

CaSA: End-to-end Quantitative Security Analysis of Randomly Mapped Caches

Thomas Bourgeat*, Jules Drean*, Yuheng Yang[†], Lillian Tsai, Joel Emer[‡], Mengjia Yan
 MIT CSAIL, [†]MIT/University of Chinese Academy of Sciences, [‡]NVIDIA/MIT
 {bthom,drean,yuhengy,tslilyai,jsemer,mengjiay}@mit.edu

*Authors contributed equally to this work.

Abstract—It is well known that there are micro-architectural vulnerabilities that enable an attacker to use caches to exfiltrate secrets from a victim. These vulnerabilities exploit the fact that the attacker can detect cache lines that were accessed by the victim. Therefore, architects have looked at different forms of randomization to thwart the attacker's ability to communicate using the cache. The security analysis of those *randomly mapped caches* is based upon the increased difficulty for the attacker to determine the addresses that touch the same cache line that the victim has accessed.

In this paper, we show that the analyses used to evaluate those schemes were incomplete in various ways. For example, they were incomplete because they only focused on one of the steps used in the exfiltration of secrets. Specifically, the step that the attacker uses to determine the set of addresses that can monitor the cache lines used by the transmitter address. Instead, we broaden the analysis of micro-architecture side channels by providing an overall view of the communication process. This allows us to identify the existence of other communication steps that can also affect the security of randomly mapped caches, but have been ignored by prior work.

We design an analysis framework, CaSA, to comprehensively and quantitatively analyze the security of these randomly mapped caches. We comprehensively consider the end-to-end communication steps and study the statistical relationship between different steps. In addition, to perform quantitative analysis, we leverage the concepts from the field of telecommunications to formulate the security analysis into a statistical problem. We use CaSA to evaluate a wide range of attack strategies and cache configurations. Our result shows that the randomization mechanisms used in the state-of-the-art randomly mapped caches are insecure.

Index Terms—Micro-architecture side channel, randomly mapped cache, security analysis.

I. INTRODUCTION

The class of attacks that exploit micro-architectural vulnerabilities to breach processor security, generally referred to as *side-channel attacks*, have become a serious security threat. Using these attacks, an attacker can steal secrets from a victim application running on the same machine by monitoring the side effects of the victim's actions on various micro-architectural states. Such attacks are effective and have been used to leak encryption keys [1], [2]. Many of these attacks employ speculative execution to modify cache states [3]–[5] to completely bypass memory isolation and leak arbitrary data.

As described in [6], there is a series of elements common to most attacks that exploit micro-architectural vulnerabilities. These include either pre-existing or attacker-generated code run

in the victim's security domain that 1) accesses secret information and 2) transmits that information over a communication channel that 3) is received by an attacker. The signal received by the receiver leaks a secret that was supposed to stay within the victim's security domain.

Focusing just on the communication phase of an attack, the *transmitter* is in the victim's code, and the *receiver* is in the attacker's code. The medium of the *communication channel* is the micro-architectural state that can be modified, i.e., *modulated*, by the activity of the transmitter. A communication channel may actually be composed of multiple *subchannels*, just as a radio transmission may use multiple frequencies.

For numerous contemporary attacks, the communication medium is the last-level cache, and each cache line can be considered a communication subchannel. In the simple case of a directly mapped cache, modulating a subchannel involves the transmitter accessing a specific address, since that will change the state of exactly one well-defined cache line. A receiver can *monitor* the state of the same cache line (subchannel) for changes (modulation) by accessing an address to occupy that cache line (subchannel) and, at a later time, measure the latency of a re-access to the same address to determine whether it is a hit or miss.

For a more complex cache, such as a set-associative cache, the receiver needs to use multiple addresses to monitor the cache set that will be used by the transmitter, i.e., multiple subchannels. In other cases, the transmitter might also use multiple addresses, i.e., modulating multiple subchannels. These sets of addresses accessed by the transmitter and receiver are referred to as the *transmitter set* and the *receiver set* respectively. The attacker generates a *receiver set* by using a so-called eviction set construction algorithm [7], [8]. Later we generalize this operation as the process of receiver *calibration*. **Randomly Mapped Caches.** Among various architectural solutions that address security vulnerabilities by disrupting communication via cache-based channels, randomly mapped caches [9]–[12] are considered highly effective with plausible security properties and low performance overhead. Randomly mapped caches aim to significantly increase an attacker's efforts to find a receiver set that monitors all the possible subchannels that might be modulated by the transmitter. They leverage one of the two ideas: make cache behavior non-deterministic by introducing some randomness into the functions used to map memory addresses to cache lines (subchannels) [10], [11], and

dynamically change these functions [9], [10], [12].

In such complex caches, the subchannels that the transmitter will modulate are not publicly known to attackers. Moreover, with non-deterministic caches, the attacker can only guess the probability of an address to be mapped to a given cache line, and modulation can only be observed probabilistically. This uncertainty requires the attacker to use complex methods to generate receiver addresses that have a high probability to monitor the cache lines modulated by the transmitter.

Unfortunately, the security claims of randomly mapped caches are quite fragile. For example, a recent secure cache design, CEASER [9], which claimed to be able to tolerate years of attack, has been broken by more advanced eviction set construction algorithms [7], [10]. Similarly, ScatterCache [11], another recently proposed randomly mapped cache design, can be broken by a new eviction set construction algorithm [13] within a few seconds.

The reason behind the failures of those designs lies in their *limited security analyses*. In fact, those defense mechanisms were designed to block very specific eviction set construction techniques. For instance, some weak security analyses [9]–[11], [13] only consider the case where the attacker tries to obtain a receiver set that monitor the subchannels used by the transmitter with high probability. Such an analysis ignores the existence of alternative strategies where the attacker could spend a modest amount of resources on constructing a receiver set that has a lower probability to monitor these subchannels. With such a weak receiver set, she would rely on repeatedly monitoring the modulations from the transmitter to ultimately leak the secret.

Communication Paradigm. In this paper, we introduce a generalized communication paradigm for micro-architecture side channels. The paradigm serves two purposes. First, it provides the overall view of the communication process and identifies the end-to-end communication steps that a comprehensive security analysis has to consider. Second, it enables us to think of micro-architecture side-channel attacks using concepts from telecommunications, so we can formulate the security analysis into a statistical problem and perform quantitative analyses.

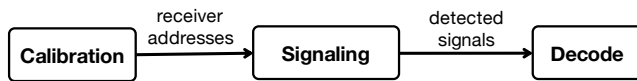


Fig. 1: Communication paradigm.

The communication paradigm, shown in Fig. 1, consists of three steps: *calibration*, *signaling* and *decode*.

First, the receiver often needs to perform a *calibration* step. Calibration is like a tuning process in a radio-based system, and aims to determine which subchannels will be modulated by the transmitter, and where to tune the receiver to monitor those subchannels. For a cache-based channel, the calibration step involves running an eviction set construction algorithm [7]. Prior analyses [9]–[11], [13] have only focused on this step.

The second step is a *signal transfer* step (*signaling* for short) where the receiver obtains the signal from the transmitter.

To obtain the signal, the receiver needs to detect the state changes of the channel caused by the transmitter. In cache-based side channels, the receiver can obtain the signal using various approaches, such as Prime+Probe [2]. The signal—made of cache hit and miss events—can then be formalized mathematically and studied with statistics and probabilities.

Finally, the receiver needs to perform a decode step to interpret the detected signals. The decode step can be straightforward if the detected signal directly corresponds to the secret value. In cache-based channels, this decode step can be complicated if it needs to cope with noise, and non-deterministic behaviors of the cache.

This Paper. We propose **Cache Security Analyzer (CaSA)** to quantitatively analyze the security of randomly mapped caches. We aim to use CaSA to comprehensively evaluate a wide range of communication strategies and cache configurations.

The design of CaSA is built from three insights. First, instead of solely focusing on the calibration step, CaSA performs an end-to-end analysis on the three communication steps in Fig. 1. Second, it leverages telecommunication concepts to formulate the security analysis into a statistical problem and quantify the security by measuring the end-to-end communication cost. Third, CaSA identifies the existence of a trade-off in distributing resources between the calibration step and the signaling step. It explores that trade-off to find the communication strategy that minimizes the overall communication cost.

We use CaSA to evaluate randomly mapped caches of different configurations and discover the *quantitative* impacts of different parameters on the communication cost (Findings 1-4 in Section VII). Furthermore, we have made new observations to better understand the limitations and benefits of randomly mapped caches that refute several common beliefs. We highlight two takeaways here:

- 1) When communicating on randomly mapped caches, spending the maximum amount of resources on calibration is neither the only nor always the best strategy. This refutes the common belief [11], [13] that an effective receiver set must be able to achieve a high eviction rate, i.e., that monitors most subchannels that the transmitter could modulate.
- 2) In the case where dynamic changes in mapping functions are used, information can be leaked and accumulated across mapping function changes. This refutes the common belief that attacks must be completed during the life of a single mapping function [9], [10].

With those insights and quantitative results, we show that the randomization mechanisms used in the state-of-the-art randomly mapped caches [9]–[11] (except for NewCache [12]) are insecure.

The contributions of this paper are:

- A three-step, end-to-end communication paradigm expanding the security analysis of cache-based side channels beyond just the calibration of the receiver.
- Formulating the security analysis into a statistical problem to enable quantitative analysis.

- CaSA, a comprehensive and quantitative security analysis framework of side-channel communication via randomly mapped caches.
- A thorough security evaluation and new observations to understand the limitations and benefits of randomly mapped caches.

II. BACKGROUND

A. Cache-based Side Channel Attacks

In a cache-based side channel attack, the transmitter and the receiver use the cache as the communication channel, and each cache line as a subchannel. Various such attacks exist [1], [8], [14]–[25], and follow the procedure described in Fig. 1.

In each attack, the receiver first performs a *calibration* step by finding a group of addresses called *receiver addresses*. The receiver uses the receiver addresses to monitor a set of subchannels, usually a cache set. Next, the receiver performs the signaling step, which consists of two substeps: *precondition* and *detection*. The receiver *preconditions* a group of subchannels into a known state in order to optimize its chances of monitoring state changes in these subchannels. The precondition generally involves accessing the receiver addresses to fill a cache set. It waits for the transmitter to modulate some of the monitored subchannels by accessing some cache lines. It then *detects* the modulation of those subchannels by either measuring the time of re-accessing the receiver addresses (Prime+Probe [1], [8]), or measuring the time of accessing the transmitter addresses (Evict+Reload [15]), or measuring the execution time of the transmitter (Evict+Time). Finally, it performs the decode step based on the measurement result.

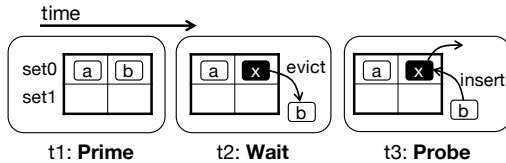


Fig. 2: An example of Prime+Probe attacks. Line *a* and *b* are receiver addresses; line *x* is the transmitter address.

Fig. 2 visualizes an example of using Prime+Probe as the signaling strategy on a two-way cache, which contains three steps: Prime, Wait, and Probe. The receiver preconditions two subchannels in set0 by accessing lines *a* and *b* (Prime). It then waits for the transmitter to modulate a subchannel in set0 by accessing line *x*, which evicts line *b* from that subchannel (Wait). At a later time, the receiver checks the state of the subchannels in set0 by re-accessing lines *a* and *b*, and measuring the access latency (Probe). Based on the long access latency, the receiver knows that line *b* missed in the cache and the state of a subchannel in set0 has been modified (modulated) by the transmitter.

B. Randomly Mapped Cache Designs

The mapping function in a cache decides how memory addresses are mapped to cache sets. Randomly mapped caches introduce randomness into the mapping functions to make

it harder for a receiver to know which subchannels will be modulated by a transmitter, and which subchannels are preconditioned or monitored by the receiver. It aims to mitigate cache attacks by significantly increasing the difficulty of the receiver’s calibration step.

There are various flavors of randomly mapped cache designs [9]–[12], each with different performance and security characteristics. To better understand their differences, we distinguish these designs based on three characteristics of the mapping function, namely whether:

- 1) It uses public or secret hash functions;
- 2) It is static or can be dynamically changed over time;
- 3) It uses a single or multiple hash functions at a point in time.

Table I categorizes each design by mapping strategy.

	Static	Dynamic
Single Hash Group	Set-associative cache* Intel sliced LLC [26]	CEASER [9] NewCache [12]
Multiple Hash Groups	ScatterCache [11]	Skewed-CEASER [10]

TABLE I: Classification of cache mapping strategies. * Uses public hash functions.

1) Public vs. Secret hash functions. Traditional set-associative caches use a public hash function, which simply extracts bits from the physical address and uses them as the set index. The other caches in Table I use secret hash functions. For example, the last-level caches in Intel processors are organized into multiple slices. The mapping function includes an undocumented slice hash function to decide which slice an address should map to. NewCache [12] uses a table-based hash function. CEASER [9] and ScatterCache [11] use encryption-based hash functions.

Even though using a secret mapping function could be thought to make calibration more difficult, it alone cannot thwart cache attacks. It has been demonstrated that there exist efficient algorithms for attackers to reverse engineer the hash function [27], [28] or even to directly construct effective receiver sets [7], [10] without needing to know anything about the mapping function.

2) Static v.s. Dynamic hash functions. To further secure the cache, researchers proposed to periodically change the hash function instead of using a static hash function. A cache with a dynamic mapping function uses one hash function in each *epoch*, and switches to a different hash function at the end of an epoch.

The length of an epoch has a significant impact on the performance and security of the design. To be secure, the epoch should be short enough so that the receiver cannot both calibrate and detect signals within one epoch. However, upon epoch switching, every line in the cache has to be remapped, and using small epochs thus incurs serious performance overhead.

NewCache [12] uses extremely small epochs—changing the hash function every time a cache conflict occurs. CEASER [9] and Skewed-CEASER [10] change the hash function when the number of cache accesses reaches a threshold. The threshold is

configured to be smaller than the number of accesses required by the state-of-the-art eviction set construction algorithm [7]. Skewed-CEASER [10] claims years of security when setting the threshold as $100 \times L$, where L is the total number of lines in the cache.

3) Single vs. Multiple hash functions. Researchers have proposed more advanced secure cache designs, namely *multi-hash caches* such as ScatterCache [11] and Skewed-CEASER [10], which use multiple hash functions at any point in time. These designs contrast with *single-hash caches*, which only ever use a single hash function.

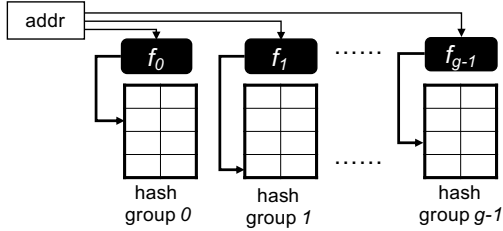


Fig. 3: A cache with multiple hash groups.

As shown in Fig. 3, a multi-hash cache is organized as multiple hash groups. Each hash group is organized as a normal set-associative cache, and uses a distinct hash function. To do a lookup in the cache, all the hash-groups are looked up, with at most one of them being a cache hit. On a cache miss, the cache first picks one of the hash groups using a uniformly *random* policy and then uses the corresponding hash function to generate the set index for that hash-group.

As a result, the mapping function becomes *non-deterministic*. An address can end up in different hash groups, i.e., modulating different subchannels, even within one epoch. It significantly increases the attacker's difficulty in obtaining a receiver set to monitor all the subchannels that will be used by the transmitter.

ScatterCache [11] uses a single-way per hash group design. Skewed-CEASER [10] makes the number of ways per hash group a configurable parameter.

III. THREAT MODEL AND SCOPE

We follow the standard threat model of cross-core cache-based side channel attacks. We assume the attacker and the victim are co-located on the same processor chip, but reside on different cores. A transmitter embedded in the victim and a receiver controlled by the attacker communicate via channels in a shared last-level cache. Even though we focus on the last-level cache, our analysis and our tool, CaSA, can be easily extended to other levels of the cache hierarchy.

The attacker can reside in a user-level process or in a malicious operating system in a secure enclave context, such as SGX [29]. Like previous work [20], we assume the receiver can use a single thread or multiple threads to control multiple cores on the chip. The transmitter may be latent in the code of the victim and execute as part of the victim's normal processing, or the attacker can leverage speculative execution [3] to provoke the execution of the transmitter.

Scope. Our analysis focuses on investigating the fundamental problems in the randomization schemes used by randomly mapped caches. Prior work [30] has shown that the mapping function used in CEASER [9] and Skewed-CEASER [10] only consists of linear functions and has a key invariant vulnerability, that is, changing the key used in the mapping function cannot change the collisions between addresses. This vulnerability can be fixed using non-linear hash functions. Note that, our analysis is independent of which hash function is used and studies new vulnerabilities that have not been explored in prior work [30]. Indeed, we focus on analyzing the fundamental problem that is intrinsic to randomly mapped caches.

Besides, we consider the analysis of the following two types of attacks orthogonal to the analysis of randomly mapped caches: flush-based attacks [14], [18] and occupancy-based attacks [31]. The reason is that randomly mapped caches are not designed to and thus are unable to mitigate these attacks. Hence, we do not analyze such attacks in this paper.

IV. MOTIVATION

Correctly reasoning about the security of randomly mapped caches is challenging. Prior security analyses have made incorrect security claims by narrowly considering two communication strategies by the attacker. We claim that a comprehensive security analysis should provide an end-to-end quantitative analysis of a broad range of communication strategies.

A. Limitations of Prior Work

Prior security analyses [10], [11], [13] only targeted specific calibration strategies that require a huge amount of resources and are unlikely to be completed within one epoch. We use the following simple example in Fig. 4 to illustrate the limitations of their analyses. Fig. 4 compares the results of three different calibration strategies on a cache with 2 hash groups and 1 way in each group. Each figure shows hash group 0 on top, hash group 1 at the bottom, and how the transmitter and receiver addresses are mapped to corresponding subchannels. The subchannels that can be used by the transmitter address are marked in grey.

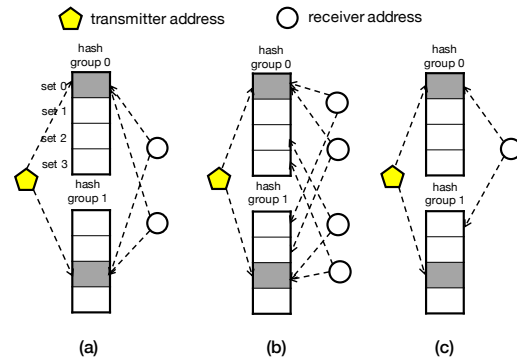


Fig. 4: An illustrative example of different calibration strategies: (a) hard-conflict receiver addresses; (b) many soft-conflict receiver addresses; (c) one soft-conflict receiver address.

The security analysis in Skewed-CEASER [10] only considered using “hard-conflict” receiver addresses for signaling. A “hard-conflict” receiver address maps to the same cache set as the transmitter address in *every* hash group, shown in Fig. 4(a). Their analysis only consider such addresses because once the attacker has enough hard-conflict addresses (2 in this example), she can perform the rest of the communication (e.g., Prime+Probe) in the same way as on a conventional cache. Randomly mapped caches are designed to make it extremely difficult to obtain hard-conflict addresses. In fact, for a given transmitter address, when there are 8 or 16 hash groups, there may not exist enough hard-conflict addresses given the limited size of the address space in the state-of-the-art systems [11]. Even though a receiver set with hard-conflict addresses is guaranteed to be functional, i.e., guarantee to monitor all the subchannels that will be used by the transmitter, it is not the only way to communicate on randomly mapped caches.

The security analysis in ScatterCache [11] considered using a large number of “soft-conflict” receiver addresses. A “soft-conflict” receiver address maps to the same set as the transmitter address in *at least one* hash group, as shown in Fig. 4(b). Soft-conflict addresses are much easier to find than hard-conflict addresses but can only be used to monitor the transmitter with some probability. The assumption behind their analysis is that the attacker needs to get a large receiver set (e.g., 256 addresses on an 8MB LLC) in order to monitor the transmitter with 99% probability. Crafting such a large receiver set is expensive and unlikely to be completed within one epoch. Consequently, strong security claims were made under such assumptions.

There exists a key problem with these analyses: they overlooked a broad range of communication strategies that are available to the attacker. In addition to the prior analysis where the receiver spends a huge amount of resources on calibration to achieve a high monitoring probability, other effective communication strategies are also possible, such as, using a small amount of resources on calibration to obtain a receiver set with low monitoring probability, and relying on repeating the signaling step to decode secrets with a high success rate. A comprehensive analysis should explore the *trade-off in distributing resources between calibration and signaling*.

Fig. 4(c) shows an example of using 1 receiver address that soft-conflicts with the transmitter on a single subchannel. Such a receiver set is fairly cheap to construct. In this example, the receiver has a probability of 0.5 to monitor that subchannel and the transmitter also has a probability of 0.5 to modulate that subchannel. As a result, when the transmitter address is accessed, the probability of the receiver observing a modulation is $0.5 \times 0.5 = 0.25$. When the transmitter is not accessed, this probability is 0. Even though the probability to observe a modulation is low, the receiver can repeat the signal transfer step to accumulate samples. Those samples are then used to infer if the transmitter was accessed (observing some modulation) or not (observing no modulation). This last phase is the decoding step and increasing the number of samples will increase the decode success rate. In our example,

by accumulating 16 samples, the receiver can know if the transmitter was accessed or not with 99% confidence, based on whether it detects modulations in at least one of the signaling samples or it detects no modulations across all samples.

Alternatively, the receiver could spend more resources on calibration to obtain two receiver addresses instead of one. In this situation, she would only need to accumulate 7 samples to decode the secret with the same level of confidence. The examples above clearly demonstrate the existence of a trade-off in distributing resources between calibration and signaling.

B. The Need for Comprehensive and Quantitative Analysis

In addition to the trade-off between calibration and signal transfer, we find it is necessary to perform a comprehensive and quantitative analysis of randomly mapped caches, since there exist multiple other factors that can affect the security of these designs. We provide the intuitions of how these factors can affect communication on randomly mapped caches below.

First, we need to consider the effects of having multiple transmitter addresses. Intuitively, having more transmitter addresses can make communication easier, because the number of subchannels associated with the transmitter increases and the communication can work as long as the receiver can successfully monitor at least one modulation from the transmitter. Note that, in practice, multi-address transmitters do occur in many security-sensitive applications. For example, the square-and-multiply exponentiation algorithm used in RSA encryption [32] acts as a multi-address transmitter: both the square and multiply functions are composed of instructions residing in multiple cache lines.

Second, we need to consider the effects of noise. Intuitively, the presence of noise can make communication more difficult, because the receiver often cannot distinguish the modulations generated by the transmitter or by the noise. CaSA quantitatively measures the impacts of noise and we discovered a new finding that noise can have a positive impact on communication.

Finally, for caches that periodically change the mapping functions, we investigate the feasibility of performing the communication across epochs. Prior work assumed that communication must complete within one epoch and no information can be carried across epochs. In this paper, we challenge this assumption. It is true that, a receiver set constructed in an epoch can only be used for the signaling steps in the same epoch. However, we observe that different receiver sets from different epochs generate signals—made of cache hit and miss events—that are similar to each other, since the signals are mainly determined by the numbers of addresses in the receiver sets. Intuitively, if the same secret bit is transmitted, the samples obtained from different epochs can be combined to increase decoding accuracy.

CaSA is designed to quantitatively analyze the impacts of the above factors on the security of randomly mapped caches. Specifically, CaSA can answer the following questions.

- Given a cache configuration, such as the one in Fig. 4, and the number of transmitter addresses, how should a receiver

distribute resources between calibration and signal transfer to exfiltrate the maximum amount of information?

- Considering background noise, how much more difficult is it for an attacker to mount a successful attack?
- Among different cache configurations (e.g., 1-way per hash group and 2-way per hash group), which one is more difficult to attack, measured by the number of attacker's cache accesses to leak one secret bit?

V. CASA OVERVIEW

The goal of CaSA is to measure the security of different configurations of randomly mapped caches. We strive to comprehensively evaluate how various communication parameters quantitatively affect the amount of information leakage on a given cache configuration. To enable quantitative analysis, we innovatively leverage concepts from the field of telecommunications (Section I) and formulate the signals in cache-based side channels into a statistical representation. In this section, we first describe the full security analysis space, and then describe the statistical representation of signals, followed by the security metric used in CaSA.

A. The Security Analysis Space

The paper strives to comprehensively evaluate the choices available with respect to the three components used in cache-based side channel communication, i.e., transmitter, receiver and channel (i.e., cache), as well as parameters related to noise. **Transmitter.** An important parameter related to the transmitter is the number of transmitter addresses. We expect program developers to set that sole transmitter parameter based on their knowledge of the applications or using program analysis tools [33]–[35].

Receiver. The receiver can choose from a wide range of calibration, signaling and decoding strategies, and it can accumulate information *across epochs* on caches periodically changing their hash functions. We investigate various possible combinations of calibration, signaling and decoding strategies, especially considering the case that the receiver spends medium to low resources on calibration.

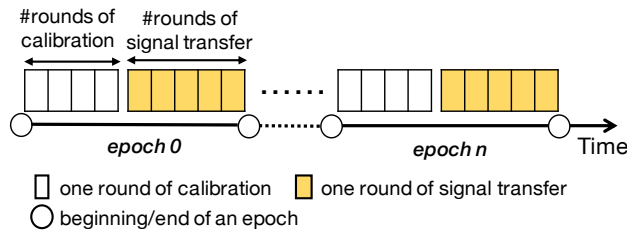


Fig. 5: Attack procedure on a multi-hash cache that periodically changes hash functions.

Fig. 5 visualizes the communication process on a multi-hash cache that periodically changes hash functions and indicates receiver parameters. The cache changes the hash functions at the end of each epoch (marked as circles). Within each epoch, the receiver generates a receiver set via multiple rounds of calibration (white boxes), and then uses the receiver set

to perform signal transfer and collect signal samples once or multiple times before the epoch ends (highlighted boxes). The receiver strategy decides how to distribute efforts between calibration and signal transfer

If enough samples have been acquired within one epoch to allow decoding of the secret value with sufficient certainty, the communication of the secret is complete. If the number of samples acquired is insufficient to decode the secret, the attacker will start a new epoch with a new calibration step, acquiring more samples.

Our analysis assumes that the epoch size is constant and epoch changes are public to the receiver. It is useful to study the security of randomly mapped caches independently of complicated epoch parameters, especially as the security of those caches is not believed to rely on hiding epoch parameters. We discuss how CaSA can be extended to analyze detecting epoch changes and handling variable-size epoch in Section VIII. **Channel (Cache).** We consider randomly mapped caches with varied configurations, in terms of the number of hash groups, the number of ways in each hash group, the number of cache sets, and the epoch length. To make our analysis tractable, in this paper, we do not consider other cache parameters such as the ones related to number of cache levels, directories, or MSHRs. However, note that CaSA can be extended to incorporate those parameters.

Noise. In a cache attack, noise can add spurious modulations and confuse the receiver. We identify and consider two types of noise. The first type is the background noise, which consists of random addresses and modulates random subchannels.

The second type is the *carrier* noise, which modulates a fixed set of subchannels independently from the secret. Taking the following victim code for example, `if (secret) {access A; access B;} else {access A;}`. Address A is a carrier address. If the receiver is calibrated on the subchannels mapped to A, it will waste resources monitoring subchannels that will not provide any useful information about the secret. As a result, carrier noise makes communication more difficult.

A summary of the parameters of all the communication components considered in this paper is shown in Table II.

Component	Parameters
transmitter	number of transmitter addresses
receiver	number of rounds of calibration in one epoch number of rounds of signal transfer in one epoch
channel (cache)	number of hash groups number of ways per group number of sets epoch length size of upper-level caches
noise	background noise carrier noise

TABLE II: The communication parameters considered in CaSA.

B. The Statistical Representation of Signals

We model the signal observed by a receiver as a random variable X , counting the number of modulations detected by the receiver during a signal transfer step. X follows a probability

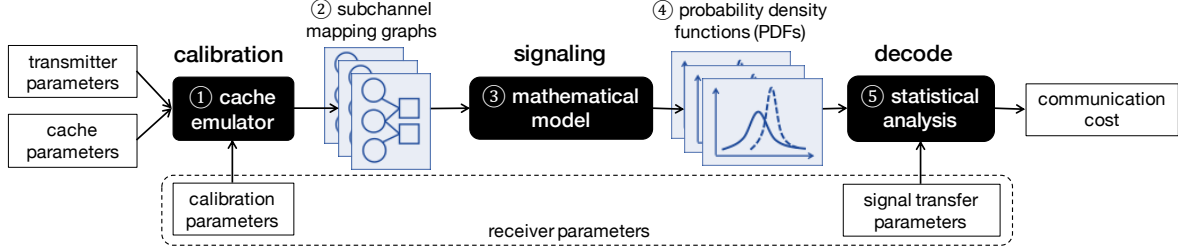


Fig. 6: CaSA: end-to-end quantitative security analysis framework.

distribution that can be characterized by a *probability density function* (PDF for short) $f(n) = P(X=n)$. We also note $F(n) = P(X \geq n)$ the *cumulative density function* (CDF for short), is sometimes more convenient to use.

To give a concrete example, let's consider a transmitter that communicates a secret bit to a receiver by modulating one subchannel to send bit "1" and doing nothing to send bit "0". We use $f_0(n)$ and $f_1(n)$ to represent the density functions for X when the bit sent is "0" or "1" respectively. To decode the signal, the receiver samples X to determine whether X follows f_0 or f_1 . Note that, one of the key tasks of CaSA is to compute the PDFs or CDFs for a given communication configuration.

When communicating on a multi-hash cache, if the receiver uses soft-conflict addresses, such as in Fig. 4(b) and (c), she will be only able to monitor the subchannels used by the transmitter with a certain probability. The corresponding PDFs are as below, with $p > 0$.

$$f_0(n) = \begin{cases} 1, & \text{if } n=0 \\ 0, & \text{otherwise} \end{cases}; \quad f_1(n) = \begin{cases} p, & \text{if } n=0 \\ 1-p, & \text{if } n=1 \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

In a noiseless environment, both $f_0(n)$ and $f_1(n)$ have non-zero values at $n=0$. Visually, the two functions partially overlap with each other. In the examples in Fig. 4, when using one soft-conflict receiver address, $p=0.25$, and when using a large number of soft-conflict receiver addresses, the value of p can decrease to 0.01 (i.e., $f_1(1)=1-p=0.99$). The smaller the value of p is, the easier the two distributions can be distinguished. In a noisy environment, even when the transmitter does nothing, the receiver can observe modulations which are generated by noise. Therefore, the corresponding PDF of the received signal, $f_0(n)$, will have non-zero values at $n \geq 1$. Moreover, if the transmitter and the receiver are composed of multiple addresses, the PDFs can have non-zero values at $n \geq 2$.

With the two PDFs partially overlapped, the decoding step becomes complicated, but still feasible. As discussed in Section IV, if the receiver can collect enough samples of the signal, she can decode the secret bit with a high success rate, e.g., 99%.

C. The Security Metric

To evaluate randomly mapped caches and compare different cache configurations, we propose to use *end-to-end communication costs* as the quantitative security metric. Recall that,

prior works [9]–[11], [13] analyze randomly mapped caches by quantifying the difficulty to perform the calibration step and have led to misleading security claims. Our end-to-end communication cost consists of the receiver's cost on the calibration step and the signaling step. Specifically, in CaSA, we use the number of cache accesses required by the receiver to decode the secret with a 99% confidence, where the cache accesses include the accesses performed by the receiver during the calibration step and the signal transfer step. Note that, other resources can also be used to express a cost, such as time or the number of times the victim is triggered.

VI. CAsA IMPLEMENTATION DETAILS

CaSA is an end-to-end quantitative security analysis framework for communication via randomly mapped caches. Note that, even though we designed the framework for randomly mapped caches, it can be easily used to analyze other simpler cache designs.

A. CaSA Work Flow

CaSA is composed of three modules analysing the three identified steps in the communication process: calibration, signaling and decode, shown in Fig. 6. It can be used to explore a large security analysis space listed in Table II and compute the *communication cost* for various communication parameters.

The first module is the calibration module (①), which uses a cache emulator to simulate the cache's behavior during calibration. It takes the transmitter parameters and the cache parameters as input, and generates the calibration result. Due to the random behavior of the cache, the module runs the calibration algorithm multiple times and each run generates a receiver set. We encode the calibration result (a group of receiver sets) as *subchannel mapping graphs* (②). The graph representation is to precisely capture the mapping relationship between addresses and subchannels (cache lines).

Fig. 7 shows an example of a subchannel mapping graph on a cache with 4 sets and 2 hash groups. The graph is a directed bipartite graph with two disjoint sets of vertices for addresses (pentagon and circle) and subchannels (square). An edge always connects an address vertex to a subchannel vertex, indicating the address can map to the subchannel. An address vertex can either be a transmitter address (pentagon) or a receiver address (circle). The graph does not include the subchannels which no address maps to, such as set 2 in hash group 1. There may exist multiple connected components in the graph, depending

on the conflict relationships between addresses, such as the two shaded areas in Fig. 7.

The second module is the signaling module (③), which takes a subchannel mapping graph as input and uses a mathematical model to compute the signal transfer result. For each value the secret can take, the signaling module outputs a signal probability density function (④), which describes the distribution of the number of modulations observed by the receiver.

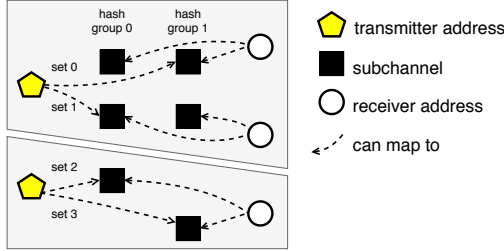


Fig. 7: An example of subchannel mapping graph.

The last module is the decode module (⑤), which takes the probability density functions (PDFs) as input and computes the end-to-end communication cost of the receiver (Section V-C). It uses a statistical analysis method to compute the number of accesses needed by the receiver to decode the secret with a given confidence value, e.g., $\geq 99\%$. Note that, it can also measure the cost for communicating across epochs.

To evaluate the security of a cache configuration, we use the above framework to compute the communication cost for different combinations of receiver parameters and find the one with the lowest cost.

We now provide details for each module.

B. The Calibration Module

The calibration module uses a cache emulator to model the state-of-the-art calibration algorithms. These algorithms are eviction set construction algorithms proposed by Qureshi et al. [10] and Purnal et al. [13]. We generalize the algorithms into three steps, shown in Fig. 8.

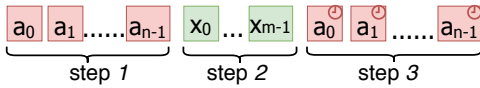


Fig. 8: The calibration algorithm for multi-address transmitters. a_0 to a_{n-1} are candidate addresses, and x_0 to x_{m-1} are transmitter addresses.

The calibration starts with a *candidate set* which is composed of many randomly chosen addresses. The candidate set should contain enough addresses so that some of them map to the subchannels used by the transmitter addresses.

- 1) The receiver accesses the addresses in the candidate set, making them all reside in the cache. This step requires multiple accesses to each candidate address to ensure every access hits in the cache. It also requires dropping some addresses if the candidate set cannot fit in the cache.
- 2) The transmitter addresses are accessed, which potentially evict some of the candidate addresses from the cache.

- 3) The receiver re-accesses the candidate set and measures the access latency of every address. Based on the latency, the algorithm decides whether a candidate address should be included in the receiver set or not.

In step 3, the algorithm can either add all the addresses that missed in the cache to the receiver set, or add only the first address that missed to the receiver set. We call the former one as a *greedy* calibration strategy and the latter one as a *non-greedy* strategy. We evaluate the cost of both calibration strategies in Section VII-A.

C. The Signaling Module

The signaling module takes the subchannel mapping graph as an input, models the cache behaviors during the signal transfer step, and computes the distributions of the signals for different secret values. Recall that, the distributions of the signals are characterized by probability density functions (PDFs) or cumulative density functions (CDFs) of the number of modulations observed by the receiver (Section V-B).

The signaling step consists of three operations: 1) the receiver performs precondition and monitors a group of subchannels; 2) the transmitter modules another group of subchannels; 3) the receiver performs detection and observes modulations on the subchannels that are *both* monitored by the receiver in step 1 and are modulated by the transmitter in step 2.

Correspondingly, given a set of subchannels, our mathematical model follows three steps to compute the probability of the receiver detecting modulations on these subchannels. Given a single subchannel s , we define the event D_s when a modulation is *detected* on the subchannel s . Similarly, the event $\bigcap_{i=1}^n D_{s_i}$ is the receiver detecting modulations on subchannels $\{s_1, \dots, s_n\}$. The three steps are as follows.

- 1) Compute the probability that the given subchannels are monitored by the receiver, noted as $P_r(s_1, \dots, s_n)$, using a Markov chain approach.
- 2) Compute the probability that these subchannels are modulated by the transmitter, noted as $P_t(s_1, \dots, s_n)$, using a simple probability calculation.
- 3) Compute the probability that the receiver detects the modulations on these subchannels by calculating the joint probability from the above two steps. Since a subchannel being monitored and being modulated are *independent* events, their joint probability is the product of their individual probabilities: $P(\bigcap_{i=1}^n D_{s_i}) = P_r(s_1, \dots, s_n) \times P_t(s_1, \dots, s_n)$.

Once we obtain the detection probability for a given set of subchannels, we can compute the cumulative density functions (CDFs) of the signals. Basically, to compute $F(n)$, we enumerate all the sets of subchannels of size n and accumulate their detection probabilities.

Note that, for simplicity, the following discussion and formulas assume each hash group has a single way. The approach is applicable to caches with multi-way hash groups.

1) *Compute Monitoring Probability:* The precondition operation generally involves accessing the receiver addresses multiple times to *ensure* all the addresses are cached. Due to

the random behavior of the cache, the precondition step is essentially a *stochastic process*.

Given a subchannel mapping graph, several methods can be used to compute the monitoring probability (P_r). For instance, we could use a cache emulator to simulate the precondition process and empirically obtain the monitoring probability. An alternative approach is to model this process as a Markov chain [36] as below.

To construct the state transition graph for the Markov chain, we enumerate all the cache states that can be reached during preconditioning and make each Markov state correspond to one of the cache states. The probability of transitioning between any two Markov states is determined by the address being accessed and the probability that the address ends up in a given subchannel, which can be known from the subchannel mapping graph. The precondition step finishes when all the receiver addresses are in the cache, which corresponds to absorbing states in the Markov chain. We can use standard approaches to compute the probability of reaching each of absorbing states. To compute the monitoring probability P_r for a given set of subchannels, we enumerate the absorbing states where these subchannels are monitored and accumulate the probabilities of reaching these states.

The Markov chain approach provides useful insights on the interactions between receiver addresses and computes the precise monitoring probability. However, it suffers from high computation complexity. The number of states in the Markov chain increases exponentially as the number of subchannels. When the subchannel mapping graph is big, using a cache emulator can be more efficient in computing the monitoring probability.

2) Compute Modulation Probability: The modulation operation involves accessing the transmitter addresses for a fixed number of times. Considering the number of transmitter addresses and assuming they do not conflict, we can compute the modulation probability (P_t) with a simple model. For each subchannel s that can be used by the transmitter as $P_t(s) = 1/g$, where g is the number of hash groups. In the case when we have a high number of transmitter addresses and these addresses share subchannels, other approaches like the Markov chain approach in Section VI-C1 or simulation can be used.

Handling Noise. We model the impact of two types of noise (Section V-A): background noise and carrier noise. Both contribute to the modulation probability (P_t).

The background noise consists of accessing randomly chosen addresses which modulate random subchannels. We model a noise access as modulating each subchannel with a probability of $1/L$, where L is the number of cache lines. The carrier noise is a part of the transmitter. Therefore, we model the carrier noise as transmitter addresses.

3) Compute Density Functions: Now that we have the probability for each given set of subchannel to be monitored by the receiver and modulated by the transmitter, we can compute the joint probability of these two events happening simultaneously i.e. the probability of detecting modulation on

these subchannels.

Next, we compute the cumulative density functions $F(n)$ by enumerating all the sets of subchannels of size n and summing their detection probabilities. Consider $F(1) = P(X \geq 1)$, which gives the probability of the event “at least one modulation being detected”. This event can be formulated as a union of the same event on each subchannel. More explicitly, for each subchannel s , the event “a modulation being detected on s ”, denoted as D_s . We note \mathcal{S} of size N the set of all subchannels. We have :

$$\text{Event “} X \geq 1 \text{”} = \bigcup_{s \in \mathcal{S}} D_s$$

We then apply the principle of inclusion-exclusion [37] to compute $F(1)$ as below.

$$F(1) = P\left(\bigcup_{s \in \mathcal{S}} D_s\right) = \sum_{k=1}^N \left((-1)^{k-1} \sum_{\{s_1, \dots, s_k\} \subseteq \mathcal{S}} P\left(\bigcap_{i=1}^k D_{s_i}\right) \right)$$

We find it extremely expensive to compute the detection probability for every possible subset of \mathcal{S} . Indeed, it is incomputable for large subchannel mapping graphs. To greatly reduce the computation complexity, we use Bonferroni’s inequalities [38] to solely compute bounds for the density functions. Specifically, to compute $F(1)$, we cut the sum at $k=1$ and $k=2$ to get the upper bound and the lower bound respectively. With the same principle, we can derive bounds of $F(n)$ for any n .

Using the bounds makes the density functions imprecise, and can affect our estimation of the communication cost. When the bound is too loose, we rely on simulation of the signaling step to empirically derive the CDFs.

D. The Decode Module

The decode module takes the density functions of the signal for each possible secret value as an input and computes the end-to-end communication cost of the receiver, which consists of the calibration cost and the signaling cost. The calibration cost can be directly derived from the calibration module by counting the number of accesses performed by the cache emulator. The signaling cost is computed by the decode module using a statistical method. Specifically, we compute the number of rounds of signal transfer that are needed by the receiver to decode the secret with 99% success rate and then convert it into number of cache accesses. We now describe how to compute the signaling cost, followed by the discussion on how to quantify the cost if communication spans across epochs.

The decode step consists in solving a statistical problem. Consider the transmitter communicates a secret bit $b \in \{0, 1\}$ to the receiver by sending a signal X . The signal X is an integer random variable that follows its CDF $F_b(n)$, which is either equal to $F_0(n)$ or $F_1(n)$. The receiver decodes the secret by sampling X and deciding which one of the two distributions X follows. Intuitively, the more samples the receiver gets, the more accurately it can decode the secret. The problem that we need to solve is how many samples are needed to distinguish the two distributions with a certain confidence. This is a standard

statistical problem with various solutions. We use an intuitive approach as follows.

To make the mathematical analysis simple, we convert the integer random variable X to a simpler signal, a Boolean variable Y . We define Y , equal to 1 if $X \geq 1$ (that is observing at least 1 modulation) and equal to 0 otherwise (observing no modulation). Hence the problem becomes: how many samples are needed to distinguish between two Boolean distributions of mean $F_0(1)$ and $F_1(1)$.

As a result, the decode strategy is straightforward. The attacker will simply perform the signal transfer several times, observe if the empirical average of Y is closer to $F_0(1)$ or $F_1(1)$, and guess the value of b accordingly.

Note that we could look at Y' equal to 1 if $X \geq 2$ (that is observing at least 2 modulation) and equal to 0 otherwise (observing 0 or 1 modulation). In some cases, Y' can be a better distinguisher than Y . In practice, we look for the value of n that maximizes the distance $|F_0(n) - F_1(n)|$, denoted as n_{\max} and define Y as equal to 1 if $X \geq n_{\max}$, 0 otherwise.

Intuitively, the more samples we get, the closer our empirical mean will get to $F_b(n_{\max})$, and the more confidence we will have for the decode result. The Chernoff-Hoeffding bound [39] provides the relationship between the number of samples and the upper bound of the decode error rate (i.e., the lower bound of the certainty) as below.

$$P\left(\left|\frac{1}{N} \sum_{i=0}^N y_i - F_b(n_{\max})\right| > \delta\right) \leq e^{-2\delta^2 N} \quad (2)$$

where N is the number of samples, y_i is a sample of the Boolean variable Y . The above formula shows that the empirical mean of Y gets closer to $F_b(n_{\max})$ exponentially fast with the number of samples N .

To compute the signaling cost, we compute the number of samples (N) needed to achieve 1% error rate, which can be obtained by setting $\delta = |F_1(n_{\max}) - F_0(n_{\max})|/2$ and $e^{-2\delta^2 N} = 1\%$ in Eq. (2).

Communication spanning across epochs. If the receiver cannot collect enough samples to achieve the required error rate within one epoch, the receiver has to decode based on samples gathered from multiple epochs. As shown in Fig. 5, since the mapping function used by the cache is changed upon switching epochs, the receiver needs to redo the calibration in each epoch to generate a different receiver set, which leads to a different distribution of the signal. This is one of the reasons that prior work [9], [10] considers it infeasible to communicate across epochs.

We claim that it is viable to communicate across epochs, because the signals from different epochs follow a global distribution that can be leveraged for the decode step. Let's use g_i to denote the subchannel mapping graph for the i th epoch and X_i for the corresponding signal. For specific calibration and cache parameters, we define the space of all subchannel mapping graphs as G and the space of the corresponding signals as \mathcal{X} . Intuitively, when the receiver performs the calibration, no matter in which epoch, it obtains a sample from the subchannel mapping graph space G . Similarly, when the receiver collects

a sample from any of the epochs, it is sampling the signal space \mathcal{X} . We call the distribution that \mathcal{X} follows the global distribution. CaSA computes the global distribution of signals and use Hoeffding bound [39] to compute the signaling cost.

We now describe how to compute the global distribution. Until now, we have computed the distribution of the signal X_i conditioned on the subchannel mapping graph g_i , denoted as $P(X_i \geq n | g_i)$. Theoretically, given the probability of generating each subchannel mapping graph $P(g)$, we can compute the global distribution of signal \mathcal{X} as below.

$$P(\mathcal{X} \geq n) = \sum_{g \in G} P(g) \times P(X \geq n | g) \quad (3)$$

With the global distribution, we can transform the signal \mathcal{X} into a Boolean random variable as before. We then apply the Hoeffding bound [39] to compute the signaling cost, i.e., the number of samples required to decode the secret with 1% error rate.

In practice, we do not directly compute Eq. (3), because we are unable to obtain the probability of generating each subchannel mapping graph $P(g)$ due to the extremely large space of G . Instead, we use an empirical approach. We have already obtained samples of subchannel mapping graphs in the calibration module and computed the conditional distributions in the signaling module. We found that if using the same calibration strategy, the conditional distributions $P(X \geq n | g)$ are fairly similar across epochs. Therefore, we can obtain a useful approximation of the global distribution using a small number of samples, e.g., using 30 subchannel mapping graph samples for 15 transmitter addresses. We show communication across epochs is feasible in Section VII-C.

VII. EVALUATION

We use CaSA to evaluate 1MB caches with 1024 sets and 16 ways. We evaluate 3 cache configurations: a) 16 hash groups with 1 way per group, b) 8 hash groups with 2 ways per group, and c) 4 hash groups with 4 ways per group. Within each hash group, if there exists multiple ways, Last Recent Used (LRU) replacement policy is used. For caches that dynamically change the hash functions, we define the length of an epoch using the number of *epoch units*. In each epoch unit, the cache is accessed for L times where L is the total number of LLC lines. We evaluate the state-of-the-art calibration strategies [10], [13] and the classical signaling strategy, i.e., Prime+Probe.

CaSA measures and compares the communication cost of different attack strategies on randomly mapped caches. In this section, we first compare different calibration strategies. Second, we show how calibration parameters, cache parameters and noise parameters quantitatively affect the signaling cost. Finally, we show evaluation results of communication spanning across epochs.

A. Comparing Calibration Strategies

We compare *calibration efficiency* of using different calibration strategies in Fig. 9. The calibration efficiency is measured by the number of receiver addresses generated per epoch unit.

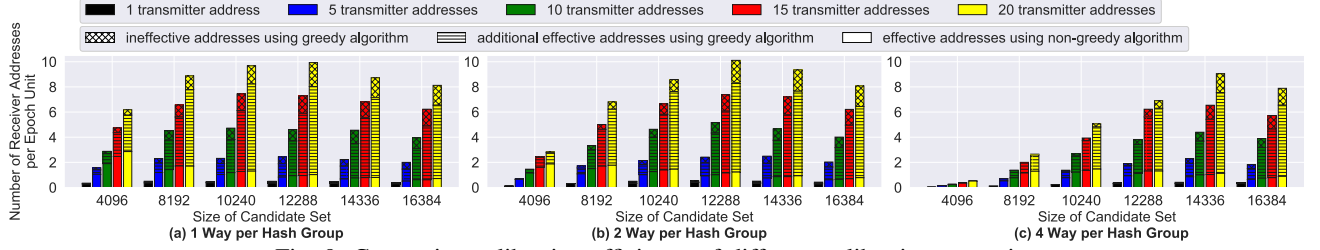


Fig. 9: Comparing calibration efficiency of different calibration strategies.

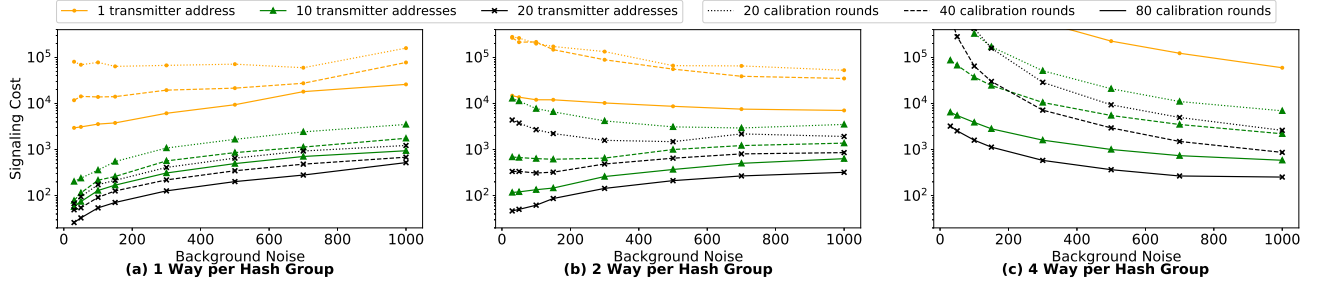


Fig. 10: Impacts of communication parameters on signaling cost.

We assume the attacker chooses the size of the candidate set that is smaller than the private caches under its control. For example, if the receiver attacker can launch two threads and use two 256KB private caches, it could use 8192 candidate addresses.

For each candidate set size, we show from left to right, the calibration efficiency when the transmitter is composed of 1, 5, 10, 15 and 20 addresses. Each bar is broken into three categories from bottom to top: the number of receiver addresses obtained using the non-greedy calibration algorithm, the additional number of effective addresses obtained using the greedy algorithm, and the number of ineffective addresses obtained by the greedy algorithm. Ineffective addresses do not map to any of the subchannels associated with the transmitters.

On the 3 cache configurations, the greedy algorithm consistently obtains more receiver addresses than the non-greedy algorithm when there is more than 1 transmitter address. However, the greedy algorithm introduces 5% to 20% ineffective addresses into the receiver sets. The calibration efficiency of the greedy algorithm increases almost *linearly* with the number of transmitter addresses. In the following evaluation, we use the greedy calibration algorithm.

Finding 1: Calibration efficiency increases almost linearly as the number of transmitter addresses increases.

B. Measuring Signaling Cost

We evaluate the impacts of communication parameters on signaling cost, including transmitter parameters, calibration parameters, cache configurations and background noise. Fig. 10 compares the signaling cost for achieving 1% error rate on 3 different cache configurations. The signaling cost is the number of samples computed using the Chernoff bound (Eq. (2)) on the empirical density functions which we obtain via sampling. In each plot, we show how the signaling cost changes with the background noise, which is modeled as accessing a certain

number of random addresses. We compare the signaling cost for different numbers of transmitters and calibration parameters.

Across the three cache configurations, more transmitter addresses and more calibration efforts both lead to lower signaling cost. On a cache with 1 way per hash group in Fig. 10(a), the signaling cost increases almost exponentially as the noise increases when the transmitter is composed of 1 address. When there are more transmitter addresses and noise is low, the signaling cost increases sub-exponentially. However, an interesting finding from our evaluation results is that noise does not always have negative impacts on communication. On caches with multiple ways per hash group in Fig. 10(b) and (c), increasing noise sometimes helps decrease signaling cost. This phenomenon can be explained intuitively using the following example. Consider a cache with 2 ways per hash group and both the transmitter and the receiver use one address to communicate. Without noise, the receiver cannot detect any modulation no matter whether the transmitter address is accessed or not. With a single noise access, the receiver has a chance to observe a modulation when both the transmitter address and the noise address modulate the subchannel that it monitors, and still no chance to observe a modulation when the transmitter is not accessed. Hence, adding noise in this example has a positive impact on communication.

Comparing the three cache configurations, we do not see a particular cache configuration has a clear advantage of security over others. When there is light background noise, the signaling cost is higher when the cache has more ways per hash group. However, communication on such caches can tolerate more noise. For example, when using 20 transmitter addresses, 80 rounds of calibration and 1000 accesses as background noise, the signaling cost on the cache with 2 ways per hash group (321 samples) is lower than the cost on the cache with 1 way per hash group (521 samples).

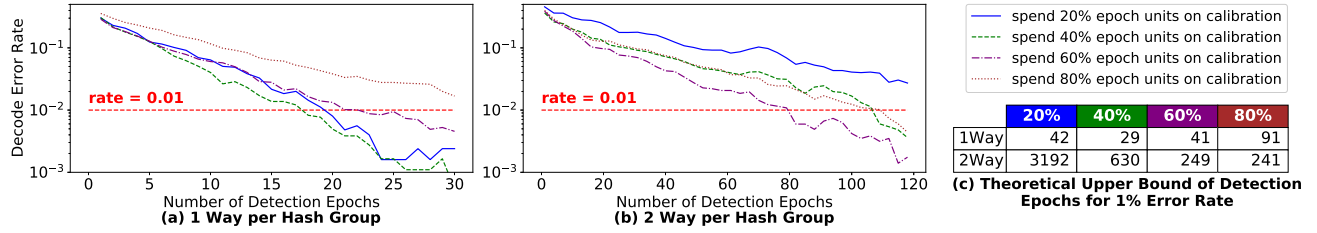


Fig. 11: Empirical decode error rate (a, b) and theoretical bound of communication cost (c) when communicating across epochs.

Finding 2: Signaling cost on caches with 1 way per hash group increases exponentially, except at low noise for multiple transmitter addresses.

Finding 3: Signaling costs on different cache configurations are mostly at the same order of magnitude.

Takeaway: Noise does not always have negative impacts on communication on randomly mapped caches.

Takeaway: There does not exist a cache configuration that has a clear advantage of security over others.

C. Communication Across Epochs when Attacking RSA

We show the feasibility of communicating across epochs using transmitter parameters from a real victim application, the square-and-multiply exponentiation function [32] in the RSA encryption algorithm. The receiver tries to distinguish whether the transmitter executes the square or multiply function. We use Pin [35] to identify 16 transmitter addresses (at cache line granularity) exclusively used by the square function and 10 carrier addresses shared by the two functions.

The end-to-end communication cost is the number of LLC accessed by the receiver during calibration and signaling. The calibration cost is directly derived from the decode module (results in Fig. 9). The signaling cost is computed by multiplying the number of signaling samples and the number of LLC accesses needed to obtain each sample. The number of LLC accesses per sample can be very different depending on whether the receiver can use the `clflush` instruction. To repeat the signaling step within one epoch, which we call *contiguous signaling*, the receiver needs to evict the receiver addresses from the cache before the next signaling round begins. This *self-eviction* operation can be completed using the `clflush` instruction with negligible cost and according to our evaluation, the communication can *always* complete within one epoch with 1% error rate. However, in the case that `clflush` is unavailable, the self-eviction operation requires accessing many random addresses. Specifically, for n receiver addresses, we additionally count $(\ln(n) + 1) \times 16k$ LLC accesses in each signaling round. Note that, this self-eviction operation on traditional set-associative caches only requires the number of accesses equal to the associativity.

In Fig. 11, we show evaluation results of using different receiver parameters to communicate on two cache configurations whose epoch sizes equal to 100 epoch units [10]. Fig. 11(a) and (b) shows the empirical number of epochs needed to achieve 1% decode error rate, and Fig. 11(c) shows the theoretical bounds.

The error rate decreases as the number of epochs increases, confirming the effectiveness of the multi-epoch strategy.

We observe that spending more resources on calibration does not always help communication. For example, on a cache with 1 way per hash group, spending 40% of the epoch on calibration achieves the highest communication efficiency, and on a cache with 2 ways per hash group, the best efficiency is achieved when spending 60% of the epoch on calibration.

Finding 4: *Contiguous signaling requires evicting receiver addresses, which increases the signaling cost by multiple thousand times on randomly mapped caches compared to traditional set-associative caches.*

Takeaway: Information can be leaked and accumulated across epochs even when mapping functions are changed.

Takeaway: Spending the maximum amount of resources on calibration is neither the only nor always the best strategy.

D. Varying Epoch Sizes

We analyze how varying epoch sizes can affect the communication cost. We repeat the experiment on attacking RSA in a cache with 1 way per hash group and vary the epoch size from 5 units to 100 units and an infinite size. The communication cost is the number of LLC accesses performed by the receiver to decode a single secret bit. For each epoch size, Fig. 12 shows the lowest communication cost and the corresponding communication strategy, which is represented using the number of calibration rounds per epoch. The signaling cost is the number of samples computed using Eq. (2) on the empirical density functions which we obtain via sampling.

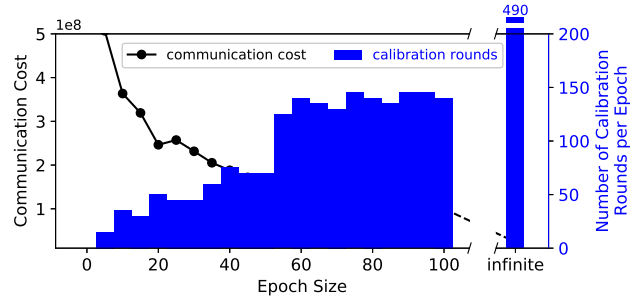


Fig. 12: The impacts of epoch sizes on communication cost.

Overall, the communication cost decreases as the epoch size increases from 5 units to 100 units. When the epoch size is equal to or below 1 unit, such as the configuration

in NewCache [12], our communication strategy becomes infeasible. However, as demonstrated in prior work [9], [10], using such a small epoch size can introduce high performance overhead. The number of calibration rounds increases as the epoch size increases from 5 to 60 units. When the epoch size is larger than 60 units, spending more resources on calibration does not help decrease the overall communication cost.

VIII. DISCUSSION

We briefly discuss how CaSA can be extended for the following cases.

1) *Handling Complex Encoding Schemes.*: We have evaluated a simple encoding scheme of one Boolean secret bit so far. When multiple transmitters are used to encode multiple secret bits using a more complex encoding scheme, we identify two potential decode strategies and the corresponding analysis that can be supported by CaSA.

First, multiple receiver sets are used for signaling. If n transmitters are used to encode n secret bits, the receiver can decode these bits using n different receiver sets. Specifically, it calibrates one receiver set for each transmitter. When performing a signaling step, the receiver uses the n receiver sets in parallel. CaSA could be easily extended to handle this decode strategy. For each receiver set, the modulations from the other transmitters and receivers are modeled as a new type of noise, as this receiver set is not calibrated for those addresses. CaSA could be easily extended to handle this type of noise.

Second, a single receiver set is used for signaling. In the case that each transmitter is composed of a different number of addresses, the receiver can calibrate to obtain a receiver set for a union of these transmitters. Accessing different combinations of the transmitters may result in 2^n different distributions of the signal. The mathematical problem that needs to be solved in the decode module is how many samples are needed to distinguish a group of distributions, instead of two. To compute bounds for this problem, CaSA will need to be extended to use more advanced mathematical methods.

2) *Computing Upper Bounds of Communication Bandwidth.*: CaSA does not compute lower bounds of communication cost (i.e., upper bounds of side-channel bandwidth). However, from the information leakage perspective, the upper bound of the communication bandwidth can be more useful, as it makes it possible to reason about the maximum number of bits that can be leaked per epoch. We think it is possible to derive a *probabilistic upper bound* using an approach similar to the one proposed by Purnal et al. [40], where they computed a probabilistic upper bound for the calibration step only.

IX. RELATED WORK

We have discussed the limitations of prior attempts to analyze randomly mapped caches [10], [11], [13] in Section IV. We now cover related work on analyzing and measuring side channel vulnerabilities.

The closest related work is the concurrently submitted paper by Purnal et al. [40]. They also aim to quantitatively analyze the security of randomly mapped caches. Similar to us, they

consider communication using smaller receiver sets and multi-address transmitters. There are two key differences. First, one of the key contributions of CaSA is to identify the existence of a trade-off between signaling and calibration and to quantify the end-to-end communication cost. However, they solely focus on the calibration step. Their contributions lie in optimizations for the attacker to reduce calibration cost. Second, one of the key findings of CaSA is that communication can happen across epochs. However, their analysis still assumes communication needs to complete within one epoch and for a given epoch size, they compute the upper bound of the success rate to obtain a receiver set with 95% eviction rate.

He et al. [41] proposed an approach to quantitatively evaluate a cache's resilience against multiple classes of attacks on traditional set-associative caches. They build a probabilistic information flow graph for steps in an attack, compute the probability of success for each step, and then compute the probability of success for the whole attack. The key difference from CaSA is that their approach only focuses on the signaling step, and assumes the calibration result is known. Without the capability to explore the calibration step, their approach cannot be used to analyze randomly mapped caches.

SVF [42] and CSV [43] are metrics used to quantitatively measure side-channel leakage in processors by computing statistical correlation between transmitter and receiver's execution traces. CaSA is different from these works. First, they use empirical approaches to obtain traces, while CaSA builds a mathematical model to obtain the communication signal. Second, they compute the metric for given attack traces, while CaSA performs space exploration to find the attack parameters that minimize communication cost.

Several tools, such as CacheAudit [44], cacheD [34] and CaSym [33], have been proposed to detect side channel vulnerabilities in software. These tools are effective in locating secret-dependent memory accesses and control flows. They generally focus on analyzing software and use simple cache models. We find these tools complementary to CaSA, and it could be promising to extend these tools to generate the transmitter parameters for CaSA.

Statistic-based analysis has been widely adopted in analyzing power side channel attacks [45]–[47]. To the best of our knowledge, we are the *first* to formulate the signals in cache-based side-channel attacks into a statistical representation.

X. CONCLUSION

In this paper, we comprehensively analyze the security of randomly mapped caches. Our result shows that the randomization mechanisms used in the state-of-the-art randomly mapped caches are insecure.

We have made key contributions in identifying the end-to-end communication procedure of microarchitecture side channels. Additionally, we leverage concepts from the field of telecommunication to formulate a security analysis of randomly mapped caches into a statistical problem. It is promising to apply our approach to analyze side channels on other types of micro-architecture structures.

REFERENCES

- [1] D. A. Osvik, A. Shamir, and E. Tromer, "Cache Attacks and Countermeasures: The Case of AES," in *Cryptographers' Track at the RSA conference*. Springer, 2006.
- [2] C. Percival, "Cache Missing for Fun and Profit," <http://www.daemonology.net/papers/htt.pdf>, 2005.
- [3] C. Canella, J. Van Bulck, M. Schwarz, M. Lipp, B. Von Berg, P. Ortner, F. Piessens, D. Evtushkin, and D. Gruss, "A Systematic Evaluation of Transient Execution Attacks and Defenses," in *Proceedings of the 28th USENIX Conference on Security Symposium*, 2019.
- [4] P. Kocher, J. Horn, A. Fogh, D. Genkin, D. Gruss, W. Haas, M. Hamburg, M. Lipp, S. Mangard, T. Prescher, M. Schwarz, and Y. Yarom, "Spectre Attacks: Exploiting Speculative Execution," in *IEEE Symposium on Security and Privacy (SP)*, 2019.
- [5] M. Lipp, M. Schwarz, D. Gruss, T. Prescher, W. Haas, A. Fogh, J. Horn, S. Mangard, P. Kocher, D. Genkin, Y. Yarom, and M. Hamburg, "Meltdown: Reading Kernel Memory from User Space," in *USENIX Security Symposium*, 2018.
- [6] V. Kiriansky, I. Lebedev, S. Amarasinghe, S. Devadas, and J. Emer, "DAWG: A Defense Against Cache Timing Attacks in Speculative Execution Processors," in *51st Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 2018.
- [7] P. Vila, B. Köpf, and J. F. Morales, "Theory and Practice of Finding Eviction Sets," in *IEEE Symposium on Security and Privacy*, 2019.
- [8] F. Liu, Y. Yarom, Q. Ge, G. Heiser, and R. B. Lee, "Last-Level Cache Side-Channel Attacks are Practical," in *Proceedings of the 2015 IEEE Symposium on Security and Privacy (SP)*, 2015.
- [9] M. K. Qureshi, "CEASER: Mitigating Conflict-Based Cache Attacks via Encrypted-Address and Remapping," in *51st Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2018.
- [10] M. K. Qureshi, "New Attacks and Defense for Encrypted-address Cache," in *Proceedings of the 46th International Symposium on Computer Architecture (ISCA)*, 2019.
- [11] M. Werner, T. Unterluggauer, L. Giner, M. Schwarz, D. Gruss, and S. Mangard, "ScatterCache: Thwarting Cache Attacks via Cache Set Randomization," in *28th USENIX Security Symposium*, 2019.
- [12] Z. Wang and R. B. Lee, "New Cache Designs for Thwarting Software Cache-Based Side Channel Attacks," *SIGARCH Comput. Archit. News*, 2007.
- [13] A. Purnal and I. Verbauwhede, "Advanced Profiling for Probabilistic Prime+Probe Attacks and Covert channels in ScatterCache," *arXiv preprint arXiv:1908.03383*, 2019.
- [14] Y. Yarom and K. Falkner, "Flush+Reload: A High Resolution, Low Noise, L3 Cache Side-channel Attack," in *USENIX Security Symposium*, 2014.
- [15] M. Lipp, D. Gruss, R. Spreitzer, C. Maurice, and S. Mangard, "AR-Mageddon: Cache Attacks on Mobile Devices," in *25th USENIX Security Symposium*, 2016.
- [16] C. Disselkoe, D. Kohlbrenner, L. Porter, and D. Tullsen, "Prime+Abort: A Timer-Free High-Precision L3 Cache Attack Using Intel TSX," in *26th USENIX Security Symposium*, 2017.
- [17] J. Bonneau and I. Mironov, "Cache-collision Timing Attacks against AES," in *International Workshop on Cryptographic Hardware and Embedded Systems*. Springer, 2006.
- [18] D. Gruss, C. Maurice, and K. Wagner, "Flush+Flush: A Stealthier Last-Level Cache Attack," in *Proceedings of the International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment*, 2016.
- [19] D. Gullasch, E. Bangerter, and S. Krenn, "Cache Games—Bringing Access-based Cache Attacks on AES to Practice," in *IEEE Symposium on Security and Privacy (SP)*. IEEE, 2011.
- [20] M. Yan, R. Sprabery, B. Gopireddy, C. W. Fletcher, R. Campbell, and J. Torrellas, "Attack Directories, Not Caches: Side Channel Attacks in a Non-Inclusive World," in *IEEE Symposium on Security and Privacy (SP)*, 2019.
- [21] D. Gruss, C. Maurice, A. Fogh, M. Lipp, and S. Mangard, "Prefetch Side-channel Attacks: Bypassing SMAP and Kernel ASLR," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS)*. ACM, 2016.
- [22] T. Hornby, "Side-Channel Attacks on Everyday Applications: Distinguishing Inputs with FLUSH+RELOAD," in *BackHat*, 2016.
- [23] T. Ristenpart, E. Tromer, H. Shacham, and S. Savage, "Hey, You, Get Off of My Cloud: Exploring Information Leakage in Third-Party Compute
- [24] Y. Xu, M. Bailey, F. Jahanian, K. Joshi, M. Hiltunen, and R. Schlichting, "An Exploration of L2 Cache Covert Channels in Virtualized Environments," in *Proceedings of the 3rd ACM workshop on Cloud computing security workshop*, 2011.
- [25] D. Gruss, R. Spreitzer, and S. Mangard, "Cache Template Attacks: Automating Attacks on Inclusive Last-level Caches," in *24th USENIX Security Symposium*, 2015.
- [26] Intel, "6th Gen Intel Core X-Series Processor Family Datasheet - 7800X, 7820X, 7900X," 2017.
- [27] C. Maurice, N. Le Scouarnec, C. Neumann, O. Heen, and A. Francillon, "Reverse Engineering Intel Last-level Cache Complex Addressing Using Performance Counters," in *Research in Attacks, Intrusions, and Defenses*. Springer, 2015.
- [28] G. Irazoqui, T. Eisenbarth, and B. Sunar, "Systematic Reverse Engineering of Cache Slice Selection in Intel Processors," in *Proceedings of the 2015 Euromicro Conference on Digital System Design (DSD)*, 2015.
- [29] Intel, "Intel Software Guard Extensions Programming Reference," <https://software.intel.com/en-us/sgx/sdk>, 2013.
- [30] R. Bodduna, V. Ganesan, P. Slpsk, C. Rebeiro, and V. Kamakoti, "BRUTUS: Refuting the Security Claims of the Cache Timing Randomization Countermeasure proposed in CEASER," *IEEE Computer Architecture Letters (CAL)*, 2020.
- [31] A. Shusterman, L. Kang, Y. Haskal, Y. Meltser, P. Mittal, Y. Oren, and Y. Yarom, "Robust Website Fingerprinting Through the Cache Occupancy Channel," in *28th USENIX Security Symposium*, 2019.
- [32] D. M. Gordon, "A Survey of Fast Exponentiation Methods," *Journal of Algorithms*, 1998.
- [33] R. Brozman, S. Liu, D. Zhang, G. Tan, and M. Kandemir, "CaSym: Cache Aware Symbolic Execution for Side Channel Detection and Mitigation," in *IEEE Symposium on Security and Privacy (SP)*. IEEE, 2019.
- [34] S. Wang, P. Wang, X. Liu, D. Zhang, and D. Wu, "CacheD: Identifying Cache-based Timing Channels in Production Software," in *26th USENIX Security Symposium*, 2017.
- [35] C.-K. Luk, R. Cohn, R. Muth, H. Patil, A. Klauser, G. Lowney, S. Wallace, V. J. Reddi, and K. Hazelwood, "Pin: Building Customized Program Analysis Tools with Dynamic Instrumentation," in *Proceedings of the 2005 ACM SIGPLAN Conference on Programming Language Design and Implementation*, 2005.
- [36] Wikipedia, "Markov Chain," https://en.wikipedia.org/wiki/Markov_chain, 2020.
- [37] Wikipedia, "Inclusion-exclusion Principle," https://en.wikipedia.org/wiki/Inclusion-exclusion_principle, 2020.
- [38] J. Galambos, "Bonferroni Inequalities," *The Annals of Probability*, 1977.
- [39] W. Hoeffding, "Probability Inequalities for Sums of Bounded Random Variables," *Journal of the American Statistical Association*, 1963. [Online]. Available: <http://www.jstor.org/stable/2282952>
- [40] A. Purnal, L. Giner, D. Gruß, and I. Verbauwhede, "Systematic analysis of randomization-based protected cache architectures," in *42th IEEE Symposium on Security and Privacy*, 5 2021.
- [41] Z. He and R. B. Lee, "How Secure is Your Cache Against Side-channel Attacks?" in *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*. ACM, 2017.
- [42] J. Demme, R. Martin, A. Waksman, and S. Sethumadhavan, "Side-channel Vulnerability Factor: A Metric for Measuring Information Leakage," in *39th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2012.
- [43] T. Zhang, F. Liu, S. Chen, and R. B. Lee, "Side Channel Vulnerability Metrics: the Promise and the Pitfalls," in *Proceedings of the 2nd International Workshop on Hardware and Architectural Support for Security and Privacy (HASP)*, 2013.
- [44] G. Doychev, B. Köpf, L. Mauborgne, and J. Reineke, "CacheAudit: A Tool for the Static Analysis of Cache Side Channels," *ACM Transactions on Information and System Security (TISSEC)*, 2015.
- [45] Y. Fei, A. A. Ding, J. Lao, and L. Zhang, "A Statistics-based Fundamental Model for Side-channel Attack Analysis," *IACR Cryptology ePrint Archive*, 2014.
- [46] S. Mangard, E. Oswald, and T. Popp, *Power Analysis Attacks: Revealing the Secrets of Smart Cards*. Springer Science & Business Media, 2008.
- [47] F.-X. Standaert, T. G. Malkin, and M. Yung, "A Unified Framework for the Analysis of Side-channel Key Recovery Attacks," in *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 2009.