

Cache Attacks and Defenses

by

William Kosasih

This thesis is submitted for the degree of

Master of Philosophy

School of Computer Science

University of Adelaide

© William Kosasih, 2023

Abstract

In the digital age, as our daily lives depend heavily on interconnected computing devices, information security has become a crucial concern. The continuous exchange of data between devices over the Internet exposes our information vulnerable to potential security breaches. Yet, even with measures in place to protect devices, computing equipment inadvertently leaks information through side-channels, which emerge as byproducts of computational activities. One particular source of such side channels is the cache, a vital component of modern processors that enhances computational speed by storing frequently accessed data from random access memory (RAM). Due to their limited capacity, caches often need to be shared among concurrently running applications, resulting in vulnerabilities. Cache side-channel attacks, which exploit such vulnerabilities, have received significant attention due to their ability to stealthily compromise information confidentiality and the challenge in detecting and countering them. Consequently, numerous defense strategies have been proposed to mitigate these attacks. This thesis explores these defense strategies against cache side-channels, assesses their effectiveness, and identifies any potential vulnerabilities that could be used to undermine the effectiveness of these defense strategies.

The first contribution of this thesis is a software framework to assess the security of secure cache designs. We show that while most secure caches are protected from eviction-set-based attacks, they are vulnerable to occupancy-based attacks, which works just as well as eviction-set-based attacks, and

therefore should be taken into account when designing and evaluating secure caches.

Our second contribution presents a method that utilizes speculative execution to enable high-resolution attacks on low-resolution timers, a common cache attack countermeasure adopted by web browsers. We demonstrate that our technique not only allows for high-resolution attacks to be performed on low-resolution timers, but is also Turing-complete and is capable of performing robust calculations on cache states. Through this research, we uncover a new attack vector on low-resolution timers. By exposing this vulnerability, we hope to prompt the necessary measures to address the issue and enhance the security of systems in the future.

Our third contribution is a survey, paired with experimental assessment of cache side-channel attack detection techniques using hardware performance counters. We show that, despite numerous claims regarding their efficacy, most detection techniques fail to perform proper evaluation of their performance, leaving them vulnerable to more advanced attacks. We identify and outline these shortcomings, and furnish experimental evidence to corroborate our findings. Furthermore, we demonstrate a new attack that is capable of compromising these detection methods. Our aim is to bring attention to these shortcomings and provide insights that can aid in the development of more robust cache side-channel attack detection techniques.

This thesis contributes to a deeper comprehension of cache side-channel attacks and their potential effects on information security. Furthermore, it offers valuable insights into the efficacy of existing mitigation approaches and detection methods, while identifying areas for future research and development to better safeguard our computing devices and data from these insidious attacks.

Declaration of Authorship

I certify that this work contains no material which has been accepted for the award of any other degree or diploma in my name in any university or other tertiary institution and, to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference has been made in the text. In addition, I certify that no part of this work will, in the future, be used in a submission in my name for any other degree or diploma in any university or other tertiary institution without the prior approval of the University of Adelaide and where applicable, any partner institution responsible for the joint award of this degree.

The author acknowledges that copyright of published works contained within this thesis resides with the copyright holder(s) of those works.

I give permission for the digital version of my thesis to be made available on the web, via the University's digital research repository, the Library Search and also through web search engines, unless permission has been granted by the University to restrict access for a period of time.

Acknowledgements

This research endeavor has been an enlightening experience for me, made possible by the collaboration, support, and warm encouragement from numerous individuals during my time at the University of Adelaide. I could not have embarked on this journey without the guidance of my supervisors, Professor Yuval Yarom and Dr. Chitchanok Chuengsatiansup. Without their assistance, this thesis would not have been achievable.

I owe immense gratitude to Dr. Cheryl Pope and Dr. Muhammad Uzair for their invaluable help and insights during some of the most difficult times in my academic journey. Their assistance allowed me to persist in my candidacy and ultimately complete this thesis.

I wish to extend my heartfelt thanks to Professor Olaf Maennel, Dr. Atanas Parashkevov, Dr. Cheryl Pope, and Dr. Muhammad Uzair for dedicating their time to review my thesis and offering valuable feedback.

Furthermore, I would like to express my appreciation to Emeritus Professor Michael Liebelt, Dr. Atanas Parashkevov, Ms. Diana Reed, and Professor Nelson Tansu for their irreplaceable counsel, support, and direction during the intricate process of my research. They have offered motivation in difficult moments and have been there to help when I faced challenges.

I am thankful for the University of Adelaide scholarship, which covered all my living expenses and allowed me to concentrate on my studies.

Lastly, I extend my heartfelt gratitude to my parents, partner, and close friends in Adelaide for their unwavering support, love, and care throughout my educational journey.

Contents

1	Introduction	1
1.1	Thesis Organization and Contributions	4
2	Background	9
2.1	Caches	9
2.2	Cache Side-Channel Attacks	13
2.3	The Processor Microarchitecture	14
2.3.1	Front End	15
2.3.2	Execution Engine	15
2.3.3	Memory Subsystem	15
2.4	Branch Prediction	16
2.5	Transient-Execution Attacks	16
2.5.1	Spectre	17
3	A Framework to Evaluate Cache Security	19
3.1	Background	20
3.2	Problem Description	21
3.3	CacheFX Design	22
3.3.1	Cache Model	22
3.3.2	Attack Model	24
3.3.3	Victim Model	25
3.3.4	The Attack Controller Function	25
3.4	Evaluation	26
3.4.1	Relative Eviction Entropy	26
3.4.2	Eviction-Set Creation	30
3.4.3	Eviction-Set Attack	34
3.4.4	Cache-Occupancy Attack	37
3.4.5	Optimal Eviction-Set Size	38
3.5	Threats to Validity and Limitations	39
3.6	Related Works	40
3.7	Conclusion	41
4	Speculative Execution Against Low-Resolution Timers	43
4.1	Gates	44
4.1.1	Computational Model	44
4.1.2	<i>NOT</i> Gate	45
4.1.3	More Complex Gates	46
4.1.4	Multiple Inputs and Outputs	47
4.1.5	Error Correction Gate	47
4.1.6	Gates With a Fixed Branch Delay	48
4.1.7	Gates Without Branch Training	50
4.1.8	Gates Evaluation	50
4.2	Circuits	51

4.2.1	ALU	51
4.2.2	SHA-1	53
4.2.3	Game of Life	54
4.3	Probe Amplification	55
4.3.1	Single-Gate Amplification	56
4.3.2	Probe Time Amplification Tree	57
4.3.3	Amplification Hyper-tree	57
4.3.4	Experimental Verification	58
4.3.5	Eviction Set Creation	60
4.4	Prime+Store: Fast Attacks with Slow Clocks	60
4.4.1	Prime+Store	61
4.4.2	Attacking ElGamal	61
4.4.3	Trace Acquisition	62
4.4.4	Trace Processing	63
4.4.5	Key Recovery	64
4.4.6	Evaluation	65
4.5	Related Work	65
4.6	Conclusions	67
5	Hardware Performance Counters in Cache Attack Detection	69
5.1	Background	70
5.1.1	Hardware Performance Counters.	70
5.1.2	HPC-Based Cache-Side Channel Attack Detection Methods	71
5.2	Survey of HPC-Based Cache Side-Channel Attack Detection Method Evaluation	73
5.2.1	Accuracy	76
5.2.2	Overhead	76
5.2.3	Detection Speed	77
5.2.4	Threat Model	78
5.3	Assessing the Quality of Attack Detection Methods	78
5.3.1	Experiment Environment	78
5.3.2	Our Method	78
5.3.3	Accuracy	80
5.3.4	Overhead	81
5.3.5	Detection Speed	83
5.3.6	Threat Model	83
5.4	Conclusions	88
6	Future Directions	90
7	Conclusion	92
	References	94
	Appendix A Chapter 4	107
A.1	Hyper-tree Amplification Implementation	107
A.2	Gates With and Without Branch Training	108
A.3	Gate Accuracy	110
A.4	Gates on Other Processors	112
A.5	Number of Cases in Gates Without Branch Training	114

List of Tables

4.1	Game of Life glider accuracy.	55
5.1	Survey result.	75
5.2	PAPI events used in our detection method.	79
5.3	Accuracy of three detection methods.	80
5.4	Accuracy of HPCache with various sampling interval.	80
5.5	Overhead of detection methods.	82
5.6	True positive rate of proof-of-concept and camouflaged attacks.	85
5.7	False positive rate of proof-of-concept and camouflaged attacks.	85
5.8	Accuracy of HPCache trained on camouflaged attacks.	86
5.9	Attack times for different Flush+Reload camouflaging aggressiveness.	88
A.1	Accuracy of gates on Intel Core(TM) i5-8250U	111
A.2	Accuracy of gates on AMD Ryzen 5 3500U	113
A.3	Accuracy of gates on Apple M1 and Samsung Exynos 2100	114
A.4	Accuracy and run time for nbt $NAND_1^{12}$ and nbt NOT_2^1 gates with different number of cases on Intel Core(TM) i5-8250U	115

List of Figures

2.1	Relationship between memory address and cache lines	9
2.2	Mapping from memory to cache entries in direct-mapped caches	10
2.3	Mapping from memory to cache lines in fully-associative caches.	11
2.4	Mapping from memory to cache lines in set-associative caches.	12
2.5	Simplified diagram of the Intel Skylake microarchitecture.	18
3.1	CacheFX design overview.	22
3.2	CacheFX overall architecture.	23
3.3	REE across cache designs with random replacement.	28
3.4	REE for CEASER-S with 2048 lines depending on ways and partitions.	29
3.5	Number of memory accesses required by eviction-set building techniques for different 2048-line caches.	31
3.6	Number of memory accesses required by eviction-set building techniques for CEASER-S depending on cache size.	32
3.7	Percentage of addresses in the constructed eviction sets that conflict with the victim’s address, using different eviction-set construction techniques and 2048-line caches.	32
3.8	Eviction set sizes found by eviction-set building techniques for different 2048-line caches.	33
3.9	Eviction success rate for the eviction sets found for different 2048-line caches.	34
3.10	Eviction set size for 2048-line caches and 90% eviction probability.	34
3.11	Eviction-set attack: Number of encryptions required to break AES and modular exponentiation with random replacement.	36
3.12	Occupancy attack: Number of encryptions required to break AES and modular exponentiation with random replacement.	37
3.13	Median number of encryptions required to break AES with different cache designs and replacement algorithms.	38
3.14	Optimal eviction set sizes for 1024-lines caches	39
4.1	<i>NOT</i> Gate.	45
4.2	A Buffer Gate with a fixed branch delay.	49
4.3	ALU accuracy.	52
4.4	SHA-1 accuracy.	54
4.5	T-tetromino heatmap.	54
4.6	One generation Game of Life accuracy.	56
4.7	Amplification tree based on NOT_Y^1 Gates.	57
4.8	Amplification Hyper-Tree in native.	59
4.9	Amplification Hyper-Tree in WebAssembly.	59
4.10	Time to find an eviction set in Chrome using 0.1 millisecond low-resolution timer	60
4.11	A segment of samples of the square operation in modular exponentiation.	63

4.12 Distribution of stitched key in relation to ground truth location. 65

Chapter 1

Introduction

Information security in computing is a very critical issue, especially in this day and age that almost all aspects of our lives are in one way or another reliant on computers as well as the network that connect them, the Internet. Used for trivial activities such as a simple web browse, to tasks that are sensitive in nature, such as banking, our computing devices are constantly interacting with all kinds of other computers all around the globe. While this digital interconnection has gifted us with communication ease and agility [Adams, 2017; Elijah et al., 2018; Levy and Strombeck, 2002], much to our unknowing (or perhaps ignorance), it has simultaneously created a risk of exposing our data essentially to anyone, or anything that we interact with on the web [Adams, 2017; Karunakaran et al., 2018; Wheatley et al., 2016].

Given the situation, it is only natural to expect people and organizations to pay utmost attention towards the security of their computing devices. Yet, day after day we continue to hear headlines about major data breaches each affecting thousands, and even millions of unfortunate users [Adams, 2017; Biddle et al., 2022; Booth, 2022; Cadwalladr and Graham-Harrison, 2018; Karunakaran et al., 2018; Pilla et al., 2023; Wheatley et al., 2016]. With these evidences in hand, it is clear that people, and especially major corporations need to allocate even more efforts towards protecting information security [Fedele and Roner, 2022; IBM, 2022].

Nevertheless, regardless of how much effort we allocate into safeguarding our devices, we cannot escape the reality that computing equipment innately leaks information as it operates [Briongos et al., 2016; Martin et al., 2012]. Computers leave traces of their computing activities as they function, in the same manner that biological organisms leave footprints of their metabolism activities as they go about living. These traces are basically side-effects of computation, called side-channels. These side-effects include acoustic [Genkin et al., 2014; O'Malley and Choo, 2014], electromagnetic [Elibol et al., 2012; Hongxin et al., 2009], power [Kocher et al., 1999; Mangard, 2003], cache [Liu et al., 2015; Osvik et al., 2006], and timing [Dhem et al., 2000; Kocher, 1996].

In this thesis, we direct our attention towards the cache side-channel. The

cache is a piece of temporary memory intended to speed up computation. Advances in technology have brought about performance enhancements to all aspects of computing. Most notably of all is the processor, which has sped up dramatically over the decades. Working in close conjunction to the processor is the random access memory (RAM). The RAM is used as a storage to provide input and store output of intermediate data to/from the processor before they are permanently recorded on disks. While the RAM has also gained significant performance improvement over the years, since the late 1990s their increase in speed has lagged considerably behind that of the processor, and the gap continues to widen as time goes by. This situation prompted the introduction of caches to processors. They are a form of smaller but much faster memory that stores copies of frequently used data from the RAM, enabling the processor to work at a much higher rate as it can now avoid the delays caused by the slow RAM.

Due its finite size, the cache often has to be shared between different applications that are simultaneously running in a computer. As mentioned above, all computing operations leave traces on the hardware they operate on, including the cache. The execution of programs alter the state of the cache. When a program performs a large number of data accesses, the cache eventually has to remove older data to make room for new ones. The cache decides what data need to be evicted typically based on their age, and is completely agnostic to whether the eviction of data is caused by the application who owns the data. Removal of other applications' data may appear harmless, after all, the processor seamlessly orders reinsertion of needed data back into the cache on an ad-hoc basis, and execution continues as usual. However, this slightly affects the data access speed of those applications whose data are temporarily removed from the cache. By closely monitoring this access time difference caused by data eviction, any application suddenly gains a potentiality to infer the execution of other programs running on the system. In other words, gaining an access to a side-channel. The act of a malicious application, seeking to gain illicit information about the system through measurements of cache-related phenomena is called cache side-channel attack.

Cache side-channel attacks have received attention both in academia and the industry in recent years given their ability to stealthily compromise information confidentiality through the exploitation of the cache hardware [Aciğmez and Seifert, 2007; Aldaya et al., 2019a; Aldaya et al., 2019b; Genkin et al., 2020; Gullasch et al., 2011; Gülmezoglu et al., 2019b; Liu et al., 2015; Oren et al., 2015; Osvik et al., 2006; Percival, 2005; Ristenpart et al., 2009; Ronen et al., 2019; Shusterman et al., 2019; Yarom and Falkner, 2014; Yarom et al., 2017; Zhang et al., 2012]. Unlike traditional malware that leaves obvious traces of their activities, cache side-channel attacks leave only microarchitectural traces and are more difficult to detect and mitigate [Alam et al., 2021; He and Lee, 2017].¹ Yet, they are sometimes seen as an improbable

¹Microarchitecture refers to the internal implementation of the processor.

threat, and a very specialized domain of computer systems security. This view may be held as a result of the high complexity involved with implementing such attacks [Apple, 2018; Easdon et al., 2022; Yarom, 2016], the “black-box” nature of microarchitectural components due to the lack of official documentations [Easdon et al., 2022], possibly leading to an incorrect viewpoint that cache attacks are more of an academic proof of concept rather than a practical attack.

In actuality, history has shown the potential for disastrous outcomes of such attacks, particularly when combined with other microarchitectural attacks such as transient-execution attacks, forming a highly effective pairing. This was evident in the 2018 discovery of Spectre [Kocher et al., 2019] and Meltdown [Lipp et al., 2018], which showcased the potential of these attacks to compromise crucial operating system information and sensitive user data. In this context, it is important to uncover and tackle potential hardware vulnerabilities and create countermeasures for existing devices.

While completely eradicating cache attacks is challenging, there are approaches designed to reduce or obscure cache side-channels, making it difficult for attackers to extract valuable information. Hardware-level mitigation such as secure caches are available [Domnitser et al., 2012; Liu et al., 2016b; Qureshi, 2018; Qureshi, 2019; Tan et al., 2020; Wang and Lee, 2007; Werner et al., 2019], along with software-level mitigation such as low-resolution timers [Chromium, n.d.; Hazen, 2018; Wagner, 2018], attack detection systems [Akram et al., 2020; Depoix and Altmeyer, 2018; Mushtaq et al., 2018a; Mushtaq et al., 2018b; Payer, 2016], and constant-time programming [Bernstein et al., 2012; Liu et al., 2015]. Despite the protection that these countermeasures presently offer, we must consistently search for new security vulnerabilities proactively, rather than reactively responding only after attacks have occurred.

The same principle applies to cache side-channel attack defense techniques; we must continuously evaluate, scrutinize, and explore ways in which they can be compromised. Therefore, this thesis aims to 1) evaluate the effectiveness of defense techniques against cache side-channel attacks, 2) identify their shortcomings and reveal new attack scenarios that may lead to their compromise, and 3) evaluate the practicality of these techniques for broad adoption.

We begin our investigation by focusing on the hardware-level, specifically on secure caches. Numerous secure cache designs have been proposed that aim to protect against data leakage. However, the diverse range of these proposals and their unique evaluation metrics make it difficult to empirically analyze their security features. While several metrics have been suggested for performing such evaluations, these tend to be limited both in terms of the potential adversaries they consider and in the applicability of the metric to real-world attacks, as opposed to attack techniques. Moreover, all existing metrics implicitly assume that a single metric can encompass the nuances of side-channel security. To fill this gap, this thesis presents CacheFX, a software-

based cache simulation framework as a way for assessing the security of secure caches, specifically by determining their effectiveness against side-channel attacks. **CacheFX** allows the evaluator to implement various cache designs, victims, and attackers, as well as to exercise them for assessing the leakage of information via the cache.

We then continue our investigation by focusing on software-level cache defense, specifically low-resolution timers. As cache attacks often require a high-resolution timing source to successfully distinguish between cache hits and misses, low-resolution timers limit the resolution of available timers, thus making attacks more difficult. As mentioned, proactive investigation of computing devices is crucial for detecting and eliminating vulnerabilities. As part of our proactive research effort, we examine the potential vulnerability of low-resolution timers as a countermeasure to attacks, specifically how this defense mechanism can be compromised by attacks capable of bypassing its safeguards. While numerous studies have looked into cache states as a “receiver” of microarchitectural outcomes from transient execution attacks [Kocher et al., 2019; Lipp et al., 2018], none have investigated the reverse perspective: the possibility of leveraging transient execution to facilitate cache attacks on systems using mitigation measures such as low-resolution timers. This thesis delves into this alternative viewpoint and evaluates the security implications of such a method against low-resolution timers as a cache attack defense technique.

We conclude our investigation with an examination of an emerging defense technique known as attack detection systems. Previously mentioned hardware-based defenses, such as secure caches, cannot be implemented on current hardware, while software-based countermeasures such as low-resolution timers may significantly disrupt functionality and degrade performance. Instead of incurring the cost of continuous protection, even during periods without attacks, an alternative approach aims to identify ongoing attacks and only activate countermeasures when an attack is detected. A widely used approach in this domain involves the utilization of hardware performance counters (HPCs). These counters monitor microarchitectural events and analyze statistical deviations to differentiate between malicious and benign software. With numerous proposals and promising reported results, we seek to investigate whether published HPC-based detection methods are evaluated in a proper setting and under the right assumptions, such that their quality can be ensured for real-world deployment against cache side-channel attacks. To achieve this goal, this thesis presents a comprehensive evaluation and scrutiny of existing literature on the subject matter in a form of a survey, accompanied by experimental evidences to support our evaluation.

1.1 Thesis Organization and Contributions

Background (Chapter 2)

This chapter provides the background for the rest of this thesis. Initially, we provide an overview of the cache and its application in computers, as well as its operational principles and the potential security risks associated with it. We then delve into the topic of cache side-channel attacks, exploring their various types. We continue by introducing the concept of the processor microarchitecture, and branch prediction which are integral parts of modern processors. We conclude with the discussion on branch prediction, and transient execution attack.

A Framework to Evaluate Cache Security (Chapter 3)

This chapter presents CacheFX, a software-based cache simulation framework as a way for assessing the security of secure caches, specifically in the context of mitigating contention-based attacks. CacheFX allows for the emulation of various cache designs, victims, and attackers, and measures the resulting leakage. Within this framework, we implement and evaluate nine cache designs: fully-associative and set-associative caches, PLCache [Wang and Lee, 2007], Newcache [Liu et al., 2016b], PhantomCache [Tan et al., 2020], ScatterCache [Werner et al., 2019], way-partitioned caches [Domnitser et al., 2012], and the two variants of CEASER [Qureshi, 2018; Qureshi, 2019]. Furthermore, we modeled five replacement algorithms: random replacement, least recently used (LRU) [Denning, 1968], two variants of pseudo-LRU, and SRRIP [Jaleel et al., 2010].

We further contribute three evaluation metrics. These not only add to the existing portfolio of metrics proposed in prior works, allowing cache designers more options for cache evaluation, but also demonstrate the flexibility of CacheFX and its ability to measure a variety of metrics. The first metric we contribute is the *Relative Eviction Entropy* (REE), which we propose to measure the information leakage from a single victim access via the cache side channel. The second evaluation metric we contribute measures the difficulty of eviction-set creation. We implement and evaluate three eviction-set building strategies: Single Holdout Method (SHM) [Liu et al., 2015], Group Elimination Method (GEM) [Vila et al., 2019] and Prime+Prune+Probe (PPP) [Purnal et al., 2021b]. The final metric we contribute evaluates the protection that the cache provides for cryptographic code. We evaluate both traditional attacks that aim to exploit eviction sets, and occupancy-based attacks [Shusterman et al., 2019; Shusterman et al., 2021a].

Through this chapter, we gain three key insights. First, we establish that different metrics highlight different aspects of the caches. In particular, we find that building eviction sets is faster in skewed caches such as ScatterCache and CEASER-S, than other randomized caches, such as fully associative caches or PhantomCache. Faster eviction-set construction reduces the effort required for mounting attacks. At the same time, our experiments show that, with the right parameters, skewed caches are not less secure when it comes to cryptographic attacks. Second, we find that the security against cryptographic attackers depends not only on the design, but also on other parameters, such

as the replacement policy and the cache associativity. We also show that *all* non-partitioned caches are vulnerable to both eviction-set and occupancy attacks. Third, we find that most non-partitioned secure cache designs offer protection against eviction-set attacks. However, cache-occupancy attacks are left unconsidered and for highly secure designs occupancy attacks are no less effective than eviction-set attacks. Our evaluation thus demonstrates that partitioning is preferable from a side-channel perspective and the resistance to eviction-set attacks of non-partitioned solutions may be tuned to match the respective complexity of cache-occupancy attacks to balance overall cache side-channel resistance and cost.

Personal Contribution

In this research, my personal contributions include various aspects. I was involved in the development of the simulator, taking on tasks such as refactoring and reorganizing its structure, as well as extending it by incorporating experiments to establish the ideal eviction set size for different cache designs. Additionally, I implemented a new replacement policy in the simulator, specifically SRRIP, and carried out the experiment to determine the optimal eviction set size, collecting the relevant data. Lastly, I wrote the sections detailing the design of the simulator and the experiments I performed.

Publication

Daniel Genkin, William Kosasih, Fangfei Liu, Anna Trikalinou, Thomas Unterluggauer, and Yuval Yarom. “CacheFX: A Framework for Evaluating Cache Security”. In: *The 18th ACM ASIA Conference on Computer and Communications Security*, ACM ASIACCS 2023.

Speculative Execution Against Low-Resolution Timers (Chapter 4)

In this chapter, we present a novel transient-execution method aimed at circumventing low-resolution timers, a common defense against cache side-channel attacks used in web-browsers. We investigate, for the first time, how transient execution can be used for improving cache attacks, specifically on systems with low-resolution timers. This chapter contributes novel logical gates whose underlying operating concept lies on speculative execution. These gates are capable of manipulating the state of the cache based on whether data is cached or not, and are both generic and versatile. They work on multiple architectures including Intel and AMD, Samsung Exynos, and Apple M1, and can even be implemented in WebAssembly.

The next contribution of this chapter is a novel cache state amplification technique using our gates. We show that our gates can be used to amplify a small timing difference and distinguish whether a memory address is cached with a low-resolution timer, and use this amplification to build eviction sets in WebAssembly, using only a low-resolution timer. In addition, this chapter contributes a novel attack, called Prime+Store that allows for high-resolution cache attacks with a low-resolution timer.

In addition to facilitating cache attacks on low-resolution timers, we show that these gates can also be leveraged to perform logical manipulation of cache states. Intrigued by this discovery, we investigate the possibility of leveraging these gates to carry out meaningful computations within the cache, without elevating the state to the architectural level. We use these gates to construct an arithmetic logic unit (ALU), the Secure Hash Algorithm 1 (SHA-1), and the Game of Life. We demonstrate that our gates are sufficiently robust to perform arbitrary calculations in the cache.

Personal Contribution

In this research, I designed the experiments for, collected data on, and analyze the results of the arithmetic logic unit (ALU), Secure Hash Algorithm 1 (SHA-1), and Game of Life implementations. Furthermore, I conducted statistical tests on all samples and adapted gates for AMD, Samsung Exynos, and Apple M1 platforms while evaluating their performance. I also devised the key-stitching technique, improved the method for converting Prime+Store traces into square and multiply sequences, and extracted the complete ElGamal key from the traces. In addition, I created an error-correction mechanism for microarchitectural computation and authored sections related to my work.

Publication

Daniel Katzman, William Kosasih, Chitchanok Chuengsatiansup, Eyal Ronen, and Yuval Yarom. “The Gates of Time: Improving Cache Attacks with Transient Execution”. In: *The 32nd USENIX Security Symposium*, USENIX Security 2023.

Hardware Performance Counters in Cache Attack Detection (Chapter 5)

This chapter presents a comprehensive survey on existing hardware performance counter (HPC)-based cache side-channel detection methods. We evaluate the settings and assumptions under which these detection methods are evaluated, and discover that the promising results reported in most of them may not be applicable in realistic scenarios. We subsequently identify four main shortcomings in many proposals: improper evaluation settings and assumptions of accuracy, overhead, and detection speed, and a weak threat model used to evaluate detection methods. Notably, we find that the threat model employed assumes that attackers use variants of published proof-of-concept attacks. To underscore the limitations of this assumption, we present a new variant of the Flush+Reload attack on GnuPG that camouflages its malicious execution behind that of a benign application, demonstrating how real-world attacks may differ significantly in terms of implementation and behavior from proof-of-concept attacks. We demonstrate the feasibility of such attacks in evading detection, while at the same time, capable of effectively stealing private encryption keys. We conclude that addressing the evaluation shortcomings we have highlighted is essential to determine the true effective-

ness of these detection methods for real-world deployment.

Personal Contribution

In this study, I established the paper’s foundation, developed its contents, conducted experiments, and completed the writing process. I performed an extensive literature review on cache side-channel attack detection and critically analyzed the subject. I identified weaknesses and drawbacks in the 48 papers examined and organized my analysis accordingly. I developed the cache side-channel attack detection method used in this chapter and evaluated its detection accuracy alongside two other publicly available detection methods. I also devised a new type of camouflaged Spectre and Flush+Reload attacks that interleave malicious activities within benign program execution. Additionally, I assessed the detection rate of these camouflaged attacks by various detection methods. I demonstrated that our camouflaged attacks successfully perform Flush+Reload attacks by stealing the ElGamal key and evaluated the duration of these attacks in relation to their camouflaging aggressiveness.

Conclusion (Chapter 7)

This chapter provides a summary of the key findings and results presented in this thesis.

Future Directions (Chapter 6)

This chapter outlines potential directions for future research that can be pursued based on the findings and results presented in this thesis.

Chapter 2

Background

2.1 Caches

Modern processors takes advantage of caches to exploit spatial and temporal locality and speed up accesses to memory data, address translation, and branch prediction information. These caches are commonly organized in hierarchies, e.g., modern Intel CPUs have dedicated L1 data and instruction caches, a unified L2 cache for each core, and a last-level cache (LLC) shared by all cores. The L1 has the lowest latency, yet is smallest in size. Conversely, the LLC is the slowest but has the largest capacity.

Cache Hits and Misses. When the processor requires memory data, it first checks in the relevant L1 cache. In a *cache hit*, when the data is found, it is served from the cache. Otherwise, in a *cache miss*, the processor retrieves the data from a lower cache level or from memory if the requested data is not present on any level in the cache hierarchy. Once the data is fetched, it is stored in the cache for future use. A similar process is used for retrieving address translation and branch prediction results.

Cache Structure. Cache memory stores data in fixed-size blocks referred to as *lines*, which serve as the smallest unit for loading and storing cached data to and from memory, for example, 64 bytes. These lines are situated in one of the available spaces, or entries, within the cache hardware. Consequently, in this thesis, the term *line* refers to the actual data block housed in the cache, while the word *entry* denotes the hardware slot within the cache that accommodates the lines.

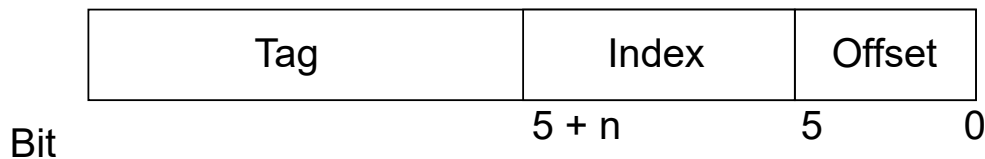


Figure 2.1: Relationship between memory address and cache lines

Figure 2.1 illustrates the relationship between memory addresses and lines

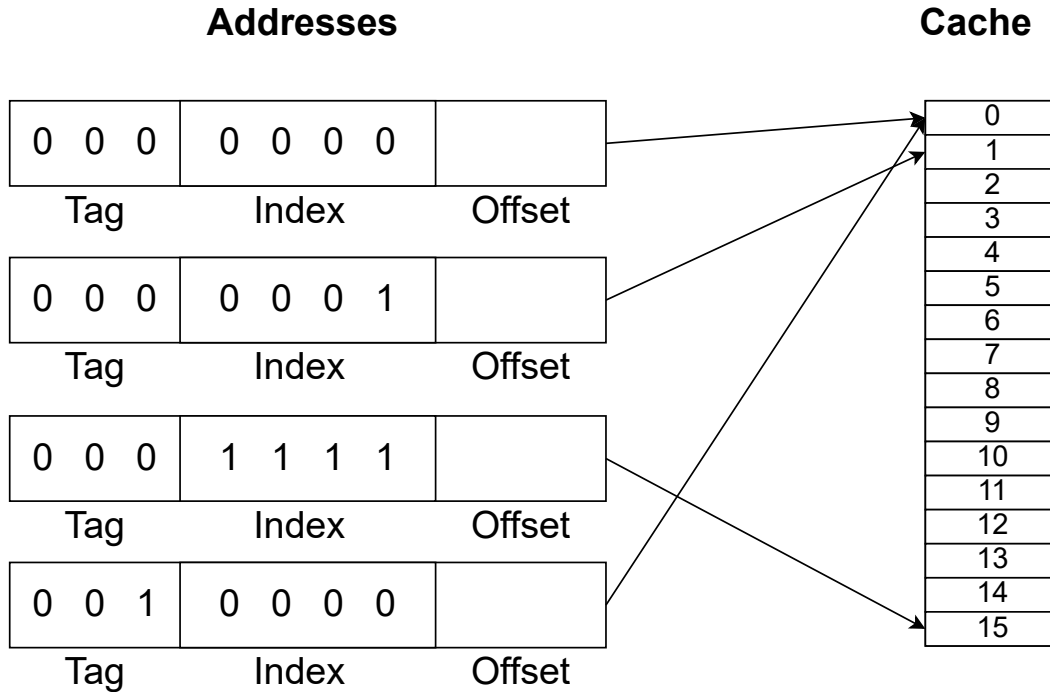


Figure 2.2: Mapping from memory to cache entries in direct-mapped caches

for a cache with 64 byte-wide line. The *offset* bits represents the position of a specific byte or word within a cache line, and is derived from the lower-six bits of the memory address. This means that access to any 64 byte within addresses with matching *tag* and *index* bits falls within a same cache line. The index bits determine which entry in the cache an address maps into. And finally, the tag is used to differentiate between addresses of different memory locations that share the same index bits, and are mapped to the same entry in the cache.

Direct-Mapped Cache A direct-mapped cache associates addresses with entries based on their index bits. This creates a one-to-one mapping between the index and the cache entry, which wraps-around when the index bits overflow. As a result, addresses sharing identical index bits are mapped to the same cache entry. To differentiate between different lines, the tag bits are used. They indicate that although different addresses are directed to the same cache entry, they are distinct cache lines if they have distinct values in their more significant bits.

Figure 2.2 illustrates the association between memory addresses and their corresponding entries in a direct-mapped cache with a 16-line capacity. Observe that the first and last address mappings align with the cache’s first entry, as their index bits are identical. However, the cache can differentiate the two due to their distinct tags.

The issue with this cache design is the potential for thrashing, which can occur when multiple memory locations with matching indexes but distinct tags

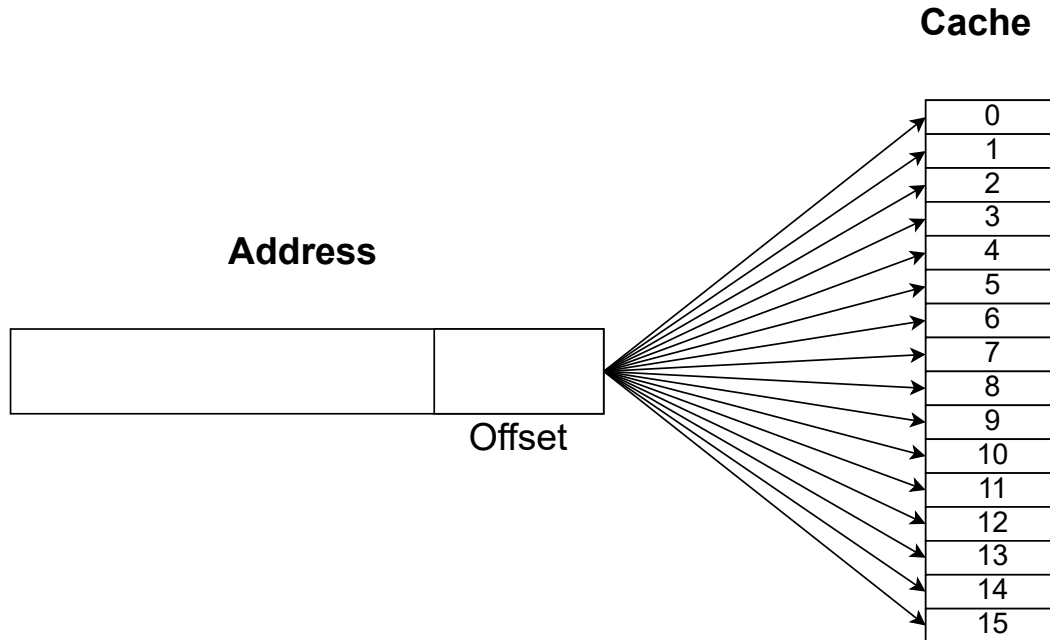


Figure 2.3: Mapping from memory to cache lines in fully-associative caches. Any address may occupy any entry in the cache.

are accessed frequently. In this situation, each access will continuously replace the others' line, irrespective of the cache's actual usage.

Fully-Associative Cache In a fully-associative cache, any memory location can be placed in any cache entry. As a result, the tag and index bits serve the sole purpose of distinguishing between cache lines and hold no significance in determining the entry where the line is to be stored. Figure 2.3 demonstrates how memory addresses map to cache entries in a fully associative cache, noting that data can reside in any of the 16 available entries. When searching for a line, fully-associative caches simultaneously searches all the cache entries for data corresponding to that line.

Set-associative Caches In this model, the cache is arranged into *sets* and *ways*. Identical to direct-mapped caches, index bits are used to determine the entry of an address in the cache. However, each entry is capable of holding more than a single line at a time. A set-associative cache capable of holding k lines in each entry is called a k -way set-associative cache. Entry in the cache as determined by the index bits is called a *set*, and each of the k entries within a set is called a *way*. Note that data can be placed in any of the k ways within a set, meaning that access within a set is associative. When searching for a line in a set-associative cache, the processor first finds the set it maps to and then uses an associative lookup in the set. Figure 2.4 illustrates the mapping of addresses in a three-way set-associative cache.

Small caches are often fully-associative, as they offer optimal hit-rate and cache utilization. However, fully-associative caches do not scale well in terms

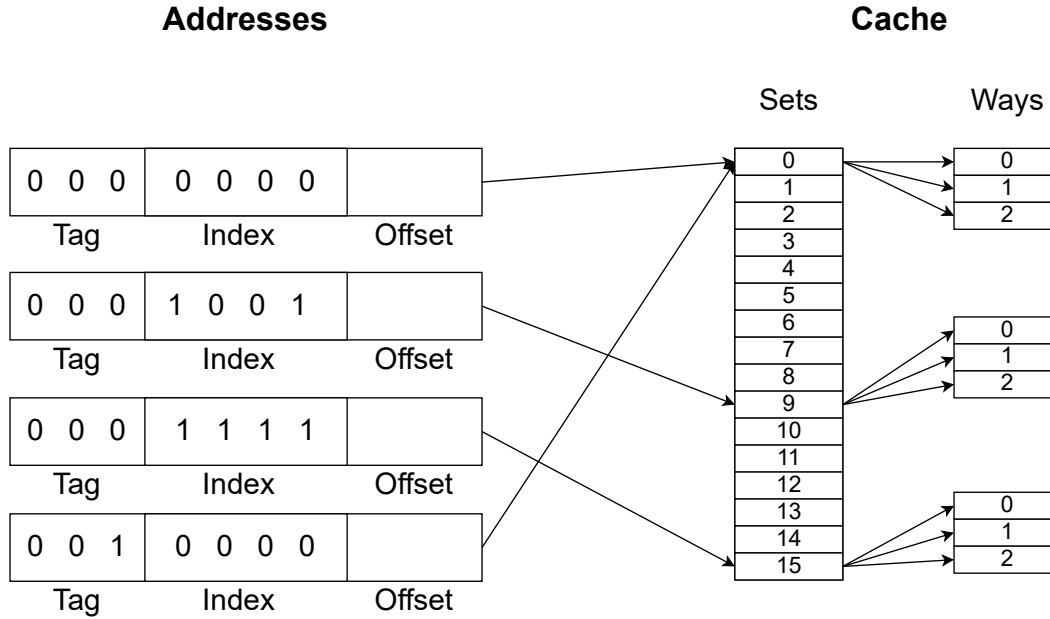


Figure 2.4: Mapping from memory to cache lines in set-associative caches. Addresses are map to the set based on their index bits. Within a set, that address may occupy any way.

of power and performance. Thus, larger caches are often set-associative, which gives a better trade-off between capacity, hit-rate, and implementation complexity.

Cache Replacement Policies. As caches eventually run out of space due to their limited size, when new data needs to be stored, the cache *evicts* older data to make room for new ones. The choice of which data to evict is dictated by *replacement policies*. One such policy is the *random* policy, which randomly chooses any cache line to evict. Another policy is the *LRU* (Least Recently Used) policy, which keeps track of the access time of data in the cache. It selects the candidate that has not been used for the longest duration. To approximate the LRU algorithm without maintaining the full access order, methods such as: *Bit-PLRU* and *Tree-PLRU* are used. Bit-PLRU utilizes a single bit for each candidate line. Initially, these bits are set to zero and are changed to one upon access. If setting a bit to one would result in all candidate bits being set, all the candidate bits are cleared before setting the bit. The replacement algorithm selects an arbitrary candidate whose bit is clear. Tree-PLRU constructs a binary tree with candidates located at the leaves. Each internal node of the tree has a direction bit that indicates the side of the most recent access. The replacement candidate is chosen by traversing the tree and selecting the half whose direction bit indicates least recently accessed data.

2.2 Cache Side-Channel Attacks

While data eviction caused due to program activity is non-malicious in the context of a single process, cross-eviction between different application contexts could potentially be exploited to break information isolation guarantees between applications. By observing timing differences of access to data, an attacker can infer whether or not the data is present in the cache, and therefore learn about the execution of a victim program.

The pioneering work of Tsunoo et al. (2002) demonstrated the ability to exploit this information for the recovery of secret cryptographic keys. Initial attacks primarily focused on the L1 data cache [Osvik et al., 2006; Percival, 2005; Tsunoo et al., 2003], but soon expanded to include other caches as well [Aciğmez, 2007; Aciğmez et al., 2007; Gras et al., 2018; Gullasch et al., 2011; Irazoqui Apecechea et al., 2015; Liu et al., 2015; Yan et al., 2019a; Yarom and Falkner, 2014]. Cache attacks have also been successfully deployed against various cryptographic systems, including symmetric cryptography [Genkin et al., 2020; Gullasch et al., 2011; Irazoqui Apecechea et al., 2015; Moghimi et al., 2017; Osvik et al., 2006; Tsunoo et al., 2002; Tsunoo et al., 2003], public-key schemes [Aciğmez et al., 2007; Dall et al., 2018; Gras et al., 2018; Liu et al., 2015; Percival, 2005; Pereida García and Brumley, 2017; Ronen et al., 2019; Yarom and Falkner, 2014], post-quantum cryptography [Groot Bruinderink et al., 2016; Pessl et al., 2017], and even non-cryptographic software [Brasser et al., 2017; Evtuyushkin et al., 2016; Gras et al., 2017; Gruss et al., 2015; Hund et al., 2013; Oren et al., 2015; Shusterman et al., 2019; Yan et al., 2020].

Reload-Based Attacks Reload-based attacks involve monitoring shared memory access [Gruss et al., 2015; Gullasch et al., 2011; Yarom and Falkner, 2014]. To prepare for an attack, an adversary evicts data from the cache either by issuing a special command [Gullasch et al., 2011; Yarom and Falkner, 2014] or by creating contention within the cache sets that contain the data [Gruss et al., 2015]. After waiting for a certain amount of time, the attacker then measures how fast its access is to the previously removed data. If the data is accessed by the victim while the attacker is waiting, the data will be brought back into the cache, allowing the attacker to also quickly fetch the data, thus learning about the cache state of said data.

Contention-Based Attacks, Contention-based attacks are attacks that exploit the limited storage capacity within the cache, particularly within individual cache sets [Aciğmez, 2007; Aciğmez et al., 2007; Gras et al., 2018; Irazoqui Apecechea et al., 2015; Liu et al., 2015; Osvik et al., 2006; Percival, 2005; Yan et al., 2019a]. The most commonly used technique in contention-based attacks is Prime+Probe. In this technique, the attacker accesses its own data, to fill some or all of the cache sets, i.e., priming the cache. After allowing the victim to execute for a certain period of time, the attacker measures the access time to the cached data. If the access is slow, it suggests that the data has been removed from the cache set as a result of victim activities.

Variations of this attack refrain from using timing information by utilizing performance counters [Bhattacharya and Mukhopadhyay, 2015; Brassler et al., 2017; Uhsadel et al., 2008] or transaction aborts [Disselkoeen et al., 2017] for contention detection, instead of relying on timing measurements.

Another contention-based attack technique is Evict+Time [Gras et al., 2018; Jain et al., 2019; Osvik et al., 2006; van Schaik et al., 2017]. In this approach, the attacker first evicts data from the cache prior to measuring the execution time of the victim. If the evicted data is accessed by the victim, it leads to a longer execution time, thereby revealing information about the cache sets accessed by the victim.

Other Types of Cache Attacks. There are certain cache attacks that do not fall under the previously mentioned categories. These attacks aim to take advantage of specific implementation characteristics of the cache itself. Examples of such attacks include leveraging port contention [Yarom et al., 2017], cache flushing time [Gruss et al., 2016], replacement policies [Purnal et al., 2021a; Vila et al., 2020; Xiong and Szefer, 2020], cache inspection operations [Brumley, 2015], or variations in power consumption [Page, 2002].

Eviction-Set Construction. In many of the aforementioned attacks, the attacker requires the ability to regularly remove certain data from the cache. This is often achieved by creating an *eviction set*, which comprises memory locations, all of which are mapped to the same set in the cache as the targeted data that needs to be evicted. Constructing an eviction set is usually straightforward when the attacker is able to gain access to the mapping information for the cache. However, in cases where the mapping function is undisclosed or the cache indexing information is not accessible to the attacker, more advanced techniques can be used to recover the absent mapping details. Previous research has demonstrated methods to reverse-engineer undocumented mapping functions [Gras et al., 2018; Inci et al., 2016; Maurice et al., 2015b; McCalpin, 2021; Yarom et al., 2015] and create eviction sets without relying on physical address information [Liu et al., 2015; Purnal et al., 2021b; Vila et al., 2019].

2.3 The Processor Microarchitecture

Contemporary processors are composed of an extensive array of components and algorithms that together execute the instruction set architecture (ISA), which is collectively referred to as the *microarchitecture*. The microarchitecture defines how program executions are carried out internally within modern processors. Program execution in these processors entails numerous stages and elements that work in tight conjunction to one another to achieve optimal performance and efficiency. A simplified representation of the Intel Skylake microarchitecture can be seen in Figure 2.5, which illustrates that three main components are responsible for this process: the front end, the execution engine, and the memory subsystem. For the purposes of this thesis, our primary focus is on gaining a deeper understanding of the Intel microarchitecture.

2.3.1 Front End

The front end is responsible for fetching instructions from memory, decoding them, and transferring them to the execution engine for execution. This involves several sub-stages, including instruction prefetching, branch prediction, and instruction decoding. Prefetching anticipates the instructions that will be needed in the near future. Branch prediction, on the other hand, attempts to predict the outcome of conditional branches, allowing the processor to speculatively execute instructions and thereby reduce the impact of control dependencies on performance. Instruction decoding is the process of translating instructions into simpler, fixed-length micro-operations (μops), which can be more efficiently processed by the execution engine.

2.3.2 Execution Engine

The execution engine is responsible for executing the μops provided by the front end. It comprises several sub-components, including the reorder buffer, reservation stations, and execution units. To maximize instruction-level parallelism, the execution engine does not strictly follow the program's specified order of instructions. Instead, it executes instructions immediately after their dependencies are resolved and suitable execution units are available. This is called *out-of-order execution*. To support out-of-order execution, the front end and execution engine utilize a shared data structure called the reorder buffer (ROB). The front end adds instructions to the tail of the ROB, and the execution engine employs a variation of the Tomasulo algorithm [Tomasulo, 1967] to execute instructions in a more optimal order. The ROB keeps track of the instructions in flight and records when execution is completed. Once the instruction at the head of the ROB finishes execution, it is removed from the ROB and retired by the front end, marking the end of its life cycle within the processor. This ensures that the retirement of instructions follows the original program order, regardless of the order in which they were executed. Reservation stations hold instructions and their operands until all dependencies are resolved, and the necessary execution units become available. The execution units consist of arithmetic logic units (ALUs), floating-point units (FPUs), and other specialized units that perform the actual computations required by the instructions.

2.3.3 Memory Subsystem

The memory subsystem is responsible for handling load and store instructions, and the ordering of memory accesses [Lipp, 2021; Schor, n.d.].

Memory instruction execution consists of three steps. The first step is the memory address generation, where addresses are computed either with absolute or segment-offset addressing mode. The second step is address translation, which does a virtual-to-physical translation of addresses computed in

the first step, if virtual memory is used. Page tables are used to facilitate this translation. To speed-up future accesses to previously translated addresses, a translation-look aside buffer (TLB) is used to cache translation information. The final step is the actual memory access. In this step, data is first retrieved from the data cache and stored in a renamed register or the reorder buffer. A cache miss occurs if the data can not be located in the data cache, in which case, the data is requested from the main memory [Lipp, 2021].

The load buffer queues load instructions that can not complete when they are dispatched by the reservation station. It is also used to listen for stores performed in other cores against completed loads to ensure memory ordering. Meanwhile, the store buffer queues all store instructions before they are dispatched to the memory in order, that is, when they have retired [Lipp, 2021].

2.4 Branch Prediction

When a branch instruction is decoded by the front end of a processor, it often encounters a situation where it cannot determine the next instruction due to pending computation of the branch condition or target address. Instead of stalling the pipeline, the front end employs a technique known as branch prediction. It predicts the branch outcome according to the history of recently executed branches and proceeds to issue instructions accordingly.

By making correct predictions, this approach optimizes the execution flow by allowing younger instructions to proceed before the branch is resolved. However, in cases where the prediction is incorrect, the execution engine invalidates all younger instructions and requests the front end to resume execution from the correct branch outcome. These squashed instructions are ignored by the front end and are never committed, meaning their results are not stored to the architectural state. The same mechanism is applied when older instructions cause exceptions, such as division by zero or memory access violations, leading to the squashing of younger instructions.

Executed instructions that are eventually squashed are called *transient* instructions. While their results are not committed to the architectural state, any effects they have on the microarchitectural state of the processor are not reversed, which can potentially introduce vulnerabilities [Canella et al., 2019a; Kocher et al., 2019; Lipp et al., 2018].

2.5 Transient-Execution Attacks

While the results of such *transient* instructions are dropped, the microarchitectural components that were affected by the computation’s side effects are not reversed. Attackers can exploit this by triggering the transient execution of instructions that access confidential data and transmitting it through

microarchitectural components such as caches. These attacks are known as transient-execution attacks and have been studied extensively [Canella et al., 2019a; Kocher et al., 2019; Lipp et al., 2018; Xiong and Szefer, 2021].

Spectre [Kocher et al., 2019] is one of the most prevalent transient-execution attacks and will be discussed further below, as it is pertinent to the subject matter of this thesis.

2.5.1 Spectre

Spectre-type attacks are a type of transient-execution attack that takes advantage of the way processors execute instructions after a control or data-flow misprediction. The attacker manipulates the branch-prediction unit, causing the processor to speculatively execute instructions that do not appear in the actual instruction stream. This is made possible by multiple prediction units that work together to determine the outcome and target of a branch. By poisoning one or more of these prediction units, Spectre-type attacks direct the processor’s execution to “gadgets” which are code sequences that allow the attacker to uncover sensitive data by exploiting microarchitectural state changes.

In this thesis, we are specifically looking at the Spectre-PHT, which takes advantage of a component in modern processors called the Pattern History Table (PHT), which is a component of the branch predictor that predicts the outcome of conditional branches [Canella et al., 2019a]. The PHT consists of saturating counter that records recent branch history [Evtvushkin et al., 2018]. The attacker repeatedly trains the PHT with valid inputs to take a certain branch direction, such that targeted instructions are executed. After training, the attacker then executes the branch with an invalid input (e.g., out-of-bounds or restricted values) that would not cause the targeted instructions to execute according to the architectural specification. However, because the PHT is trained with valid inputs prior to this execution, it mispredicts the branch direction, causing processor speculatively executes the targeted instructions with the invalid input. The attacker can use this technique to perform architecturally-invalid reads, and forward the illicit data into a temporary buffer such as the cache for the purpose of retrieval into the architectural state later.

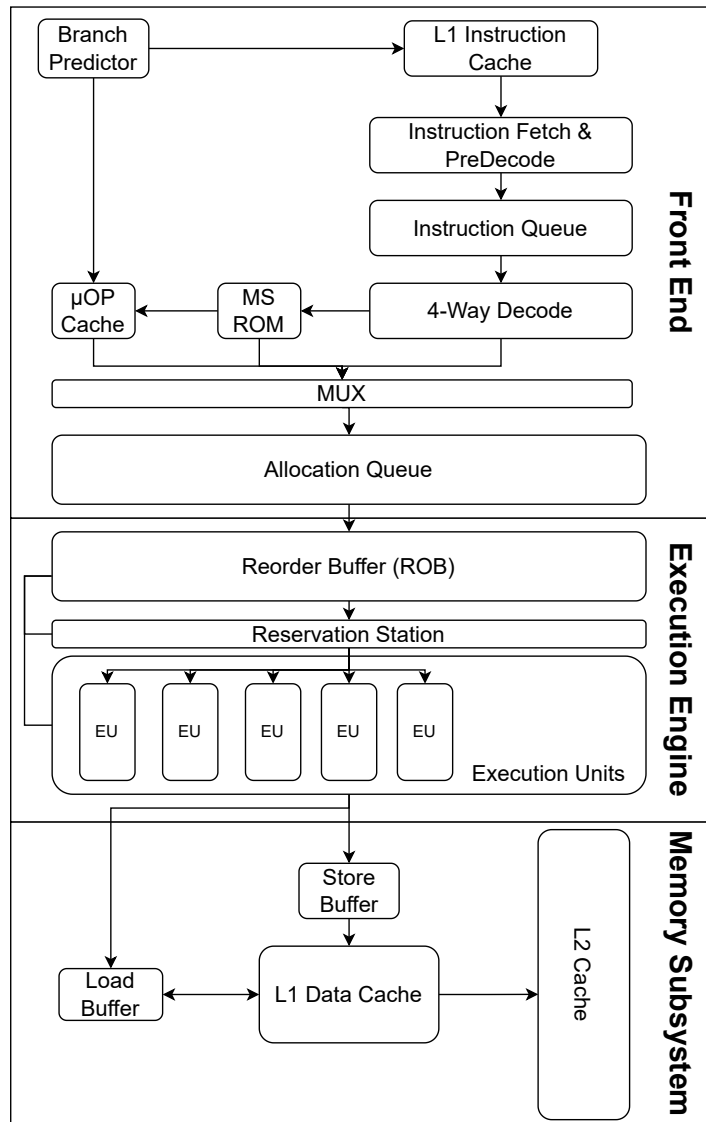


Figure 2.5: Simplified diagram of the Intel Skylake microarchitecture [Schor, n.d.].

Chapter 3

A Framework to Evaluate Cache Security

As we commence this thesis endeavor, in this chapter, our attention is directed towards the first cache side-channel attack defense approach, particularly safeguarding against contention-based attacks through secure cache designs. Various designs have been proposed to address contention-based cache attacks. For example, partitioned caches aim to prevent contention [Domnitser et al., 2012; Wang and Lee, 2007], while randomized caches create noise in the signal, obstructing attackers from analyzing the side-channel [Liu et al., 2016b; Qureshi, 2018; Qureshi, 2019; Tan et al., 2020; Werner et al., 2019]. Randomized caches generally prevent attackers from linking addresses to predicted cache line indexes, an essential element for the attack. Additionally, some suggestions aim to counter cross-core attacks by modifying the inclusion properties of shared cache levels in modern processors [Green et al., 2017; Kayaalp et al., 2017; Yan et al., 2017; Yan et al., 2019b].

As various proposals exist for defending against cache side-channel attacks, a technique to assess their security is essential. Several methods for evaluating secure caches have been proposed [Cock et al., 2014; Demme et al., 2012; Demme et al., 2013; Deng et al., 2019; Deng et al., 2020; Doychev et al., 2013; Ge et al., 2018; Ge et al., 2019; Ghasempouri et al., 2020; He and Lee, 2017; Köpf et al., 2012; Wang et al., 2019; Zhang and Lee, 2014; Zhang et al., 2013]. Nonetheless, these approaches have certain drawbacks as they either work only with basic cache models, focus on theoretical analysis, lack automation, or do not address the entire spectrum of cache attacks. Moreover, evaluating multiple metrics is crucial to bolster confidence in cache design security. Consequently, we explore the following question in this chapter:

How can we evaluate the security that cache designs offer against contention-based cache attacks?

To address this question, we introduce in this chapter a software framework, called CacheFX, to assess cache security. CacheFX allows for the emulation

of various cache designs, victims, and attackers, and measures the resulting leakage using various metrics.

3.1 Background

In this section we present specific background on secure caches. Several proposed cache designs aim to mitigate contention-based attacks. Their mitigation strategies are either based on partitioning [Domnitser et al., 2012; Wang and Lee, 2007] or randomization [Liu et al., 2016b; Qureshi, 2018; Qureshi, 2019; Werner et al., 2019].

Partitioned Caches Way-partitioned caches [Domnitser et al., 2012] enforce a strong partitioning between security domains¹ by letting each security domain use a different subset of the cache ways. Hence, domains not sharing cache ways will not see any interference. Alternatively, Partition-Locked (PL) [Wang and Lee, 2007] caches share the whole cache among all security domains, but offer to pin cache lines in the cache. These pinned lines cannot be evicted by other security domains, preventing contention-based attacks. However, aggressive pinning can starve other domains and severely degrade their performance.

CEASER. The CEASER cache [Qureshi, 2018] is based on an ordinary set-associative cache but uses encryption to randomize the mapping of addresses to cache sets. As a result, attackers need to first profile the victim’s accesses of interest to find a suitable eviction set before they can perform contention-based attacks. To limit the attacker’s time for finding such eviction set, CEASER regularly changes the encryption key. However, in this chapter we only measure information leakage in each key epoch (i.e., during the time period the cache uses the same key) and do not model re-keying. This allows assessing the security of pure cache-set randomization as it is needed to appropriately tune the re-keying interval for long-term security.

CEASER-S and ScatterCache. With improving eviction set building techniques [Qureshi, 2019; Vila et al., 2019], CEASER required higher key refresh rates to maintain security, resulting in increased overheads. CEASER-S [Qureshi, 2019] and ScatterCache [Werner et al., 2019] thus use a skewed cache [Seznec, 1993] to improve cache randomization. These skewed caches split the cache into partitions along its ways and use a different key to encrypt the address to index into each partition. The partition count can vary between 1 (CEASER) and the number of ways (ScatterCache) and allows to control the degree of randomization. As before, we omit re-keying to assess the security of pure cache randomization with skewing.

PhantomCache. PhantomCache [Tan et al., 2020] builds upon set-associative caches and maps each address to multiple sets via multiple hash functions, i.e.,

¹Security domain refers to code and data regions in memory that have the same level of security. (e.g., sensitive vs. normal).

it looks in multiple sets for a cache hit. On a cache miss, PhantomCache randomly selects one of the sets the address maps to and inserts the cache line into the chosen set. The number of cache sets looked up in parallel determines the degree of randomization and the cost of lookup. As before, we evaluate PhantomCache without re-keying.

NewCache. Rather than randomizing the cache mapping, NewCache [Liu et al., 2016b] is a more efficient implementation of a fully-associative cache. NewCache allows every cache line to be stored in any of the entries of the cache. Compared to a standard associative design, it optimizes power using a two-step look-up procedure: For a cache that can hold 2^n entries, NewCache first looks up $n + k$ index bits of the cache line address in a 2^n -element Content-Addressable Memory (CAM), which has a 1:1 mapping to the actual cache lines. Only if these $n + k$ bits match, this *index hit* is secondly followed by checking the tag for the respective entry. If the *tag hits*, the cache line is found and returned. If there is a *tag miss* for the same security domain in the second step, the tag and cache line are simply replaced. If there is an *index miss*, any of the 2^n cache lines in the cache is randomly replaced. While for large k NewCache resembles a traditional fully associative cache, a smaller k significantly reduces power and implementation cost.

3.2 Problem Description

With the abundance of secure cache designs, there is a clear need for systematically evaluating the security of caches to ensure that emerging cache architectures deliver the promised protection. Tackling this task, previous works [Cock et al., 2014; Demme and Sethumadhavan, 2014; Demme et al., 2012; Demme et al., 2013; Deng et al., 2019; Deng et al., 2020; Domnitzer et al., 2010; He and Lee, 2017; Wang et al., 2019; Zhang and Lee, 2014; Zhang et al., 2013] have suggested several metrics. However, all of these tend to suffer from some limitations to their practicality. For example, measuring the amount of information that can be transferred via a cache side channel [Cock et al., 2014] or the correlation between a specific victim’s activity and attacker observation [Demme et al., 2012] may not translate easily to cryptographic attack scenarios. Possibly the most common limitation of these approaches is the attempt to provide a single metric that somehow represents the security of the cache.

A General Cache Evaluation Framework. Instead of focusing on a single metric, this chapter proposes **CacheFX**, a framework for evaluating the security of cache designs. The main design aim of **CacheFX** is flexibility: **CacheFX** is extensible and allows evaluating various combinations of victims, cache designs, and attack strategies. As a proof of **CacheFX**’s generality, this chapter implements and evaluates three security metrics on nine different cache designs.

A Leakage Upper Bound. While we try to evaluate in realistic scenarios, CacheFX aims to provide an upper bound on the amount of leakage an attacker can obtain from a cache design. We thus assume an attacker who has significant control over the victim and is tightly synchronized with the victim’s execution. We further assume that the attacker has access to victim’s memory layout and code and thus knows the position of “interesting” data in the cache (e.g., cache lines containing secrets). This allows the attacker to craft inputs to the victim that may cause specific cache footprints. Unless required by the cache design, we assume a noise-free environment without any system activity besides the attacker and victim.

Using such strong assumptions allows CacheFX to properly evaluate the security offered by the cache design, as opposed to being misled by security guarantees stemming from other components’ noise. We note that previous works have demonstrated that attackers can find interesting cache lines [Liu et al., 2015] and overcome noise [Cock et al., 2014].

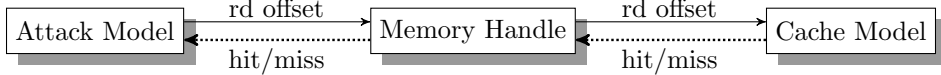


Figure 3.1: CacheFX design overview.

3.3 CacheFX Design

As mentioned above, CacheFX is designed to provide an easily extendable framework for evaluating 1) the security of emerging cache designs and 2) the applicability and complexity of new attack strategies to both deployed and emerging caches. To facilitate these goals, CacheFX is split into three major components as depicted in Figure 3.1. First, the *attack model* provides a set of interfaces and their implementations to model different attack and security evaluation strategies. Attack models use a *memory handle* to request reads, writes, and cache line invalidations to the memory system by specifying a certain offset into a memory region that is associated with the memory handle. The memory handle translates the requests to cache line addresses and queries the *cache model* correspondingly. The cache model returns whether the request hit or missed in the cache via the memory handle to the attack model, which then proceeds with the attack accordingly. Finally, the cache model provides a generic cache interface allowing for multiple different cache implementations. In the following, we give additional details about each of these components.

3.3.1 Cache Model

CacheFX’s cache model offers a generic cache interface that the memory handle and the attacker model can use to issue read, write, and invalidation requests

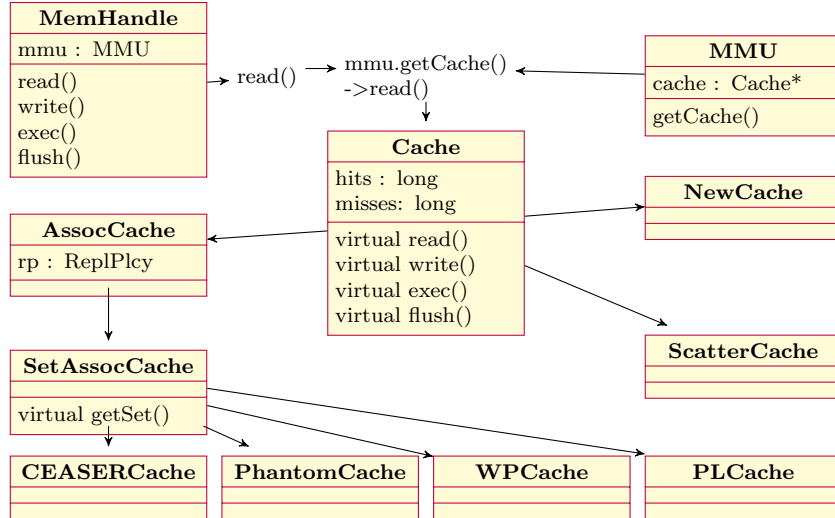


Figure 3.2: CacheFX overall architecture.

to the cache under test. For each of these requests, the cache responds with whether the request hit or missed. This indication removes the need to distinguish between hits and misses using (potentially noisy) timing measurements, providing an upper bound on the amount of leakage available to the attacker and consequently lower bounding the attack’s complexity.

Supported Cache Designs. The `cache model` currently provides multiple implementations of security-oriented cache designs: fully associative cache, set-associative cache, way-partitioned cache, partition-locked cache, CEASER and CEASER-S [Qureshi, 2019], ScatterCache [Werner et al., 2019], NewCache [Liu et al., 2016b], and PhantomCache [Tan et al., 2020]. These cache implementations are parameterized by the number of sets, ways, replacement policy, and cache-specific parameters. Unless a cache design mandates a specific replacement policy, all the implementations support LRU, Bit-PLRU, Tree-PLRU, and random replacement.

The Cache Interface. The internal interface of the cache model is defined by the `Cache` virtual class which acts as an interface with mainly four functions, namely: `read`, `write`, `exec`, and `flush` all of which embody actual requests that are sent to the cache hardware. Being an interface, these functions’ specific implementations are delegated to the derived classes of the `Cache` virtual class, such as, `NewCache`, `ScatterCache`, `SetAssocCache`, as shown in Figure 3.2. Each inherited class devises distinct mechanisms for dealing with the four requests which are carefully modeled based on the actual cache hardware functionality. This design promotes effective encapsulation and abstraction along with ease of implementation when adding new models as CacheFX’s simulation logic interfaces only with these top level functions. To support generic set-associative classes, CacheFX supports the `SetAssocCache` class, where each set is implemented as an `AssocCaches`. This allows for easy creation of the different set-based caches and reuses the code of the associative cache class,

e.g., to support multiple replacement algorithms.

The mechanism for selecting a set is implemented by the function *getSet()*. For the *SetAssocCache*, this is a simple modulus operation of the cache line and the number of cache sets. For *CEASERCache* and *PhantomCache* the set selection mechanism is based on a hashing algorithm. For both of these cases, the *getSet()* function is simply overloaded while the underlying implementation of *SetAssocCache* remains unchanged. Similarly, *WPCache* separates its context into two partitions, one for sensitive, and another for general data. This model simply consists of two *SetAssocCache* instances that are chosen based on the security context of a data access. Finally, *PLCache*'s design is based on the *SetAssocCache*, but with a minor change to the replacement algorithm to facilitate pinning of specific lines.

Statistics Generation. The abstract cache model automatically tracks the number of cache hits and misses for each security domain. In addition, the cache model can return the evicted address, if a cache access causes an eviction. While attackers usually do not have direct access to such information, providing the address allows us to apply novel and efficient techniques, such as the *Relative Eviction Entropy* (REE) in Section 3.4.1, for analyzing cache security.

3.3.2 Attack Model

CacheFX's `attack` model implements the actual adversarial strategy and evaluates the cache design under test. Currently, CacheFX supports three security evaluation strategies:

The Attacker. CacheFX allows to model synchronized pairs of victims and attackers, aiming to evaluate the security of cache designs with respect to realistic attacks, such as cache attacks against cryptographic block ciphers. CacheFX supports two types of attackers, *EvictionAttacker*, and *OccupancyAttacker*, both are subclasses of a generic *Attacker* class that manages the attack and collects the success statistics.

Information Leakage Assessment. CacheFX supports entropy-based security metrics that quantify information leakage during cache attacks (e.g., mutual information analysis). Most noteworthy, CacheFX implements a novel technique for evaluating information leakage in cache designs via the REE, by efficiently analyzing the statistical properties of a cache's cross-domain eviction behavior.

Eviction-Set Profiling. CacheFX provides an environment that allows for the evaluation of strategies to construct eviction sets for different cache designs.

Experiment Randomization and Automation. CacheFX allows to conduct each of these experiments multiple times with randomized address ranges to automatically obtain statistical data like maximum, minimum, etc. CacheFX hereby collects data such as cache statistics and attack success rates.

3.3.3 Victim Model

The purpose of this model is to simulate the behavior of victim applications within CacheFX’s simulation. CacheFX implements a number of victim models including:

- **SingleAccessVictim** is the simplest victim model that repeatedly accesses a single address in the memory.
- **AESVictim** simulates the behavior of AES encryption.
- **SquareMultVictim** imitates the square-and-multiply routine used in popular algorithms such as RSA.

All models are carefully crafted to resemble their actual attack characteristics. Take for example the *AESVictim* model whose code is taken directly from the original AES implementation, but is adapted to call into CacheFX’s API on each T-table access. In other words, all memory operations of the victim are redirected to CacheFX for further simulation. A similar approach is taken for *SquareMultVictim* where the cache line containing the multiplication code is executed conditionally based on the exponent. In this context, a call to CacheFX’s cache line read function is invoked at the beginning of the multiplication basic block to notify the simulator of the instruction cache read. With this methodology, we ensure high precision of CacheFX’s simulated model characteristic in comparison to the actual implementation.

Note that CacheFX victims currently focus on cryptographic code as its properties are well understood and are well suited to analyze the properties of the underlying cache design. However, CacheFX is generic enough to similarly model other leaking code, e.g., (de-)compression algorithms, data en-/decoders [Sieck et al., 2021] or neural networks [Yan et al., 2020].

3.3.4 The Attack Controller Function

As its name suggests, the purpose of this function is to moderate interactions between *Attack Model* and *Victim Model* (both of which subsequently interact with *Cache Model*).

Listing 3.1: Attack Controller main loop.

```
// Eviction or Occupancy Attacker
Attacker* a = createAttacker(attackerModel);

// Single, AES, or SquareMult Victim
Victim* v = createVictim(victimModel);

while (controller_run) {
    a->prime();
    v->cipher();
    a->probe();
}
```

Listing 3.1 outlines the main workings of the attack controller function. At the outset, pointers to both the *Attacker* and *Victim* classes are instantiated

to their desired model. At the heart of this function is a loop that interleaves the execution of the attacker and the victim, i.e., the prime, cipher, and probe methods. Note that the controller is agnostic of the specific implementations of these victim and attack functions and that the choice is left up to polymorphism.

The controller consolidates all major simulation components and can be thought of as the main driver of CacheFX.

3.4 Evaluation

Recognizing that no single metric is sufficient for measuring the resilience of caches to side-channel attacks, we evaluate emerging cache designs with respect to multiple metrics using our framework. First, the *Relative Eviction Entropy* (REE) metric measures the amount of information (in bits) that an attacker can deduce following a single memory access performed by the victim. Our second metric measures the complexity of creating eviction sets in randomized caches. Our third metric measures the complexity of performing cache attacks on cryptographic implementations. It evaluates both traditional attacks that seek to exploit eviction sets and cache-occupancy attacks [Cock et al., 2014; Shusterman et al., 2019], which do not require eviction sets.

We now discuss each metric in detail and compare different designs according to each of the measurement metrics.

3.4.1 Relative Eviction Entropy

In this section we introduce our *Relative Eviction Entropy* (REE) technique for effectively measuring the amount of information available to an attacker following a single memory access performed by the victim. We begin by observing that traditional mutual information analysis [Cock et al., 2014; Zhang and Lee, 2014] achieves such estimation for general side channels by computing a 2-dimensional joint probability distribution, which describes the likelihood of each victim activity (side channel input) to be mapped to an effect observable by an attacker (side channel output). For the case of caches, this means that for any address i accessed by the victim, and for all cached addresses a , we need to compute $p_e(a, i)$ which is the probability that a is evicted from the cache assuming that the victim accesses address i . We note that mutual information techniques typically measure average leakage across accesses, and thus do not capture the worst-case leakage.

Avoiding Quadratic Overheads. To avoid the quadratic overhead of computing the 2-dimensional joint probability distribution, we start by observing that natural cache designs typically do not have different eviction behavior between cache line addresses, and instead use the same replacement policy constantly across all cache lines. In addition to simplifying cache designs, this

property implies that all cache line addresses exhibit the same leakage behavior. Leveraging this fact, we can thus fix an arbitrary address i to be accessed by the victim, and simply sample $p_e(a, i)$ for all other addresses a . This avoids iterating over all possible values of i and thus makes the evaluation time of our metric linear in the size of the victim and attacker address spaces. When the value of i is fixed and clear from the context, we will simply omit i from the notation.

Quantifying Information Leakage. To capture the amount of leakage available to the attacker (in bits), we start from the intuition that fully associative caches with a random replacement policy leak the least amount of information among all cache designs that share cache lines between security domains, i.e., without consideration of partitioned caches. We argue that this assumption is reasonable, since fully associative caches with uniformly-random replacement only leak whether an address a was evicted or not, and do not reveal any information about which address i accessed by the victim caused the eviction of a . To evaluate leakage of new cache designs, we thus measure the REE as the statistical distance (in bits) of the eviction behavior of the tested cache design from the eviction behavior of ideal fully associative caches with random replacement.

Computing Relative Eviction Entropy. Our strategy for computing a cache design’s REE is as follows. First, we allocate a chunk of memory in the adversary’s address space, typically a small multiple of the cache size. We denote the set of cache line addresses within that adversary’s memory as $a \in [0, \dots, N - 1]$. Second, for a single victim access to some fixed address i , we estimate the eviction probability $p_e(a)$ for each cache line $a \in [0, \dots, N - 1]$ in the adversary’s memory, using our implementation of the cache design under test. The distribution $p_e(a)$ will reflect the cache’s placement policy: e.g., if the single victim access can evict every adversary address a , as in a fully associative cache with random replacement, $p_e(a)$ will be uniform among all adversary addresses a . If the single victim access can only evict adversary addresses a mapping to the same cache set in a set-associative cache, $p_e(a)$ will be uniform among those addresses a mapping to the same set as the victim address i and zero otherwise. The reference eviction distribution of a fully associative cache with random replacement is set to $p_u(a) = 1/N$ for all addresses a , reflecting that every adversary address is equally likely to get evicted. Finally, we compute the REE as the statistical distance in bits between the eviction probability distributions $p_e(a)$ and $p_u(a)$ using the Kullback-Leibler (KL) divergence to measure,

$$D_{KL}(p_e||p_u) = \sum_{a \in [0, \dots, N-1]} p_e(a) \log_2 \frac{p_e(a)}{p_u(a)}. \quad (3.1)$$

Note that the KL divergence does not fulfill the requirements of a metric and is asymmetric. Nevertheless, $D_{KL}(p_e||p_u)$ describes the relative entropy of $p_e(a)$ with respect to $p_u(a)$ and is a measure of the information lost if $p_u(a)$ was used

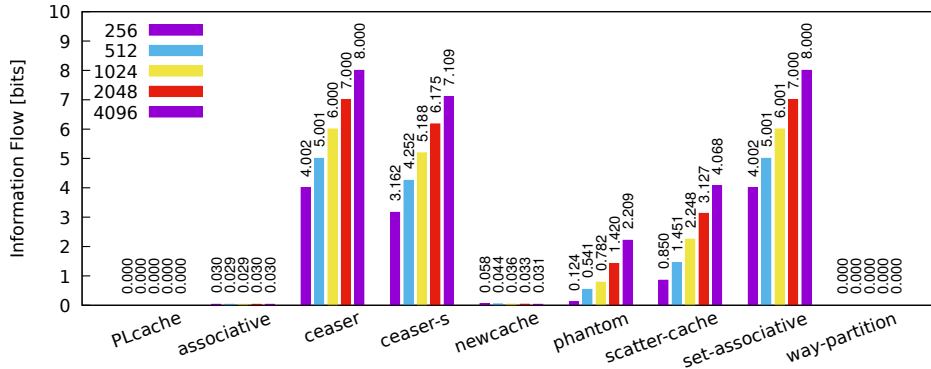


Figure 3.3: REE across cache designs with random replacement. All but NewCache and the fully associative cache use 16 ways.

to approximate $p_e(a)$. Mapped to cache side channels, the KL divergence thus nicely characterizes the leakage of a cache design with an eviction probability distribution $p_e(a)$ relative to the distribution $p_u(a)$ in a fully associative cache design.

Sampling $p_e(a)$. As $p_e(a)$ is generally unknown, we sample $p_e(a)$ and use the plug-in estimator [Zhang and Grabchak, 2014] for the KL divergence to estimate the REE: we simply count the number of evictions for the attacker’s cache lines when the victim repeatedly accesses a fixed, randomly chosen address. More specifically, we first fill the cache by randomly accessing cache line addresses from the memory chunk corresponding to the attacker’s security domain. To keep track of self-evictions and hence the attacker’s lines that are actually cached, we utilize our cache model’s capability to return which cache line is evicted with each access, as described in Section 3.3.1. We note that this is an over-approximation of the attacker’s capabilities, as on real systems this translates to an attacker who can perfectly monitor cache evictions and accurately determine address collisions in the cache. Once the cache is entirely filled with the attacker’s data, we access a fixed secret address from the victim’s security domain, forcing an eviction of one of the attacker’s addresses. We then increment the eviction counter for the attacker address that is being reported as evicted from the cache. We repeat this sampling step multiple times and finally divide the per-address eviction counts by the total number of observed evictions, thereby obtaining $p_e(a)$. The repeated sampling procedure reduces the error of the sampled eviction probabilities proportional to $\sqrt{(r)}$, where r is the number of samples collected.

Definition *Relative Eviction Entropy* is evaluated as follows:

1. Select a victim address v and initialize a pool of memory \mathcal{P} .
2. Sample $p_e(a) \forall a \in \mathcal{P}$, which is the probability that v evicts address a from the cache.
3. Compute the REE via Equation 3.1, where $p_u(a)$ denotes a uniform distribution over all addresses $a \in \mathcal{P}$.

Evaluation Results. Figure 3.3 depicts the information leakage in the

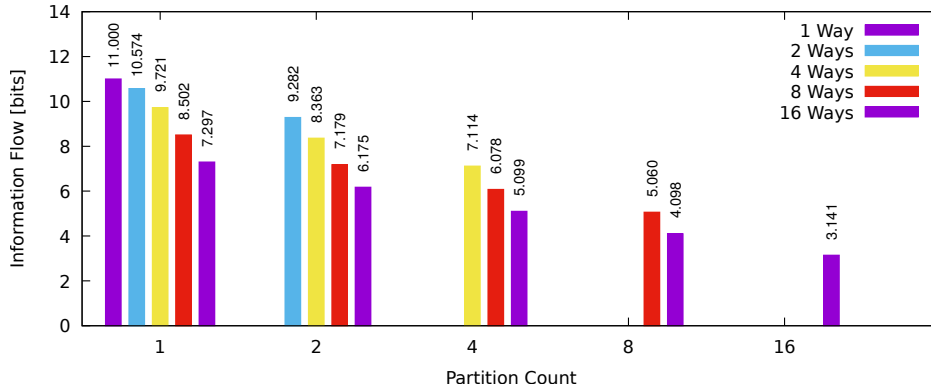


Figure 3.4: REE for CEASER-S with 2048 lines depending on ways and partitions.

analyzed cache designs for various cache sizes and using random replacement. While the partitioned cache designs exhibit zero leakage, the leakages for CEASER and set-associative caches is the number of sets, i.e., $\log_2(\#sets)$ bits, thereby confirming the validity of our results. Next, we attribute the slightly above-zero leakages in NewCache and the fully associative cache to statistical noise. Note that CEASER-S and ScatterCache (with 2 and 16 partitions, respectively) show considerably lower leakage than standard set-associative caches. Moreover, as PhantomCache is looking up 8 sets, i.e., 128 lines, in parallel, PhantomCache stands out with significantly lower leakage per access than other designs, but also hurts chip area and power consumption.

Figure 3.4 analyzes the leakage in skewed caches like CEASER-S depending on way and partition count. Figure 3.4 clearly shows that increasing the number of ways and partitions effectively reduces leakage, with the difference between the best and worst configuration being 8 bits per access.

Supporting More General Cache Designs. We note that our Relative Eviction Entropy method can be computed in linear time, allowing us to evaluate different cache designs within minutes. However, we do assume some properties of the replacement policy of the cache being tested, namely that every line in the considered cache design exhibits the same leakage behavior, which in turn is independent from the specific address accessed by the victim. We rely on this assumption in our procedure for sampling $p_e(a)$, evaluating the eviction distribution using only a single fixed address accessed by the victim. We argue that this assumption is natural and holds for most cache designs, including all the caches considered in this chapter, as typical replacement policies do not differentiate between cache line addresses. While a single access does not reflect practical attack scenarios, it gives strong insight into the theoretical leakage caused by the caches’s structural mapping of addresses to cache lines. However, while the REE is a highly efficient tool to approximate leakage, we recognize that its underlying assumptions may be limiting its use in various corner cases, e.g., when replacement decisions are based on the actual

address.

To better understand the practical exploitability of leakage determined via the REE, we conduct application-specific tests using cryptographic routines later in Section 3.4.3. However, note that the REE metric can be easily adapted to other cases as well, by simply testing multiple victim addresses and reporting the range of the occurring leakage as a function of victim’s address.

3.4.2 Eviction-Set Creation

To perform contention-based cache attacks, attackers first construct suitable eviction sets, i.e., minimal sets of addresses in their own address space that collide with the victim’s accesses of interest. Due to its perceived importance, multiple cache designs aim at randomizing the cache to prevent efficient eviction-set creation and thus contention-based attacks.

Definition Eviction-set creation is evaluated as follows:

1. Select an eviction-set construction algorithm \mathcal{A} , an address a , a target eviction set size T , and number of repetitions R
2. Select a pool \mathcal{P} of candidate addresses.
3. Run \mathcal{A} on pool \mathcal{P} to find an eviction set $\mathcal{E} \subset \mathcal{P}$ for address a until the target size $|\mathcal{E}| = T$ is reached and repeat R times.
4. Evaluate minimum/maximum/median/average for R samples of key metrics, e.g., number of attacker accesses and final set size.

Constructing Eviction Sets on Randomized Caches. Previous works proposed a range of methods for finding eviction sets in randomized caches. Taking a top-down approach, the Single Holdout Method (SHM) and the Group Elimination Method (GEM) [Qureshi, 2019; Vila et al., 2019] both start from a large set of attacker addresses that evicts a certain victim address and then shrink this conflict set to a minimal eviction set by trying to remove (groups of) addresses while continuously verifying that the cache conflict remains. Taking a bottom-up approach, the Prime+Prune+Probe (PPP) method [Purnal et al., 2021b] pre-fills the cache with a set of candidate addresses, and subsequently triggers the victim access of interest. PPP then tests for cache misses in its candidate set, thereby locating conflicting addresses. Note that all of these approaches allow for optimizations specific to the cache replacement strategy in use.

Evaluating Difficulty of Eviction Set Construction. As protecting against eviction set construction is a major design goal for randomized caches, CacheFX allows to evaluate the effectiveness of SHM, GEM, and PPP on a candidate cache design. In particular, CacheFX quantifies the number of memory accesses required by an attacker, the number of conflicting addresses found, and the success rate of using the found addresses for evicting the victim address. These figures eventually allow to configure cache re-keying intervals, e.g., for CEASER and CEASER-S. To set a level playing field and support an equal comparison across cache designs, we use the same implementations

of eviction-set construction techniques for all evaluated cache designs. We intentionally avoided cache-specific optimizations, opting for comparable results rather than for optimal strategies. Specifically, all of our implementations iterate until they find (or shrink a conflict set to) the minimum number of addresses required for an eviction set, or until a predefined maximum iteration count is reached. The latter is a necessity to perform bulk testing as some algorithms do not terminate for every cache design. We leave the question of identifying optimal strategies and evaluating them to future work.

Measurement Setup. To measure the success rate, we set up a clean cache environment 1000 times and count the number of successful evictions of the cached victim address given the found eviction set. We extracted the cache hit/miss statistics to evaluate the number of attacker accesses needed for eviction-set creation. We determine the number of true conflicts in the eviction set by testing every found address for a collision with the victim address in the cache. While this is not directly possible on real systems, CacheFX provides this feature to assess how well each algorithm works for every cache design.

In our experiments, we used random replacement, 2048 lines and 16 ways where applicable, i.e., except for NewCache and the fully associative cache, which only have one set. We operated CEASER-S with 2 partitions, NewCache with $k = 2$, and PhantomCache with 8 parallel set lookups. We set up the algorithms to look for as many addresses as there are cache ways. For PhantomCache, however, we require 8x the number of ways, because it can place lines in 8 different sets.

Evaluating the Number of Memory Accesses for Eviction Set Construction. Figure 3.5 shows the number of memory operations done by SHM, GEM, and PPP for different cache designs. As L1 cache accesses take about five CPU cycles, these results give an indication about the execution time of each technique when used against a specific cache design.

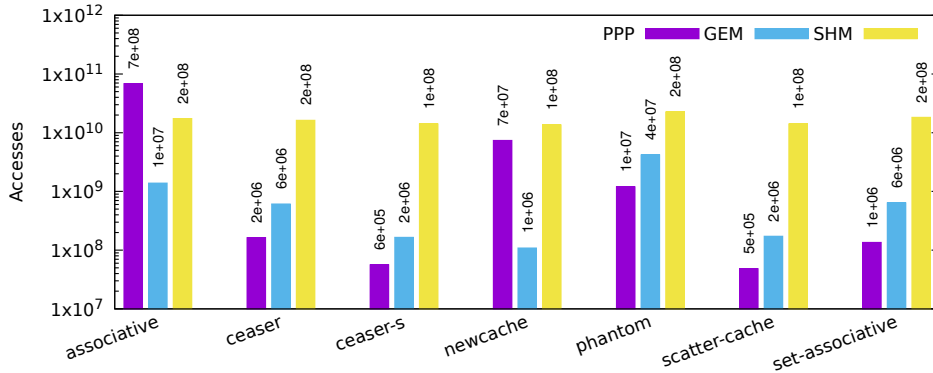


Figure 3.5: Number of memory accesses required by eviction-set building techniques for different 2048-line caches.

As the figure shows, the number of memory accesses for SHM is the highest, and in the same order of magnitude for all designs. In contrast, the complexity of PPP scales with the eviction set size, e.g., PPP is two orders of magnitude

faster for ScatterCache than for NewCache. PPP also tends to be more efficient for skewed caches, as it is 3x faster for CEASER-S than for CEASER. The performance of GEM is mostly in between PPP and SHM, but tends to be faster than PPP in the case of large eviction sets (e.g., for NewCache). Figure 3.6 gives further performance figures for CEASER-S and shows the linear increase in complexity with the number of cache lines.

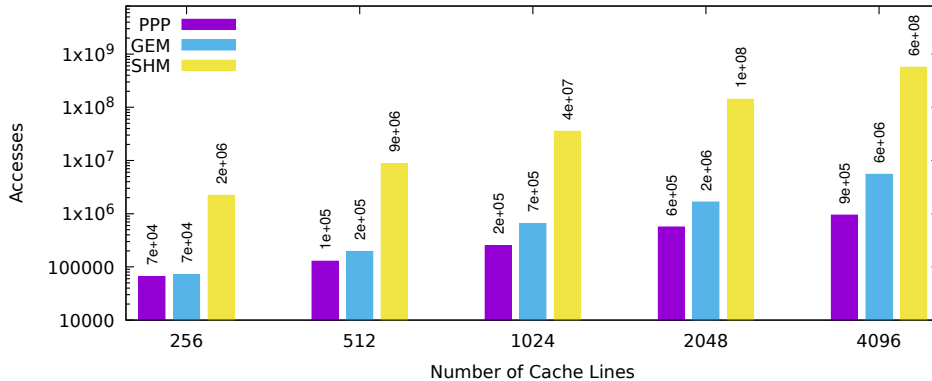


Figure 3.6: Number of memory accesses required by eviction-set building techniques for CEASER-S depending on cache size.

Evaluating Eviction Coverage. Different eviction set construction techniques can also produce eviction sets of different quality. Figure 3.7 thus shows the percentage of addresses in the found eviction sets that truly conflict with the victim address: PPP works best for all of the tested cache designs, producing eviction sets where all of its addresses truly conflict with the victim address. In contrast, SHM and GEM are less reliable, producing eviction sets where many of the addresses do not conflict with the victim address. The main reason for this is that SHM and GEM are highly susceptible to noise, which stems from both random replacement and cache skewing.

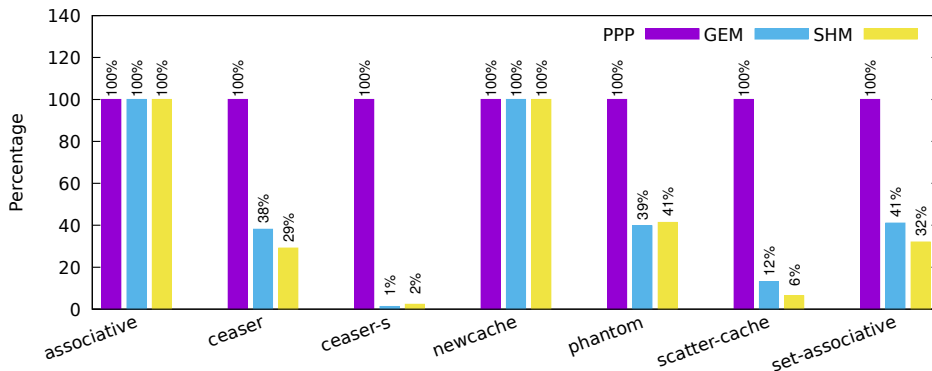


Figure 3.7: Percentage of addresses in the constructed eviction sets that conflict with the victim's address, using different eviction-set construction techniques and 2048-line caches.

To verify this, Figure 3.8 shows the constructed eviction sets’ sizes for SHM, GEM and PPP. Except for NewCache and fully associative caches, both SHM and GEM stop shrinking the conflict set before it becomes minimal, which results in eviction sets where many addresses do not conflict with the victim address. This effect is particularly strong for the skewed cache designs CEASER-S and ScatterCache. Moreover, SHM and GEM also fail on PhantomCache, where both algorithms terminate with 10x as many addresses as needed. Finally, for NewCache and fully associative caches every address is equally suitable for an eviction set, which automatically results in 100% of the addresses conflicting with the victim address.

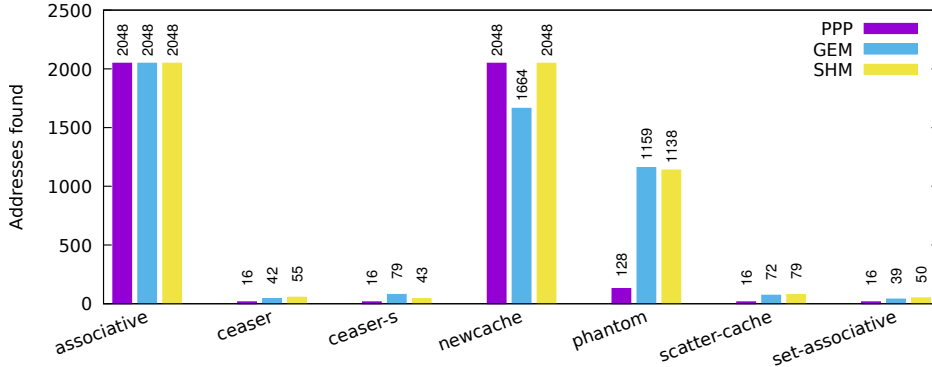


Figure 3.8: Eviction set sizes found by eviction-set building techniques for different 2048-line caches.

Evaluating Eviction Success Rate. We also evaluate the constructed eviction sets for their ability to effectively evict the victim address of interest. As Figure 3.9 shows, the eviction sets found by all three eviction set construction techniques perform equally well for CEASER, NewCache, set- and fully associative caches. For CEASER-S and ScatterCache, PPP yields better eviction success rates than SHM and GEM, because PPP is generally more accurate (cf. Figure 3.7). For PhantomCache, however, GEM and SHM yielded better eviction rates as the found eviction set makes up roughly 50% of the cache. As skewed caches exhibit a significantly smaller probability of successful eviction (e.g., 2-4% for ScatterCache), eviction sets might be chosen larger to obtain high eviction probabilities and and Prime+Probe observability.

Obtaining a Specific Eviction Probability. To learn how many addresses would be needed to yield a certain eviction probability α , we start with an empty eviction set and successively add conflicting addresses until the eviction probability reaches α . Figure 3.10 presents the results of this routine for $\alpha = 90\%$, across different caches and replacement policies. It shows that LRU and Tree-PLRU allow for smaller eviction sets than Bit-PLRU and random replacement. In addition, skewing significantly increases the number of conflicting addresses needed, e.g., ScatterCache requires 10x more addresses than CEASER with equal sets and ways.

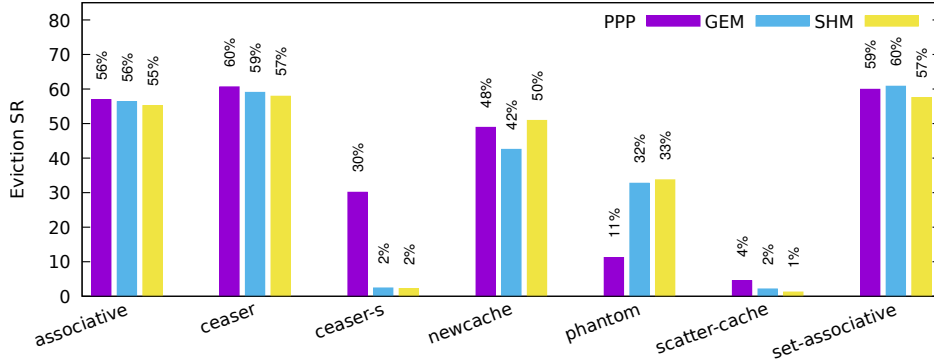


Figure 3.9: Eviction success rate for the eviction sets found for different 2048-line caches.

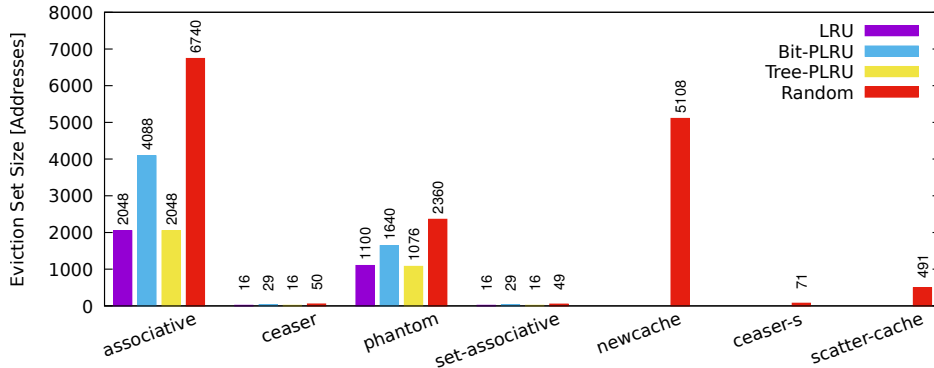


Figure 3.10: Eviction set size for 2048-line caches and 90% eviction probability.

3.4.3 Eviction-Set Attack

This section focuses on measuring the security offered by cache designs when performing attacks on cryptographic implementations. To that aim, we simulate victims that use a cryptographic algorithm while the attacker tries to learn enough information to distinguish between two keys used by the victim. We use two cryptographic algorithms, each representing a different type of cache attack.

The AES Victim. Our AES victim is based on code from OpenSSL, which uses a set of tables, called T-tables, implemented as arrays. The attack focuses on the first four accesses made to the first T-table during the encryption. The two keys are selected such that, when encrypting some *vulnerable* plaintexts with the first key, all of these accesses fall in the first cache line of the T-table. Conversely, when encrypting vulnerable plaintexts with the second key, each of the four accesses falls in a different cache line. Finally, to further facilitate the attack, we allow the attacker to choose as many random vulnerable plaintexts as required for the attack.

In a more realistic scenario, the attacker can guess the characteristics of the vulnerable plaintexts. Specifically, the attacker can fix the first byte of the

plaintext, and test every combination of plaintext values for the other three bytes that affect access to the first T-table. For each such combination, the attacker then performs the attack. If any of the combinations show statistical difference between the keys, the attack succeeds. With T-tables that span 16 cache lines, there are 16 possible values for each of these bytes. Thus, allowing to select vulnerable plaintexts represents a constant factor of $16^3 = 4096$ improvement in attack complexity.

Modular Exponentiation Victim Our second victim implements modular exponentiation, a core operation in multiple public-key schemes, e.g., RSA. Our modular exponentiation victim gets a 2048-bit base b , a 2048-bit modulus m , and a 32-bit exponent e . The victim then uses the square-and-multiply algorithm [Gordon, 1998] to calculate $b^e \bmod m$. The square-and-multiply algorithm maintains an accumulator a that is initialized to 1. For each bit of the exponent e , the algorithm squares a , and if the bit is set the algorithm also multiplies a by b , reducing a modulo m as necessary. Thus, the multiplication code is only executed when the exponent bit is 1, and the effect on the cache is that when the bit is 1, more cache lines are accessed.

The keys are selected so that the value of a bit at a specific index (7 in our tests) of the exponent is 0 in the first key and 1 in the second. The other bits of each exponent are randomly chosen. We simulate an attacker that runs concurrently with the victim. The attacker can manipulate the cache whenever the victim finishes processing an exponent bit to distinguish between the number of cache lines accessed depending on the exponent bit.

Attacker Setup In the attack setup phase, the attacker is provided with an eviction set that evicts a monitored victim cache line with a probability 90%. We construct this eviction set by successively adding conflicting addresses to an initially empty set as outlined in Section 3.4.2. See there for an analysis of eviction-set construction.

Attacker Procedure. The attack proceeds as a sequence of rounds. In each round, the attacker asks the victim to encrypt a plaintext with the two selected keys, randomizing the order of using the keys in each round to avoid cache effects that depend on the order of the use of keys. Before each encryption, the attacker accesses the eviction set three times to prime the cache. After each encryption, the attacker accesses the eviction set, counting the number of cache misses during these accesses. Finally, the attacker calculates the average number of cache misses for each key, and stops when achieving a 95% confidence that the averages differ, or when hitting a predefined number of rounds. (10^3 for the modular exponentiation and 10^5 for AES.) To overcome the case where the eviction set and the victim all fit in the cache, the attacker accesses some arbitrary memory when no cache evictions are observed during a round.

Selecting Cache Designs for Evaluation. We perform the attack on a sample of the cache designs considered in this chapter. First, we do not test partitioned caches, because these do not leak information as there is no

resource contention between the attacker and the victim. Secondly, to ensure that results are comparable, we limit our experiments to a cache size of 256 lines. Where applicable, we vary the associativity, testing all powers of two between 1 and 16. For each configuration, we run the attack 1,000 times and report the median of the number of encryptions required for distinguishing the keys. We use the median rather than the mean because in some cases the distribution has a long tail, skewing the mean towards a small number of cases where many encryptions are required.

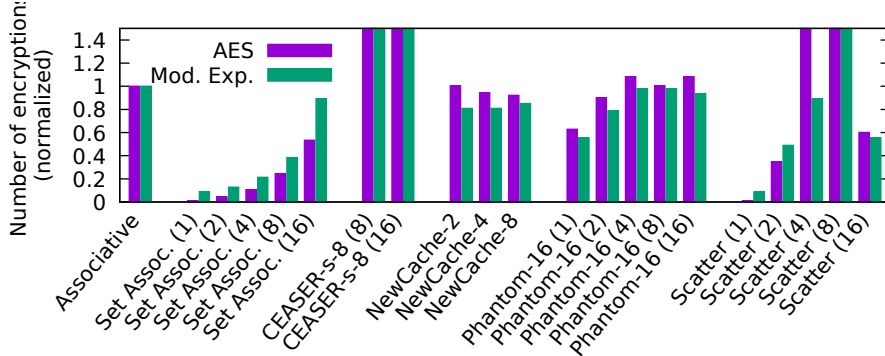


Figure 3.11: Eviction-set attack: Number of encryptions required to break AES and modular exponentiation with random replacement. CEASER, CEASER-s-1, and Phantom-1, which show behavior similar to set associative caches, have been omitted from the figure. (Normalized to a random-replacement associative cache.)

Observing Key Leakage. Figure 3.11 shows the median number of encryptions required for the attacks when using random replacement. We normalize the results to the number of encryptions required for the fully associative cache. (10,590 and 94 for AES and modular exponentiation, respectively.) For brevity, we omit the results of CEASER, CEASER-S with one partition, and PhantomCache with one set lookup, all of which do not seem to offer any advantage over a set-associative cache with the same associativity.

The figure shows that all NewCache variants and PhantomCache with 16 set lookups are mostly equivalent to the fully associative cache. CEASER-S with 8 partitions provides a stronger protection: the majority of the AES attacks on CEASER-S with 8 ways and of the modular exponentiation attacks with 16 ways were not successful.

The results with ScatterCache are a mixed bag. When the associativity is four or eight, the design provides a good protection, equivalent or surpassing the fully associative cache. (In particular, the AES attack fails in most cases on an 8-way cache.) However, the protection is lower for the other cases.

3.4.4 Cache-Occupancy Attack

We now turn our attention to an emerging cache attack strategy that ignores spatial information and instead only utilizes the victim’s overall cache usage [Maurice et al., 2015a; Shusterman et al., 2019; Shusterman et al., 2021a]. To measure resistance against so called *cache-occupancy attacks*, we use the same cryptographic victims as in Section 3.4.3. The attacker is still tasked with distinguishing between two keys, but instead of using an eviction set targeting a specific cache line, the attacker uses a cache-size buffer and counts the number of cache misses when scanning the buffer. (A different sized buffer may also work [Shusterman et al., 2021b], but requires further investigation.) Most other aspects of the attack are the same as in our eviction-set attack. We do not, however, handle failed eviction because using a cache-size buffer ensures contention on the cache.

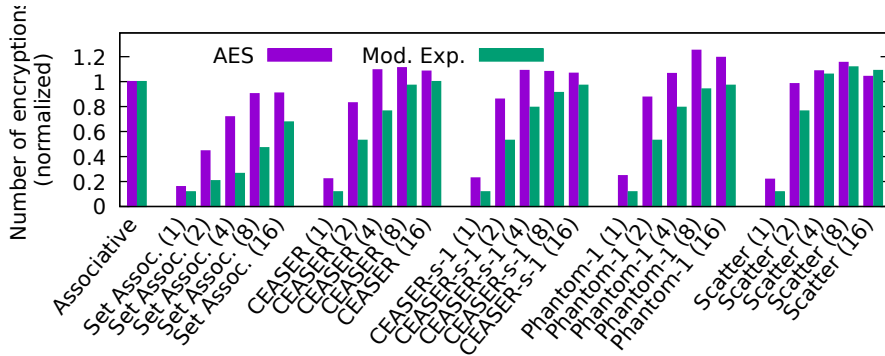


Figure 3.12: Occupancy attack: Number of encryptions required to break AES and modular exponentiation with random replacement. CEASER-S-8, NewCache, and Phantom-16, which show behavior similar to fully associative caches, have been omitted from the figure. (Normalized to a random-replacement associative cache.)

Observing Key Leakage. Figure 3.12 shows the median number of encryptions required for the cache occupancy attacks when using a random replacement strategy. As in Figure 3.11, we normalize the results to the number of encryptions required for the fully associative cache. Similar to the eviction-set attack, NewCache, CEASER-S with 8 partitions, and PhantomCache with 16 set lookup achieve a protection similar to that of fully associative cache. (We omitted these three from the figure for brevity.) Most other configurations achieve a protection level which is significantly better than set-associative caches, in particular for the attack on AES.

Due to normalization, Figures 3.11 and 3.12 do not show that occupancy attacks on the fully associative cache require significantly less encryptions than eviction-set attacks. (5664 and 68 for AES and modular exponentiation, compared to 10590 and 94.) The cause is that the eviction set algorithm targets 90% eviction rate, which for fully associative caches leads to eviction sets that

are larger than the cache-sized buffer used in the occupancy attack and thus more self evictions.

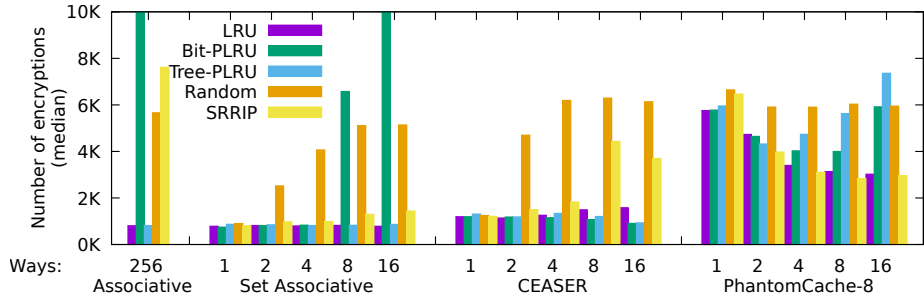


Figure 3.13: Median number of encryptions required to break AES with different cache designs and replacement algorithms. Fully associative and 16-way set associative caches are not fully represented, requiring 16,984 and 22,116 encryptions for Bit-PLRU, respectively.

Comparing Different Replacement Algorithms. Figure 3.13 shows the effect of changing the replacement policy on the attack complexity. As the figure demonstrates, in most cases, caches with a random replacement policy offer significantly better protection than those with deterministic replacement.

For deterministic replacement policies, we observe that CEASER only provides marginal benefit over set-associative caches, whereas PhantomCache provides a significantly better protection than other cache designs. We believe that the reason is that PhantomCache is inherently non-deterministic, hence, even with deterministic replacement algorithms, PhantomCache can reduce the correlation between the victim’s access and the attacker’s observation.

Attacks on deterministic cache designs that use bit-based pseudo-LRU replacement exhibit an anomaly that increases the number of encryptions required for statistical confidence. The cause is that the algorithm experiences some rare cases where a single cache miss causes cascading evictions of the eviction set. These rare cases increase the variance of the number of evictions observed, and with it the number of samples required. Modifying the attack to ignore outliers will eliminate these rare cases and significantly improve the attack. Hence, the results do not indicate that bit-based pseudo-LRU is more secure than random replacement.

3.4.5 Optimal Eviction-Set Size

In Section 3.4.2 we evaluate eviction sets based on the probability of evicting a victim cache line from the cache. However, as discussed in Section 3.4.4, larger eviction sets can result in lower attack efficiency, apparently due to self evictions. Specifically, increasing the size of the eviction set increases the probability of cache conflicts between elements of the eviction set. These self evictions introduce measurement noise that increases the variance in the

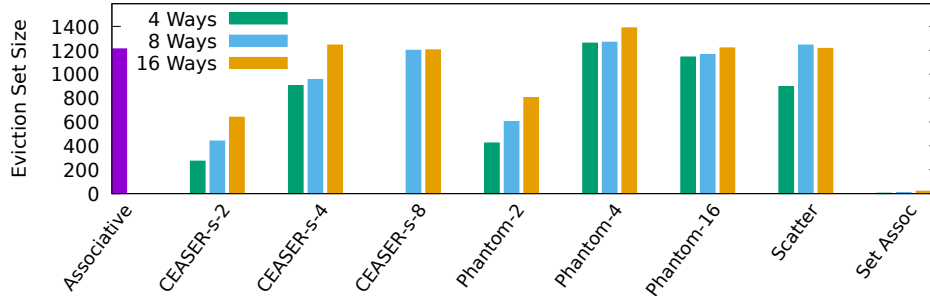


Figure 3.14: Optimal eviction set sizes for 1024-lines caches

measurements and consequently the number of samples the attacker needs to observe to distinguish the keys. As a secondary effect, larger eviction sets require more memory accesses for both the prime and the probe steps of the attack, reducing attack efficiency.

As a final example of the flexibility of CacheFX, we now use it to find the eviction-set size that allows for the most efficient attack. Specifically, we experiment with various cache designs, all with size 1024, our AES victim, and our eviction-set attacker. We vary the eviction-set size between 1 and 2048, and measure the median number of encryptions required for distinguishing the keys. Figure 3.14 reports the eviction-set size that allows the attack with the minimal median. As we can see, a lower associativity allows for smaller eviction-set sizes. However, when the associativity grows to 16, in most cache designs the best eviction-set size is similar to that of a fully associative cache, indicating that occupancy-based attacks are as effective as eviction-set attacks.

3.5 Threats to Validity and Limitations

At the moment, CacheFX does not support the evaluation of cache hierarchies. Consequently, designs that rely on the hierarchy for defense are outside the scope for this chapter. Moreover, evaluations using CacheFX currently assume a noise-free scenario, which provides a conservative security estimate as the absence of noise is the best case for attackers. However, practical cache attacks also face systematic and random noise stemming from other system activity.

For our cryptographic attack evaluations, CacheFX models a strong, synchronized attacker and an artificial victim that computes (and leaks) upon the attacker’s request. As for noise, this is a very strong attack model that allows to obtain a lower bound for security. However, real-world attacks involve various kinds of complexities that cannot be assessed with a simple model like CacheFX.

Another aspect of secure caches is their performance. CacheFX currently does not support the evaluation of cache performance.

CacheFX simplifies cache models to efficiently analyze the security of caches

against contention-based attacks. As a result, **CacheFX** does not lend itself to model cache-based attacks that relate to speculative execution [Katzman et al., 2023] or other microarchitectural structures, e.g., cache ports [Yarom et al., 2017] and fill buffers [Van Schaik et al., 2021]. Note that **CacheFX** allows to model **Flush+Reload** attacks, but as **Flush+Reload** is well understood we do not expect new insights from doing so.

3.6 Related Works

Past work on evaluating the security of caches against side channel attacks mainly focused on three aspects: 1) formal model of cache and theoretical analysis of information leakage, 2) metrics for empirical quantification of information leakage, 3) modeling of cache side-channel attacks.

Formal Cache Model and Theoretical Analysis. This line of research [Doychev et al., 2013; Köpf et al., 2012] tries to formally model the state change of the cache and extend the program execution semantics to include cache state changes by leveraging prior work on formal analysis of cache miss rates. Eventually they can estimate the number of reachable cache states and give an upper bound on the leakage in terms of channel capacity, for a given program under analysis. Similarly, Ghasempouri et al. (2020) models caches and cache attacks as automata to verify cache security using model checking. Due to the restrictions of formal methods, these works are limited to simple cache models (e.g., set-associative cache with LRU replacement) and can only give a very loose upper bound of leakage. Hence, they are not suitable for comparing the security of various complex secure cache designs. In contrast, **CacheFX** empirically evaluates a number of metrics to quantify side-channel leakage in software cache models and evaluates the exploitability of cache leakage for programs such as cryptographic algorithms.

Metrics for Empirical Quantification of Information Leakage. Another line of research introduces metrics to empirically evaluate the security of cache designs and implementations, such as by using mutual information and min-leakage [Cock et al., 2014], by using a linear correlation coefficient between oracle traces and the attacker’s observations [Demme and Sethumadhavan, 2014; Demme et al., 2012; Demme et al., 2013; Zhang et al., 2013], by measuring the accuracy of deep learning models trained to learn the relationship between victim accesses and the attacker’s cache observations [Zhang et al., 2018], or by modeling and statistically analyzing cache side channels using communication theory [Bourgeat et al., 2020]. **CacheFX** as well tries to empirically characterize the leakage of cache designs. However, as we point out, a single metric is insufficient to entirely capture cache security. Moreover, none of these works looks at cache occupancy channels or tries to assess security by using well-studied cryptographic targets.

Modeling of Cache Side-Channel Attacks. Some works tried to model caches and cache attacks such as to detect and quantify cache leakage. For

instance, Zhang and Lee (2014) model the cache as a finite state machine to identify interference and determine the mutual information. He and Lee (2017) model cache attacks as a Probabilistic Information Flow Graph (PIFG) to derive for each cache and attack an overall probability of success. Wang et al. (2019) derive a risk score from modeling attacks using Petri nets and calculating the success probabilities of concrete attacks. Deng et al. [Deng et al., 2019; Deng et al., 2020; Deng et al., 2021] model cache attacks as a series of three consecutive read/invalidation steps, identify vulnerable three-step patterns using a simulator, and use the model for evaluating the security of the caches in multiple ARM devices. In addition, their work introduces a Cache Timing Vulnerability Score (CTVS) from running vulnerable patterns on real machines. While these prior works greatly improve the understanding of cache attacks, many are based on simple cache models. CacheFX thus takes another step forward and automatically evaluates arbitrary software models of cache designs w.r.t. to a number of different metrics and attack complexity to provide a comprehensive security report.

3.7 Conclusion

As the threat of cache side-channel attacks continues to grow, numerous measures have been implemented to protect systems from such attacks. One such measure is the development of secure caches, which aim to prevent contention-based cache attacks at their source—the cache itself. Given the various existing proposals for defending against cache side-channel attacks, it is crucial to have a technique for assessing their security. This chapter introduces a flexible framework called CacheFX, which addresses the need to evaluate cache designs for security using multiple metrics.

Using CacheFX, we conducted experiments with three related metrics to assess and compare multiple secure cache designs. Our observations revealed that all non-partitioned caches leak information, with the leakage being significant enough to enable cryptographic attacks. However, partitioned caches are likely not practical for many use cases. We also demonstrated that a single metric may not fully capture all complexities. For instance, the REE of CEASER-S indicates reduced leakage as the number of ways or partitions increases (see Figure 3.4). This is consistent with the intuition that leakage in set associative caches correlates with the number of cache sets, which decreases as associativity increases, given a constant cache size. However, caches with 4 or 8 partitions provide better resistance to eviction-set attacks than those with 1 or 16 partitions. Thus, the main conclusion of this chapter is that there is no “best” cache design. Instead, we believe that caches need to be designed for the anticipated use cases.

The flexibility of CacheFX also enables the comparison of attack strategies against existing caches. Specifically, we demonstrate that the Prime+Prune+Probe method for eviction set construction yields more accurate results than the Sin-

gle Holdout and Group Elimination Methods. Additionally, we found that for caches with low randomization, constructing an eviction set is an effective strategy for cryptographic attacks. However, in highly randomized designs, the cache-occupancy attack emerges as a more efficient approach. As a result, we recommend that secure cache designers take this attack into consideration.

Chapter 4

Speculative Execution Against Low-Resolution Timers

In this chapter, we shift our focus to the second countermeasure against cache side-channel attacks, which primarily aims to reduce the accuracy of available timers. Cache attacks depend on distinguishing cache hits from misses by calculating the time required to access the cache, with the time difference between these events being incredibly small (less than a hundred nanoseconds). As a result, high-resolution timers are beneficial for carrying out cache attacks and are often crucial for successful attacks. To counter this, browsers decreased the resolution of timers they offer and eliminate certain methods for generating artificial timers as a defense against cache side-channel attacks [Chromium, n.d.; Hazen, 2018; Schwarz et al., 2017; Wagner, 2018]. Besides impeding attackers’ ability to differentiate between cache hits and misses, lowering timer resolution also restricts the attacker’s capability to construct eviction sets—groups of congruent addresses in the cache that correspond to the same cache set. Identifying eviction sets is a crucial step in enabling attacks such as Prime+Probe and Evict+Time [Osvik et al., 2006]. Since techniques for discovering eviction sets also require the ability to discern cache hits from misses, low-resolution timers obstruct this critical step in executing cache attacks.

While some cache attacks have been designed to use only low-resolution timers [Cronin et al., 2021; Hadad and Afek, 2018; McIlroy et al., 2019; Röttger and Janc, 2021; Schwarzl et al., 2021; Shusterman et al., 2019; Shusterman et al., 2021a], to the best of our knowledge, none have been specifically aimed at finding eviction sets, and the issue of discovering eviction sets with low-resolution timers remains open. Furthermore, in all reported attacks, the timer resolution constrains the sampling rate. In particular, no high-resolution tracing attacks, for example, against modular exponentiation [Bernstein et al., 2017; Liu et al., 2015; Yarom and Falkner, 2014; Zhang et al., 2012], have been demonstrated using only such timers.

Additionally, there are numerous studies that investigate the utilization of cache states to implement transient-execution attacks [Agarwal et al., 2022;

Behnia et al., 2021; Canella et al., 2019b; Chen et al., 2019; Hadad and Afek, 2018; Kirzner and Morrison, 2021; Koschel et al., 2020; Lipp et al., 2018; Maisuradze and Rossow, 2018; Ragab et al., 2021a; Ragab et al., 2021b; Röttger and Janc, 2021; Van Schaik et al., 2019; Van Schaik et al., 2020; Van Schaik et al., 2021; Schwarz et al., 2019; Stecklina and Prescher, 2018; Van Bulck et al., 2018; Van Bulck et al., 2020], yet the inverse question has not been explored thus far. This realization, coupled with the interest in uncovering potential vulnerabilities of low-resolution timers as a cache side-channel attack defense mechanism, lead us to pose the following question in this chapter:

Can transient execution improve cache attacks in a low-resolution timer environment?

In this chapter, we affirmatively answer the question. We show that the effects of transient execution can significantly improve cache attacks on systems with low-resolution timers. The central concept is that cache states can influence the duration of speculative execution and how it modifies future states. We devise logical gates that enable us to manipulate cache states, amplify them, store information in them, and even compute on them. We then use these gates to demonstrate how we can conduct high-resolution cache attacks using only a low-resolution timer.

4.1 Gates

This section describes the main ideas behind our implementations of logical gates based on cache states and speculative execution. We explain the computational model, assumptions, and design rationale.

4.1.1 Computational Model

We use logical gates to implement computation on the microarchitectural cache states of memory addresses. We define the logical value of “uncached” addresses as ‘0’, and of “cached” addresses (either in L1, L2, or LLC) as ‘1’. Changing the logical state from ‘0’ to ‘1’ is straightforward: we simply need to read the data stored in an address, and it will be fetched into the cache. However, we do not assume access to low-level instructions (e.g., `clflush`) that can set the cache state to ‘0’ directly. Instead, we assume that the initial value for hitherto unused addresses is ‘0’. Our gates compute a logical function of their inputs and store the result in the output. Because testing an input is done by accessing it, our gate evaluation is destructive. That is, computing over the state of an address brings it to the cache, setting its logical value to ‘1’. The main implication of destructive gates is that we cannot reuse the same address as an input to multiple gates.

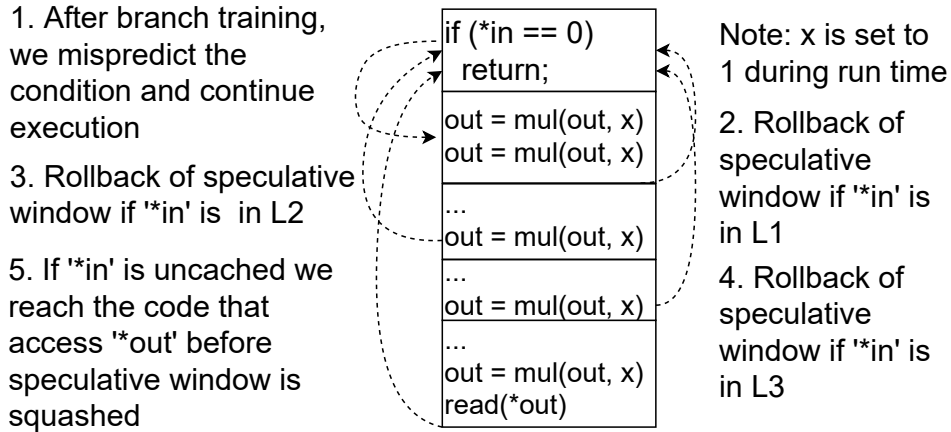


Figure 4.1: *NOT* Gate. `read(*out)` is transiently executed only when `*in` is not cached.

4.1.2 *NOT* Gate

The main insight that motivates our design is that the length of the speculation window of a mispredicted branch is not fixed. Instead, it depends on the time it takes the processor to evaluate the condition of the mispredicted branch. The length of the speculative window determines the number of instructions that are speculatively executed after the branch. Hence, the number of instructions that are executed speculatively varies with the time it takes to resolve the branch condition.

Figure 4.1 shows how we use these variable-length speculative windows to create a NOT gate. We use a branch conditioned on `*in`. Architecturally, we have `*in=0`. Hence the code should return without executing any of the subsequent instructions. However, the value of `*in` is not immediately available to the CPU as the pointer needs to be dereferenced

Before using the gate, we train the branch predictor to assume `*in` is not zero. Consequently, while waiting for the branch condition to be evaluated, the processor continues to speculatively execute instructions from the predicted branch. The length of this speculative window depends on the time it takes the processor to evaluate the branch condition, which is dominated by the time it takes to retrieve the value of `*in`.

In turn, the time to retrieve the value of `*in` depends on where that value is stored. If the variable `*in` is cached in the L1 data cache, accessing it will be quick (≈ 4 cycles). The time will be longer if `*in` is retrieved from the L2 cache, and even longer if it needs to be retrieved from the LLC. Finally, if the value of `*in` is not cached, it will need to be retrieved from memory, which would take a few hundreds of cycles.

The mispredicted branch contains a sequence of dummy operations (we use `imul` instructions to repeatedly multiply the pointer `out` by 1) followed by a memory access to the output variable `*out`. Tying the value of `out` to the

multiplication ensures that the processor does not execute the memory access before all of the dummy instructions complete execution.

The number of dummy operations is carefully chosen such that if `*in` is cached in any of the cache levels, the speculative window will terminate before the memory access to `*out` is issued and `*out` will not be accessed. However, if the value of `*in` is not cached, the speculative window is long enough and the access to `*out` executes speculatively. Eventually, the processor retrieves the value of `*in` and squashes all instructions on the mispredicted branch. However, because memory accesses execute asynchronously, the memory access to `*out` will complete even if the instruction is squashed.

If `*in` is cached ('1') then `*out` is not accessed and maintains its original logical value. Conversely, if `*in` is not cached, the memory access to `*out` executes transiently, bringing the value of `*out` to the cache, which sets the logical value to '1'. As we assume the initial state of `*out` is uncached ('0'), the end result is that, after executing the gate, the logical value of `*out` is the inverse of the original logical value of `*in`. Hence, the gate computes the logical NOT function.

4.1.3 More Complex Gates

The technique used for implementing the *NOT* in Section 4.1.2 can be extended to implement more complex logical gates. We now demonstrate how we can combine inputs to create a *NAND* gate and add branches to create a *NOR*.

NAND Gate

To create a *NAND* gate, we take our *NOT* gate and replace the if statement `if (*in == 0)` with:

```
if (*in1 + *in2 == 0)
```

Similar to the *NOT* gate, after we train the branch predictor, a speculative execution window is opened. It continues to run until the values of both `*in1` and `*in2` are made available to the CPU. The CPU processes the two read requests in parallel. Thus, the length of the speculative window is approximately the longer of the two access times.

If either of the input addresses is uncached, the processor needs to wait until the contents is retrieved from memory, resulting in a long speculation window. Consequently, in such a case, speculative execution reaches the `read(*out)` instruction, setting the value of the output to '1'. On the other hand, if both addresses are cached, the length of the speculative window is shorter. Consequently, the misspeculation is squashed before it reaches the read code, and the state of the output address remains '0'. To summarize, if the state of both input addresses is '1', the output value remains '0'. Otherwise, the output is set to '1'. Hence, the code computes the logical *NAND* function.

NOR Gate

For a *NOR* gate, we replace the single if statement of the *NOT*, i.e., `if (*in == 0)`, with two consecutive if statements:

```
if (*in1 == 0) {return;}
if (*in2 == 0) {return;}
```

If either input addresses is cached, the speculative window of the corresponding if statement is short and speculation is squashed before the processor executes the read. This leaves the state of the output address at '0'. However, if both input addresses are not in the cache, the processor needs to retrieve both values from memory before it can squash the speculation of any of the branches. This allows a long speculation window, which would execute the read command, setting the state of the output to '1'. To summarize, only if the state of both inputs values is '0', we get an output value of '1'. Otherwise, the output value is '0'. This is exactly the logical value of a *NOR* function.

4.1.4 Multiple Inputs and Outputs

By repeating the patterns in Section 4.1.3, we can increase the number of inputs in the *NAND* and *NOR*. For example, for a four input *NAND* we use the following if statement:

```
if (*in1 + *in2 + *in3 + *in4 == 0)
```

Similarly, we can replicate the output of the gates into multiple output variables by adding read statements to the misspeculated branch. We use the notation $GATE_{out}^{in}$ for gate *GATE* with *in* inputs and *out* outputs.

The processor uses a structure called *line fill buffer* (LFB) to track memory loads that miss on the L1 cache. Consequently, the number of LFBs limits the number of reads that can be processed concurrently, and the fan-in and fan-out of our gates. Specifically, when the total of the fan-in and fan-out exceeds the number of LFB entries (12 in the processors we use), the gates may fail.

4.1.5 Error Correction Gate

We further expand the usefulness of our gates, particularly in performing error correction of cache states. Given that cache states are volatile, it is important to have a technique for error correction to ensure robust calculations. With our gates having the capability to operate on multiple inputs and outputs, we can store the outputs of each gate to multiple locations, creating copies of their results. This allows the output states of each gate to tolerate errors of up to half of their copies. However, before further computation can be performed on the copied states, an additional gate is needed to compute the majority value of the copies, i.e., the most commonly occurring cache states of the copies. We design a gate that takes these copies as its inputs, computes their majority state, and outputs the majority state of these copies ('0' or '1').

To create a majority (out of five) gate, we take a five input *NAND*, and replace the if statement with the following:

```

if (*in1 + *in2 + *in3 == 0) return;
if (*in1 + *in2 + *in4 == 0) return;
if (*in1 + *in2 + *in5 == 0) return;
if (*in1 + *in3 + *in4 == 0) return;
if (*in1 + *in3 + *in5 == 0) return;
if (*in1 + *in4 + *in5 == 0) return;
if (*in2 + *in3 + *in4 == 0) return;
if (*in2 + *in3 + *in5 == 0) return;
if (*in2 + *in4 + *in5 == 0) return;
if (*in3 + *in4 + *in5 == 0) return;

```

If either three of the five addresses are cached, the speculative window of the if statements above is short, and the speculation is terminated before the processor executes the read statement. This means that the state of the output address is '0' if three or more of the input addresses are cached (the majority is '1'). On the other hand, if less than three of the input addresses are cached, the processor needs to retrieve at least three of the addresses from memory before it can squash the speculation of any of the if statements above. This causes a long speculation window, which would allow for the read statement to be executed by the processor. This means that if less than three addresses are cached (the majority is '0'), the state of the output address is '1'.

The behavior described above behaves the opposite of the majority functionality we are aiming to implement. In fact, the gate outputs the “minority” states of the five addresses. Because we are working with binary states, to get the majority, we can simply negate the output using a *NOT* gate. Therefore, we define the sequence of the gate described followed by a *NOT* gate as a *MAJORITY* gate.

4.1.6 Gates With a Fixed Branch Delay

All of the gates we use operate by creating a race between the time a branch is resolved and the time the outputs are accessed speculatively. In all of the gates we have seen so far, the timing of resolving the branch varies depending on the logical state of the inputs, whereas the timing of the memory access is fixed by the computation of the *imul* instructions. In this section we create additional gate types by swapping the fixed and the variable paths of execution.

Figure 4.2 shows an example of a *BUFFER* gate, which copies the logical state of the input to the output. The gate operates by first calculating a sequence of *sqrt* operations, which are set to return the (architectural) value 0. It then branches on the result, returning if the value is indeed 0.

During the computation of the *sqrt* operations, the processor cannot predict the final result. Consequently, the processor predicts the outcome of the branch, which we exploit by setting the branch predictor to mispredict

1. After branch training, we mispredict the condition and continue execution
2. Rollback of speculative window if '*in' is uncached
3. Only if '*in' is cached we reach the code that access '*out' and only then rollback the speculative window

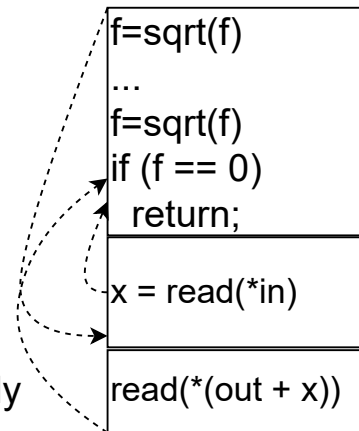


Figure 4.2: A Buffer Gate with a fixed branch delay.

the branch outcome. During the ensuing speculation window, the processor proceeds to execute the remaining code of the function.

The misspeculated code first reads the input `*in`. It then adds this value to the pointer `out`, and reads from the resulting address. Before using the gate, we ensure that the memory that `in` points to contains the value 0. Thus, the second `read` operation accesses the location pointed by `out`. However, the processor has to read `*in` before it can read `*(out+x)`. Consequently, the timing of the read from `out` depends on whether `*in` is in the cache or not, i.e., whether its state is '0' or '1'. If `*in` is in the cache (state '1'), the read from `out` will be executed before the speculative window is squashed and its state will be set to '1'. However, if `*in` is not the cache (state '0'), the speculative window will be squashed before its value will be made available to the CPU. In that case the second read instruction will not be executed speculatively, and the state of `out` will remain '0'. This is exactly the truth table for a *BUFFER* gate.

As before, by reading more output addresses (e.g., `read(*(out2 + x))`, `read(*(out3 + x))`, etc.) we can increase the fan-out of the gate and extend it to *BUFFER*_Y¹.

Similar to the extension of *NOT*_Y¹ to *NAND*_Y^X gate, we can extend *BUFFER*_Y¹ to *AND*_Y^X gate. This is done by making the reading of `*out` dependent on the sum of multiple input addresses instead of a single input address (e.g., `x = read(*in1) + read(*in2)`, `x = read(*in1) + read(*in2) + read(*in3)`, etc.). Only if all input addresses are cached, their sum is available to the CPU, and the `read` instructions are executed before the speculative window is squashed.

4.1.7 Gates Without Branch Training

In order to open the speculative window, we need to train the branch predictor to mispredict the initial branch we are speculating on. We use the training method introduced in [Google, 2021]. It starts with two “dry runs” of the gate that train the branch predictor, followed by the actual “wet run”. Moreover, each run of the gate (either “dry” or “wet”) starts with an empty `for` loop that creates a consistent branch history. As a consequence, the training of the gates is relatively long, increasing the overall run time of our gates.

To reduce training time and allow faster gates, we use a novel approach to cause the branch predictor to reliably mispredict a branch without the need to “retrain” it. The main idea is to replace the single `if` condition with a large `switch` statement. The correct case is determined by combining the value of the input address with a counter that is incremented after each evaluation of the gate. As the value of the input address is not available to the CPU, the branch predictor mispredicts by jumping to the previously chosen case based on the counter’s previous value. Listing A.1 and Listing A.2 in Appendix A.2 show the code for both types of gates. According to our experimental evaluation of the gates (see Appendix A.3), switching from the branch training-based approach (“bt”) to our non-branch training-based approach (“nbt”) can significantly reduce the gates’ run time at the cost of a slight reduction in the gates’ accuracy. Thus, the “nbt” gates present a tradeoff between performance and accuracy.

Compilers offer multiple implementations for `switch` statements. For our technique to work, we need the compiler to use an indirect branch with a jump table. The choice of implementation method depends on the number of cases and their values. In our experiments we find that in both the native and the WebAssembly compiler we use, the compiler chooses a jump-table-based implementation when there are at least eight cases. Once the implementation is chosen, the number of cases has little impact on the gates’ accuracy or performance. Hence we use the required minimum of eight cases for implementing our gates.

4.1.8 Gates Evaluation

In this chapter, we focused our experiments on Intel’s range of processors where we optimized the implementation of our gates to work on those processors. See Appendix A.3 for complete experimental results with different gate types including a discussion on the various values for fan-in and fan-out.

Note that many of the gates require tailoring parameters where we perform the tuning on a case-by-case basis (e.g., number of dummy instructions). We can fully optimize many of the gates achieving an accuracy of approximately 99.9% or above. All branch training-based gates achieve an accuracy of over $\approx 99.5\%$, while the slightly less accurate non-branch training-based gates still achieve an accuracy of over $\approx 95.8\%$. We can also see that our non-branch

training-based gates are indeed significantly faster than our branch training-based variants. The non-branch training-based gates are approximately 300 cycles faster. Specifically, they are around twice as fast when the output is '1' and three times faster when the output is '0'.

We could also reproduce similar results on an AMD Ryzen 5 3500U (see Appendix A.4). We further implemented some of the gate types on ARM processors. Appendix A.4 also presents experimental results on a Macbook Air laptop with Apple M1 processor and a Galaxy S21 phone with a Samsung Exynos 2100 SoC. We leave the optimization of gates for these platforms to future work.

4.2 Circuits

To demonstrate the versatility of our gates, we now show how they can be used to build complex logical circuits. We investigate three circuits: the arithmetic logic unit (ALU) from Nisan and Schocken (2021), the SHA-1 hash function [National Institute of Standards and Technology, 2015], and finally Conway’s Game of Life [Gardner, 1970].

Experimental Setup. We carry out the experiments in this section on an Intel NUC 9 Extreme Kit equipped with an Intel Core i7-9750H CPU running Xubuntu 21.10. Due to the load-sensitive nature of the speculative gates and the cache states, we isolate the core used to run the experiment using the `isolcpus` kernel parameter and enable huge pages.

Our experiments explore the effects of the prefetcher (enable vs. disable) on the accuracy of our circuits [Viswanathan, 2014]. We have also tested the effects of CPU frequency scaling (fixed vs. variable). We find that gates need to be tuned for the processor frequency. However, once tuned, the accuracy of the gates is similar. Moreover, we find that when setting a variable frequency, the processor mostly executes at the highest frequency. Therefore, we only report the results of variable frequency to reflect a more realistic scenario.

4.2.1 ALU

We first demonstrate the viability of our speculative gates through a construction of a four-bit ALU adapted from Nisan and Schocken (2021).

The ALU takes two four-bit numbers, x and y , along with six control values, as input. Then it produces a four-bit value output. The control values are used to direct the ALU to which operations to perform on each operand. Specifically, by using different combinations of the control values, the ALU performs different operations such as increment, decrement, addition, subtraction, negation, binary AND and binary OR.

To initialize the state of x , y , and the control values, we either read the corresponding address to signify '1' or flush the address to signify '0'. We

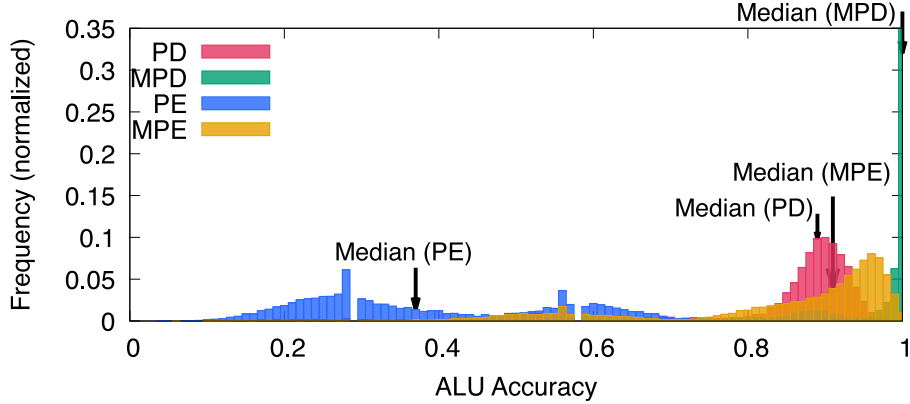


Figure 4.3: ALU accuracy. PD denotes accuracy with prefetcher disabled, MPD denotes accuracy with error correction (majority gate) with prefetcher disabled. PE signifies accuracy with prefetcher enabled, while MPE signifies accuracy with majority gate and prefetcher enabled. For visibility, the Y axis is trimmed at a frequency of 0.35.

then let the ALU perform its computation. Finally, we read the results by measuring the access time of the output addresses.

We utilize error correction techniques to improve the accuracy of our ALU. Specifically, we run the ALU five times with five copies of the same initial values. We then use the *majority-out-of-5* gate (described in Section 4.1.5) to compute the final value. Performing such redundant calculation and following a majority vote is a classic error correction approach as used in, for example, redundant coding to detect and correct errors from bit flips [Kim et al., 2009; Stroud and Barbour, 1989; Von Neumann, 2016]. Our ALU is built from 250 logic gates without majority (1 258 logic gates with majority). It consists of 336 intermediate states (1 688 with majority), of which only the four output bits, or 1.19% (0.24% with majority) are exposed architecturally.

We perform 10 000 sets of experiments to measure the performance of our ALU where we focus on the correctness of the four-bit output. Each set of experiments contains 100 runs of the ALU, where the inputs are selected at random. For each run, the accuracy is one if all the four bits are correct. If any of the four bits are incorrect, the accuracy is zero. The average of the accuracy of 100 runs represents the accuracy of that set of experiments. Figure 4.3 illustrates the result of our experiments in histogram. Specifically, it shows the accuracy of the ALU calculation with/without error correction and with/without prefetcher enabled. The histogram clearly highlights a significant increase in the accuracy when using the majority gates, namely, from a median of 89.0% to 100% (82.0% to 95.4% average) when disabling the prefetcher and from a median of 37.0% to 91.0% (43.7% to 84.1% average) when enabling the prefetcher. On average, performing an ALU instruction takes 106 microseconds without error correction and 529 microseconds with error correction.

4.2.2 SHA-1

Our second example of a circuit is an implementation of a cryptographic hash function SHA-1 [National Institute of Standards and Technology, 2015]. Note that in contrast to Evtvyushkin et al. (2021) our entire round of SHA-1 calculations are performed in microarchitectural states where we interact with them only for the initial state setting and the final output reading. Generally, SHA-1 consists of a loop with 80 iterations to produce a 160-bit message digest output. In our experiment, we perform one round of SHA-1; Listing 4.1 shows the pseudocode.

Listing 4.1: Pseudocode for the first round of SHA-1

```
void sha1_round(A, B, C, D, E, W) {
    temp = circular_shift(5, A) + ((B & C) |
        (~B) & D) + E + W + 0x5A827999;
    E = D; D = C;
    C = circular_shift(30, B);
    B = A; A = temp;
}
```

As the listing shows, one round of SHA-1 consists of two circular left shifts, four additions, two binary AND, one binary OR, and one binary NOT, each operating on 32-bit words. The calculation of a round of SHA-1 requires 32-bit adder, AND, OR, and NOR. In total, the circuit consists of 2 208 logic gate primitives and exposes 1.07% of its 2,976 total microarchitectural intermediate logical states to the architectural state.

We evaluate our implementation by running the SHA-1 and testing the rate of which the full 160-bit result is correctly computed. Specifically, we take measurements of 10 000 experiments, with each experiment performing 100 runs of SHA-1 with random inputs.

Figure 4.4 shows the distribution of the accuracy for a single round of SHA-1. With prefetcher disabled, we obtain an average and median of 94.95% and 99.00% respectively. When the prefetcher is enabled, we obtain an average and median of 66.55% and 58% respectively. Each round of SHA-1 takes only 969 microseconds to run.

We further evaluate the robustness of our circuit by instrumenting it to compute two blocks of SHA-1, each consisting of 80 rounds. For a complete implementation, we also use four round functions, as specified in the SHA-1 standard [National Institute of Standards and Technology, 2015]. To perform the calculation, we chain consecutive invocation of SHA-1 round circuits. That is, after performing one round, we sample the result and copy it to the architectural state of the processor. We then use the sampled data to set up the cache state for the following round. We further increase the accuracy by computing each round ten times, and using the per-bit majority to determine the output of each round.

Repeating the full computation 1 000 times, we observe a 95.1% probability that the output from our two-block SHA-1 is correct. Each run involves

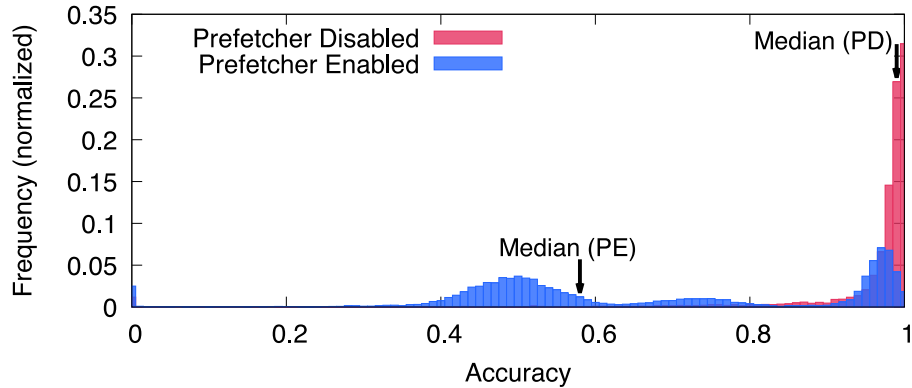


Figure 4.4: SHA-1 accuracy.

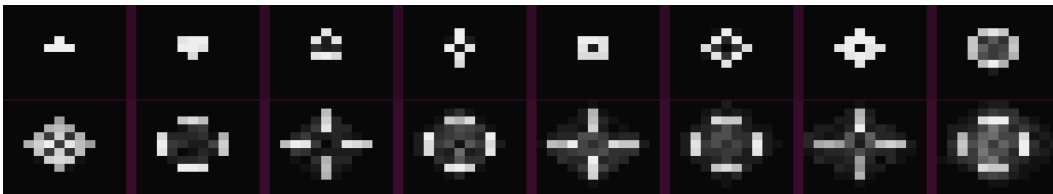


Figure 4.5: T-tetromino heatmap. (calculated from 300 repetitions, the brighter the cell the higher accuracy.)

3 737 600 logic gates and 5 068 800 intermediate values, 1.01% of which is exposed architecturally.

Evyushkin et al. (2021) implement SHA-1 using their “weird gates”. Their implementation relies on Intel Transactional Synchronization Extensions (TSX) [Intel, 2021a], a feature of Intel processor but has been mostly disabled due to security issues [Intel, 2021b]. Also, their implementation exposes a significant part (41.9%) of the logical state of SHA-1 to the architectural state of the program. In contrast, our implementation uses generic processor features, which are available across multiple architectures, and exposes only 1.01% of the logical states.

4.2.3 Game of Life

As a third example for complex logical circuits based on our gates, we implement Conway’s game of life [Gardner, 1970] for a universe up to size 12×12 . Recall that the game is a cellular automaton, consisting of a grid of cells. Each cell has a state which can be either ‘live’ or ‘dead’. Each generation, the state of a cell is updated based on the values of the cell and of its eight neighbors, using the following rules: 1) a live cell that has two or three live neighbors remains live; 2) a dead cell with exactly three live neighbors becomes live; 3) other live cells become dead, and other dead cells remain dead. In our implementation, we denote a live cell by ‘1’ and a dead cell by ‘0’.

According to the rules, calculating the next state of a cell requires evalu-

Generation	Prefetcher Disabled		Prefetcher Enabled	
	Average (percent)	Median (percent)	Average (percent)	Median (percent)
1	59.10	69.00	62.76	73.00
10	48.74	48.00	46.99	48.00
20	22.76	22.00	25.58	25.00
30	15.28	13.00	15.53	11.00
40	17.09	16.00	11.10	9.00
50	10.99	9.00	4.70	3.00

Table 4.1: Game of Life glider accuracy.

ating the state of that particular cell and its eight neighbors. Since evaluating a cell changes its value, it becomes crucial to not destroy the state of the cells that are still needed for future evaluations. To tackle this challenge, we microarchitecturally copy the value of a cell into two locations. The first is used to perform an actual calculation while the second is used to restore the original state of the cell.

We implement games over multiple generations with initial states such as T-tetromino [*Game of Life Wiki* n.d.] and glider. Figure 4.5 shows 16 generations of the game, starting from a T-tetromino pattern, in a 12×12 universe. The brightness of a cell shows the probability that our circuits calculates it as live. Table 4.1 summarizes the accuracy of a glider across 50 generations for an 12×12 universe. We achieve a high accuracy for the first generation; however, due to error propagations onto future generations, the accuracy drops as the game progresses.

Each generation requires 7 808 logic gates to perform its calculation with the total of 11 456 intermediate microarchitectural states. This means that running 10 generations computes 114 560 intermediate microarchitectural states. As we expose only the final output of the circuit to the microarchitectural states, the number of states exposed constantly remains 64 regardless of the number of generations. Specifically, for one generation we expose 0.56% ($64/11\,456$), whereas for 10 generations we expose 0.06% ($64/114\,560$). One generation of the game takes 3.19 millisecond to run.

4.3 Probe Amplification

This section demonstrates how we use our gates for side-channel probe amplification. A fundamental ability that most microarchitectural side-channel and speculative execution based attacks require is to determine or “probe” if a specific cache line is in the cache or not. Because the difference in access time is only in the order of approximately 100 clock cycles, this requires access to

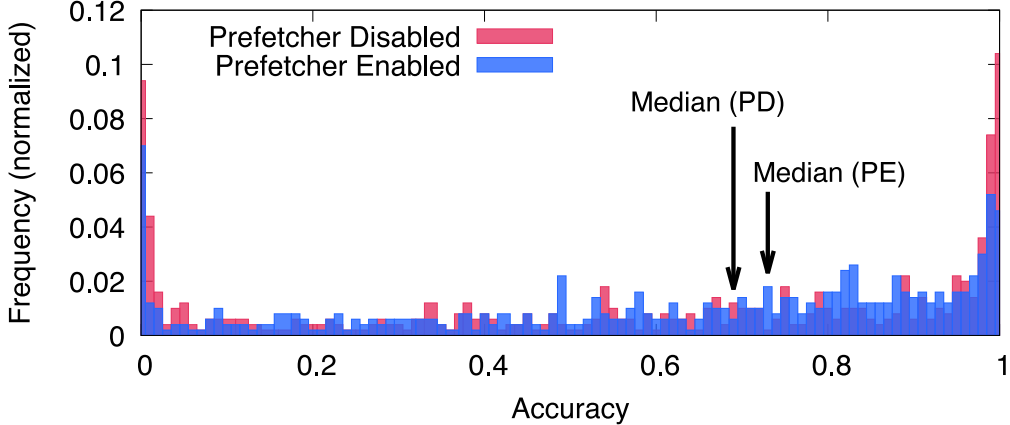


Figure 4.6: One generation Game of Life accuracy.

high-resolution timers. In many settings, e.g., WebAssembly and JavaScript code in modern browsers, access to such timers is actively blocked to prevent this type of attacks. We propose an amplification approach, using our gates, that can amplify the minute timing difference between cache hit and cache miss, allowing the use of timers with arbitrarily low resolution to distinguish the two.

The scheme consists of three amplification steps. In the first step, we use a single gate to achieve a small amplification. In the second step, we create a tree-like structure achieving a theoretical timing difference of up to 4 milliseconds. Finally, in the third step, we combine multiple trees to achieve an arbitrarily long timing difference.

4.3.1 Single-Gate Amplification

The first step in our proposed scheme is using NOT_Y^1 , a NOT gate with a fan-out of Y (e.g., $Y = 4$), to gain a small amplification.¹ We denote the access time to an address cached in the LLC by t_{in} , the access time to an uncached address in the main memory by t_{RAM} , and their difference by $\Delta_{cache} = t_{RAM} - t_{in}$. Assume we want to test if $addr_{in}$ is cached or not. Instead of directly measuring the access time to $addr_{in}$, we use $addr_{in}$ as an input for the NOT_Y^1 gate. Then we measure the total time it takes to sequentially access all Y output addresses. If $addr_{in}$ was uncached, the Y output addresses will be cached and the total access time will be $Y \cdot t_{in}$. Otherwise, the Y output addresses will be uncached and the total access time will be $Y \cdot t_{RAM}$. This amplifies the timing difference by a factor of Y , from Δ_{cache} to $\Delta_{gate} = Y \cdot \Delta_{cache}$.

¹We use the NOT_Y^1 and not the $BUFFER_Y^1$ because it is easier to implement for multiple environments (e.g., native and WebAssembly).

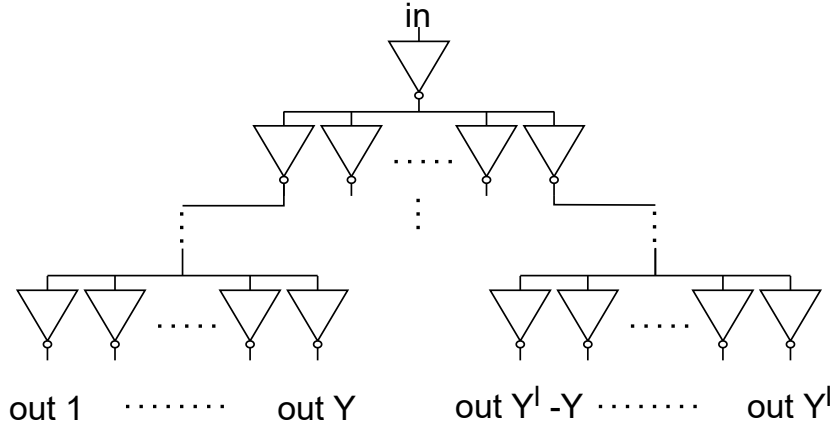


Figure 4.7: Amplification tree based on NOT_Y^1 Gates.

4.3.2 Probe Time Amplification Tree

As mentioned in Section 4.1.4, the fan-out is limited by the size of the LFB, allowing only a small constant amplification. To support a larger amplification factor, our next step is to use a tree structure with tree depth l as shown in Figure 4.7. Again, we use $addr_{in}$ as an input to a NOT_Y^1 gate. However, instead of simply accessing the Y output addresses, we now use each of them as inputs to NOT_Y^1 gates. This gives us a total of Y^2 output addresses. We then continue in the same manner for the full l layers, resulting in a total of Y^l output addresses. If the number of layers l is even, we expect all output addresses to be cached if $addr_{in}$ was cached, and uncached otherwise. If l is odd, the cache state of the output addresses is negated. In either case, measuring the total time of sequentially accessing all of the Y^l output addresses allows us to amplify the timing difference to $Y^l \cdot \Delta_{cache}$. Note that the tree is generated in a breadth-first order, i.e., generating each layer before continuing to the next one.

We note that the amplification factor of this tree structure is limited by the cache size as the number of possible output addresses is limited by the number of cache lines. For example, if we assume an LLC of size 8 MB (2^{17} cache lines) and $\Delta_{cache} = 100$ clock cycles (0.033 microseconds on a 3 GHz CPU) then $\Delta_{tree} \leq Y^l \cdot \Delta_{cache} \approx 4$ milliseconds. Moreover, due to practical considerations (e.g., memory prefetcher, risk of self eviction between layers of the tree, etc.), the maximal number of output addresses we can reliably use is much lower.

4.3.3 Amplification Hyper-tree

To overcome the limitations of tree amplification method, we use a hyper-tree structure. We start our hyper-tree amplification by using an l -layer amplification tree to copy (or negate if l is odd) the cache state of $addr_{in}$ to a bank of Y^l

output addresses. We then continue to iterate over all addresses in the bank. In each iteration, we use the address from the bank as an input to a new tree, then sequentially access all the Y^l output addresses of the resulting tree. We measure the time it takes to generate the Y^l sub-trees and sequentially access all the leaves of each sub-tree. Such a two-level hypertree produces a total of $Y^{2 \cdot l}$ output addresses. However, it only has at most $2 \cdot Y^l$ “live” memory addresses at each time, a significant improvement over the single tree case. If needed, we can extend this hyper-tree structure to d levels of sub-trees and a total of $Y^{d \cdot l}$ output addresses at a space cost of $d \cdot Y^l$.

Note that in contrast to the simple tree amplification, we cannot store all of the $Y^{d \cdot l}$ output addresses in the cache at the same time because the cache might not be large enough. This means that we cannot merely measure the access time to the output addresses but need to measure the entire process of the amplification. This also means that regardless of the cache state of $addr_{in}$, we measure the time it takes to access all of the addresses in all the nodes and leaves of the hyper-tree. However, if the addresses in the layer before the last are uncached, the output addresses are accessed from inside the speculative window in parallel. In such a case, it takes t_{RAM} time to access all Y addresses inside the speculative window and then $Y \cdot t_{in}$ time to access them sequentially. However, if the addresses in the layer before the last are cached, they are only accessed sequentially at the end, with a total time of $Y \cdot t_{RAM}$. Hence, theoretically we obtain an overall amplified timing difference of approximately

$$\Delta_{hypertree} \approx Y^{2 \cdot l - 1} \cdot ((Y - 1) \cdot t_{RAM} - Y \cdot t_{in}) \approx Y^{2 \cdot l} \cdot \Delta_{cache}.$$

Note that the actual difference is lower due to the access time in the intermediate layers. Appendix A.1 provides details on the hyper-tree amplification implementation.

4.3.4 Experimental Verification

To demonstrate our amplification hyper-tree scheme, we implement and test it both in native code and in WebAssembly code. We run the experiments on a Dynabook TECRA A50-EC, with an Intel Core i5-8250U CPU running Ubuntu 20.04.3 LTS. In our experiments, we set the frequency governance to performance. The WebAssembly code was tested under Chromium 99.0.4843.0 (Developer Build). Figure 4.8 shows the results of running an amplification hyper-tree in native. The hyper-tree is composed of three tree layers. The topmost layer is an amplification tree from 1 to 16. The bottom two layers are amplification trees from 1 to 512. In total the tree amplifies the access time of a single address to $16 \cdot 512 \cdot 512 = 4\,194\,304$ memory accesses. For each run, we measure the total time it takes to generate the whole tree and access all the leaves. We run the amplification for a total of 1 000 times—500 runs with the root address we amplify cached, and 500 runs with it uncached. The

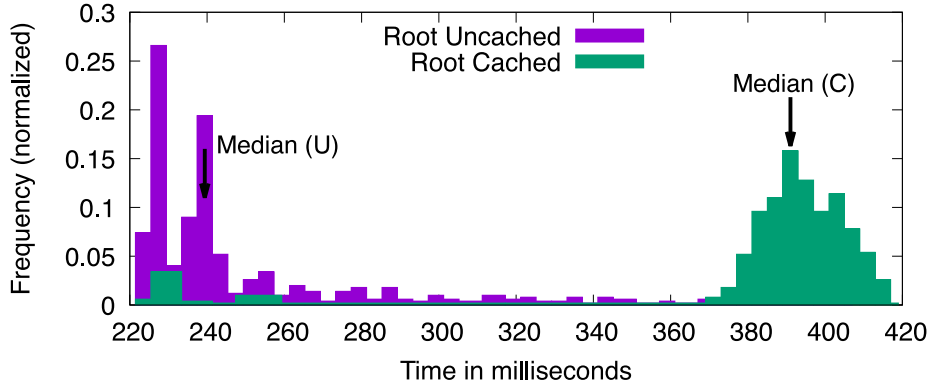


Figure 4.8: Amplification Hyper-Tree in native.

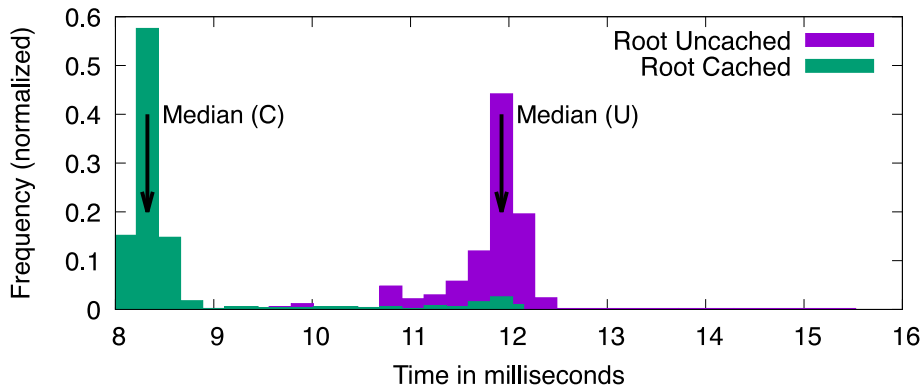


Figure 4.9: Amplification Hyper-Tree in WebAssembly.

difference between the median values of the two distributions is more than 100 milliseconds. Our statistical t -test analysis of the two distributions (root cached vs. uncached) yield the p-value (two-tailed) of 0.036, which strongly confirms that we can, indeed, distinguish between the two scenarios.

Figure 4.9 shows the results of running an amplification hyper-tree implemented in WebAssembly, after discarding measurements that takes longer than 20 milliseconds as they are too noisy to use ($\approx 3\%$ of the measurements). The hyper-tree is composed of two tree layers. The topmost layer is an amplification tree from 1 to 192. The second layer is an amplification tree from 1 to 512. In total the tree amplifies the access time of a single address to $192 \cdot 512 = 98\,304$ memory accesses. For each run, we measure the total time it takes to generate the whole tree and access all the leaves. Similar to the case of native code, we run 1000 experiments—500 with a cached root and 500 with an uncached root. The difference between the median values of the two distributions is more than 2 milliseconds, which is the current timer resolution provided in the Firefox browser. Similarly, we perform a statistical t -test analysis of the two distributions (root cached vs. uncached) yield the p-value (two-tailed) of 0.010, which, again, strongly confirms that we can, indeed, distinguish between the two scenarios.

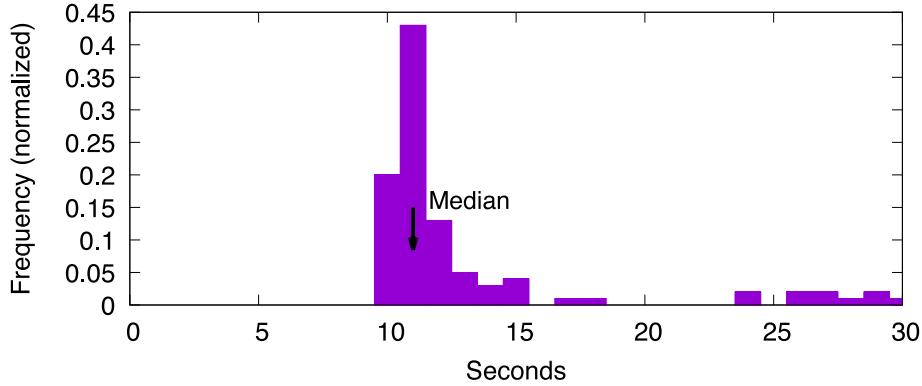


Figure 4.10: Time to find an eviction set in Chrome using 0.1 millisecond low-resolution timer

4.3.5 Eviction Set Creation

We now show how we use the probe time hyper-tree amplification scheme from Section 4.3 to create eviction sets using only low-resolution timers available to JavaScript and WebAssembly code running inside a browser. We implement the eviction set creation algorithm from Vila et al. (2019), while using our probe time amplification to support the low-resolution timers provided by Google Chrome.

We run our experiment on an unmodified Chrome 102.0.5005.61 (Official Build; 64-bit) on the same setup as before. Our WebAssembly code chooses a memory address and tries to find its congruent eviction set that includes 12 addresses. We start with a set of 3000 addresses that, with a very high probability, contain all 128 possible eviction sets with the same page offset as the target address. In 83 out of 100 runs, we are able to find the correct eviction set. Another 8 runs are “close”, meaning that only one address returned is not in the real eviction set. We are not able to find an eviction set for the remaining 9 runs. Failed runs are due to excessive noise. The algorithm detects such failures, and re-running typically finds an eviction set.

Figure 4.10 shows the measured running time of the algorithm. To summarize, our algorithm is able to find the correct eviction set in 78% of the runs. We only use the 0.1 millisecond resolution timer provide by Chrome, and the median run time of algorithm is approximately 11 seconds.

4.4 Prime+Store: Fast Attacks with Slow Clocks

Slow clocks introduce two problems for microarchitectural side-channel attacks. The first issue is that it is hard to distinguish microarchitectural events with a slow clock; we address this problem in Section 4.3. The second issue is that the clock limits the rate at which we can measure events; each measurement takes at least one clock tick. This section presents the Prime+Store

attack, which overcomes this limitation. We first describe the attack and then demonstrate how we use it against a vulnerable version of ElGamal.

4.4.1 Prime+Store

Our Prime+Store attack is a variant of Prime+Probe. Recall that a Prime+Probe attack consists of two main actions. In the prime step, the attacker accesses all of the members of an eviction set, bringing them to the cache. In the probe step, the attacker accesses the members of the eviction set again to measure the access time and detect if any of the members of the eviction set has been evicted from the cache. Specifically, the probe step is a function that takes the cache state of the eviction set members and returns `false` if all of them are in the cache and `true` if some of them are not in the cache. We note that under our computational model, if eviction set members are in the cache, they represent the logical value '1'. Hence, the probe function effectively calculates the *NAND* of the logical values of the addresses in the eviction set.

Based on this observation, we design our Prime+Store attack using a $NAND_1^x$ gate, where x is the associativity of the cache. For the attack, we use an eviction set as the input to the $NAND_1^x$ gate. This stores the probe result as the cache state of the output of the gate. To perform multiple probes of the same cache set, we repeatedly invoke the $NAND_1^x$ gate with the eviction set as input, but each invocation we set the output to a different memory address. After we finish sampling, we can then test each of the memory addresses to determine the outputs of the gates. Thus, using this technique, we decouple the cache measurements from the sampling, allowing us to perform repeated samples at a high rate.

Recall that the total fan-in and fan-out of our gates is limited by the size of the LFB, but the fan in of the *NAND* gate used for the probe operation is the associativity of the LLC. Hence, if the associativity of the LLC is larger than the size of the LFB, the *NAND* gate may fail to work. We note, however, that in most cases, the victim only evicts one entry from the cache. Consequently, most of the eviction set remains cached. Accesses to cached memory free the LFB fast, allowing the attack to operate even though the fan-in is larger than the size of the LFB. The attack may still fail when several entries of the eviction set are evicted from the cache. We ignore this case, considering the failure as noise. If such noise is not acceptable, the attacker can use a more complex circuit to compute the *NAND* function using multiple gates.

4.4.2 Attacking ElGamal

To demonstrate the effectiveness of Prime+Store, we use it to recover the private key from a vulnerable implementation of the ElGamal public-key encryption scheme [El Gamal, 1984]. Specifically, we target the modular exponentiation operation, which raises a base b to the power e modulo some modulus

m , i.e., calculating $b^e \bmod m$. During ElGamal decryption, the private key is used as the exponent e . Hence, our attack aims to recover the exponent.

The attack itself consists of three steps. We first collect traces of memory activity that correspond to segments of the modular exponentiation operation. We then process these traces to recover the operations performed during the observed segments. Finally, we “stitch” the segments to recover the private key. In this section we describe the attack setup and the victim we target. Following subsections describe the steps of the attack.

Attack Setup. We run the experiment on Dynabook TECRA A50-EC, with an Intel Core i5-8250U CPU running Ubuntu 20.04.3 LTS with LFB size 12. The frequency governance is set to performance. We run two processes, a victim and a spy. The victim uses GnuPG 1.4.13 to repeatedly perform ElGamal decryption with a 4096-bit modulus. With this parameter, the GnuPG private key is of length 457 bits. The spy performs the attack, collecting traces as described below in Section 4.4.3.

Victim Implementation. To calculate the modular exponentiation, GnuPG 1.4.13 uses the square-and-multiply algorithm. The algorithm consists of three main operations, square, multiply, and modular reduction. It scans the bits of the exponents from the most to the least significant. For a bit value zero, it performs square followed by a modular reduction. For a bit value one, it performs a sequence of square, modular reduction, multiply, and modular reduction. The algorithm is known to be vulnerable to side-channel attacks and has been attacked multiple times [Liu et al., 2015; Yarom and Falkner, 2014; Zhang et al., 2012]. Specifically, by recovering the sequence of square and multiply operations that the algorithm performs, the attacker can recover the exponent.

4.4.3 Trace Acquisition

In our attack, we find an eviction set for the code of the square operation and use Prime+Store to repeatedly sample cache usage in the cache set. Once we have collected several samples, we use our amplification technique from Section 4.3 to amplify each sample to allow recovery with our clock.

Because measurement takes a long time, data we collect may decay, e.g., due to spurious cache evictions. To overcome the decay, we use two techniques. First, we observe that decays tend to change values only in one direction (‘1’ to ‘0’). We therefore oversample and then coalesce three consecutive samples using a *NOR* gate before measuring the result. Hence, unless all three samples decay, we do not miss a ‘1’ sample.

However, when measurement time is long, oversampling is not sufficient to avoid decay. Consequently, as a second measure, we limit the number of samples we collect in each trace, so that the trace only correspond to a small segment of the exponentiation operation.

For the attack on ElGamal, we collect a total of 100 000 traces, each con-

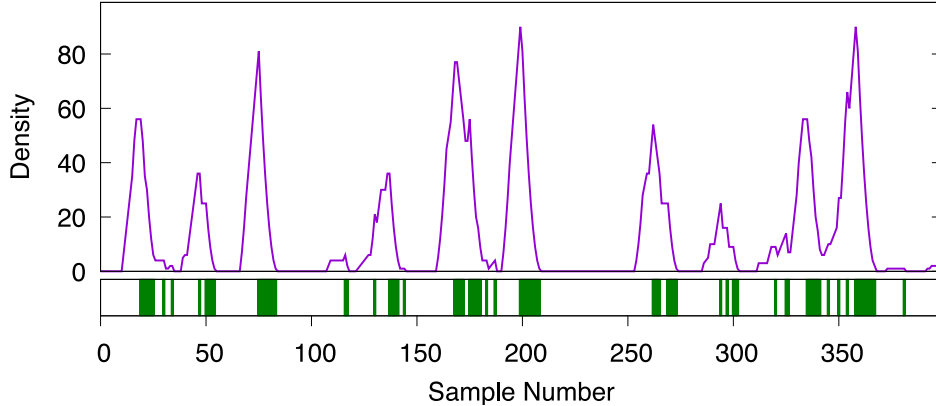


Figure 4.11: A segment of samples of the square operation in modular exponentiation. The bottom shows samples in which we detect eviction. The top shows the sample density. Peaks with density above 15 correspond to a square operation.

sisting of 2793 samples, at a rate of 0.33 microsecond per sample. We then coalesce groups of three consecutive samples, and amplify for measurement with a 0.1 millisecond clock. Measurement of a trace takes approximately 0.34 seconds, resulting in a trace of 931 coalesced samples, which we further process. Overall, collecting 100 000 traces takes 9 hours and 40 minutes.

The limited number of samples means that we can only observe a single segment of the exponent at a time. As in past works [Van Schaik et al., 2019; Schwarz et al., 2019; Zhang et al., 2012], our aim is to collect a large number of segments and then stitch them together.

4.4.4 Trace Processing

After collecting the traces, we process them to recover the sequence of square (S) and multiply (M) operations in the segments of the modular exponentiation that they cover. The bottom part of Figure 4.11 shows an example of a trace. (Trimmed to 400 samples for clarity.) Shaded areas indicate that our Prime+Store attack detected activity in the set in the corresponding coalesced sample. We clearly observe blocks of cache activity that indicate a square operation. However, the samples are noisy, with both gaps during square operations and spurious activity between squares.

To detect the blocks, we measure the density of evictions in an area. For each sample, we count the number of evictions in the subsequent 9 and 15 samples, and multiply the counts to obtain a measure of the density. The top part of Figure 4.11 shows the density measure for the displayed trace. We then perform peak detection to identify the samples at which a square operation starts. Specifically, we define a peak as a sample that has a higher density than any of the preceding and subsequent eight samples. Peaks with density above a threshold of 12 indicates a start of a square operation.

After recovering the positions of the square operations we use the distance between consecutive square operations. Specifically, we find that when the exponent bit is 0, the distance between consecutive squares is about 30 samples, and when the bit is 1, the distance is around 60 samples. As such, we assume that peaks at a distance of 15–45 samples are consecutive square operations, whereas peaks at a distance of 46–75 samples are a square followed by a multiply. Thus, the sequence of operations in the segment covered in Figure 4.11 is SSSMSSSMSSSSM. We ignore peaks that are closer than 15 samples to the previous peak, and treat peak distances of over 75 samples as unknown operations.

4.4.5 Key Recovery

The next step is to “stitch” the segments and recover the key. For that we draw on an algorithm from DNA sequencing [Wilkinson et al., 2017], adapted to the binary case. Our stitching algorithm relies on the observation that long enough sequences of square and multiply operations are unlikely to appear more than once within the exponent. Thus, the algorithm iteratively extends a guess of the sequence of square and multiply operations used during exponentiation with the key. For the sake of exposition we first explain a naive algorithm that assumes no errors in the traces.

Naive Algorithm. Our naive algorithm starts from the longest segment, which it uses as the current guess, and iteratively extends it to recover the full key. If there are multiple longest segments, it just picks one arbitrarily. To extend the guess, the algorithm searches all captured segment, looking for the matching segment with the largest overlap with the current guess. That is, it looks for a segment that when aligned at some position, has matches on all the positions that overlap with the guess, and out of those it picks the one with the longest overlap. It then merges the segment into the guess, extending the guess in the case that the segment extends beyond the guess. The algorithm stops when running out of segments or when the complete key is recovered.

Handling Trace Errors. The main problem with the naive algorithm is that traces do contain errors. To handle errors we modify the algorithm slightly. In the modified algorithm, instead of just associating an operation with a position, we track the likelihood for both operations, guessing the more likely for each position. To generate the initial guess we search for a segment that repeats the largest number of times in the collected traces, and use it for initial guess.

Instead of keeping a single guess for each position, we track the support for a square and a multiply operation. If support for square is larger than support for multiply, we predict that the operation is square. Otherwise, we predict that the operation is multiply.

To extend the guess, the algorithm searches the collected segments for the segment that has the largest match with the prediction in the guess. It

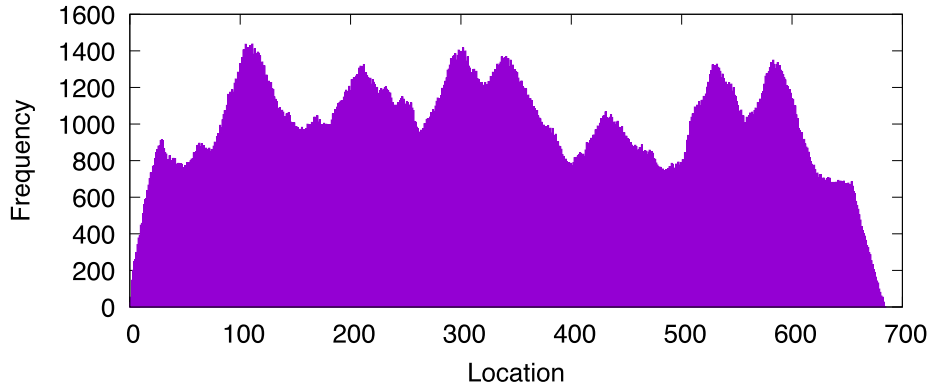


Figure 4.12: Distribution of stitched key in relation to ground truth location.

then calculates the weight of the segment, which is the number of positions in which the operation in the segment agrees with the operation predicted for the matching location in the guess. Finally, the algorithm updates the prediction by adding the weight of the segment to the guess support in each position.

4.4.6 Evaluation

As discussed in Section 4.4.3, we collect 100 000 traces over a period of almost 10 hours. Among those, 613 attempts fail, leaving us with 99 387 traces.

We first filter collected segments to remove apparent trace errors. Specifically, we ignore segments that contain 11 or more consecutive square operations. We find that such sequences often appear due to capture errors, and their inclusion confuses our stitching algorithm. We note that with a 457-bit ElGamal private key, we expect about one in five keys to include such a sequence. Omitting these sequences means that for such keys we will need an extra step of adding the missing segments. Overall, there are 12 568 such segments.

To test the coverage of the trace collection, we compare each collected segment with the ground truth. Figure 4.12 shows the distribution of the positions at which the segments best match the ground truth. Focusing on segments with length of 20 or more operations, we find that 3 572 segments out of the remaining 86 819, completely match the private key. Furthermore, 9 467 long segments match the ground truth with one operation error, and 14 007 match with two errors. Running the stitching algorithm, we find that after merging 22 863 we recover the full key. The process takes 14 minutes and 31 seconds.

4.5 Related Work

Early Attacks from the Browser. Genkin et al. (2018) demonstrate an LLC attack in WebAssembly, recovering the keys of multiple cryptographic

schemes. Oren et al. (2015) implement a cache attack in JavaScript, performing website fingerprinting. Gras et al. (2017) use a cache attack to break address space layout randomization. All use high-resolution timers that were available to JavaScript and WebAssembly at the time.

Alternative Timers. Several works investigate alternative timers for use in web browsers [Kohlbrenner and Shacham, 2016; Rokicki et al., 2021; Schwarz et al., 2017]. Following Wray (1991), Schwarz et al. (2017) propose a timer based on a shared counter, which has been used for Spook.js [Agarwal et al., 2022]. Mainstream browsers have eliminated the `SharedArrayBuffer` feature which is used for implementing the shared counter [Chromium, n.d.; Hazen, 2018; Wagner, 2018]. While some browsers have since partially re-enabled the feature, it is not fully enabled [MDN Web Docs, 2022].

Cache Occupancy Attacks. The cache occupancy attack [Shusterman et al., 2019; Shusterman et al., 2021a] measures the access time to a cache-size buffer. The timing differences between many and few cache misses can be detected with a low-resolution clock. Moreover, the papers show that by counting the number of times the buffer can be accessed between clock ticks reveals information even when the timer resolution is as low as 100 milliseconds. The downside of the attack is that it has a low spatial resolution, i.e., it provides a proxy of the amount of memory activity, but does not reveal information on which memory addresses are accessed. Under ideal conditions, the attack can be used for cryptography [Genkin et al., 2022], but it is not clear how this translates to real environments. Moreover, due to the time it takes to access the whole buffer, the temporal resolution of the attack is low.

Amplifying through Repetition. When the event that sets the cache state can be repeated, repeating it multiple times can amplify the difference between hits and misses. Such amplification has been used for a Spectre attack, where the attacker can repeat the attack as long as the leaked contents do not change [Hadad and Afek, 2018; McIlroy et al., 2019; Schwarz et al., 2019]. Theoretically, this approach could also be used to find eviction sets, but we are not familiar with any implementation of such approach and it would seem that a naive implementation will be too noisy to be effective.

PLRU Attack. Röttger and Janc (2021) propose an L1 cache attack that exploits the Pseudo-LRU replacement algorithm used in the Intel L1 caches. L1 attacks do not apply across cores, but it may be possible, in combination with the Prime+Scope attack [Purnal et al., 2021a], to apply cross-core attacks. However, it is not clear if and how the technique can be used for finding eviction sets with low-resolution timers.

Weird Gates. Evtvyushkin et al. (2021) describe “weird” gates that use transient execution to compute over microarchitectural state. They propose two types of gates. BP gates compute a logical function of the state of the branch predictor and the memory location that contains the code for the gate. They store the result as a state of a location in the data cache. Due to the differences between the types of states used for input and output, BP gates

are not composable and it is not clear how to create gates from them.

The second type of gates, TSX gates, use Intel’s Transactional Synchronization Extensions (TSX) [Intel, 2021a] to implement the gates. The feature is Intel-specific and is not supported by other processors. Due to security issues [Van Schaik et al., 2019; Van Schaik et al., 2021; Schwarz et al., 2019], Intel disabled the feature by default. Thus, TSX gates cannot work on newer and patched processors. Moreover, JavaScript and WebAssembly do not support TSX, hence TSX gates cannot be implemented in these languages.

While TSX gates are composable like ours, their accuracy (about 99%) is significantly lower than some of our gates’ (over 99.9%). Consequently, Evtyushkin et al. (2021) frequently transfer gates’ output to the architectural state. For example, their SHA-1 implementation exposes 41% of the intermediate values to the architectural state. In contrast, our SHA-1 implementation performs a full round without exposing any intermediate values. Moreover, to increase robustness, the software needs to execute each gate multiple times and perform statistical analysis to decide the likely correct outcome. While our circuits also use redundant computation to increase robustness, we use our majority gates to select the output without exposing the intermediate values to the software.

Evtyushkin et al. (2021) propose to use their gates to implement stealthy computation that is harder to detect and analyze (e.g., a form of program obfuscation to prevent reverse engineering). Similarly, our gates evaluation is also dependent on the microarchitectural layer, which also stores their state. Thus, the same analysis on stealthiness and program obfuscation apply to our gates.

Concurrent Work. In a concurrent and independent work, Kaplan (2023) also shows how to use speculative execution to create logical gates. The work mentions the possibility of using branch prediction (which we exploit), but also shows gates based on return address prediction, which require no branch training. It also demonstrates how to combine gates to create logical circuits, to speed up cache attacks, and to amplify cache measurement up to 600 millisecond.

4.6 Conclusions

Conducting proactive research on potential attack vectors of computing devices is a crucial practice in ensuring their security. In this chapter, our focus is on exploring vulnerabilities of low-resolution timers as a countermeasure against cache side-channel attacks. Specifically, we investigate how transient execution can be used to enhance cache attacks on systems that utilize low-resolution timers.

We present three types of logical gates that operate on the state of the cache. Our gates are sufficiently robust to perform Turing complete calculations, and are versatile enough to work on a range of environments, including

in browsers and across multiple processor architectures. We demonstrate that using our gates we can amplify the timing signal of cache miss vs. hit by six orders of magnitude, achieving a timing difference of 100 milliseconds. Our amplification strategy works well within a browser, and we demonstrate its use for building eviction sets in Chrome, using no timing source other than the JavaScript timer, whose resolution is 0.1 milliseconds.

We further present the Prime+Store attack, a variant of Prime+Probe that decouples cache sampling from the timing measurement. We demonstrate the power of Prime+Store by using it to attack the modular exponentiation implementation in a version of GnuPG. We show that we can sample at a rate that is more than 100 times faster than our clock rate, allowing us to obtain (secret) exponent bits from GnuPG.

We believe that our gates pave the way for implementing other primitives that manipulate microarchitectural state. Furthermore, by exposing the security ramifications that our gates pose, we hope to prompt the necessary measures to address these issues and enhance the security of systems in the future.

Chapter 5

Hardware Performance Counters in Cache Attack Detection

In this chapter, we focus on the final countermeasure against cache side-channel attacks discussed in this thesis, specifically software-oriented cache side-channel attack detection techniques that rely on hardware performance counters (HPC) for data. Owing to the serious implications of cache attacks, numerous defenses in software [Barthe et al., 2014; Carruth, 2018; Liu et al., 2016a; Shi et al., 2011; Zhang et al., 2023] and hardware [Feng et al., 2021; Khasawneh et al., 2019; Liu et al., 2016b; Loughlin et al., 2021; Qureshi, 2018; Wang and Lee, 2007; Werner et al., 2019; Yan et al., 2019b] have been suggested to safeguard against such threats. However, implementing these defense strategies can be complicated. Hardware-based defenses, such as secure caches discussed in Chapter 3, cannot be applied to existing hardware, and software countermeasures, such as low-resolution timer in browsers discussed in Chapter 4, might cause significant disruption to functionality and performance degradation. Instead of perpetually incurring the cost of protection, even when no attacks are happening, an alternate strategy seeks to identify ongoing attacks and only use countermeasures when attacks are detected [Akram et al., 2020]. A prevalent method in this area involves the use of HPCs, a collection of machine-specific registers that observe microarchitectural events, searching for statistical deviations that differentiate malicious from benign software. The statistical tests employed range from basic threshold techniques [Carnà et al., 2022; Payer, 2016] to sophisticated machine learning [Depoix and Altmeyer, 2018; Mushtaq et al., 2018b]. Given the large number of proposals, and the high quality of the reported results, in this chapter we ask the following question:

Are published HPC-based detection methods properly evaluated, such that their quality can be ensured for real-world deployment against cache side-channel attacks?

In this chapter, we find that the answer is an unfortunate negative.

We evaluate the correctness of performance evaluation of 48 relevant side channel detection literature, and find four commonly occurring problems in their settings and assumptions, particularly in the case of accuracy, overhead, detection speed and threat model. We subsequently present our findings in a form of a survey. We further demonstrate how these improper evaluation settings, and especially weak threat models assumed by publicly available detection methods leave them vulnerable to sophisticated attacks that do not make such weak assumptions. To illustrate this weakness, we develop new camouflaged attacks that mask their malicious execution patterns behind benign program execution, allowing them to evade detection effectively. We demonstrate the success of these attacks in stealing sensitive data, while simultaneously remaining undetected by two publicly available detection methods, that otherwise detect proof-of-concept attacks with reasonable accuracy. To further support our case, we construct our own implementation of a detection method, named **HPCache**, that can detect proof-of-concept cache attacks with perfect accuracy, and yet we show that it is still fallible to our camouflaged attacks. This emphasizes the importance of assessing detection methods against more advanced attack models. We show that merely claiming defense against a specific attack type, particularly when based on proof-of-concept implementations, without providing comprehensive insight into its implementation and threat model, falls short of ensuring reliable protection.

By conducting a survey and our own experimentation on the vulnerability of HPC-based cache side-channel attack detection methods, we have found that issues persist in the performance evaluation of current detection approaches, as mentioned earlier, which hamper their practical implementation on a wider scale. We conclude that without addressing the aforementioned evaluation shortcomings, it is uncertain whether real-time cache side-channel attack detection systems can truly be deemed effective for practical use in real-world scenarios.

5.1 Background

5.1.1 Hardware Performance Counters.

Modern CPUs include hardware performance counters (HPCs), initially introduced for the purpose of debugging [Das et al., 2019]. This is done through recording CPU events such as number of cycles, and branch misses. In Intel processors, this functionality is implemented under the name *Performance Monitoring Units* (PMUs). They consist of individual counters called *Performance Monitoring Counters* (PMCs). HPCs can be programmed by setting specific *Machine Specific Registers* (MSRs) in the processor.

5.1.2 HPC-Based Cache-Side Channel Attack Detection Methods

HPCs have been used to facilitate debugging and dynamic profiling [Das et al., 2019]. However, recent research has uncovered another area where these counters may be advantageous, namely in the realm of cache side-channel attack detection [Akram et al., 2020; Mushtaq et al., 2018a]. In this section, we examine the challenges of using HPCs for security-related purposes. We then explore the classification of HPC-based detection methods for cache side-channel attacks.

Challenges of Using Performance Counter for Security Practices

Das et al. (2019) have brought to light issues of non-determinism and contamination when utilizing hardware performance counters for security purposes. This is due to the fact that the usage of non-architectural events, which are specific to the microarchitecture of a processor (e.g., cache accesses, branch prediction, and TLB accesses) for security applications of hardware performance counters can be problematic. These events differ across processor architectures and may also change with processor enhancements. The second issue, contamination, arises because performance monitoring units operate at the hardware level and are application agnostic. Therefore, when an interrupt is configured to notify of performance monitoring events, the PMU generate interrupts for all processes running on a given processor core. To obtain an accurate profile of an application, it is essential to filter the performance counter data relevant only to the process of interest, as performance data can be contaminated by the events of other processes. The authors note that while these issues may not have significant consequences for certain applications, they can have a significant impact on approaches whose security depends on having accurate and consistent hardware performance counter measurements. For instance, malware exploit defenses are vulnerable to non-deterministic effects and contamination of events, as security applications rely on small variations in performance counter data to distinguish between suspicious and benign behaviors. Even minor variations of 1–5% in counter values can cause these models to perform poorly. Therefore, it is especially important to address these challenges in security-related applications.

In line with what was suggested by Das et al. (2019), Zhou et al. (2018) conducted an experimental study on traditional malware that demonstrate how the use of microarchitectural level information obtained from hardware performance counters (HPCs) cannot differentiate between benignware and malware. Previous HPC-based malware detectors rely on the assumption that malicious behavior affects measured HPC values differently than benign behavior. However, it is debatable and counter-intuitive as to why the semantically high-level distinction between benign and malicious behavior would manifest itself in the microarchitectural events that are measured by HPCs. The authors do not be-

lieve that there is a causal relationship between low-level microarchitectural events and high-level software behavior. They argue that the positive results in previous research are due to a series of optimistic assumptions and unrealistic experimental setups.

Furthermore, Jiang et al. (2022) evaluate existing detection tools for cache attacks on Secure Guard Extension (SGX), Intel’s implementation of a trusted execution environment in x86_64 processors. They identify how these tools fail to accommodate various subtleties in the use of HPCs in the case of exploit prevention and malware detection. They also demonstrate how an adversary can manipulate HPCs to bypass certain security defenses, making detection tools less effective in detecting side-channel attacks on SGX enclaves. Existing detection mechanisms are geared towards an adversary that interferes with the victim’s execution to extract the most secret bits, causing significant performance degradation that can signal an attack. However, they show that an adversary leaking smaller portions of secret, as small as a single bit at each execution of the victim, can remain undetected. They specifically demonstrate that an adversary can profile a victim enclave to identify the precise moment during execution when a specific part of the secret can be leaked via a side-channel attack. By running the victim multiple times and leaking a different part of the secret each time, their technique can recover the whole secret while remaining undetected. They adapt known attacks that leverage page tables, L3 cache, or a combination of the two and evaluate their performance on routines on libcrypto, used by cryptographic algorithms like ElGamal, RSA, and EdDSA. They show that an adversary using their attack technique cannot be detected by existing detection tools unless they tolerate a large number of false positives. They also provide evidence that any detection tool that monitors the performance of the victim is equally likely to fail.

Classification of Cache Attack Detection Methods Using HPCs

Despite the difficulties and literature advising against the use of HPCs in a security context, a considerable number of papers support the application of HPCs to detect cache attacks. We gathered 48 papers related to cache side-channel attack detection methods that use HPCs. To gain a better understanding of these detection techniques, we first present the categories into which they can be grouped. For this purpose, we refer to the study by Akram et al. (2020), which classifies academic papers on cache attack detection methods promoting HPC use into two primary groups based on their detection pattern: *signature-based* and *anomaly-based*. Moreover, these detection approaches can be sorted by their classifier type, either as *machine learning-based* or *threshold-based*.

Signature-based detection methods analyze the status of microarchitectural components to identify any patterns that may be indicative of an attack. This technique looks for similarities to known attack patterns. If the HPC

readings of an application reach a certain similarity threshold with a known attack pattern, the detection mechanism is triggered.

Anomaly-based detection methods continuously scan microarchitectural patterns to search for similarities with a benign application. These methods identify potential attacks by comparing the behavior of monitored applications by reading their HPC values and comparing them with the expected values of a benign application. When an application’s readings deviate from what is expected of a harmless application, the detection method flags the application as a possible threat. This is based on the premise that benign applications usually generate a modest number of microarchitectural HPC readings, and any number that exceed a certain threshold are considered anomalous and may indicate the presence of cache attacks.

Signature-based detection is more accurate in detecting known attacks, but is more prone to false negatives when faced with new attacks, whereas anomaly-based can possibly detect new attacks but at the same time may experience false positives when execution of benign program unexpectedly changes [Alam et al., 2017].

The detection techniques for cache side-channel attacks can be further classified based on their classifiers, which are the methods used to determine the likelihood of an attack based on the collected data. There are two main methods of classification:

Threshold-based detection methods use a simple limit-based classification method, this is done in such a way that when HPC values are above certain threshold, the detection method associates that trace with an attack.

Machine learning-based detection methods use a more advanced way of classifying data, i.e., with the help of machine learning classifiers. The idea behind this approach is to let machine learning algorithms learn features in the collected data. This is done with hopes that the classifier can generalize better over the data, to improve detection accuracy and detect new attacks better.

Mushtaq et al. (2022) note that in controlled settings with minimal background noise, threshold-based techniques may be adequate for identifying attacks. However, in more realistic and noisy environments, these methods have difficulty distinguishing between benign features and attacks. Therefore, they suggest that machine learning-based approaches are better suited for use in these types of settings.

5.2 Survey of HPC-Based Cache Side-Channel Attack Detection Method Evaluation

In Section 5.1.2, we briefly discuss a previous survey by Akram et al. (2020) that classified various works on cache side-channel attack detection methods, revealing that 20 of them rely on hardware performance counters (out of 23 in

total). Given the growing popularity of this approach, we set out to investigate whether the use of performance counters is a viable option for detecting cache attacks and whether it is an appropriate and effective approach.

As part of our research, we searched scholarly works related to cache side channel detection using HPCs and identified 48 articles on the subject. After conducting a survey, we identified several recurring issues across these works that we consolidated into four categories of methodological shortcomings. These categories consist of improper measurement of accuracy, overhead, and detection speed, as well as a weak threat model used for assessing the effectiveness of detection methods.

Our findings are summarized in Table 5.1, where each circle represents the conformance of the proposed method to the criteria described above.

Ref	ACC	OV	DS	TM	Ref	ACC	OV	DS	TM
[Ahmad, 2019]	●	●	●	●	[Mushtaq et al., 2018b]	●	○	○	○
[Ahmad, 2020]	●	●	●	●	[Mushtaq et al., 2018c]	●	●	●	○
[Alam et al., 2017]	●	○	○	○	[Mushtaq et al., 2018d]	●	●	●	○
[Alam et al., 2021]	●	●	●	●	[Mushtaq et al., 2020]	●	●	●	●
[Albalawi et al., 2022]	●	●	○	○	[Mushtaq et al., 2021]	●	●	●	○
[Allaf et al., 2019]	●	●	○	○	[Mushtaq et al., 2022]	●	●	●	○
[Bazm et al., 2018]	●	●	○	○	[Payer, 2016]	○	●	○	●
[Briongos et al., 2018]	●	○	●	●	[Polychronou et al., 2021]	●	●	○	●
[Carnà et al., 2022]	●	●	●	○	[Prada et al., 2019]	○	○	○	○
[Chiappetta et al., 2016]	●	●	○	●	[Sabbagh et al., 2018]	●	○	○	○
[Cho et al., 2020]	●	●	●	○	[Singh and Rebeiro, 2021]	○	●	○	●
[Choudhari et al., 2022]	●	○	●	○	[Tao et al., 2021]	○	○	○	○
[Chouhan and Halabi, 2016]	○	○	○	○	[Tong et al., 2020]	●	○	○	○
[Depoix and Altmeyer, 2018]	●	●	●	○	[Tong et al., 2022]	●	○	○	●
[Dutta and Sinha, 2019]	○	○	○	○	[Vanathi and Chokkalingam, 2018]	○	○	○	○
[Ferracci, 2019]	○	●	●	○	[Wang et al., 2020a]	●	●	○	●
[Gregory and Harini, 2021]	●	○	○	●	[Wang et al., 2020b]	●	○	○	○
[Gülmezoglu et al., 2019a]	●	●	●	●	[Wang et al., 2021a]	●	○	○	○
[Hamza et al., 2021]	●	●	●	○	[Wang et al., 2021b]	●	●	○	○
[Kim et al., 2021]	●	●	○	○	[Wang et al., 2022]	●	●	○	○
[Kulah et al., 2019]	●	●	○	●	[Wu et al., 2022]	●	●	●	●
[Lantz, 2021]	○	○	●	○	[Yan and Cui, 2022]	●	○	○	○
[Li and Gaudiot, 2018]	●	○	○	○	[Zhang et al., 2016]	●	●	●	○
[Li and Gaudiot, 2019]	●	○	○	○	[Zheng et al., 2022]	●	●	●	○

Table 5.1: Survey result. ACC represents accuracy, OV represents overhead, DS represents detection speed, TM represents evaluation against stronger threat models. An empty circle ○ signifies that the criterion has not been assessed, nor mentioned, while a semi-filled circle ◐ implies that the literature mentions the specific criterion but lacks certain attributes required for accurate assessment. A fully-filled circle ● demonstrates that the method conducts its assessment of the criterion appropriately.

5.2.1 Accuracy

Assessing the accuracy of cache attack detection methods entails examining both false positive and false negative rates. Keeping a low false negative rate is essential, as detection techniques strive to detect potential attacks on the systems they safeguard, thus reducing the likelihood of overlooking malicious efforts. While having a false negative rate of 0% is the ideal situation, it may be challenging due to the continuous emergence of new threats and the intrinsic disadvantage defenders experience in a “cat and mouse” context with malicious actors perpetually avoiding detection. As such, we do not mandate a specific false negative rate for detection methods to reach. Nevertheless, we expect that research papers evaluate the accuracy of their detection method.

Out of all the papers we examined, the majority (39 out of 48) evaluate the detection accuracy of their methods. Nine exceptions do not perform this assessment. We denote these nine papers with empty circles in Table 5.1 ¹.

False positive rate is another important aspect to consider. A detection method with a high false positive rate can be detrimental to the protected system, as legitimate programs may be incorrectly flagged and terminated.

Because of this, we expect detection systems to have an absolute minimum false positive rate, considering that they may be run on a system for a long period of time and the potentially large number of applications running on the protected system. False positive rate, even very small, will result in a cumulatively large number of wrongly flagged applications.

We found that the majority of evaluated papers (32 out of the 39 that evaluate their accuracy) report false positives at a rate higher than zero, for which we assign a semi-filled circle in Table 5.1. The remaining six meets our criteria of absolute zero false positive rate, to which we assign a fully-filled circle.

5.2.2 Overhead

The term “overhead” refers to the decrease in system performance that occurs when attack detection methods are executed.

In order to ensure a fair assessment of the overhead, we recommend that benchmark applications be either pinned to the same core as the detection method or run on all cores. This guarantees that the detection method influences the benchmark by ensuring that the benchmark is scheduled on the same core as the detection method. This approach prevents the misconception of the detection method being “overhead-free” due to them being scheduled on different cores. Only three papers, specifically [Carnà et al., 2022; Kulah

¹It is important to note that there are situations where assessing accuracy may not be necessary, such as in the case of [Singh and Rebeiro, 2021]. Their work focuses on developing a mitigation that minimizes performance penalty in the case of false positive. We only take note of whether accuracy assessment is present. An absence of such evaluation does not necessarily imply any incompleteness in the research.

et al., 2019; Singh and Rebeiro, 2021] appropriately examine this criterion. Their detection techniques operate during every context switch (which occurs on each core of the system), signifying that their benchmark application and the detection method run on the same core, thus allowing for a fair overhead assessment.

Out of the 48 papers reviewed, 31 assess their detection overhead, while 17 do not perform any overhead evaluation. For papers that do not perform any overhead evaluation, we assign an empty circle in Table 5.1. For papers that evaluate their detection method’s overhead but do not meet the aforementioned criterion, we assign a semi-filled circle. For the three papers that meet the criterion, we assign a fully-filled circle.

Out of the 17 papers that do not evaluate overhead, five [Chouhan and Halabi, 2016; Gregory and Harini, 2021; Li and Gaudiot, 2018; Li and Gaudiot, 2019; Tong et al., 2020] claim that the detection overhead associated with collecting HPC data is low, without evaluating their overhead. We believe that this claim is inadequate to demonstrate the low overhead of their detection system. This is because evaluation of system overhead should include not only the cost of reading HPC data but also the overhead associated with performing attack classification, and potentially scanning the processes running on the system.

Despite measuring their overhead with the benchmark application and data collection process pinned to the same core, we consider the overhead evaluation of Chiappetta et al. (2016) to partially cover the true overhead of their detection method. This is because their evaluation only encompasses the overhead of their HPC collection module and disregards the classification step, which is typically more resource-intensive of the two. Therefore, we assign a semi-filled circle.

5.2.3 Detection Speed

We assess research papers based on whether they examine the detection speed of their detection techniques. A thorough evaluation of detection speed involves conducting such an assessment and presenting the results in terms of either the time required to identify attacks or the percentage of attack completion when the attack is detected. We assign fully-filled circles to papers that meet these criteria. Out of 48 papers analyzed, six [Cho et al., 2020; Choudhari et al., 2022; Chouhan and Halabi, 2016; Gülmezoglu et al., 2019a; Hamza et al., 2021; Zhang et al., 2016] present their detection speed using the former metric, while ten [Briongos et al., 2018; Carnà et al., 2022; Lantz, 2021; Mushtaq et al., 2018c; Mushtaq et al., 2018d; Mushtaq et al., 2020; Mushtaq et al., 2021; Mushtaq et al., 2022; Wu et al., 2022; Zheng et al., 2022] use the latter. We consider the latter metric to be more informative because it signifies the maximum amount of key leakage from a cryptographic algorithm. Nevertheless, both approaches are arguably valid, and therefore we assign a

fully-filled circle in Table 5.1 for these 16 papers.

Among the remaining 32 papers, 15 provide the detection method’s HPC sampling interval but do not evaluate the system’s detection speed in detecting attacks, for which we assign a semi-filled circle, while 17 do not mention this criterion at all, to which we assign an empty circle.

5.2.4 Threat Model

HPC-based detection methods face a significant problem in that many of them are developed with the assumption that attackers will only use naive, proof-of-concept implementations of attacks. However, this assumption is inaccurate because in reality, attackers are likely to use advanced techniques to evade detection. Therefore, it is critical to determine whether current detection systems can effectively detect these evasive attacks.

Overall, 22 out of 48 papers acknowledge the possibility of evasive attacks. However, out of these 22, only eight evaluate their detection method against any sort of attack modification efforts. For these eight papers, we assign a fully-filled circle, whereas for the remaining 14 papers, we assign a semi-filled circle. For the 26 papers that do not acknowledge this issue, we assign an empty circle.

5.3 Assessing the Quality of Attack Detection Methods

In Section 5.2, we highlight several performance assessment issues with recent proposals for HPC-based cache side-channel attack detection methods. In this section, we apply the evaluation criteria proposed in Section 5.2 to publicly accessible detection methods, as well as our own method, to assess their performance.

5.3.1 Experiment Environment

The experiments in this section are conducted on an Intel NUC 9 Extreme Kit that comes with an Intel Core i7-9750H CPU. The system runs on Ubuntu 22.04.

5.3.2 Our Method

We made efforts to acquire the implementation code from the authors, but we were only able to obtain a limited number of solutions. Out of the total number of papers (48), we found two available online and contacted the authors of the remaining 46 papers via email. We received responses from 21 of them, out of

which 13 provided us with the code. However, only two were functional, which means that they were able to compile and perform the detection as intended.

Due to the unavailability of reliable implementations and our inability to verify the quality of the detection methods we find, we supplement our experiments with our own cache attack detection solution that uses comparable technique to other proposed methods. This detection method is called **HP-Cache**.

It is important to note that our detection method does not aim to offer flawless detection accuracy, minimum performance overhead, nor the ability to detect advanced threat models. Instead, it serves as an illustration of how to apply the evaluation criteria outlined in Section 5.2.

It consists of three modules: the Process Checker, the Data Collector, and the Classifier. The Process Checker module scans the system for running processes, tracks started and killed processes, and sends process information to the Data Collector module. The Data Collector module uses the process information from the Process Checker to collect HPC data from each running process in the system every 100 milliseconds. The collected data is associated with the process from which it was sampled and then passed on to the Classifier module. The Classifier module processes the HPC data using a classifier algorithm chosen by the user (in the experiments in this paper, a neural network classifier is used) to determine whether the HPC data is indicative of cache attacks.

PAPI Event Name	Intel Mnemonic
PAPIL1_DCM	L1D.REPLACEMENT
PAPIL1_ICM	L2_RQSTS.ALL_CODE_RD
PAPIL1_TCM	L1D.REPLACEMENT, L2_RQSTS.ALL_CODE_RD
PAPIL2_ICM	L2_RQSTS.CODE_RD_MISS
PAPIL2_TCA	L2_RQSTS.ALL_CODE_RD, L2_RQSTS.ALL_DEMAND_REFERENCES

Table 5.2: PAPI events used in our detection method, their description and native performance counter events in x86_64.

To gather performance counter data for our detection, we utilize the PAPI library [Terpstra et al., 2009], which offers a consistent interface and approach for accessing the performance counter hardware present in most major microprocessors. The HPC events we use as a data source for our detection method are listed in Table 5.2. This table includes the names of PAPI events, and their corresponding Intel performance counter mnemonics. These specific events are chosen because they encompass a wide range of cache-related events in the L1 and L2 caches, considering the limitations of the available performance counters. Since we are dealing with non-inclusive caches, it is important to note

that LLC misses also result in misses in the L1 and L2 caches. Consequently, these counters take into consideration such misses.

5.3.3 Accuracy

First, to assess the accuracy criterion, we test the accuracy of these detection methods in detecting standard proof-of-concept attacks. To this end, we have selected an implementation of Spectre-PHT [Crozone, 2023], and an implementation of the Flush+Reload attack on GnuPG from the Mastik library [Yarom, 2016]. These attacks are chosen due to their significant security implications [Kocher et al., 2019; Yarom and Falkner, 2014]. In particular, the risk posed by Flush+Reload is noteworthy as it is capable of facilitating the theft of cryptographic keys [Yarom and Falkner, 2014].

	[Payer, 2016]	[Depoix and Altmeyer, 2018]	HPCache
Criteria			
Number of Datapoints	2000	2000	2000
Number of true positive	861	843	1000
Number of false positive	0	500	0
Number of true negative	1000	500	1000
Number of false negative	139	157	0
False negative rate	13.9%	15.7%	0%
False positive rate	0%	50%	0%

Table 5.3: Accuracy of three detection methods.

	HPCache 100ms	HPCache 10ms	HPCache 1ms
Criteria			
Number of Datapoints	2000	2000	2000
Number of true positive	1000	570	577
Number of false positive	0	0	28
Number of true negative	1000	1000	972
Number of false negative	0	430	423
False negative rate	0%	43%	42.3%
False positive rate	0%	0%	2.8%

Table 5.4: Accuracy of HPCache with 100, 10, and 1 millisecond sampling interval.

We assess the accuracy of our detection technique, along with two other detection methods [Depoix and Altmeyer, 2018; Payer, 2016] across an 8-hour time frame, during which we execute multiple benign programs, encompassing

the CPU stress-testing application `stress-ng`, the `gcc` compiler for compilation, GnuPG for decryption, and the SPEC CPU 2017 `gcc_r` benchmark. Additionally, we run malicious applications such as Spectre and Flush+Reload. Our detection method is put to the test against 1,000 benign and 1,000 malicious applications.

Table 5.3 presents our detection method’s accuracy compared to that of Payer (2016) and Depoix and Altmeyer (2018). To make a fair comparison, we test all detection methods against the same collection of benign and malicious samples. We find that Depoix and Altmeyer (2018) only scans for processes running before starting their detection method, therefore we conduct our evaluation by initiating the malicious attack, then launching their detection method, and repeating this sequence for each subsequent experiment. As for Payer (2016), we encounter some memory errors after a few minutes of running their detection method and had to modify our testing approach to start a new instance of the detection method for each sample being tested.

Furthermore, we conduct experiments to evaluate the effect of different sampling intervals on the accuracy of our detection method. Table 5.4 shows that using a sampling interval of 1 or 10 milliseconds leads to a decline in accuracy compared to using a 100 millisecond sampling interval. Under a 1 millisecond sampling interval, both the false negative and false positive rates increase. Similarly, under a 10 millisecond sampling interval, the false negative rate increases, while the false positive rate remains at zero. We conclude that this outcome results from inadequate amount of data being collected within both 1 and 10 millisecond intervals, as both benign and malicious applications have execution periods where the cache miss rate is exceptionally high or low. Collecting HPC data within these intervals fails to capture a comprehensive view of program execution. Therefore, we determine that the sampling interval of 100 milliseconds is optimal for our detection method.

Our tool’s detection capabilities and functionality are comparable to others in the field [Choudhari et al., 2022; Mushtaq et al., 2018a; Mushtaq et al., 2018b; Mushtaq et al., 2018c; Mushtaq et al., 2018d; Mushtaq et al., 2020; Mushtaq et al., 2021; Mushtaq et al., 2022; Tong et al., 2020; Tong et al., 2022; Wu et al., 2022; Yan and Cui, 2022], as demonstrated by the low false positive and false negative rate of 0%. In conclusion, it is important to conduct accurate evaluations of detection methods’ accuracy. A low false negative rate ensures effective protection, and a low false negative rate prevents excessive false positive rate that could render them impractical.

5.3.4 Overhead

Second, we evaluate the overhead of these detection methods using the technique we recommend for overhead evaluation in Section 5.2.2.

To evaluate overhead, we use the 7zip benchmark to measure CPU performance during file compression and decompression. This benchmark is chosen

Detection Method	[Payer, 2016]			[Depoix and Altmeyer, 2018]			HPCache		
	ST	AT	PT	ST	AT	PT	ST	AT	PT
Not Running	3075	24162	3096	3075	24162	3096	3075	24162	3096
Running	3091	24051	3062	2978	22827	1804	2970	20013	1494
Overhead (%)	-0.5	0.5	1.1	3.2	5.5	41.7	3.4	17.2	51.7

Table 5.5: Overhead of detection methods. ST signifies single-threaded 7zip benchmark. AT signifies 12-threaded of the benchmark, while PT signifies the a single-threaded benchmark, pinned on the same core along with the detection method tested.

because of its high CPU usage, which makes it an ideal test to determine the influence of the detection methods on CPU performance.

We perform three experiments to evaluate overhead. Initially, we execute a single-thread benchmark followed by a 12-thread benchmark, which aligns with the number of cores in our CPU. The third experiment is to run the benchmark (single-thread) on the same core as the detection method using the `taskset` command. We repeat 80 compression and decompression tasks in each scenario, with and without the detector running. The rationale behind running the 12-thread benchmark is to ensure the detection method has an effect on the benchmark, by ensuring that at least one of the benchmark threads is scheduled on the same core as the detection method. Note that the benchmark program monitors the performance of all the threads simultaneously. This approach prevents an illusion of the detection method being “overhead-free” due to them being scheduled on different cores. Similarly, allocating both the benchmark and detection methods to the same core also avoid underestimating the detection method’s overhead.

The results of our experiment are shown in Table 5.5, with scores expressed in million instructions per second (MIPS). As mentioned, running a single-thread benchmark and detection method without ensuring they are scheduled on the same core results in a low overhead reading. Interestingly, Payer (2016) showed 0.5% higher performance when running their detection method than when it was not running. Depoix and Altmeyer (2018) and HPCache showed around 3% overhead.

Running the benchmark on all cores results in higher detection overhead in all three detection methods. Pinning a single-thread instance of the benchmark with the detection method results in an even higher overhead evaluation for all detection methods. We note that Depoix and Altmeyer (2018) and HPCache show significantly higher overhead in this configuration as they are both multi-threaded, causing a sharp decrease of the benchmark’s performance when pinned to a single core with the detection method. The results highlight the importance of correctly configuring the benchmark to avoid unfairly low overhead results, as demonstrated in the single-threaded benchmark experiment of the detection methods.

In conclusion, it is crucial to use a proper overhead benchmark setting to ensure fair evaluation. Precise reporting of overhead is important as unexpectedly high overhead can hinder the adoption and practicality of detection methods.

5.3.5 Detection Speed

Third, we assess the detection speed of detection methods by measuring the time it takes for them to identify ongoing attacks.

With a sampling interval of 100 milliseconds, **HPCache** can identify attacks within 300 milliseconds of the execution of a malicious program. By increasing the sampling interval to 10 milliseconds, the tool can detect attacks in just 100 milliseconds. However, we have found that using a sampling interval of 1 millisecond leads to a longer detection time. This is because the amount of information collected during this period is insufficient, as discussed in Section 5.3.3. Consequently, the detection method can only recognize an attack trace as malicious after it has been running for a longer period, resulting in longer detection times.

For comparison, Payer (2016) uses a sampling interval of 1000 milliseconds, and detects attacks within 1100 milliseconds. Regarding Depoix and Altmeyer (2018), we were unable to test their detection speed because their method only detects attacks that were executed prior to its start-up, meaning that it cannot detect new attacks.

These findings highlight the discrepancy between the sampling interval of HPC and the detection speed of a detection method. They emphasize the importance of accurately evaluating the detection speed rather than solely stating sampling interval used (as seen in numerous papers in Section 5.2.3). For example, collecting HPC-data every 100 milliseconds does not guarantee that attacks are detected within such time-frame. Such precise reporting of detection speed enables users to make informed decisions regarding the suitability of detection methods in safeguarding their systems against specific threats.

5.3.6 Threat Model

Last, we assess the accuracy of the detection method against a stronger attack model. We are interested in determining how effective these detection methods are in identifying camouflaged attacks, which hide their malicious activities within benign code. We also compare the accuracy of these detection methods when detecting proof-of-concept attacks to understand any discrepancies between these two setups, and to understand whether a weak threat model assumption leads to overestimation of detection methods' capability. To achieve this, we select the SPEC CPU 2017 `gcc_r` benchmark program [SPEC, 2020] as our benign program target and inject Flush+Reload and Spectre attacks within the program's normal execution. The Flush+Reload attack specifically

aims to recover the private key from a vulnerable implementation of the ElGamal encryption algorithm [El Gamal, 1984] used in the cryptographic software GnuPG version 1.4.13. The core of the attack lies on the modular exponentiation operation, which involves raising a base b to the power e modulo some modulus m , i.e., calculating $b^e \bmod m$. In the context of ElGamal decryption, the private key serves as the exponent e . Consequently, the attack aims to retrieve the exponent.

To create our camouflaged attacks, we inject the Flush+Reload attack and the Spectre attack into the SPEC CPU 2017 `gcc.r` benchmark. The `gcc.r` benchmark is a C compiler that tests the optimization and code generation capabilities of the CPU. Attack injection is done by inserting an injection code that triggers those attacks to the method which is called most often in the benchmark application. The attacks are executed with a small probability, making their executions infrequently interleaved between benign, benchmark code.

Listing 5.1: Injection Code

```
void inject_attack() {
    if (rand() < PROB)
        do_fr_or_spectre();
}

static inline bitmap_element *
bitmap_find_bit (bitmap head, unsigned int bit)
{
    // Inject attack at the beginning of this function.
    inject_attack();

    bitmap_element *element;

    ...
}
```

We choose the `bitmap_find_bit` function which is called the most number of times during the benchmark execution, as reported by the profiling tool GProf. In Listing 5.1 we show the insertion of `inject_attack` function at the very beginning of the `bitmap_find_bit` function. Inside `inject_attack`, we set with a low probability (as listed in Table 5.9) that an actual Flush+Reload and Spectre attacks are executed. With these, we essentially interleave the execution of Flush+Reload and Spectre alongside the actual benchmark. Note that the probability of running the attack code is set to be very small, and therefore, for the majority of the injected program, its execution largely resembles that of the actual benchmark.

In order to minimize the impact on the benchmark execution and reduce the differences in HPC readings (plain benchmark vs. attack-injected benchmark) caused by the attack activities, we schedule the attack infrequently between actual benchmark procedures and ensure that the execution of the attack at each iteration is brief. For Flush+Reload, this results in shorter traces that do not capture the complete key. Furthermore, since we do not assume

Name	Classifier	Spectre		Flush+Reload	
		POC	CMF	POC	CMF
[Payer, 2016]	Threshold	73%	0%	100%	0%
[Depoix and Altmeyer, 2018]	NN ¹	100%	100%*	62%	100%*
HPCache	NN ¹	100%	0%	100%	0%

* Depoix and Altmeyer (2018) detected camouflaged attacks with 100% due to false positive in the original SPEC CPU 2017 gcc_r benchmark

¹ Neural-Network

Table 5.6: True positive rate of proof-of-concept (POC) and camouflaged attacks (CMF).

Name	Classifier	Spectre		Flush+Reload	
		POC	CMF	POC	CMF
[Payer, 2016]	Threshold	0%	0%	0%	0%
[Depoix and Altmeyer, 2018]	NN	100%	100%	100%	100%
HPCache	NN	0%	0%	0%	0%

Table 5.7: False positive rate of proof-of-concept (POC) and camouflaged attacks (CMF).

any synchronization between our attack and the victim’s ElGamal encryption algorithm, the traces may begin at any stage of the encryption algorithm. To address these issues, we use the approach described in Section 4.4.5 to recover the complete key from partial traces.

Detection Discrepancies Between POC and Camouflaged Attacks

We conduct experiments to gauge the effectiveness of Payer (2016), Depoix and Altmeyer (2018), and HPCache in detecting attacks, in particular, we test these detection methods capability in detecting both regular and camouflaged attacks. We run both Spectre and Flush+Reload attack applications 1000 times and allow the methods a maximum of ten seconds to identify each scenario. If the detection method is able to detect an attack within ten seconds, we consider it a true positive. Otherwise, we consider it a false negative. We also test the detection method against benign SPEC CPU 2017 gcc_r benchmark for 1000 times. If the benchmark is detected as malicious, we consider this a false positive, otherwise we consider it a true negative.

We summarize the results of our experiments in Table 5.6 and Table 5.7, which present the true positive rate and false positive of each detection method in identifying both proof-of-concept and camouflaged Spectre and Flush+

Reload attacks.

HPCache achieves perfect accuracy in detecting proof-of-concept Spectre and Flush+Reload attacks on GnuPG, however it fails to detect any of the camouflaged Spectre and Flush+Reload attacks.

Depoix and Altmeyer (2018) detect proof-of-concept Spectre and Flush+Reload attacks with 100% and 62% accuracy respectively, they also detect 100% of the camouflaged as malicious. At the same time, they falsely detect the benchmark without any attack with 100% false positive rate, while HPCache and Payer (2016) do not falsely detect the benchmark application.

Payer (2016) detect proof-of-concept Spectre and Flush+Reload attacks with 73% and 100% accuracy respectively, however they fail to detect any of the camouflaged Spectre and Flush+Reload attacks.

The findings show the discrepancy in accuracy when detecting proof-of-concept attacks compared to camouflaged attacks, emphasizing the need for evaluating detection methods against a stronger threat model. It is evident that simply stating defense against a particular attack without offering details of its implementation and threat model can lead to an overestimation of the effectiveness of detection methods.

Re-training Model with Camouflaged Attacks

At first glance, the failure of the detection methods may appear to be caused by the training data not being well-suited to such camouflaged attacks, and a simple re-training of the classifier or adjusting threshold values could solve the problem. As suggested by [Depoix and Altmeyer, 2018], retraining of classifiers is needed when deploying a detector in a new environment or when supporting detection of new attacks.

Total number of datapoints	2000
Number of true positives (TP)	1000
Number of false positives (FP)	1000
number of true negatives (TN)	0
number of false negatives (FN)	0
False negative rate	0%
False positive rate	100%

Table 5.8: Accuracy of HPCache when trained with camouflaged attacks labeled as malicious.

Our analysis indicates that the root cause of the problem goes beyond inadequate training data. Retraining HPCache to include camouflaged attacks proved ineffective, as it resulted in a sharp increase in false positives, resulting in 100% false positive rate. Table 5.8 shows the result of the detection method trained with camouflaged attack labeled as malicious.

This shows that the detection method in discerning between genuinely benign program execution and the execution of a benign program injected with malicious attacks. Since the detection method is unable to distinguish between the two, it labels genuinely benign applications as malicious. This is because during the training of the detection method, the training data for the execution of injected attacks is marked as malicious. Consequently, the detection method becomes confused and starts classifying benign applications as malicious.

Cost of Evading Detection

Despite the fact that camouflaging attack behind the execution of a benign application can be an effective means of sidestepping detection, adversaries do not employ this strategy without incurring costs. Attackers must expend valuable CPU time executing decoy code, which, in the case of Flush+Reload, ultimately does not contribute to their primary goal of attacking cryptographic keys. The more aggressively an attacker tries to camouflage its attack by mimicking benign behavior, the less time is available to execute its malicious payload. Hence, we are interested to determine the extent of the trade-off in terms of attack time that attackers must accept when executing camouflaged attacks.

In Table 5.9, we present the results of our experiment on the camouflaged Flush+Reload attack, including the frequency of attack injection, the time required for data collection to allow complete key recovery, and detection rate. The frequency column represents the probability of executing the malicious attack injection within the *bitmap_find_bit* function listed in Listing 5.1.

As the table shows, when the camouflaged attack is executed with an injection probability of one in ten million, the time needed to recover the full ElGamal key is approximately 18 hours. However, when the attack is run more aggressively, with an injection probability of one in ten thousand, the time required to recover the key drops to around six minutes. It is worth noting that even at the highest level of attack injection aggressiveness, none of the attacks are detected by HPCache (trained with normal benign and malicious applications, i.e., not trained with injected attacks), which is in contrast to the detection rate of the Flush+Reload attack without camouflaging, which is detected 100% of the time. Consequently, attacks with less injection aggressiveness also remains undetected. We correctly recover 436 out of 459 bits private ElGamal decryption key.

These findings reveal that our camouflaged attack can stealthily extract secret keys and evade detection methods, even when executed under the most aggressive settings. This highlights the practicality and feasibility of the attack for adversaries, underscoring the importance for detection methods to be evaluated under a stronger threat model to avoid overestimating their detection capabilities. In cases where a detection method fails to identify attacks under

Frequency	Time needed	Detected
1/10,000,000	18:13:15	0%
1/1,000,000	1:50:34	0%
1/100,000	00:15:08	0%
1/10,000	00:05:41	0%

Table 5.9: Attack times for different Flush+Reload camouflaging aggressiveness.

such a threat model, it is crucial for authors to acknowledge this limitation. Transparency in acknowledging such limitations is preferable, as it provides users with valuable information and helps them understand any potential unexpected compromises they may encounter in their system’s protection.

5.4 Conclusions

Prevention and mitigation techniques against cache side-channel attacks have been proposed to counter the ever-increasing threat of these attacks. However, the high cost of hardware solutions has prompted researchers to explore cheaper software-based alternatives, such as HPC-based attack detection methods. In this chapter, we reveal that the performance evaluation of current proposed methods are insufficiently conducted to ensure effective protection in practical real-world scenarios. We analyzed 48 papers and found that none of them performed proper evaluation of all the necessary criteria of accuracy, overhead, detection speed, and threat model evaluation.

We highlight how the inadequate evaluation of these criteria compromises the protection provided by detection methods. Initially, we show the importance of conducting accurate evaluation of detection accuracy to ensure effective protection. Additionally, attention should be given to prevent excessive false positives, which can render detection methods impractical and diminish their adoption.

Furthermore, we demonstrate the importance of appropriately evaluating the overhead. We demonstrate how an improper setup of benchmark applications can result in unfairly low overhead evaluations. Precise reporting of overhead is crucial since unexpectedly high overhead can impede the adoption and practicality of detection methods.

Additionally, we underscore the importance of properly evaluating detection speed and highlight the discrepancy between the sampling interval of HPC data collection and the detection speed of a detection method. Providing accurate information about detection speed enables users to make informed decisions about the suitability of detection methods in safeguarding their systems against specific threats.

Finally, we illustrate how a weak threat model can lead to an overestima-

tion of the effectiveness of detection methods. To illustrate this, we performed an assessment of three cache side-channel attack detection methods. During our evaluation, we tested their ability to detect proof-of-concept Spectre and Flush+Reload attacks and found that they successfully identified these attacks with high accuracy. However, when we subjected the detection methods to our camouflaged attacks, their effectiveness was significantly compromised. Specifically, our camouflaged Flush+Reload attack successfully extracted keys from GnuPG, while remaining undetected. Based on these findings, we propose that authors should acknowledge the limitations of their detection methods when they fail to identify attacks under a stronger threat model. Transparent disclosure of such limitations is crucial for users to avoid unexpected compromises in their protection due to a lack of information.

In conclusion, we find that HPC-based cache side-channel attack detection methods still have a long way to go before they can be considered practical and widely applicable. We conclude that without addressing the aforementioned evaluation shortcomings, it remains uncertain whether these detection methods can truly be deemed effective for deployment in real-world scenarios.

Chapter 6

Future Directions

The findings of this thesis provide valuable insights into cache side-channel attacks, their mitigation techniques, and detection methods. However, there are several avenues for further research and development to enhance the security of computing devices against such attacks. In this section, we outline potential future works that could be explored to address the remaining challenges and limitations identified in our study.

In order to build upon the work presented in Chapter 3, various enhancements can be explored. In that chapter, we show the efficacy of `CacheFX` for evaluating cache designs' security and shared the insightful discoveries made using `CacheFX`. Nonetheless, we recognize that the simulation model used in this chapter is relatively simplistic and could be refined in several ways. First, integrating cache hierarchies into the analysis would allow for a more thorough examination of diverse cache designs. Second, the study could benefit from the inclusion of noise models that account for both systematic and random noise generated by system activity, resulting in a more realistic assessment since real cache attacks often involve such noises. Furthermore, considering the complexity of real-world attacks, it would be advantageous to model more intricate scenarios in future research. Finally, expanding the scope of the investigation to encompass cache performance evaluation, in addition to cache security evaluation, would be a valuable supplement to the existing framework.

The work in Chapter 4 can be enhanced in several ways. In that chapter, we demonstrate the feasibility of constructing intricate circuits using our speculative gates mainly on Intel architecture. While theoretically possible to create these circuits on other architectures, such as ARM, this issue remains open for future research. However, we can highlight some challenges we encountered when attempting to implement these circuits on the ARM architecture. The first challenge was the lack of user-accessible high-resolution timers on the Macbook Air system we performed our experimentation on. A potential solution could involve using alternative techniques to obtain more accurate timing measurements, such as setting specific MSRs to enable more precise timing. Although ARM devices possess these MSRs, we could not use them

due to limitations imposed by the operating system (MacOS). An alternative OS might permit access to this feature. The second issue was the absence of a precise cache flushing mechanism that operates at cache-line granularity, crucial for achieving high circuit execution rates. In the absence of this feature, the entire cache must be cleared before each circuit execution. Another option is to locate an eviction set that maps to the desired cache line. However, this method lacks precision and may result in the eviction of other critical cache states. We note that both approaches are relatively time consuming, and therefore constrain the circuit's execution rate. A potential solution could involve modifying the operating system to enable such functionality, but this remains a topic for future exploration.

In order to enhance the work presented in Chapter 5, several refinements can be considered. In that chapter, we identify the shortcomings of cache side-channel attack detection methods and provided experimental evidence to support our observations. Through experimentation, we demonstrate that these detection methods cannot effectively differentiate between camouflaged attacks and benign application execution when malicious behavior is embedded. This conclusion was reached by examining three detection methods and their inability to accurately identify the attacks. However, this argument could be further supported with statistical evidence, such as gathering HPC data and applying equivalence tests such as the two one-sided t -test (TOST) to compare benign traces and attack-injected benign traces, proving their similarity. Additionally, the detection method developed in this chapter could be improved by incorporating time-series classifiers such as Pearson correlation and dynamic time warping (DTW).

By addressing these future research directions, the security community can continue to advance the understanding of cache side-channel attacks and develop more effective mitigation and detection techniques. These efforts will ultimately contribute to the creation of more secure computing environments, protecting sensitive data and ensuring the integrity of our digital infrastructure.

Chapter 7

Conclusion

Throughout this thesis, we explore the critical topic of cache side-channel attacks, which pose a significant threat to the confidentiality and integrity of information in modern computing devices. Our primary focus is on investigating and evaluating the effectiveness of secure cache designs, mitigation techniques, and detection methods for these attacks. We demonstrate that while some progress has been made in securing computing devices against cache side-channel attacks, there are still numerous challenges and limitations that need to be addressed.

Our first contribution is a software framework to evaluate the security of secure cache designs. We show that secure caches can effectively mitigate eviction-set based attacks but are vulnerable to occupancy-based attacks. This finding highlights the need for further research into novel cache designs that provide better resistance to both types of attacks.

Our second contribution is a technique that leverages speculative execution to bypass low-resolution timers, a common mitigation method in browsers. Our technique demonstrates that speculative execution not only can be used for enabling high-resolution cache side-channel attacks on low-resolution timers, but also possesses Turing-complete capabilities for performing robust calculations on cache states. By revealing this new attack vector on low-resolution timers, we aim to raise awareness of this vulnerability and prompt the implementation of appropriate measures to address the issue and enhance the security of systems in the future.

Our final contribution is a survey of existing hardware performance counter-based cache side-channel attack detection methods. We highlight limitations in their evaluation of accuracy, overhead, detection speed, and threat model, and provided experimental results to support our verdict. Furthermore, we demonstrate a new camouflaged Flush+Reload that is able to bypass these detection methods and leak GnuPG's ElGamal private key within a reasonable timeframe without triggering any detection. Based on our findings, we conclude that HPC-based cache side-channel attack detection methods are still far from being practical and widely applicable. We suggest that unless

the evaluation shortcomings we highlighted are addressed, it remains uncertain whether these detection methods can genuinely be considered effective for real-world deployment.

In conclusion, this thesis provides valuable insights into the current state of cache side-channel attack research and identifies key areas for future work to enhance the security of computing devices against these stealthy and hard-to-detect attacks. By further exploring novel cache designs, improving existing mitigation techniques, and developing advanced detection methods, the security research community can continue to make strides in safeguarding our digital infrastructure and protecting sensitive information from cache side-channel attacks.

References

- Onur Aciıçmez (2007). “Yet Another Microarchitectural Attack: Exploiting I-Cache.” In: *CSAW*, pp. 11–18.
- Onur Aciıçmez, Çetin Kaya Koç, and Jean-Pierre Seifert (2007). “Predicting Secret Keys Via Branch Prediction.” In: *CT-RSA*, pp. 225–242.
- Onur Aciıçmez and Jean-Pierre Seifert (2007). “Cheap Hardware Parallelism Implies Cheap Security.” In: *FDTC*, pp. 80–91.
- Mackenzie Adams (2017). “Big Data and Individual Privacy in the Age of The Internet of Things.” In: *Technology Innovation Management Review 7.4*.
- Ayush Agarwal et al. (2022). “Spook.js: Attacking Chrome Strict Site Isolation via Speculative Execution.” In: *IEEE SP*.
- Bilal A. Ahmad (2019). “Detecting Spectre and Meltdown Attacks Using Hardware Performance Counters and Machine Learning.” PhD thesis. University of the Punjab.
- (2020). “Real Time Detection of Spectre and Meltdown Attacks Using Machine Learning.” In: *arXiv preprint arXiv:2006.01442*.
- Ayaz Akram et al. (2020). “Meet the Sherlock Holmes’ of Side Channel Leakage: A Survey of Cache SCA Detection Techniques.” In: *IEEE Access* 8, pp. 70836–70860.
- Manaar Alam, Sarani Bhattacharya, and Debdeep Mukhopadhyay (2021). “Victims Can Be Saviors: A Machine Learning-Based Detection for Micro-Architectural Side-Channel Attacks.” In: *ACM J. Emerg. Technol. Comput. Syst.* 17.2, 14:1–14:31.
- Manaar Alam et al. (2017). “Performance Counters to Rescue: A Machine Learning Based Safeguard Against Micro-Architectural Side-Channel-Attacks.” In: *IACR Cryptol. ePrint Arch.*, p. 564.
- Abdullah Albalawi, Vassilios G. Vassilakis, and Radu Calinescu (2022). “Protecting Shared Virtualized Environments Against Cache Side-channel Attacks.” In: *ICISSP*, pp. 507–514.
- Alejandro Cabrera Aldaya et al. (2019a). “Cache-Timing Attacks on RSA Key Generation.” In: *CHES 2019.4*, pp. 213–242.
- Alejandro Cabrera Aldaya et al. (2019b). “Port Contention for Fun and Profit.” In: *IEEE SP*, pp. 870–887.
- Zirak Allaf, Mo Adda, and Alexander E. Gegov (2019). “Malicious Loop Detection Using Support Vector Machine.” In: *INISTA*, pp. 1–6.

- Apple (2018). *About Speculative Execution Vulnerabilities in Arm-Based and Intel CPUs*. URL: <https://support.apple.com/en-us/HT208394>.
- Gilles Barthe et al. (2014). “System-Level non-Interference for Constant-Time Cryptography.” In: *CCS*, pp. 1267–1279.
- Mohammad-Mahdi Bazm et al. (2018). “Cache-based Side-Channel Attacks Detection through Intel Cache Monitoring Technology and Hardware Performance Counters.” In: *FMEC*, pp. 7–12.
- Mohammad Behnia et al. (2021). “Speculative Interference Attacks: Breaking Invisible Speculation Schemes.” In: *ASPLOS*, pp. 1046–1060. DOI: 10.1145/3445814.3446708.
- Daniel J. Bernstein, Tanja Lange, and Peter Schwabe (2012). “The Security Impact of a New Cryptographic Library.” In: *LATINCRYPT*. Vol. 7533. Springer, pp. 159–176. DOI: 10.1007/978-3-642-33481-8_9.
- Daniel J. Bernstein et al. (2017). “Sliding Right into Disaster: Left-to-Right Sliding Windows Leak.” In: *CHES*, pp. 555–576.
- Sarani Bhattacharya and Debdeep Mukhopadhyay (2015). “Who Watches the Watchmen?: Utilizing Performance Monitors for Compromising Keys of RSA on Intel Platforms.” In: *CHES*, pp. 248–266.
- Nicholas Biddle, Matthew Gray, and Steven McEachern (2022). “Public Exposure and Responses to Data Breaches in Australia: October 2022.” In: Emily Booth (2022). “What Comes First, a Breach of The Law or a Data Breach? In the Wake of Optus, Medibank and Other Recent High-Profile Data Breaches, What do In-House Lawyers Need to Know About How to Comply With Privacy Obligations?” In: *PRIVACY LAW BULLETIN* 19.8, pp. 160–162.
- Thomas Bourgeat et al. (2020). “CaSA: End-to-end Quantitative Security Analysis of Randomly Mapped Caches.” In: *MICRO 2020*, pp. 1110–1123.
- Ferdinand Brasser et al. (2017). “Software Grand Exposure: SGX Cache Attacks Are Practical.” In: *WOOT*.
- Samira Briongos et al. (2016). “Modeling Side-Channel Cache Attacks on AES.” In: *Proceedings of the Summer Computer Simulation Conference*, pp. 1–8.
- Samira Briongos et al. (2018). “CacheShield: Detecting Cache Attacks through Self-Observation.” In: *CODASPY*, pp. 224–235.
- Billy Bob Brumley (2015). “Cache Storage Attacks.” In: *CT-RSA*, pp. 22–34.
- Carole Cadwalladr and Emma Graham-Harrison (2018). “Revealed: 50 Million Facebook Profiles Harvested for Cambridge Analytica in Major Data Breach.” In: *The guardian* 17.1, p. 22.
- Claudio Canella et al. (2019a). “A Systematic Evaluation of Transient Execution Attacks and Defenses.” In: *USENIX Security*, pp. 249–266.
- Claudio Canella et al. (2019b). “Fallout: Leaking Data on Meltdown-resistant CPUs.” In: *CCS*, pp. 769–784.
- Stefano Carnà et al. (2022). “Fight Hardware with Hardware: System-Wide Detection and Mitigation of Side-Channel Attacks Using Performance Counters.” In: *Digital Threats: Research and Practice*.

- Chandler Carruth (2018). *Speculative Load Hardening*. <https://l1vm.org/docs/SpeculativeLoadHardening.html>.
- Guoxing Chen et al. (2019). “SgxPectre: Stealing Intel Secrets from SGX Enclaves Via Speculative Execution.” In: *IEEE EuroSP*, pp. 142–157.
- Marco Chiappetta, Erkay Savas, and Cemal Yilmaz (2016). “Real Time Detection of Cache-Based Side-Channel Attacks Using Hardware Performance Counters.” In: *Appl. Soft Comput.* 49, pp. 1162–1174.
- Jonghyeon Cho et al. (2020). “Real-Time Detection for Cache Side Channel Attack Using Performance Counter Monitor.” In: *Applied Sciences* 10.3, p. 984.
- Amit Choudhari, Sylvain Guilley, and Khaled Karray (2022). “SpecDefender: Transient Execution Attack Defender using Performance Counters.” In: *ASHES*, pp. 15–24. DOI: 10.1145/3560834.3563830. URL: <https://doi.org/10.1145/3560834.3563830>.
- Munish Chouhan and Hasbullah Halabi (2016). “Adaptive Detection Technique for Cache-Based Side Channel Attack using Bloom Filter for Secure Cloud.” In: *ICCOINS*, pp. 293–297.
- Chromium (n.d.). *Mitigating Side-Channel Attacks*. <https://www.chromium.org/Home/chromium-security/ssca/>. Accessed: 2022-01-25.
- David Cock et al. (2014). “The Last Mile: An Empirical Study of Timing Channels on seL4.” In: *CCS*, pp. 570–581.
- Patrick Cronin et al. (2021). “An Exploration of ARM System-Level Cache and GPU Side Channels.” In: *ACSAC*, pp. 784–795.
- Crozone (2023). *Proof of Concept Code for The Spectre CPU Exploit*. URL: <https://github.com/crozone/SpectrePoC>.
- Fergus Dall et al. (2018). “CacheQuote: Efficiently Recovering Long-Term Secrets of SGX EPID via Cache Attacks.” In: *IACR Trans. Cryptogr. Hardw. Embed. Syst.* 2018.2, pp. 171–191.
- Sanjeev Das et al. (2019). “SoK: The Challenges, Pitfalls, and Perils of Using Hardware Performance Counters for Security.” In: *IEEE (SP)*, pp. 20–38.
- John Demme and Simha Sethumadhavan (2014). “Side-Channel Vulnerability Metrics: SVF vs. CSV.” In: *Workshop on Duplicating, Deconstructing and Debunking*.
- John Demme et al. (2012). “Side-Channel Vulnerability Factor: A Metric for Measuring Information Leakage.” In: *ISCA*, pp. 106–117.
- (2013). “A Quantitative, Experimental Approach to Measuring Processor Side-Channel Security.” In: *IEEE Micro* 33.3, pp. 68–77.
- Shuwen Deng, Wenjie Xiong, and Jakub Szefer (2019). “Analysis of Secure Caches Using a Three-Step Model for Timing-Based Attacks.” In: *J. Hardware and Systems Security* 3.4, pp. 397–425.
- (2020). “A Benchmark Suite for Evaluating Caches’ Vulnerability to Timing Attacks.” In: *ASPLOS*. ACM, pp. 683–697.
- Shuwen Deng et al. (2021). *Evaluation of Cache Attacks on ARM Processors and Secure Caches*. arXiv 2106.14054.

- Peter J. Denning (May 1968). “The working set model for program behavior.” In: *Communications of the ACM* 11.5, pp. 323–333.
- Jonas Depoix and Philipp Altmeyer (2018). “Detecting Spectre Attacks by Identifying Cache Side-Channel Attacks Using Machine Learning.” In: *Advanced Microkernel Operating Systems* 75.
- Jean-Francois Dhem et al. (2000). “A Practical Implementation of the Timing Attack.” In: *Smart Card Research and Applications: Third International Conference*. Springer, pp. 167–182.
- Craig Disselkoen et al. (2017). “Prime+Abort: A Timer-Free High-Precision L3 Cache Attack using Intel TSX.” In: *USENIX Security*, pp. 51–67.
- Leonid Domnitsler, Nael B. Abu-Ghazaleh, and Dmitry Ponomarev (2010). “A Predictive Model for Cache-Based Side Channels in Multicore and Multi-threaded Microprocessors.” In: *MMM-ACNS*, pp. 70–85.
- Leonid Domnitsler et al. (2012). “Non-Monopolizable Caches: Low-Complexity Mitigation of Cache Side Channel Attacks.” In: *TACO* 8.4, 35:1–35:21.
- Goran Doychev et al. (2013). “CacheAudit: A Tool for the Static Analysis of Cache Side Channels.” In: *USENIX Security*, pp. 431–446.
- Swastika Dutta and Sayan Sinha (2019). “Performance Statistics and Learning Based Detection of Exploitative Speculative Attacks.” In: *CF*, pp. 206–210.
- Catherine Easdon et al. (2022). “Rapid Prototyping for Microarchitectural Attacks.” In: *USENIX Security*, pp. 3861–3877.
- Taher El Gamal (1984). “A Public Key Cryptosystem and a Signature Scheme Based on Discrete Logarithms.” In: *CRYPTO*, pp. 10–18.
- Fürkan Elibol, Uğur Sarac, and Işin Erer (2012). “Realistic Eavesdropping Attacks on Computer Displays with Low-Cost and Mobile Receiver System.” In: *European Signal Processing Conference (EUSIPCO)*. IEEE, pp. 1767–1771.
- Olakunle Elijah et al. (2018). “An Overview of Internet of Things (IoT) and Data Analytics in Agriculture: Benefits and Challenges.” In: *IEEE Internet of things Journal* 5.5, pp. 3758–3773.
- Dmitry Evtuyushkin, Dmitry V. Ponomarev, and Nael B. Abu-Ghazaleh (2016). “Jump Over ASLR: Attacking Branch Predictors to Bypass ASLR.” In: *MICRO*, 40:1–40:13.
- Dmitry Evtuyushkin et al. (2018). “BranchScope: A New Side-Channel Attack on Directional Branch Predictor.” In: *ASPLOS*, pp. 693–707.
- Dmitry Evtuyushkin et al. (2021). “Computing with Time: Microarchitectural Weird Machines.” In: *ASPLOS*, pp. 758–772.
- Alessandro Fedele and Cristian Roner (2022). “Dangerous Games: A Literature Review on Cybersecurity Investments.” In: *Journal of Economic Surveys* 36.1, pp. 157–187.
- Yusi Feng et al. (2021). “Constant-Time Loading: Modifying CPU Pipeline to Defeat Cache Side-Channel Attacks.” In: *TrustCom*, pp. 1132–1140.
- Serena Ferracci (2019). “Detecting Cache-based Side Channel Attacks using Hardware Performance Counters.” PhD thesis. Sapienza, University of Rome.

- Game of Life Wiki* (n.d.). URL: <https://conwaylife.com/wiki/T-tetromino>.
- Martin Gardner (1970). “Mathematical Games: the Fantastic Combinations of John Conway’s New Solitaire Game “life”.” In: *Sci. Am.* 223, pp. 120–123.
- Qian Ge, Yuval Yarom, and Gernot Heiser (2018). “No Security Without Time Protection: We Need a New Hardware-Software Contract.” In: *APSys*, 1:1–1:9.
- Qian Ge et al. (2019). “Time Protection: The Missing OS Abstraction.” In: *EuroSys*, 1:1–1:17.
- Daniel Genkin, Adi Shamir, and Eran Tromer (2014). “RSA Key Extraction via Low-Bandwidth Acoustic Cryptanalysis.” In: *CRYPTO*. Springer, pp. 444–461.
- Daniel Genkin et al. (2018). “Drive-By Key-Extraction Cache Attacks from Portable Code.” In: *ACNS*, pp. 83–102.
- Daniel Genkin et al. (2020). “Cache vs. Key-Dependency: Side Channeling an Implementation of Pilsung.” In: *IACR Trans. Cryptogr. Hardw. Embed. Syst.* 2020.1, pp. 231–255.
- Daniel Genkin et al. (2022). *CacheFX: A Framework for Evaluating Cache Security*. arXiv/2201.11377.
- Tara Ghasempouri et al. (2020). “A Security Verification Template to Assess Cache Architecture Vulnerabilities.” In: *DDECS*, pp. 1–6.
- Google (2021). *Spectre*. <https://leaky.page>.
- Daniel M. Gordon (1998). “A Survey of Fast Exponentiation Methods.” In: *Journal of Algorithms* 27.1, pp. 129–146.
- Ben Gras et al. (2017). “ASLR on the Line: Practical Cache Attacks on the MMU.” In: *NDSS*.
- Ben Gras et al. (2018). “Translation Leak-Aside Buffer: Defeating Cache Side-Channel Protections with TLB Attacks.” In: *USENIX Security*, pp. 955–972.
- Marc Green et al. (2017). “AutoLock: Why Cache Attacks on ARM Are Harder Than You Think.” In: *USENIX Security*, pp. 1075–1091.
- Nick Gregory and Kannan Harini (2021). “Using Undocumented Hardware Performance Counters to Detect Spectre-Style Attacks.” In.
- Leon Groot Bruinderink et al. (2016). “Flush, Gauss, and Reload - A Cache Attack on the BLISS Lattice-Based Signature Scheme.” In: *CHES*, pp. 323–345.
- Daniel Gruss, Raphael Spreitzer, and Stefan Mangard (2015). “Cache Template Attacks: Automating Attacks on Inclusive Last-Level Caches.” In: *USENIX Security*, pp. 897–912.
- Daniel Gruss et al. (2016). “Flush+Flush: A Fast and Stealthy Cache Attack.” In: *DIMVA*, pp. 279–299.
- David Gullasch, Endre Bangerter, and Stephan Krenn (2011). “Cache Games - Bringing Access-Based Cache Attacks on AES to Practice.” In: *IEEE SP*, pp. 490–505.
- Berk Gülmezoglu et al. (2019a). “FortuneTeller: Predicting Microarchitectural Attacks via Unsupervised Deep Learning.” In: *CoRR* abs/1907.03651.

- Berk Gülmezoglu et al. (2019b). “Undermining User Privacy on Mobile Devices Using AI.” In: *AsiaCCS*, pp. 214–227.
- Noam Hadad and Jonathan Afek (2018). *Overcoming (some) Spectre Browser Mitigations*. <https://alephsecurity.com/2018/06/26/spectre-browser-query-cache/>. Accessed: 2022-01-25.
- Ameer Hamza et al. (2021). “Diminisher: A Linux Kernel Based Countermeasure for TAA Vulnerability.” In: *ESORICS*, pp. 477–495.
- John Hazen (2018). *Mitigating Speculative Execution Side-Channel Attacks in Microsoft Edge and Internet Explorer*. <https://blogs.windows.com/msedgedev/2018/01/03/speculative-execution-mitigations-microsoft-edge-internet-explorer/>. Accessed: 2022-01-25.
- Zecheng He and Ruby B. Lee (2017). “How Secure is Your Cache Against Side-Channel Attacks?” In: *MICRO*, pp. 341–353.
- Zhang Hongxin et al. (2009). “Recognition of Electro-Magnetic Leakage Information from Computer Radiation with SVM.” In: *Computers & Security* 28.1-2, pp. 72–76.
- Ralf Hund, Carsten Willems, and Thorsten Holz (2013). “Practical Timing Side Channel Attacks Against Kernel Space ASLR.” In: *NDSS*.
- IBM (2022). *Cost of a Data Breach Report 2022*. URL: <https://www.ibm.com/downloads/cas/3R8N1DZJ>.
- Mehmet Sinan Inci et al. (2016). “Cache Attacks Enable Bulk Key Recovery on the Cloud.” In: *CHES*, pp. 368–388.
- Intel (Dec. 2021a). *Intel 64 and IA-32 Architectures Software Developer’s Manual Volume 1: Basic Architecture*. <https://cdrdv2.intel.com/v1/dl/getContent/671436>.
- (2021b). *Performance Monitoring Impact of Intel Transactional Synchronization Extension Memory Ordering Issue*. <https://www.intel.com/content/dam/support/us/en/documents/processors/Performance-Monitoring-Impact-of-TSX-Memory-Ordering-Issue-604224.pdf>.
- Gorka Irazoqui Apecechea, Thomas Eisenbarth, and Berk Sunar (2015). “S\$A: A Shared Cache Attack That Works across Cores and Defies VM Sandboxing - and Its Application to AES.” In: *IEEE SP*, pp. 591–604.
- Himanshi Jain, D. Anthony Balaraju, and Chester Rebeiro (2019). “Spy Cartel: Parallelizing Evict+Time-Based Cache Attacks on Last-Level Caches.” In: *J. Hardw. Syst. Secur.* 3.2, pp. 147–163.
- Aamer Jaleel et al. (2010). “High Performance Cache Replacement Using Reference Interval Prediction (RRIP).” In: *ISCA*. ACM, pp. 60–71. DOI: 10.1145/1815961.1815971.
- Jianyu Jiang, Claudio Soriente, and Ghassan Karame (2022). “On the Challenges of Detecting Side-Channel Attacks in SGX.” In: *RAID*, pp. 86–98.
- David A. Kaplan (2023). *Optimization and Amplification of Cache Side Channel Signals*. arXiv/2303.00122. AMD. DOI: 10.48550/arXiv.2303.00122.
- Sowmya Karunakaran et al. (2018). “Data Breaches: User Comprehension, Expectations, and Concerns with Handling Exposed Data.” In: *SOUPS*, pp. 217–234.

- Daniel Katzman et al. (2023). “The Gates of Time: Improving Cache Attacks with Transient Execution.” In: *USENIX Security*.
- Mehmet Kayaalp et al. (2017). “RIC: Relaxed Inclusion Caches for Mitigating LLC Side-Channel Attacks.” In: *DAC*, 7:1–7:6.
- Khaled N. Khasawneh et al. (2019). “SafeSpec: Banishing the Spectre of a Meltdown with Leakage-Free Speculation.” In: *DAC*, p. 60.
- Hodong Kim, Changhee Hahn, and Junbeom Hur (2021). “Real-Time Detection of Cache Side-channel Attack Using Non-cache Hardware Events.” In: *ICOIN*, pp. 28–31.
- Man Ho Kim, Suk Lee, and Kyung Chang Lee (2009). “Kalman Predictive Redundancy System for Fault Tolerance of Safety-Critical Systems.” In: *IEEE Transactions on Industrial Informatics* 6.1, pp. 46–53.
- Ofek Kirzner and Adam Morrison (2021). “An Analysis of Speculative Type Confusion Vulnerabilities in the Wild.” In: *USENIX Security*, pp. 2399–2416.
- Paul Kocher, Joshua Jaffe, and Benjamin Jun (1999). “Differential Power Analysis.” In: *CRYPTO*. Springer, pp. 388–397.
- Paul Kocher et al. (2019). “Spectre Attacks: Exploiting Speculative Execution.” In: *IEEE SP*, pp. 1–19.
- Paul C Kocher (1996). “Timing Attacks on Implementations of Diffie-Hellman, RSA, DSS, and other systems.” In: *CRYPTO*. Springer, pp. 104–113.
- David Kohlbrenner and Hovav Shacham (2016). “Trusted Browsers for Uncertain Times.” In: *USENIX Security*, pp. 463–480.
- Boris Köpf, Laurent Mauborgne, and Martín Ochoa (2012). “Automatic Quantification of Cache Side-Channels.” In: *Computer Aided Verification*.
- Jakob Koschel et al. (2020). “TagBleed: Breaking KASLR on the Isolated Kernel Address Space using Tagged TLBs.” In: *EuroS&P*, pp. 309–321.
- Yusuf Kulah et al. (2019). “SpyDetector: An Approach for Detecting Side-Channel Attacks at Runtime.” In: *Int. J. Inf. Sec.*, pp. 393–422.
- David Lantz (2021). *Detection of Side-Channel Attacks Targeting Intel SGX*.
- Judith A Levy and Rita Strombeck (2002). “Health Benefits and Risks of the Internet.” In: *Journal of medical systems* 26, pp. 495–510.
- Congmiao Li and Jean-Luc Gaudiot (2018). “Online Detection of Spectre Attacks Using Microarchitectural Traces from Performance Counters.” In: *SBAC-PAD*, pp. 25–28.
- (2019). “Detecting Malicious Attacks Exploiting Hardware Vulnerabilities Using Performance Counters.” In: *COMPSAC*, pp. 588–597.
- Moritz Lipp (2021). “Exploiting Microarchitectural Optimizations from Software.” In: *Diss., Graz University of Technology*.
- Moritz Lipp et al. (2018). “Meltdown: Reading Kernel Memory from User Space.” In: *USENIX Security*, pp. 973–990.
- Fangfei Liu et al. (2015). “Last-Level Cache Side-Channel Attacks are Practical.” In: *IEEE SP*, pp. 605–622.
- Fangfei Liu et al. (2016a). “CATalyst: Defeating Last-Level Cache Side Channel Attacks in Cloud computing.” In: *HPCA*, pp. 406–418.

- Fangfei Liu et al. (2016b). “Newcache: Secure Cache Architecture Thwarting Cache Side-Channel Attacks.” In: *IEEE Micro* 36.5, pp. 8–16.
- Kevin Loughlin et al. (2021). “Dolma: Securing Speculation with the Principle of Transient Non-Observability.” In: *USENIX Security*, pp. 1397–1414.
- Giorgi Maisuradze and Christian Rossow (2018). “ret2spec: Speculative Execution Using Return Stack Buffers.” In: *CCS*, pp. 2109–2122.
- Stefan Mangard (2003). “A Simple Power-Analysis (SPA) Attack on Implementations of the AES Key Expansion.” In: *Information Security and Cryptology—ICISC 2002: 5th International Conference Seoul, Korea, November 28–29, 2002 Revised Papers 5*. Springer, pp. 343–358.
- Robert Martin, John Demme, and Simha Sethumadhavan (2012). “Timewarp: Rethinking Timekeeping and Performance Monitoring Mechanisms to Mitigate Side-Channel Attacks.” In: *ISCA*, pp. 118–129.
- Clémentine Maurice et al. (2015a). “C5: Cross-Cores Cache Covert Channel.” In: *DIMVA*, pp. 46–64.
- Clémentine Maurice et al. (2015b). “Reverse Engineering Intel Last-Level Cache Complex Addressing Using Performance Counters.” In: *RAID*, pp. 48–65.
- John D. McCalpin (2021). *Mapping Addresses to L3/CHA Slices in Intel Processors*. Tech. rep. TR-2021-03. ACELab, The University of Texas at Austin.
- Ross McElroy et al. (2019). *Spectre is Here to Stay: An Analysis of Side-Channels and Speculative Execution*. arXiv/1902.05178.
- MDN Web Docs (2022). *Planned Changes to Shared Memory*. https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference/Global_Objects/SharedArrayBuffer/Planned_changes. Accessed: 2022-01-30.
- Ahmad Moghimi, Gorka Irazoqui, and Thomas Eisenbarth (2017). “CacheZoom: How SGX Amplifies the Power of Cache Attacks.” In: *CHES*, pp. 69–90.
- Maria Mushtaq et al. (2018a). “Cache-Based Side-Channel Intrusion Detection using Hardware Performance Counters.” In: *CryptArchi*.
- Maria Mushtaq et al. (2018b). “Machine Learning for Security: The Case of Side-Channel Attack Detection at Run-Time.” In: *ICECS*, pp. 485–488.
- Maria Mushtaq et al. (2018c). “NIGHTs-WATCH: A Cache-based Side-Channel Intrusion Detector Using Hardware Performance Counters.” In: *HASP*, 1:1–1:8.
- Maria Mushtaq et al. (2018d). “Run-Time Detection of Prime+Probe Side-Channel Attack on AES Encryption Algorithm.” In: *GIIS*, pp. 1–5.
- Maria Mushtaq et al. (2020). “WHISPER: A Tool for Run-Time Detection of Side-Channel Attacks.” In: *IEEE Access* 8, pp. 83871–83900.
- Maria Mushtaq et al. (2021). “Transit-Guard: An OS-Based Defense Mechanism Against Transient Execution Attacks.” In: *ETS*, pp. 1–2.
- Maria Mushtaq et al. (2022). “The Kingsguard OS-Level Mitigation Against Cache Side-Channel Attacks Using Runtime Detection.” In: *Ann. des Télécommunications*, pp. 731–747.

- National Institute of Standards and Technology (2015). *FIPS 180-4: Secure Hash Standard (SHS)*. <https://nvlpubs.nist.gov/nistpubs/fips/nist.fips.180-4.pdf>.
- Donald J. Newman (1960). “The Double Dixie Cup Problem.” In: *The American Mathematical Monthly* 67.1, pp. 58–61.
- Noam Nisan and Shimon Schocken (2021). *The Elements of Computing Systems: Building a Modern Computer from First Principles*. MIT press. ISBN: 9780262539807.
- Yossef Oren et al. (2015). “The Spy in the Sandbox: Practical Cache Attacks in JavaScript and their Implications.” In: *CCS*, pp. 1406–1418.
- Dag Arne Osvik, Adi Shamir, and Eran Tromer (2006). “Cache Attacks and Countermeasures: The Case of AES.” In: *CT-RSA*, pp. 1–20.
- Samuel O’Malley and Kim-Kwang Raymond Choo (2014). “Bridging the Air Gap: Inaudible Data Exfiltration by Insiders.” In: *Americas Conference on Information Systems*, pp. 7–10.
- Dan Page (2002). *Theoretical Use of Cache Memory as a Cryptanalytic Side-Channel*. Cryptology ePrint Archive 2002/169.
- Mathias Payer (2016). “HexPADS: A Platform to Detect “Stealth” Attacks.” In: *ESSoS*, pp. 138–154.
- Colin Percival (2005). “Cache Missing for Fun and Profit.” In: *Proceedings of BSDCan*. URL: <https://www.daemonology.net/papers/htt.pdf>.
- Cesar Pereida García and Billy Bob Brumley (2017). “Constant-Time Callees with Variable-Time Callers.” In: *USENIX Security*, pp. 83–98.
- Peter Pessl, Leon Groot Bruinderink, and Yuval Yarom (2017). “To BLISS-B or not to be: Attacking strongSwan’s Implementation of Post-Quantum Signatures.” In: *CCS*, pp. 1843–1855.
- Ravi Pilla, Taiwo Oseni, and Andrew Stranieri (2023). “A Study Into the Impact of Data Breaches of Electronic Health Records.” In: *Proceedings of the 2023 Australasian Computer Science Week*, pp. 252–254.
- Nikolaos Foivos Polychronou et al. (2021). “MaDMAN: Detection of Software Attacks Targeting Hardware Vulnerabilities.” In: *DSD*, pp. 355–362.
- Iván Prada, Francisco D. Igual, and Katalin Olcoz (2019). “Detecting Time-Fragmented Cache Attacks Against AES Using Performance Monitoring Counters.” In: *JCC&BD*, pp. 3–15.
- Antoon Purnal, Furkan Turan, and Ingrid Verbauwhede (2021a). “Prime+Scope: Overcoming the Observer Effect for High-Precision Cache Contention Attacks.” In: *CCS*, pp. 2906–2920.
- Antoon Purnal et al. (2021b). “Systematic Analysis of Randomization-based Protected Cache Architectures.” In: *IEEE SP*, pp. 987–1002.
- Moinuddin K. Qureshi (2018). “CEASER: Mitigating Conflict-Based Cache Attacks via Encrypted-Address and Remapping.” In: *MICRO*, pp. 775–787.
- (2019). “New Attacks and Defense for Encrypted-Address Cache.” In: *ISCA*, pp. 360–371.

- Hany Ragab et al. (2021a). “CrossTalk: Speculative Data Leaks Across Cores Are Real.” In: *IEEE S&P*, pp. 1852–1867.
- Hany Ragab et al. (2021b). “Rage Against the Machine Clear: A Systematic Analysis of Machine Clears and Their Implications for Transient Execution Attacks.” In: *USENIX Security*, pp. 1451–1468.
- Thomas Ristenpart et al. (2009). “Hey, You, Get Off of My Cloud: Exploring Information Leakage in Third-Party compute clouds.” In: *ACM CCS*, pp. 199–212.
- Thomas Rokicki, Clémentine Maurice, and Pierre Laperdrix (2021). “SoK: In Search of Lost Time: A Review of JavaScript Timers in Browsers.” In: *EuroS&P*, pp. 472–486. DOI: 10.1109/EuroSP51992.2021.00039.
- Eyal Ronen et al. (2019). “The 9 Lives of Bleichenbacher’s CAT: New Cache Attacks on TLS Implementations.” In: *IEEE SP*, pp. 435–452.
- Stephen Röttger and Artur Janc (2021). *A Spectre Proof-of-Concept for A Spectre-Proof Web*. <https://security.googleblog.com/2021/03/a-spectre-proof-of-concept-for-spectre.html>.
- Majid Sabbagh et al. (2018). “SCADET: A Side-Channel Attack Detection Tool for Tracking Prime+Probe.” In: *ICCAD*, p. 107.
- Stephan Van Schaik et al. (2019). “RIDL: Rogue In-Flight Data Load.” In: *IEEE S&P*, pp. 88–105.
- Stephan Van Schaik et al. (2020). *SGAxe: How SGX Fails in Practice*. <https://sgaxeattack.com/>.
- Stephan Van Schaik et al. (2021). “CacheOut: Leaking Data on Intel CPUs via Cache Evictions.” In: *IEEE S&P*, pp. 339–354.
- David Schor (n.d.). *Skylake (client)*. URL: [https://en.wikichip.org/wiki/intel/microarchitectures/skylake_\(client\)](https://en.wikichip.org/wiki/intel/microarchitectures/skylake_(client)).
- Michael Schwarz et al. (2017). “Fantastic Timers and Where to Find Them: High-Resolution Microarchitectural Attacks in JavaScript.” In: *Financial Cryptography*, pp. 247–267.
- Michael Schwarz et al. (2019). “ZombieLoad: Cross-Privilege-Boundary Data Sampling.” In: *CCS*, pp. 753–768.
- Martin Schwarzl et al. (2021). *Dynamic Process Isolation*. arXiv/2110.04751.
- André Seznec (1993). “A Case for Two-Way Skewed-Associative Caches.” In: *ISCA*, pp. 169–178.
- Jicheng Shi et al. (2011). “Limiting Cache-Based Side-Channel in Multi-Tenant Cloud Using Dynamic Page Coloring.” In: *DSN Workshops*, pp. 194–199.
- Anatoly Shusterman et al. (2019). “Robust Website Fingerprinting through the Cache Occupancy Channel.” In: *USENIX Security*, pp. 639–656.
- Anatoly Shusterman et al. (2021a). “Prime+Probe 1, JavaScript 0: Overcoming Browser-based Side-Channel Defenses.” In: *USENIX Security*, pp. 2863–2880.
- Anatoly Shusterman et al. (2021b). “Website Fingerprinting through the Cache Occupancy Channel and its Real World Practicality.” In: *IEEE Trans. Dependable Secur. Comput.* 18.5, pp. 2042–2060.

- Florian Sieck et al. (2021). “Util: : Lookup: Exploiting Key Decoding in Cryptographic Libraries.” In: *CCS*. ACM, pp. 2456–2473.
- Nikhilesh Singh and Chester Rebeiro (2021). “LEASH: Enhancing Micro-Architectural Attack Detection With A Reactive Process Scheduler.” In: *CoRR* abs/2109.03998.
- SPEC (2020). *gcc_r SPEC CPU®2017 Benchmark*. URL: https://www.spec.org/cpu2017/Docs/benchmarks/502.gcc_r.html.
- Julian Stecklina and Thomas Prescher (2018). *LazyFP: Leaking FPU Register State using Microarchitectural Side-Channels*. arXiv/1806.07480.
- Charles E Stroud and Ahmed E Barbour (1989). “Design for Testability and Test Generation for Static Redundancy System Level Fault-Tolerant Circuits.” In: *Proceedings. Meeting the Tests of Time’, International Test Conference*. IEEE, pp. 812–818.
- Qinhan Tan et al. (2020). “PhantomCache: Obfuscating Cache Conflicts with Localized Randomization.” In: *NDSS*.
- Xiaojie Tao et al. (2021). “SCAMS: A Novel Side-Channel Attack Mitigation System in IaaS Cloud.” In: *MILCOM*, pp. 329–334.
- Daniel Terpstra et al. (2009). “Collecting Performance Data with PAPI-C.” In: *International Workshop on Parallel Tools for High Performance Computing*. Springer, pp. 157–173. DOI: 10.1007/978-3-642-11261-4_11. URL: https://doi.org/10.1007/978-3-642-11261-4_11.
- Robert M Tomasulo (1967). “An Efficient Algorithm for Exploiting Multiple Arithmetic Units.” In: *IBM Journal of research and Development* 11.1, pp. 25–33.
- Zhongkai Tong et al. (2020). “Cache Side-channel Attacks Detection Based on Machine Learning.” In: *TrustCom*, pp. 919–926.
- Zhongkai Tong et al. (2022). “Attack Detection Based on Machine Learning Algorithms for Different Variants of Spectre Attacks and Different Melt-down Attack Implementations.” In: *CoRR* abs/2208.14062.
- Yukiyasu Tsunoo et al. (2002). “Cryptanalysis of Block Ciphers Implemented on Computers with Cache.” In: *ISITA*.
- Yukiyasu Tsunoo et al. (2003). “Cryptanalysis of DES Implemented on Computers with Cache.” In: *CHES*, pp. 62–76.
- Leif Uhsadel, Andy Georges, and Ingrid Verbauwhede (2008). “Exploiting Hardware Performance Counters.” In: *FDTC*, pp. 59–67.
- Jo Van Bulck et al. (2018). “Foreshadow: Extracting the Keys to the Intel SGX Kingdom with Transient Out-of-Order Execution.” In: *USENIX Security*, pp. 991–1008.
- Jo Van Bulck et al. (2020). “LVI: Hijacking Transient Execution through Microarchitectural Load Value Injection.” In: *IEEE S&P*, pp. 54–72.
- Stephan van Schaik et al. (2017). “RevAnC: A Framework for Reverse Engineering Hardware Page Table Caches.” In: *EUROSEC*, 3:1–3:6.
- R. Vanathi and Sp. Chokkalingam (2018). “Cache-Based Side Channel attack Discovery using Intelligent-Detection Algorithm for Securing the Cloud Computing Environment.” In.

- Pepe Vila, Boris Köpf, and José F. Morales (2019). “Theory and Practice of Finding Eviction Sets.” In: *IEEE SP*, pp. 39–54.
- Pepe Vila et al. (2020). *Flushgeist: Cache Leaks from Beyond the Flush*. arXiv 2005.13853.
- Krishnaswamy Viswanathan (2014). *Disclosure of Hardware Prefetcher Control on Some Intel® Processors*. <https://www.intel.com/content/www/us/en/develop/articles/disclosure-of-hw-prefetcher-control-on-some-intel-processors.html>.
- John Von Neumann (2016). “Probabilistic Logics and the Synthesis of Reliable Organisms from Unreliable Components.” In: *Automata Studies.(AM-34), Volume 34*. Princeton University Press, pp. 43–98.
- Luke Wagner (2018). *Mitigations Landing for New Class of Timing Attack*. <https://blog.mozilla.org/security/2018/01/03/mitigations-landing-new-class-timing-attack/>. Accessed: 2022-01-25.
- Han Wang et al. (2020a). “Phased-Guard: Multi-Phase Machine Learning Framework for Detection and Identification of Zero-Day Microarchitectural Side-Channel Attacks.” In: *ICCD*, pp. 648–655.
- Han Wang et al. (2020b). “SCARF: Detecting Side-Channel Attacks at Real-time using Low-level Hardware Features.” In: *IOLTS*, pp. 1–6.
- Han Wang et al. (2021a). “Evaluation of Machine Learning-Based Detection Against Side-Channel Attacks on Autonomous Vehicle.” In: *AICAS*. IEEE, pp. 1–4.
- Limin Wang, Lei Bu, and Fu Song (2022). “Locality Based Cache Side-Channel Attack Detection.” In: *International Workshop 87*.
- Limin Wang et al. (2019). “Colored Petri Net Based Cache Side Channel Vulnerability Evaluation.” In: *IEEE Access* 7, pp. 169825–169843.
- Wubing Wang et al. (2021b). “Specularizer: Detecting Speculative Execution Attacks via Performance Tracing.” In: *DIMVA*, pp. 151–172.
- Zhenghong Wang and Ruby B. Lee (2007). “New Cache Designs for Thwarting Software Cache-Based Side Channel Attacks.” In: *ISCA*, pp. 494–505.
- Mario Werner et al. (2019). “ScatterCache: Thwarting Cache Attacks via Cache Set Randomization.” In: *USENIX Security*, pp. 675–692.
- Spencer Wheatley, Thomas Maillart, and Didier Sornette (2016). “The Extreme Risk of Personal Data Breaches and The Erosion of Privacy.” In: *The European Physical Journal B* 89, pp. 1–12.
- Mike J. Wilkinson et al. (2017). “Replacing Sanger with Next Generation Sequencing to improve coverage and quality of reference DNA barcodes for plants.” In: *Scientific Reports* 7.1, p. 46040. DOI: 10.1038/srep46040.
- John C. Wray (1991). “An Analysis of Covert Timing Channels.” In: *IEEE SP*, pp. 2–7. DOI: 10.1109/RISP.1991.130767.
- Minjun Wu et al. (2022). “PREDATOR: A Cache Side-Channel Attack Detector Based on Precise Event Monitoring.” In: *IEEE SEED*, pp. 25–36.
- Wenjie Xiong and Jakub Szefer (2020). “Leaking Information through Cache LRU States.” In: *HPCA*, pp. 139–152.

- Wenjie Xiong and Jakub Szefer (2021). “Survey of Transient Execution Attacks and Their Mitigations.” In: *ACM Comput. Surv.* 54.3, 54:1–54:36.
- Hui Yan and Chaoyuan Cui (2022). “CacheHawkeye: Detecting Cache Side Channel Attacks Based on Memory Events.” In: *Future Internet* 14.1, p. 24.
- Mengjia Yan, Christopher W. Fletcher, and Josep Torrellas (2020). “Cache Telepathy: Leveraging Shared Resource Attacks to Learn DNN Architectures.” In: *USENIX Security*, pp. 2003–2020.
- Mengjia Yan et al. (2017). “Secure Hierarchy-Aware Cache Replacement Policy (SHARP): Defending Against Cache-Based Side Channel Attacks.” In: *ISCA*, pp. 347–360.
- Mengjia Yan et al. (2019a). “Attack Directories, Not Caches: Side Channel Attacks in a Non-Inclusive World.” In: *IEEE SP*, pp. 888–904.
- Mengjia Yan et al. (2019b). “SecDir: A Secure Directory to Defeat Directory Side-Channel Attacks.” In: *ISCA*, pp. 332–345.
- Yuval Yarom (2016). *Mastik: A Micro-Architectural Side-Channel Toolkit*.
- Yuval Yarom and Katrina Falkner (2014). “Flush+Reload: A High Resolution, Low Noise, L3 Cache Side-Channel Attack.” In: *USENIX Security*, pp. 719–732.
- Yuval Yarom, Daniel Genkin, and Nadia Heninger (2017). “CacheBleed: a Timing Attack on OpenSSL Constant-Time RSA.” In: *J. Cryptographic Engineering* 7.2, pp. 99–112.
- Yuval Yarom et al. (2015). *Mapping the Intel Last-Level Cache*. IACR Cryptology ePrint Archive, Report 2015/905.
- Tianwei Zhang and Ruby B. Lee (2014). “New Models of Cache Architectures Characterizing Information Leakage from cache side channels.” In: *ACSAC*. Ed. by Charles N. Payne Jr. et al., pp. 96–105.
- Tianwei Zhang, Yinqian Zhang, and Ruby B. Lee (2016). “CloudRadar: A Real-Time Side-Channel Attack Detection System in Clouds.” In: *RAID*, pp. 118–140.
- (2018). “Analyzing Cache Side Channels Using Deep Neural Networks.” In: *ACSAC*, pp. 174–186.
- Tianwei Zhang et al. (2013). “Side Channel Vulnerability Metrics: The Promise and the Pitfalls.” In: *HASP@ISCA*, p. 2.
- Yinqian Zhang et al. (2012). “Cross-VM Side Channels and Their Use to Extract Private Keys.” In: *CCS*, pp. 305–316.
- Zhiyi Zhang and Michael Grabchak (2014). “Nonparametric Estimation of Küllback-Leibler Divergence.” In: *Neural Comput.* 26.11, pp. 2570–2593.
- Zhiyuan Zhang et al. (2023). “Ultimate SLH: Taking Speculative Load Hardening to the Next Level.” In: *USENIX Security*.
- Beilei Zheng et al. (2022). “CBA-Detector: A Self-Feedback Detector Against Cache-Based Attacks.” In: *IEEE TDSC* 19.5, pp. 3231–3243.
- Boyoun Zhou et al. (2018). “Hardware Performance Counters can Detect Malware: Myth or Fact?” In: *AsiaCCS*, pp. 457–468.

Appendix A

Chapter 4

A.1 Hyper-tree Amplification Implementation

Because our hyper-tree amplification scheme is aimed for restricted environments, we needed to address the following main constraints:

1. We do not assume access to low-level instructions that are able to manipulate and set the cache state to uncached directly (e.g., `clflush`).
2. We have limited amount of memory that we can allocate for the amplification process (e.g., total allocated memory is smaller than Y^{2^l} cache lines but is larger than the LLC cache).
3. Our access pattern must handle the CPU's memory prefetcher mechanism.
4. The amplification process does not assume that we know how to find LLC eviction sets.

To accommodate these constraints, we develop an algorithm that allocates the addresses for each node in the hyper-tree from a large buffer. Our allocation algorithm results in an access pattern that fulfills the following conditions:

1. As we build each layer in the sub-tree in a breadth-first order, we ensure that the addresses we access in layer $i + 1$ do not evict any addresses in layer i from the cache. We achieve this through ensuring that all addresses in odd layers of a tree are mapped to different cache sets from the addresses in even layers of the tree.
2. As we generate the sub-trees in a depth-first manner, we ensure that none of the addresses in any of the lower sub-trees can evict addresses in the bank generated by the top tree.

3. Consecutive generations of a fixed number of sub-trees evict all addresses used in previous sub-trees from the cache with a very high probability. This *self-eviction* property allows us to reuse the same addresses for the next sub-trees we generate and limit the size of the buffer.
4. Any two consecutive accesses to memory are in different memory pages. This helps us not to trigger the memory prefetcher.

Although we do not know the mapping of the addresses in the buffer to LLC eviction sets, we use the fact that two addresses can be in the same eviction set only if their virtual page offset is the same. We look at each memory page in the buffer and partition possible offsets into three parts. The first part is used only for addresses in the bank generated by the top tree. The second part is used only by even layers in the lower sub-tree, and the third part for the odd layers.

We further assume we have a LLC cache with associativity X , and a total of N_{sets} cache sets or $N_{sets}/64$ cache sets for each 64-byte offset in a memory page. For a single offset, if we access $C \cdot X \cdot N_{sets}/64$ cache lines for some small constant C , with very a high probability, we obtain at least X addresses in each of the eviction sets [Newman, 1960] and evict all addresses previously accessed in this offset. We arrange the cache access pattern to ensure that for both odd and even layers of the tree, after a fixed number of sub-tree generated we access enough addresses in each offset to create such an offset eviction set. Now we can switch between two groups of addresses. After finish generating sub-trees using one group, with a very high probability, we evict all of the addresses in the second group and can reuse them.

A.2 Gates With and Without Branch Training

Listing A.1: Code for Not Gate (Need training)

```

_not_gate(int wet_run, void* addr_in, void* addr_out) {
    // "Stabilize" the branch predictor.
    for (int i = 0; i < 128; i++) {asm volatile("");}
    //Open Speculative window on addr_in
    if (wet_run == read_addr(addr_in)) {

```

```

    return 0;
}
//Return if we are in dry run/
if (!wet_run) (*@\label{f:NotGate:wet_run}@*)
    return 0;
//Slowdown to run longer than a
//speculative window on an address in cache
for (int i = 0; i < SPEC_SLOW_PARAM; i++) {
    nop;}
//Access output address only if
//speculative window is long enough
read_addr(addr_out);
return 0;
}

run_not_gate(void* addr_in, void* addr_out) {
    int train_in = 1;
//Run two dry runs to train branch prediction.
    _not_gate(0, &train_in, addr_out);
    _not_gate(0, &train_in, addr_out);
//Run gate
    _not_gate(1, addr_in, addr_out);
}

```

Listing A.2: Code for Not Gate (No training)

```

_nbt_not_gate(int counter, void* addr_in, void* addr_out) {
//Open Speculative window on switch case
    switch ((counter & 7) + *addr_in) {
        case 0x0: if(counter == 0x0) return 0; break;
        case 0x1: if(counter == 0x1) return 0; break;
        case 0x2: if(counter == 0x2) return 0; break;
        case 0x3: if(counter == 0x3) return 0; break;
        case 0x4: if(counter == 0x4) return 0; break;
        case 0x5: if(counter == 0x5) return 0; break;
        case 0x6: if(counter == 0x6) return 0; break;
        case 0x7: if(counter == 0x7) return 0; break;
    }
    volatile dummy = 0;
//Slowdown to run longer than a
//speculative window on an address in cache
    for (int i = 0; i < SPEC_SLOW_PARAM; i++)
        dummy *= dummy;
//Access output address only if
//speculative window is long enough
    read_addr(addr_out);
    return 0;
}

run_nbt_not_gate(void* addr_in, void* addr_out) {
    static int counter = 0;
//Run gate
    _nbt_not_gate(counter, addr_in, addr_out);
    counter++;
}

```

}

A.3 Gate Accuracy

We ran extensive test to validate the design and accuracy of our proposed gates. We ran the experiments on a Dynabook TECRA A50-EC, with a Intel Core i5-8250U CPU running Ubuntu 20.04.3 LTS. This CPU's LFB size and LLC associativity is 12. In our experiments we fixed the CPU's frequency to 1.6 GHz. Table A.1 shows the results of our experiments. "bt" stands for branch training, "nbt" for no branch training, and "fbd" for fixed branch delay. Most gates provide accuracy that is very close to 100%. We consider the output to be correct if *all* the state of all output addresses is correct.

For each gate, we ran more than 100 000 tests for each of the two possible output value '1' (cached) and '0' (uncached). For the *NAND* and *AND* gates, either all inputs were cached or all but one. For the *NOR* gate, either all inputs were uncached, or all but one. For the *NOT* and *BUFFER* gates, the single input was either cached or not.

Gate	Accuracy (percent)		Avg. run time (cycles)		Accuracy (percent)		Avg. run time (cycles)	
	'1'	'0'	'1'	'0'	'1'	'0'	'1'	'0'
btNAND ₁ ²	99.996	100.000	420	260	97.665	98.116	451	278
btNAND ₁ ¹¹	92.740	96.112	463	291	96.687	96.164	465	292
btNAND ₂ ²	100.000	100.000	417	261	99.988	98.101	439	281
btNAND ₂ ¹¹	99.984	95.801	466	291	99.998	96.111	448	289
btNAND ₂ ¹¹	99.209	99.996	411	249	98.034	98.104	449	279
btNAND ₁₁ ¹¹	94.996	95.697	456	291	94.194	95.861	460	293
btNOR ₁ ²	99.998	100.000	409	327	99.982	100.000	423	403
btNOR ₁ ¹¹	99.960	100.000	420	407	99.996	100.000	422	339
btNOR ₂ ⁸	99.994	100.000	415	395	99.998	100.000	400	242
btNOT ₁ ²	99.998	100.000	417	259	99.284	100.000	407	248
nbtNAND ₁ ²	99.952	100.000	223	64	99.974	98.210	251	88
nbtNAND ₁ ¹¹	99.516	95.626	268	94	99.861	96.161	284	98
nbtNAND ₂ ²	99.944	100.000	231	66	99.936	98.266	269	88
nbtNAND ₂ ¹¹	99.767	95.784	278	95	99.869	95.912	283	99
nbtNOT ₁ ¹	100.000	100.000	220	58	99.998	100.000	220	59
fbdAND ₁ ²	100.000	100.000	580	580	98.888	99.258	580	580
fbdAND ₁ ¹¹	96.995	99.202	583	580	96.204	100.000	579	579
fbdAND ₂ ²	100.000	99.998	580	580	98.946	99.293	579	579
fbdAND ₂ ¹¹	96.688	99.196	579	579	96.066	99.998	580	580
fbdAND ₁₁ ¹¹	99.996	100.000	580	580	98.966	99.283	583	579
fbdAND ₁₁ ¹¹	96.850	99.258	580	581	95.886	100.000	579	579
fbdBUFFER ₁ ¹	100.000	99.998	579	579	99.998	100.000	580	580
fbdBUFFER ₁₁ ¹	99.996	99.998	580	580	99.998	100.000	580	580

Table A.1: Accuracy of gates on Intel Core(TM) i5-8250U

A.4 Gates on Other Processors

We measure the accuracy of our gates on various devices including 1) HP laptop with AMD Ryzen 5 3500U running an up-to-date version of WSL and Windows 11 as shown in Table A.2. Table A.3 shows the accuracy of our gates on 2) 2020 Macbook Air with Apple M1 running MacOS Monterey 12.2 and 3) Samsung Galaxy S21 with Samsung Exynos 2100 running Android 11.

Some challenges associated with the ARM processors are 1) the lack of cacheline flushing mechanism. To overcome this problem, we were able to reliably flush the whole cache by accessing a large number of memory addresses. This method was sufficient to test the functionality of our gates. However, a better cache flushing mechanism is needed to build a running circuit to avoid corrupting wire states and to reduce runtime. 2) The second problem is the unavailability of high-resolution timer available in userspace. Our testing discovered that the timestamp counter accessible on both machines are running approximately 20-50 times slower than the CPU clock at maximum frequency. This would still allow for differentiating between cache hits and misses, albeit in lower resolution.

We perform the accuracy measurement of both the Exynos and M1 processors over a smaller number of experiments, mainly due to the slow throughput of the tests as a side effect of flushing the whole cache to reset wire states. Measurements for both processors are collected from 10 000 runs of each gates.

Gate	Accuracy (percent)		Avg. run time (cycles)		Accuracy (percent)		Avg. run time (cycles)	
	'1'	'0'	'1'	'0'	'1'	'0'	'1'	'0'
btNAND ₁ ²	99.610	98.447	817	522	96.821	98.341	853	555
btNAND ₁ ¹¹	99.613	98.246	824	524	99.459	100.000	872	572
btNAND ₂ ²	99.511	98.484	832	517	99.595	98.322	848	531
btNAND ₂ ¹¹	99.428	98.324	863	543	99.496	100.000	852	538
btNAND ₂ ¹¹	99.297	98.498	931	513	99.389	98.171	957	535
btNAND ₁₁ ¹¹	99.441	98.178	962	542	99.326	99.979	978	548
btNOR ₁ ²	99.119	98.880	813	785	97.476	99.508	895	862
btNOR ₁ ¹¹	98.458	99.634	931	909	96.782	98.704	831	786
btNOR ₂ ⁸	98.762	99.478	879	848	99.589	99.387	820	516
btNOT ₂ ¹	99.482	98.435	834	513	99.088	98.175	948	515
nbtNAND ₁ ²	95.201	98.466	486	167	94.280	98.427	521	179
nbtNAND ₁ ¹¹	96.965	98.218	526	176	95.231	99.993	529	177
nbtNAND ₂ ²	93.242	98.608	508	168	93.927	98.304	529	176
nbtNAND ₂ ¹¹	94.956	98.304	540	173	94.009	99.993	552	183
nbtNOT ₁ ¹	97.233	99.442	469	165	92.506	98.518	489	170
fbdAND ₁ ²	99.229	99.214	881	866	99.526	98.677	894	883
fbdAND ₁ ¹¹	99.421	98.742	904	883	98.693	99.945	909	892
fbdAND ₂ ²	99.443	99.026	904	878	99.173	98.783	911	903
fbdAND ₂ ¹¹	99.311	98.601	921	892	99.434	99.939	910	883
fbdAND ₁₁ ¹¹	99.192	99.110	1043	880	99.086	98.726	1040	876
fbdBUFFER ₁ ¹	99.334	98.719	1018	879	99.005	99.966	1012	883
fbdBUFFER ₁ ¹¹	99.411	99.485	871	873	99.289	98.857	899	868
fbdBUFFER ₁₁ ¹¹	99.475	98.908	1011	882				

Table A.2: Accuracy of gates on AMD Ryzen 5 3500U

Gate	Apple M1		Exynos 2100	
	Average (percent)	Median (percent)	Average (percent)	Median (percent)
bt $NAND_1^2$	87.19	91.00	91.70	94.00
bt NOR_1^2	80.01	83.00	81.10	89.00
bt NOT_1^1	85.23	86.00	91.47	96.00

Table A.3: Accuracy of gates on Apple M1 and Samsung Exynos 2100

A.5 Number of Cases in Gates Without Branch Training

As mentioned in Section 4.1.7, we test the accuracy and run time of gates without branch training using a different number of gates in each experiment. Recall that our compiler will not generate the required jump table if we use a small number of cases. To test gates with a small number of cases we compiled all gates with 16 cases but only used the required number of cases during run time. For example, to test a gate with 2 cases, we incremented our counter modulo 2, so only the case 0 and 1 were used.

As expected, gates with only 1 case do not work as this case is always used, and there is no misprediction that can open a speculative window. However, the rest of the gates seem to have similar accuracy and run time, so we conclude that the number of cases does not have a significant effect on the gates.

Table A.4 shows the accuracy and run time for nbt $NAND_1^{12}$ and nbt NOT_2^1 gates with different number of cases on Intel Core(TM) i5-8250U. Note that the run time is slightly longer than the results in Table A.1 due to the different code bases required to support the variable number of cases.

Cases	Gate	Accuracy (percent)		Avg. run time (cycles)		Accuracy (percent)		Avg. run time (cycles)	
		'1'	'0'	'1'	'0'	'1'	'0'	'1'	'0'
1	nbtNAND ₁ ¹²	0.0%	100.0%	278	95	0.0%	100.0%	216	61
2	nbtNAND ₁ ¹²	99.4%	95.6%	306	122	100.0%	100.0%	232	72
3	nbtNAND ₁ ¹²	98.8%	96.7%	302	119	100.0%	100.0%	236	78
4	nbtNAND ₁ ¹²	99.6%	97.0%	299	120	100.0%	100.0%	238	77
5	nbtNAND ₁ ¹²	99.7%	96.8%	297	117	100.0%	100.0%	237	79
6	nbtNAND ₁ ¹²	99.5%	96.7%	300	117	100.0%	100.0%	240	82
7	nbtNAND ₁ ¹²	99.7%	96.9%	308	120	100.0%	100.0%	237	81
8	nbtNAND ₁ ¹²	99.9%	96.7%	305	120	100.0%	100.0%	238	77
9	nbtNAND ₁ ¹²	99.6%	96.6%	305	116	100.0%	100.0%	241	79
10	nbtNAND ₁ ¹²	99.8%	96.8%	307	117	100.0%	100.0%	240	83
11	nbtNAND ₁ ¹²	99.8%	96.6%	307	119	100.0%	100.0%	239	80
12	nbtNAND ₁ ¹²	99.5%	95.9%	308	124	100.0%	100.0%	237	81
13	nbtNAND ₁ ¹²	99.8%	96.1%	305	118	100.0%	100.0%	236	76
14	nbtNAND ₁ ¹²	92.4%	96.7%	302	115	92.7%	100.0%	239	82
15	nbtNAND ₁ ¹²	93.0%	97.1%	299	115	93.3%	100.0%	240	81
16	nbtNAND ₁ ¹²	99.9%	96.8%	303	114	100.0%	100.0%	237	78

Table A.4: Accuracy and run time for nbtNAND₁¹² and nbtNOT₂¹ gates with different number of cases on Intel Core(TM) i5-8250U