any cache line if the mapping is unknown. Thus, we focused on the leakage analysis when memory sharing is enabled.

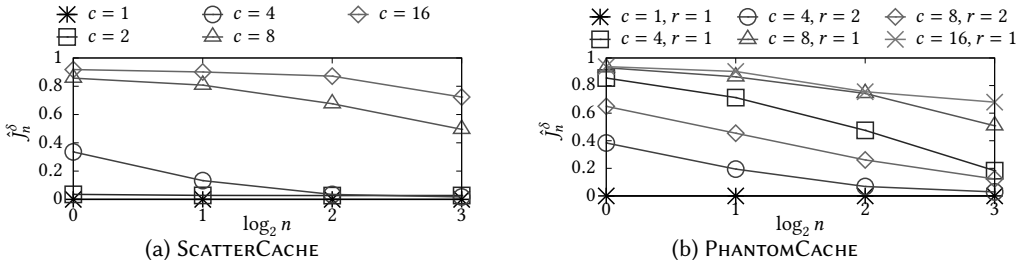
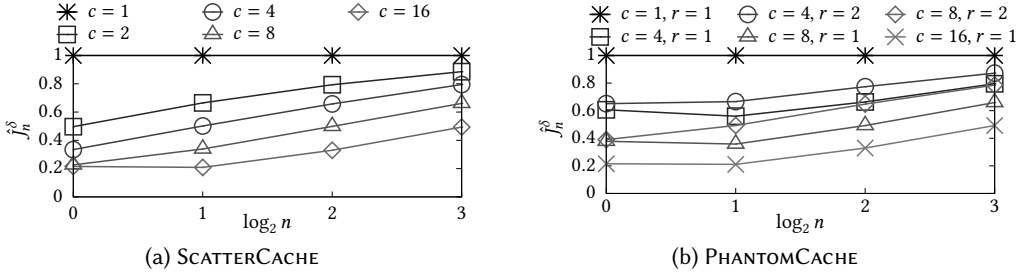
Intuitively, FLUSH+RELOAD is the best attacker strategy for a normal cache design when memory sharing is enabled. However, for a new cache design, it may not be clear that it is still the best. Our leakage rules provide some insight for SCATTERCACHE and PHANTOMCACHE. For example, one top-ranking rule for SCATTERCACHE, with precision ≥ 0.80 and recall of ≈ 0.02 , is:

$$\begin{aligned} & \bar{s}(\text{'secret'})[3] \geq 1 \wedge \bar{s}(\text{'secret'})[2] < 1 \wedge \bar{s}(\text{'secret'})[1] < 1 \wedge \bar{s}(\text{'secret'})[0] < 1 \\ & \wedge \bar{i}(M_0\{8\}\{1\}) \geq 5 \wedge \bar{i}(M_1\{8\}\{1\}) \geq 5 \wedge \bar{c}(\text{'load'})[8] < 1 \end{aligned} \quad (21)$$

This rule is similar to (20) but with some additional predicates about M_0 . Specifically, (21) adds $\bar{i}(M_0\{8\}\{1\}) \geq 5 \wedge \bar{i}(M_1\{8\}\{1\}) \geq 5$ to the rule when setting $\bar{c}(\text{'load'})[8] = 0$ (i.e., attacker FLUSHES $block_8$) and $\bar{s}(\text{'secret'}) = 8$, which indicates that the $block_8$ should occupy line $k = 1$ in set $j = 5$ in both the victim's and attacker's domains to ensure leakage about whether $\bar{s}(\text{'secret'}) = 8$ when the attacker RELOADS $block_8$.

Thus, an attacker should FLUSH+RELOAD all blocks that could share cache lines between victim's and attacker's domain to cause more leakage. Since the memory-to-cache mapping is unknown, an attacker may FLUSH+RELOAD all shared memory blocks. The resulting \hat{J}_n is shown in Fig. 11(a) for SCATTERCACHE and Fig. 11(b) for PHANTOMCACHE. \hat{J}_n is high when n is large, indicating the attacker can precisely determine $\bar{s}(\text{'secret'})$ when leakage occurs. Our results indicate that lower cache set granularity leaks more: In Fig. 11(a), $c = 1$ leaks the most, which is similar to the normal cache. When $c > 1$, the leakage is reduced.

⁵In contrast to the original paper [Tan et al. 2020], we do not force each memory block to map to r unique cache sets, i.e., we do not constrain $M'_{dom}\{baddr\}\{k\} \neq M'_{dom}\{baddr\}\{k'\}$ for $k \neq k'$.

Fig. 12. Memory sharing disabled (PRIME+PROBE attack), $\vec{\Delta}$ ('info') $\leftarrow \vec{\Gamma}(M)$ (or $\vec{\Gamma}(M^r)$)Fig. 13. Memory sharing enabled (FLUSH+RELOAD attack), $\vec{\Delta}$ ('info') $\leftarrow \vec{\Gamma}(M)$ (or $\vec{\Gamma}(M^r)$)

Overall, with same cache set granularity, \hat{J}_n is higher with PHANTOMCACHE with $r = 2$ than PHANTOMCACHE with $r = 1$ and SCATTERCACHE. This is because setting $r = 2$ allows one physical address to be mapped to more cache sets and so gains more chance to share cache lines across domains.

We also see that \hat{J}_n for ' $c = 8, r = 2$ ' is close to that for ' $c = 4, r = 1$ ', as randomly mapping to 2 out of 8 sets is similar to mapping to 1 out of 4 cache sets. Our evaluation results suggests that SCATTERCACHE and PHANTOMCACHE eliminate side-channel leakage when there is no shared memory and largely restrict it when there is shared memory, if the address-to-cache mapping is random and remains unknown to the attacker.

6.4.2 Declassifying the Memory-to-Cache Mapping. When $\vec{\Gamma}(M)$ is unknown to the attacker, our previous analysis shows that cache-based side channels are mitigated. Werner et al. [2019] also discussed the possibility of this mapping being disclosed to the attacker, however, through a profiling procedure. If we declassify $\vec{\Gamma}(M)$, the interference \hat{J}_n^δ will increase: Fig. 12(a) shows \hat{J}_n^δ due to PRIME+PROBE attacks in this case, and Fig. 13(a) shows the impact of this declassification on FLUSH+RELOAD attacks.

Similarly, using $\vec{\Delta}$ ('info') $\leftarrow \vec{\Gamma}(M^r)$, we evaluate PHANTOMCACHE's leakage when the random mapping is declassified; results are shown in Fig. 12(b) and Fig. 13(b). Comparing Fig. 12(b) and Fig. 12(a), PHANTOMCACHE's leakage (measured by \hat{J}_n^δ) for unshared memory is higher than SCATTERCACHE's when $r = 1$. The strength of PHANTOMCACHE is revealed when r increases, since it allows memory blocks to map to more than one cache set. Specifically, the leakage for SCATTERCACHE's ' $c = 4$ ' is much less than PHANTOMCACHE's ' $c = 4, r = 1$ ' but is similar to PHANTOMCACHE's ' $c = 4, r = 2$ '. However, PHANTOMCACHE with $r = 2$ provides weaker protection for FLUSH+RELOAD than PHANTOMCACHE with $r = 1$ and SCATTERCACHE.

6.5 Leaking Exponent in Modular Exponentiation

The evaluations in Sec. 6.3 and Sec. 6.4 focused on whether the adversary could detect the victim's access to a particular memory block, which is a well-known vector of information leakage. To further demonstrate the utility of our framework in measuring this type of leakage, here we consider a classic example whereby the secret is not a memory address, but rather is a cryptographic secret that, due to the algorithm in use, can influence the victim's cache footprint.

The particular example we evaluate here is modular exponentiation as used in algorithms such as RSA. A textbook implementation of modular exponentiation uses a sliding-window method that is known to leak information in caches [Bernstein et al. 2017; Zhang et al. 2012]. As shown in Fig. 14(a), the algorithm leverages some small powers $b[k]$ of a base b (where $k < 2^W - 1$) to compute a larger power. Accesses to those precomputed powers is determined by the window-sized segment d_i of the private key d in each loop iteration i . First, this procedure will leak via the cache whether d_i is zero. Second, since the precomputed elements are addressed by d_i , an attacker may identify up to $\log_2 c$ bits about d_i if those precomputed powers map to different cache sets.

To evaluate the one-round leakage of Fig. 14(a), we used the RISC-V assembly shown in Fig. 14(b) in BOOM with a 2-way, 8-set cache ($c = 8$). The \hat{J}_n measure shown in Fig. 15(a) indicates that the amount of leakage for one loop iteration i is limited, when $W \leq 4$ and so the precomputed b only uses up to $4 \times 2^4 = 64$ bytes (i.e., one cache line). When $4 < W < 8$, the side channel will leak more about d_i when W increases. Thus, choosing $W = 4$ is the best choice to protect the secret in our cache configuration.

To further diagnose the cause of leakage, we generated the interference rules for $W = 1$, $W = 4$, and $W = 8$. When $W = 1$, we obtain a single rule with precision and recall of 1.0, namely

$$\bar{c}(\text{'load'})[0] \geq 1 \wedge \bar{c}(\text{'load'})[8] \geq 1$$

This has no \bar{s} or \bar{s}' related conjuncts, indicating that the 1-bit secret d_i is fully leaked when an attacker PRIMES one cache set. In contrast, when $W = 4$, the top rules (precision of 1.0, recall ≥ 0.5) include some \bar{s} or \bar{s}' related conjuncts, constraining the secret value to be zero, e.g.,

$$\bar{s}(d_i) < 1 \wedge \bar{c}(\text{'load'})[0] \geq 1 \wedge \bar{c}(\text{'load'})[8] \geq 1$$

That is, it only leaks whether it is zero or not for a 4-bit secret.

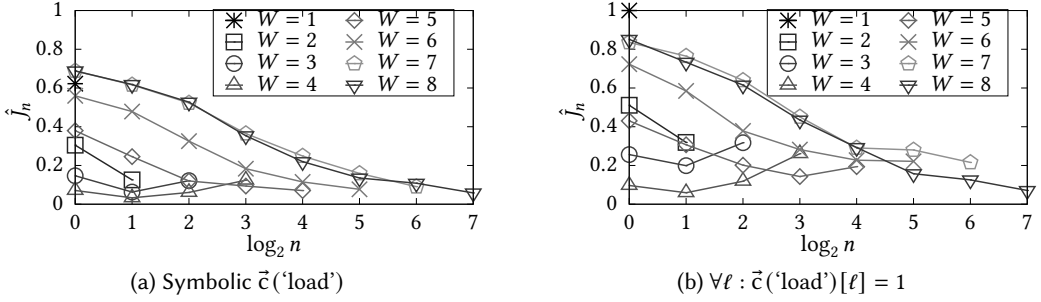
When $W > 4$, however, the most important cause of leakage changes from whether a memory access happens to which cache set is used by d_i . For example, when $W = 8$, one highly ranked rule (precision of 1.0, recall ≥ 0.04) is

$$\begin{aligned} & \bar{s}'(d_i)[6] < 1 \wedge \bar{s}'(d_i)[5] \geq 1 \wedge \bar{s}'(d_i)[4] < 1 \wedge \bar{s}(d_i)[4] \geq 1 \\ & \wedge \bar{c}(\text{'load'})[10] \geq 1 \wedge \bar{c}(\text{'load'})[2] \geq 1 \end{aligned} \quad (22)$$

which indicates that the attacker can distinguish an $\bar{s}'(d_i)$ with $\bar{s}'(d_i)[4 : 6] = 2$ from an $\bar{s}(d_i)$ with $\bar{s}(d_i)[4 : 6] \in \{1, 3, 5, 7\}$ if the attacker PRIMES cache set 2. Similar to the analysis in Sec. 6.3.1, rules

<pre> 1: function MODEXP(b,d) 2: e ← 1 3: for i ← n to 1 do 4: e ← e × e mod M 5: if d_i ≠ 0 then 6: e ← e × b[d_i] 7: end if 8: end for 9: return e 10: end function </pre>	<pre> proc(c̄, ī, s̄) li sp, 0x80000400 li a0, 1 li a2, M li a3, s̄(d_i) oneIteration: mulw a0, a0, a0 remw a0, a0, a2 beqz a3, .NextIteration sll a5, a3, 2 add a5, sp, a5 lw a5, 0(a5) mulw a0, a0, a5 remw a0, a0, a2 </pre>
(a) Algorithm	(b) Assembly for one iteration

Fig. 14. Sliding window modular exponentiation. d is the private key where each d_i ($i = 1, \dots, n$) is a W -bit value.

Fig. 15. \hat{J}_n for MODEXP in 2-way, 8-set cache

```
conditionalAccess(offset, arr1.size)
```

```
if (offset < arr1.size)
  tmp ← arr2[(arr1[offset] × 64) & 1023]
  declassify(arr1[offset])
```

(a) Conditional memory access

```
victimFunc(offset, secret, arr1.size)
```

```
arr1[offset] ← secret
arr1.size ← (arr1.size × 257) mod 256
arr1.size ← (arr1.size × 257) mod 256
conditionalAccess(offset, arr1.size)
```

(b) Bounds check with long dependency

```
victimFunc(offset, secret)
```

```
arr1[offset] ← secret
read arr1.size from memory;
conditionalAccess(offset, arr1.size)
```

(c) Bounds check with short dependency

```
1 .shortDependency:
2 lbu a0, 0x100(t3)
```

(d) Short speculation

```
1 proc(c, i, s)
2 .prepareData:
3 li a0, i('arr1.size')
4 li a1, c('offset')
5 li a2, s('secret')
6 //t3 ← arr1.addr
7 //t4 ← arr2.addr
8 add a3, t3, a1
9 sb s2, 0(a3)
10 .complexDependency:
11 li t1, 0x101
12 li t2, 0x100
13 mul a4, a0, t1
14 remuw a4, a4, t2
15 mul a4, a4, t1
16 remuw a0, a4, t2
17 .conditionalAccess:
18 bleu a0, a1, .end
19 add t3, t3, a1
20 lbu a3, 0x0(t3)
21 sll a3, a3, 6
22 and a3, a3, 0x3ff
23 add a3, t4, a3
24 lbu a4, 0(a3)
```

(e) Long speculation

Fig. 16. Speculative execution example. Assembly in (e) is snippet from compilation of pseudocode in (b). Replacing lines 10–16 with (d) gives the analogous assembly for the pseudocode in (c).

for $W = 8$ illustrate that an attacker can reveal the cache set used by the victim (e.g., secret bits 4-6) when priming all cache sets.

6.6 Cache-Based Side Channels in Speculative Execution

SPECTRE and its variants have received widespread attention in recent years. In a SPECTRE attack, a CPU predicts the outcome of a conditional branch and executes instructions based on that prediction to reduce delays incurred by those instructions if its prediction was correct. However, even if the prediction is incorrect, then some changes to the hardware state caused by speculative execution will persist even after the mispredicted computations have been discarded. These changes propagate information to exploitable cache-based side channels, allowing the attacker to steal it.

To explore such leaks using our framework, we used the software pseudocode in Fig. 16(b) and Fig. 16(c), each of which accesses an element of array `arr2` at a secret index `arr1[offset]`. The bounds check on `offset` is dependent on a complex sequence of computations in Fig. 16(b) and on reading `arr1.size` from memory in Fig. 16(c). Theoretically, speculative execution may leak

arr1[offset] through cache-based side channels in both cases if the dependency is not resolved before speculative execution, i.e., by bringing arr2[(arr1[offset] × 64) & 1023] into cache. Fig. 16(e) shows an important snippet of RISC-V assembly for Fig. 16(b) running on BOOM with a 2-way, 8-set cache. To evaluate the software snippet in Fig. 16(c), we change the block denoted by .complexDependency (Lines 10–16) with the .shortDependency in Fig. 16(d). Furthermore, we evaluated a mitigation similar to **lfence** [Int 2018], by adding a RISC-V instruction ‘fence r, r’ just after Line 18 in Fig. 16(e).

We assume the attacker can control the offset value \vec{c} (‘offset’), train the *GShare* branch predictor \vec{c} (‘bpd’) shown in Fig. 7, and use FLUSH+RELOAD to observe \vec{o} (‘hit’). The attacker can use the FLUSH+RELOAD-style attacks to precisely determine the index into arr2 if arr2 is shared and thus four bits of arr1[offset]. Note that the secret value \vec{s} (‘secret’) is assigned to arr1[offset] as the first step of Fig. 16(c) and Fig. 16(b). We presume that \vec{i} (‘arr1.size’) is an attacker-known but not controlled variable; thus, we include it as one output parameters as well, i.e., \vec{o} (‘arr1.size’) ← \vec{i} (‘arr1.size’).

As shown in Fig. 17, the \hat{J}_n measures for ‘ShortSpec’ (denoting Fig. 16(d)) and ‘Fence’ are somewhat similar to that for ‘LongSpec’ (denoting Fig. 16(e))—contrary to what intuition would suggest. This counterintuitive result is due to the fact that leakage from *in-bounds* array accesses is also being counted. By declassifying in-bounds array elements (i.e., declassifying arr1[offset] if \vec{c} (‘offset’) < \vec{i} (‘arr1.size’)), we obtain a better picture of when leakage occurs. Specifically, when measuring the leakage with declassification of in-bounds array elements, \hat{J}_n^δ indicates that both *proc* with the short dependency (‘ShortSpec+ δ ’) and *proc* with the **fence** mitigation (‘Fence+ δ ’) do not leak out-of-boundary memory contents, while the *proc* with the longer dependency (‘LongSpec+ δ ’) continues to leak secret data and indeed, is just slightly lower than ‘complexDepend’.

In generating interference rules for *proc* with a long speculation (Fig. 16(e)), the linear feature

$$L_0 = 0.005 \times \vec{s}(\text{‘secret’}) - 0.003 \times \vec{s}'(\text{‘secret’}) - 0.494 \times \vec{c}(\text{‘offset’}) + 0.496 \times \vec{i}(\text{‘arr1.size’}) \quad (23)$$

$$\approx 0.5 \times \vec{i}(\text{‘arr1.size’}) - 0.5 \times \vec{c}(\text{‘offset’}) \quad (24)$$

and specifically the conjunct $L_0 < 1$ appears in many of the top ranked rules. Using the approximation of L_0 above, $L_0 < 1$ implies that $\vec{i}(\text{‘arr1.size’}) < \vec{c}(\text{‘offset’}) + 2$, and so the offset is indeed out-of-bounds.

An example rule with precision 1.0 and recall 0.30 is

$$L_0 < 1 \wedge \vec{c}(\text{‘bpd}\{0\}.\text{state}') [1] < 1 \wedge |\vec{s}(\text{‘secret’}) [2] - \vec{s}'(\text{‘secret’}) [2]| \geq 1 \quad (25)$$

This rule indicates that an attacker can determine the third bit of the secret when the second bit of the state of the prediction entry \vec{c} (‘bpd{0}.state’) is 0 (‘strongly untaken’) or 1 (‘weakly untaken’). Analogous rules appear in the list for each of bits 0-2 and 4 of the secret. Other highly ranked rules (also with precision 1.0 and recall 0.30) are

$$L_0 < 1 \wedge \vec{c}(\text{‘bpd}\{0\}.\text{CFI}') [0] \geq 1 \wedge |\vec{s}(\text{‘secret’}) [0] - \vec{s}'(\text{‘secret’}) [0]| \geq 1 \quad (26)$$

$$L_0 < 1 \wedge \vec{c}(\text{‘bpd}\{0\}.\text{CFI}') [1] < 1 \wedge |\vec{s}(\text{‘secret’}) [3] - \vec{s}'(\text{‘secret’}) [3]| \geq 1 \quad (27)$$

Rule (26) leaks the first bit of the secret when the ‘CFI’ value (i.e., \vec{c} (‘bpd{0}.CFI’)) in the prediction entry is 1 or 3, and (27) leaks the fourth bit when the ‘CFI’ value is 0 or 1. In these cases, the ‘CFI’

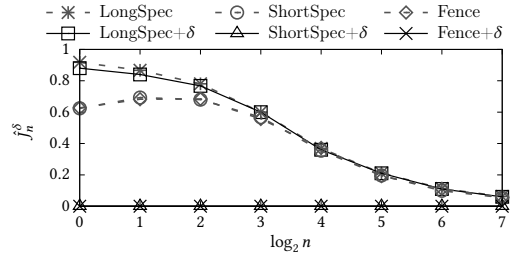


Fig. 17. \hat{J}_n^δ for SPECTRE in different procedures

value does not match the CFI portion of the instruction address (i.e., the address of Line 18 in Fig. 16(e)), which was $0x80000800 + 0x44$ ($= 0b01000100$), yielding a CFI portion of $0b10$ and $bidx$ of $0b00$. Because of the mismatch on CFI value, \vec{c} ('bpd{0}.state') is ignored and so speculation will not execute Lines 19–24. Though (26) and (27) are specific to the first or fourth bit of the secret, respectively, analogous rules appear for each of bits 0-3.

The simplicity of these rules stands in stark contrast to the complexity of the Yosys-generated per-cycle transition logic $\tau_{proc}(\vec{H}^{t-1}, \vec{H}^t)$, which includes 459,170 bit variables and 1,922,229 clauses in CNF, or the postcondition Π_{proc} , which still includes 5,413 bit variables and 41,940 clauses. Clearly, our interpretation rules are vastly simpler for the analyst to consider than these alternatives.

Fig. 18(a) shows the cumulative precision and recall for all leakage rules in this case study. However, we do not need to use all rules for interpretation, since most rules do not help much with the cumulative recall. For example, considering only rules that improve cumulative recall by $\geq 1\%$ gives 12 rules that achieve 0.97 precision and 0.98 recall (Fig. 18(b)).

We have performed this evaluation using earlier BOOM versions and noticed that the out-of-bounds leakage was partially eliminated in version 2.2.3.⁶ Since version 2.2.1, the miss handling (MSHR) module of the L1 cache tracks branch prediction results and discards the pending cache refill request if a misprediction is detected before the refill commit.

7 PERFORMANCE

In this section, we discuss the runtime performance of DINOme on the case studies described in Sec. 6. In DINOme, we have four important components: an automated logical formula generator (Sec. 5.1), a model counter (Sec. 5.2), a sampler (Sec. 5.3), and a rule learner (Sec. 4.3). This section reports the time costs in the first three stages for all case studies we have evaluated. We performed those experiments on a DELL PowerEdge R815 server with 2.3GHz AMD Opteron 6376 processors and 128GB memory.

⁶In BOOM version 2.2.1, the victim program described in Fig. 16(c) also suffers the out-of-bounds leakage and thus has 'ShortSpec' close to 'LongSpec' and 'ShortSpec+ δ ' close to 'LongSpec+ δ '.

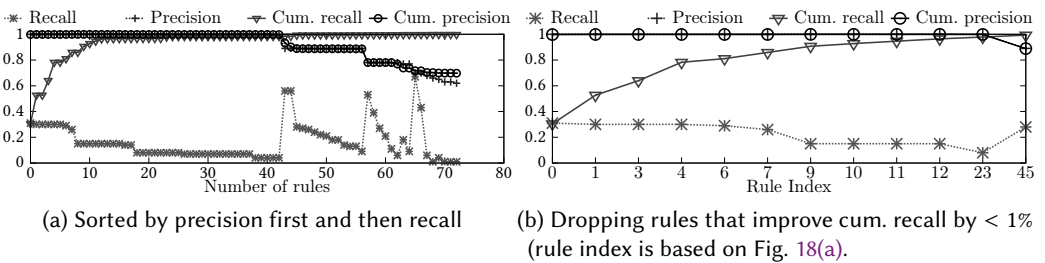


Fig. 18. Cumulative precision and recall vs. rules

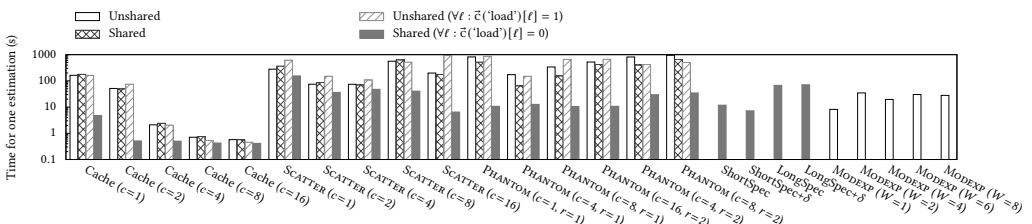


Fig. 19. Time used in one estimation of $\hat{J}^\delta(S, S')$

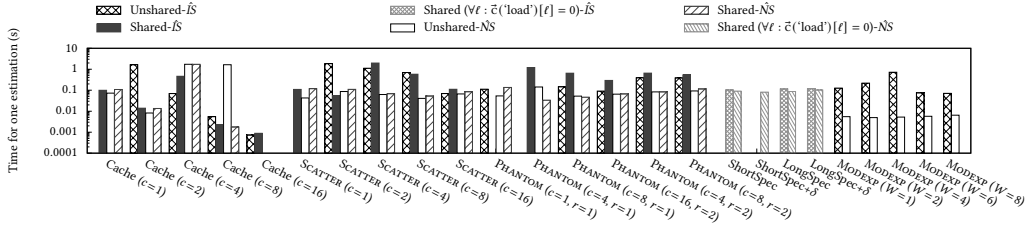


Fig. 20. Time used in generating one tuple in $\hat{N}S$ or $\hat{I}S$

The time to generate and simplify the logical postcondition is primarily influenced by the number of RISC-V BOOM cycles represented by that postcondition, as we incrementally compose the formula cycle by cycle. Computing Π_{proc} required 20-40 minutes for the memory accessing experiments (100 cycles) in Sec. 6.3 and Sec. 6.4; 45 minutes for the modular exponentiation experiments (120 cycles) in Sec. 6.5; and around 2 hours for the SPECTRE experiments (150 cycles) in Sec. 6.6. Different from Zhou et al. [2018], DINoME assembles the postcondition without path splitting per branch (and so avoids path explosion) and defers its solving task to a simplification step and final cycle, which reduces the complexity dramatically.

Fig. 19 shows the runtime to compute *one* estimate of $\hat{J}(S, S')$ or $\hat{J}^\delta(S, S')$ in the model counting process; note the logarithmic y-axis. Specifically, counting for cache-based side channels in SCATTERCACHE and PHANTOMCACHE are much more expensive than others, where one estimate requires up to 16 minutes. The difficulty in counting for SCATTERCACHE (denoted by ‘SCATTER’) and PHANTOMCACHE (denoted by ‘PHANTOM’) is due to the large size of their counting variables. For SCATTERCACHE, the memory-to-cache mapping uses $\log_2(c) \times w$ bits per domain per memory block for 32 memory blocks. Specifically, the 8-way 2-set SCATTERCACHE (denoted by ‘SCATTER ($c = 2$)’), uses 512 bits to represent $\vec{\tau}(M)$, which means the counting process would add hundreds of XOR constraints to compute one estimate, which greatly increases the difficulty to find a feasible solution. To obtain the sample sets $\hat{I}S$ and $\hat{N}S$, the sampling process generates a tuple in $\hat{I}S$ or $\hat{N}S$ within seconds, as illustrated in Fig. 20.

Our reported results reflect estimations of $\hat{J}(S, S')$ or $\hat{J}^\delta(S, S')$ for at least 100 S, S' pairs per n , and we sampled up to 100,000 tuples in $\hat{I}S$ and $\hat{N}S$. These estimations and samplings are trivially parallelizable and so, with horizontal scaling, can be performed in total times approaching those in Fig. 19 and Fig. 20 to the extent budget allows.

8 LIMITATIONS

Despite the scalability represented by DINoME specifically for analyzing processor designs, it still has limitations. First, due to the complexity of hardware logic, generating the postcondition $\Pi_{proc}(\vec{c}, \vec{o}, \vec{\tau}, \vec{s})$ for a *proc* representing both the OS and the application would require more CPU cycles than the number to which we have been able to scale DINoME thus far. The DINoME workloads described in this paper represent a tradeoff, using a sequence of opcodes with concretized operations and selected symbolic operands above a partially symbolic hardware specification. To evaluate with more complicated software, a possible solution is to highly concretize the initial hardware state (especially for the memory and cache states) or highly concretize the software, at the cost of possibly missing some potential leakage that remains hidden due to this concretization.

A second limitation of DINoME, and specifically of its generation of interpretation rules to explain leakage, is that the interpretation rules may not be complete, for two reasons. First, the interpretation rules might skip a rule that covers few leakage samples (i.e., with low recall). A possible way to address this source of incompleteness is to declassify the sources of leakage exposed in the inference rules that *are* learned, and then rerun the learning process again. Second, the

conditions that result in leakage might be more complicated than can be learned using decision trees built using local linear classifiers. Alternative learning methods might be tried, though doing so while retaining interpretability will be a challenge.

9 CONCLUSION

Scaling high-fidelity, static noninterference measurement to complex computations has been a challenge since the introduction of noninterference in the 1980s [Goguen and Meseguer 1982]. We believe that we have advanced the state-of-the-art in this area both generally and specifically for its application to processor designs. Certain innovations in our DINoME framework, such as the cycle-by-cycle construction of the logical postcondition for processor execution, are specific to processor designs. Others, such as our methods for declassification and interpreting leakage results, are not. Together, however, they permit the measurement of leakage in complex scenarios, as we demonstrated through using DINoME to analyze leakage due to speculative execution in the BOOM core and of published defenses to mitigate it. Our analysis enables comparisons between defenses to discover, e.g., the processor and defense parameterizations where one defense outperforms the other. Though the performance of DINoME suggests that static measurement of noninterference for processors is still too time-intensive for highly interactive use, it is fast enough to permit multiple analysis iterations per day in many cases. And through its improvements in declassification and interpretability, it substantially facilitates human understanding of its measurement results.

ACKNOWLEDGMENTS

We are grateful to our shepherd, Prof. Andrew Myers, and to the anonymous reviewers for numerous constructive suggestions for improving this paper. This work was supported in part by grant 2113345 from the National Science Foundation and by a gift from Intel.

REFERENCES

2018. *Intel Analysis of Speculative Execution Side Channels*. Technical Report. Intel Corp. <https://www.intel.com/content/www/us/en/architecture-and-technology/intel-analysis-of-speculative-execution-side-channels-paper.html>
- O. Aciğmez. 2007. Yet another microarchitectural attack: Exploiting I-cache. In *ACM Workshop on Computer Security Architecture*. 11–18. <https://doi.org/10.1145/1314466.1314469>
- J. B. Almeida, M. Barbosa, G. Barthe, F. Dupressoir, and M. Emmi. 2016. Verifying constant-time implementations. In *25th USENIX Security Symposium*. 53–70.
- R. A. Aziz, G. Chu, C. Muise, and P. Stuckey. 2015. # \exists SAT: Projected Model Counting. In *18th International Conference on Theory and Applications of Satisfiability Testing (LNCS)*. 121–137. https://doi.org/10.1007/978-3-319-24318-4_10
- M. Backes, B. Kopf, and A. Rybalchenko. 2009. Automatic discovery and quantification of information leaks. In *30th IEEE Symposium on Security and Privacy*. 141–153. <https://doi.org/10.1109/SP.2009.18>
- T. Ball, B. Cook, V. Levin, and S. K. Rajamani. 2004. SLAM and Static Driver Verifier: Technology transfer of formal methods inside Microsoft. In *4th International Conference on Integrated Formal Methods (LNCS)*, Vol. 2999. 1–20. https://doi.org/10.1007/978-3-540-24756-2_1
- A. Banerjee, D. A. Naumann, and S. Rosenberg. 2008. Expressive declassification policies and modular static enforcement. In *29th IEEE Symposium on Security and Privacy*. 339–353. <https://doi.org/10.1109/SP.2008.20>
- G. Barthe, G. Betarte, J. Campo, C. Luna, and D. Pichardie. 2014. System-level non-interference for constant-time cryptography. In *21st ACM Conference on Computer and Communications Security*. 1267–1279. <https://doi.org/10.1145/2660267.2660283>
- D. J. Bernstein, J. Breitner, D. Genkin, L. G. Bruinderink, N. Heninger, T. Lange, C. V. Vredendaal, and T. Yarom. 2017. Sliding right into disaster: Left-to-right sliding windows leak. In *19th International Conference on Cryptographic Hardware and Embedded Systems (LNCS)*, Vol. 10529. 555–576. https://doi.org/10.1007/978-3-319-66787-4_27
- S. Blazy, D. Pichardie, and A. Trieu. 2019. Verifying constant-time implementations by abstract interpretation. *Journal of Computer Security* 27, 1 (2019), 137–163. https://doi.org/10.1007/978-3-319-66402-6_16
- C. Celio, P. Chiu, B. Nikolic, D. A. Patterson, and K. Asanovic. 2017. BOOMv2: an open-source out-of-order RISC-V core. In *1st Workshop on Computer Architecture Research with RISC-V (CARRV)*.
- S. Chakraborty, K. S. Meel, and M. Y. Vardi. 2013. A scalable approximate model counter. In *Principles and Practice of Constraint Programming (LNCS)*, Vol. 8124. 200–216. https://doi.org/10.1007/978-3-642-40627-0_18

- P. Chapman and D. Evans. 2011. Automated black-box detection of side-channel vulnerabilities in web applications. In *18th ACM Conference on Computer and Communications Security*. 263–274. <https://doi.org/10.1145/2046707.2046737>
- S. Chattopadhyay, M. Beck, A. Rezine, and A. Zeller. 2017. Quantifying the Information Leak in Cache Attacks via Symbolic Execution. In *15th ACM International Conference on Formal Methods and Models for System Design* (Vienna, Austria). New York, NY, USA, 25–35. <https://doi.org/10.1145/3288758>
- S. Chattopadhyay and A. Roychoudhury. 2018. Symbolic verification of cache side-channel freedom. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 37, 11 (2018), 2812–2823. <https://doi.org/10.1109/TCAD.2018.2858402>
- C. Chen, K. Lin, C. Rudin, Y. Shaposhnik, S. Wang, and T. Wang. 2018. An Interpretable Model with Globally Consistent Explanations for Credit Risk. arXiv:cs.LG/1811.12615
- T. Chen and C. Guestrin. 2016. Xgboost: A scalable tree boosting system. In *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 785–794. <https://doi.org/10.1145/2939672.2939785>
- S. Chong and A. C. Myers. 2004. Security policies for downgrading. In *11th ACM conference on Computer and communications security*. 198–209. <https://doi.org/10.1145/1030083.1030110>
- W. W. Cohen and Y. Singer. 1999. A simple, fast, and effective rule learner. *16th AAAI Conference on Artificial Intelligence 99* (1999), 335–342. <https://doi.org/10.5555/315149.315320>
- G. Doychev, B. Köpf, L. Mauborgne, and J. Reineke. 2013. CacheAudit: A tool for the static analysis of cache side channels. In *22nd USENIX Security Symposium*. 431–446.
- B. Dutertre. 2015. Solving exists/forall problems with yices. In *Workshop on satisfiability modulo theories*.
- C. Dwork, F. McSherry, K. Nissim, and A. Smith. 2006. Calibrating Noise to Sensitivity in Private Data Analysis. In *3rd Theory of Cryptography Conference (LNCS)*, Vol. 3876. 265–284. https://doi.org/10.1007/11681878_14
- J. Fan. 1993. Local linear regression smoothers and their minimax efficiencies. *The Annals of Statistics* (1993), 196–216. <https://doi.org/10.1214/aos/1176349022>
- R. Fan, K. Chang, C. Hsieh, X. Wang, and C. Lin. 2008. LIBLINEAR: A library for large linear classification. *Journal of Machine Learning Research* 9, Aug (2008), 1871–1874. <https://doi.org/10.5555/1390681.1442794>
- A. Ferraiuolo, R. Xu, D. Zhang, A. C. Myers, and G.E. Suh. 2017. Verification of a practical hardware security architecture through static information flow analysis. In *22nd International Conference on Architectural Support for Programming Languages and Operating Systems*. 555–568. <https://doi.org/10.1145/3093337.3037739>
- M. Fokkema. 2020. Fitting Prediction Rule Ensembles with R Package pre. *Journal of Statistical Software* 92, 12 (2020), 1–30. <https://doi.org/10.18637/jss.v092.i12>
- J. H. Friedman. 2001. Greedy function approximation: a gradient boosting machine. *Annals of Statistics* (2001), 1189–1232.
- J. H. Friedman and B. E. Popescu. 2008. Predictive learning via rule ensembles. *The Annals of Applied Statistics* 2, 3 (2008), 916–954. <https://doi.org/10.1214/07-AOAS148>
- R. Giacobazzi and I. Mastroeni. 2004. Abstract non-interference: Parameterizing non-interference by abstract interpretation. *ACM SIGPLAN Notices* 39, 1 (2004), 186–197. <https://doi.org/10.1145/982962.964017>
- R. Giacobazzi and I. Mastroeni. 2018. Abstract non-interference: a unifying framework for weakening information-flow. *ACM Transactions on Privacy and Security (TOPS)* 21, 2 (2018), 1–31. <https://doi.org/10.1145/3175660>
- K. V. Gleissenthall, R. G. Kıcı, D. Stefan, and R. Jhala. 2019. IODINE: Verifying Constant-Time Execution of Hardware. In *28th USENIX Security Symposium*. 1411–1428.
- P. Godefroid, M. Y. Levin, and D. Molnar. 2012. SAGE: Whitebox Fuzzing for Security Testing. *Queue* 10, 1 (2012), 20–27. <https://doi.org/10.1145/2090147.2094081>
- J. A. Goguen and J. Meseguer. 1982. Security policies and security models. In *3rd IEEE Symposium on Security and Privacy*. 11–20. <https://doi.org/10.1109/SP.1982.10014>
- J. W. Gray. 1991. Toward a mathematical foundation for information flow security. In *12nd IEEE Symposium on Security and Privacy*. 21–34. <https://doi.org/10.1109/RISP.1991.130769>
- X. Guo, R. G. Dutta, J. He, M. M. Tehranipoor, and Y. Jin. 2019. QIF-Verilog: Quantitative Information-Flow based Hardware Description Languages for Pre-Silicon Security Assessment. In *IEEE International Symposium on Hardware Oriented Security and Trust*. 91–100. <https://doi.org/10.1109/HST.2019.8740840>
- J. Kelsey. 2002. Compression and information leakage of plaintext. In *9th International Workshop on Fast Software Encryption*. 263–276. https://doi.org/10.1007/3-540-45661-9_21
- P. Kocher, J. Horn, A. Fogh, D. Genkin, D. Gruss, W. Haas, M. Hamburg, M. Lipp, S. Mangard, and T. Prescher. 2019. Spectre attacks: Exploiting speculative execution. In *40th IEEE Symposium on Security and Privacy*. 1–19. <https://doi.org/10.1109/SP.2019.00002>
- M. Lipp, M. Schwarz, D. Gruss, T. Prescher, W. Haas, A. Fogh, J. Horn, S. Mangard, P. Kocher, S. Genkin, Y. Yarom, and M. Hamburg. 2018. Meltdown: Reading Kernel Memory from User Space. In *27th USENIX Security Symposium*. 973–990.
- P. Malacaria, MHR. Khouzani, C. S. Pasareanu, Q. Phan, and K. Luckow. 2018. Symbolic Side-Channel Analysis for Probabilistic Programs. In *31st IEEE Computer Security Foundations Symposium*. 313–327. <https://doi.org/10.1109/CSF.2018.00030>

- M. McCall, H. Zhang, and L. Jia. 2018. Knowledge-Based Security of Dynamic Secrets for Reactive Programs. In *31st IEEE Computer Security Foundations Symposium*. 175–188. <https://doi.org/10.1109/CSF.2018.00020>
- C. Molnar. 2019. *Interpretable Machine Learning*. <https://christophm.github.io/interpretable-ml-book/>.
- S. Nilzadeh, Y. Noller, and C. S. Păsăreanu. 2019. DiffFuzz: Differential Fuzzing for Side-Channel Analysis. In *41st International Conference on Software Engineering*. 176–187. <https://doi.org/10.1109/ICSE.2019.00034>
- O. Oleksii, T. Bohdan, S. Mark, and F. Christof. 2020. SpecFuzz: Bringing Spectre-type vulnerabilities to the surface. In *29th USENIX Security Symposium*.
- D. A. Osvik, A. Shamir, and E. Tromer. 2006. Cache attacks and countermeasures: The case of AES. In *Topics in Cryptology – CT-RSA (LNCS)*, Vol. 3860. 1–20. https://doi.org/10.1007/11605805_1
- C. Percival. 2005. Cache missing for fun and profit. In *BSDCan 2005*. <https://doi.org/10.1.1.187.8383>
- Q. Phan and P. Malacaria. 2014. Abstract model counting: A novel approach for quantification of information leaks. In *9th ACM Symposium on Information, Computer and Communications Security*. 283–292. <https://doi.org/10.1145/2590296.2590328>
- A. Pnueli, Y. Rodeh, O. Strichman, and M. Siegel. 2002. The small model property: How small can it be? *Information and Computation* 178, 1 (2002), 279–293. [https://doi.org/10.1016/S0890-5401\(02\)93175-5](https://doi.org/10.1016/S0890-5401(02)93175-5)
- M. T. Ribeiro, S. Singh, and C. Guestrin. 2016. “Why should I trust you?”: Explaining the predictions of any classifier. In *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1135–1144. <https://doi.org/10.1145/2939672.2939778>
- M. T. Ribeiro, S. Singh, and C. Guestrin. 2018. Anchors: High-precision model-agnostic explanations. In *32nd AAAI Conference on Artificial Intelligence*. 1527–1535. <https://ojs.aaai.org/index.php/AAAI/article/view/11491>
- A. Sabelfeld and A. C. Myers. 2003. A model for delimited information release. In *2nd International Symposium on Software Security – Theories and Systems (LNCS)*, Vol. 3233. 174–191. https://doi.org/10.1007/978-3-540-37621-7_9
- A. Sabelfeld and D. Sands. 2009. Declassification: Dimensions and Principles. *Journal of Computer Security* (2009), 517–548. <https://doi.org/10.5555/1662658.1662659>
- S. Sahai, P. Subramanyan, and R. Sinha. 2020. Verification of Quantitative Hyperproperties Using Trace Enumeration Relations. In *32nd International Conference on Computer Aided Verification (LNCS)*, Vol. 12224. 201–224. https://doi.org/10.1007/978-3-030-53288-8_11
- T. Seidenfeld. 1986. Entropy and uncertainty. *Philosophy of Science* 53, 4 (1986), 467–491. <https://doi.org/10.1086/289336>
- G. Smith. 2009. On the Foundations of Quantitative Information Flow. In *12th International Conference on Foundations of Software Science and Computational Structures (LNCS)*, Vol. 5504. 288–302. https://doi.org/10.1007/978-3-642-00596-1_21
- G. Smith. 2011. Quantifying Information Flow Using Min-Entropy. In *8th International Conference on Quantitative Evaluation of Systems*. 159–167. <https://doi.org/10.1109/QEST.2011.31>
- Dawn Xiaodong Song, David Wagner, and Xuqing Tian. 2001. Timing Analysis of Keystrokes and Timing Attacks on SSH. In *10th USENIX Security Symposium*.
- M. Soos and K. S. Meel. 2019. BIRD: Engineering an Efficient CNF-XOR SAT Solver and its Applications to Approximate Model Counting. In *36th AAAI Conference on Artificial Intelligence*. 1592–1599. https://doi.org/10.1007/978-3-030-80223-3_37
- Qinhan Tan, Zhihua Zeng, Kai Bu, and Kui Ren. 2020. PhantomCache: Obfuscating Cache Conflicts with Localized Randomization. In *27th Network and Distributed System Security Symposium*. <https://doi.org/10.14722/ndss.2020.24086>
- T. Wang, T. Wei, Lin Z, and W. Zou. 2009. IntScope: Automatically Detecting Integer Overflow Vulnerability in x86 Binary Using Symbolic Execution. In *16th Network and Distributed System Security Symposium*. https://doi.org/10.1007/978-3-642-15497-3_5
- Z. Wang and R. B. Lee. 2007. New cache designs for thwarting software cache-based side channel attacks. In *34th International Symposium on Computer Architecture*. 494–505. <https://doi.org/10.1145/1273440.1250723>
- M. Werner, T. Unterluggauer, L. Giner, M. Schwarz, D. Gruss, and S. Mangard. 2019. ScatterCache: Thwarting Cache Attacks via Cache Set Randomization. In *28th USENIX Security Symposium*. Santa Clara, CA, 675–692.
- C. Wolf. [n.d.]. Yosys Open SYnthesis Suite. <http://www.clifford.at/yosys/>.
- Y. Xiao, Y. Zhang, and R. Teodorescu. 2020. SPEECHMINER: A Framework for Investigating and Measuring Speculative Execution Vulnerabilities. In *27th Network and Distributed System Security Symposium*. <https://doi.org/10.14722/ndss.2020.23105>
- Y. Yarom and K. E. Falkner. 2014. FLUSH+RELOAD: A high resolution, low noise, L3 cache side-channel attack. In *23rd USENIX Security Symposium*. 719–732.
- H. Yasuoka and T. Terauchi. 2014. Quantitative information flow as safety and liveness hyperproperties. *Theoretical Computer Science* 538 (2014), 167–182. <https://doi.org/10.4204/EPTCS.85.6>
- D. Zhang, Y. Wang, G. E. Suh, and A. C. Myers. 2015. A Hardware Design Language for Timing-Sensitive Information-Flow Security. In *20th International Conference on Architectural Support for Programming Languages and Operating Systems* (Istanbul, Turkey). Association for Computing Machinery, New York, NY, USA, 503–516. <https://doi.org/10.1145/2694344.2694372>

- K. Zhang, Z. Li, R. Wang, X. Wang, and S. Chen. 2010. Sidebuster: Automated detection and quantification of side-channel leaks in web application development. In *17th ACM Conference on Computer and Communications Security*. 595–606. <https://doi.org/10.1145/1866307.1866374>
- R. Zhang, C. Deutschbein, P. Huang, and C. Sturton. 2018. End-to-End Automated Exploit Generation for Validating the Security of Processor Designs. In *51st IEEE/ACM International Symposium on Microarchitecture*. 815–827. <https://doi.org/10.1109/MICRO.2018.00071>
- Y. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart. 2012. Cross-VM side channels and their use to extract private keys. In *19th ACM Conference on Computer and Communications Security*. 305–316. <https://doi.org/10.1145/2382196.2382230>
- Z. Zhou. 2020. *Evaluating Information Leakage by Quantitative and Interpretable Measurements*. Ph.D. Dissertation. The University of North Carolina at Chapel Hill.
- Z. Zhou, Z.Y Qian, M. K. Reiter, and Y. Zhang. 2018. Static Evaluation of Noninterference using Approximate Model Counting. In *39th IEEE Symposium on Security and Privacy*. 514–528. <https://doi.org/10.1109/SP.2018.00052>
- Z. Zhou, M. K. Reiter, and Y. Zhang. 2016. A Software Approach to Defeating Side Channels in Last-Level Caches. In *23rd ACM Conference on Computer and Communications Security*. 871–882. <https://doi.org/10.1145/2976749.2978324>