

# Mitigating Rowhammer Attacks with Software Diversity

DIPLOMARBEIT

zur Erlangung des akademischen Grades

**Diplom-Ingenieur**

im Rahmen des Studiums

**Software Engineering & Internet Computing**

eingereicht von

**Manuel Wiesinger, BSc**

Matrikelnummer 00825632

an der Fakultät für Informatik

der Technischen Universität Wien

Betreuung: Privatdoz. Mag.rer.soc.oec. Dipl.-Ing. Dr. techn Edgar Weippl

Wien, 30. August 2018

---

Manuel Wiesinger

---

Edgar Weippl



# Mitigating Rowhammer Attacks with Software Diversity

DIPLOMA THESIS

submitted in partial fulfillment of the requirements for the degree of

**Diplom-Ingenieur**

in

**Software Engineering & Internet Computing**

by

**Manuel Wiesinger, BSc**

Registration Number 00825632

to the Faculty of Informatics

at the TU Wien

Advisor: Privatdoz. Mag.rer.soc.oec. Dipl.-Ing. Dr. techn Edgar Weippl

Vienna, 30<sup>th</sup> August, 2018

---

Manuel Wiesinger

---

Edgar Weippl



# Erklärung zur Verfassung der Arbeit

Manuel Wiesinger, BSc  
Goldschlagstraße 112/48  
1150 Wien

Hiermit erkläre ich, dass ich diese Arbeit selbständig verfasst habe, dass ich die verwendeten Quellen und Hilfsmittel vollständig angegeben habe und dass ich die Stellen der Arbeit – einschließlich Tabellen, Karten und Abbildungen –, die anderen Werken oder dem Internet im Wortlaut oder dem Sinn nach entnommen sind, auf jeden Fall unter Angabe der Quelle als Entlehnung kenntlich gemacht habe.

Wien, 30. August 2018

---

Manuel Wiesinger



# Acknowledgements

I would like to thank my advisor Univ.-Prof. Dr. Stefan Brunthaler for initiating the idea of this thesis, for constantly giving helpful advice and motivation whenever I needed it. He always supported me generously and encouraged me to attempt a scientific career. Moreover, I want to thank Privatdoz. Mag.rer.soc.oec. Dipl.-Ing. Dr. techn. Edgar Weippl for his advice and supervision.

Furthermore, I want to express my gratitude to SBA Research, the company that has employed me for working on operation system security and is giving me the opportunity to continue my research as employee.

My gratitude goes to my parents Brigitte and Thomas, as their generosity allowed me to study without pressure. They always kept supporting me during my studies. Especially I want to thank my father who awakened my interests for natural science, technology and computer science. Furthermore, I want to thank my sister who always listened to me and made me laugh. I also want to thank my grandmother Gertrud for showing interest, supporting me and constantly believing in me.

I want to thank my girlfriend Klara for her interminable support throughout my studies, for always listening to me and giving me essential advice. I cannot imagine finishing this thesis without her.

Finally, I would like to thank my university friends Florian Schweikert, Simon Wolfsteiner, Kristoffer Kleine and Albert Dengg with whom I attended many classes, had lunch and coffee during my studies. Additionally, my gratitude goes to Christoph Roschger who took the part of a mentor since my first semester.





# Kurzfassung

Die Rowhammer-Schwachstelle stellt ein drastisches Sicherheitsrisiko für moderne Computersysteme dar. Bössartige Angreifer können damit gezielt Daten manipulieren bzw. Kontrolle über ganze Computersysteme erlangen, indem sie Daten ohne jegliche Autorisierung modifizieren. Die vorliegende Arbeit analysiert die öffentlich bekannten Rowhammer-Angriffe auf Betriebssysteme und präsentiert erstmals eine alternative, auf Softwarediversität beruhende Lösung für dieses Problem. Da die bekannten Angriffe stets auf der Vorhersagbarkeit von in modernen Betriebssystemen eingesetzten Algorithmen basieren, beabsichtigen wir, diese Vorhersagbarkeit mittels Softwarediversität zu unterbinden. In concreto stellen wir PAGE SACRIFICE vor: eine leicht umzusetzende, effiziente Modifikation bestehender Mechanismen, die Rowhammer Angriffe verhindert, indem bei jeder Speicherzuweisung physisch benachbarte Speicherblöcke zufällig freigelassen und somit Speicherzuweisungen unvorhersehbar werden. Der Speicher wird dadurch allerdings nicht dauerhaft belegt — sobald das Betriebssystem einen Speicherblock wieder frei gibt, werden auch die freigelassenen Blöcke wieder verfügbar. Um unsere Verteidigungsmaßnahme zu evaluieren, wurde ein Prototyp für Linux entwickelt, der anhand weitverbreiteter Messmethoden untersucht wurde. PAGE SACRIFICE stellt eine neue Maßnahme gegen Rowhammer-Angriffe dar.



# Abstract

The Rowhammer vulnerability allows the modification of arbitrary data without authorization. This poses a dramatic security risk for modern computing systems, as it allows malicious attackers to manipulate data or even gain the control over entire systems. This thesis analyzes publicly available attacks based on Rowhammer as well as known defenses against them. Finally, it suggests a novel solution to the problem based on software diversity. As existing attacks are typically based on the predictability of algorithms used in modern operating systems, we aim to prevent these predictability using software diversity. In concreto we present PAGE SACRIFICE, an easy-to-implement, efficient modification of existing operating system mechanisms which prevents Rowhammer attacks by making the locations of memory blocks unpredictable. This is done by randomly leaving neighboring blocks free when memory is allocated. Leaving memory blocks free does not mean that they remain occupied: As soon as the operating system frees a memory block, it also frees the blocks which were left empty during the memory allocation. To evaluate our defense, we implemented a prototype for Linux which we evaluated using widely known benchmarks. PAGE SACRIFICE is a new defense against Rowhammer based attacks.



# Contents

<b>Kurzfassung</b>	<b>ix</b>
<b>Abstract</b>	<b>xi</b>
<b>Contents</b>	<b>xiii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Structure of the Thesis . . . . .	2
<b>2 Background</b>	<b>3</b>
2.1 Hardware Architecture . . . . .	3
2.2 CPU Cache and Main Memory Structure . . . . .	5
2.3 Memory Management . . . . .	6
2.4 Paging . . . . .	6
2.5 Huge Pages . . . . .	9
2.6 Page Cache . . . . .	10
2.7 Kernel Space And User Space . . . . .	10
2.8 Rowhammer . . . . .	10
2.9 CPU Cache Flushing and Eviction Strategies . . . . .	15
<b>3 State of the Art</b>	<b>19</b>
3.1 B-CATT . . . . .	19
3.2 G-CATT . . . . .	20
3.3 Kernel Page-Table Isolation . . . . .	20
3.4 Static Binary Analysis . . . . .	21
3.5 Rowhammer Defenses based on CPU Performance Counters . . . . .	22
<b>4 Design</b>	<b>25</b>
4.1 The Buddy System . . . . .	25
4.2 Rowhammer-based Attacks . . . . .	27
4.3 Phys Feng Shui . . . . .	29
4.4 Software Diversity . . . . .	36
<b>5 Design and Implementation of Page Sacrifice</b>	<b>37</b>
	xiii

5.1	Overview . . . . .	37
5.2	Preliminaries . . . . .	38
<b>6</b>	<b>Evaluation</b>	<b>47</b>
6.1	Memory Usage . . . . .	47
6.2	CPU Benchmarks . . . . .	48
6.3	Security evaluation . . . . .	49
<b>7</b>	<b>Conclusion</b>	<b>51</b>
7.1	Contribution . . . . .	51
7.2	Limitations . . . . .	51
7.3	Future Work . . . . .	52
<b>A</b>	<b>SPEC CPU2006 Benchmarks</b>	<b>53</b>
	<b>Bibliography</b>	<b>55</b>

# Introduction

Modern society relies heavily on secure computing infrastructure. This infrastructure must be protected against any kind of vulnerability, because cybercriminals and intelligence agencies use vulnerabilities of digital infrastructure to acquire or manipulate information. Therefore, recent hardware vulnerabilities such as Meltdown [1] and the Spectre-family [2] [3] have attracted great media interest, as they affect almost any computer. Another serious hardware vulnerability is the so called *Rowhammer* bug [4]. It affects most memory modules and consequently most computers. Rowhammer can be used to modify arbitrary data if attackers succeed to execute code on their victim's machine. Similarly, to Meltdown and the Spectre family, it may be sufficient if victims are tricked to visit vulnerable websites. Avoiding these hardware vulnerabilities is difficult though. Unlike software vulnerabilities hardware vulnerabilities cannot be easily fixed by installing software updates, instead hardware vulnerabilities require hardware replacement. If critical vulnerabilities affect widely used hardware, it is not feasible to deploy fixes in time. Instead of replacing hardware, software often can be used to circumvent hardware bugs. Circumvention makes it possible to timely deploy software-based fixes for critical hardware problems. Proposed defenses against circumventing the Meltdown, the Spectre-family and Rowhammer vulnerabilities introduce significant performance loss, though [1] [2] [5].

This thesis aims to approach the problem of Rowhammer attacks by analyzing existing attacks and defenses. Based on our analysis we propose PAGE SACRIFICE, a novel mitigation strategy inspired by the ideas of software diversity (which we discuss in section 4.4). To the best of our knowledge all publicly available attacks based Rowhammer take advantage of deterministic operating system behavior. Applying software diversity allows us to eliminate these determinisms. To evaluate our ideas, we provide a proof of concept implementation of non-deterministic memory allocation for the Linux kernel. Finally, we analyze the effects of our proposed solution on operating system performance. This analysis covers memory consumption benchmarks, CPU benchmarks and a security

evaluation. We hope that this work will lead to further research questions for the applicability of software diversity in modern operating systems.

### 1.1 Structure of the Thesis

This thesis is structured as following: Chapter 2 gives preliminary background information on hardware architectures and fundamental operating system tasks relevant for understanding Rowhammer attacks and PAGE SACRIFICE. Chapter 3 gives a comprehensive overview of proposed software and hardware defenses against Rowhammer attacks. In chapter 4 we explain the mechanisms necessary to understand the design of our proposed defense. In chapter 5 we extensively discuss how we implemented a prototype of our defense for Linux. Finally, we evaluate our implementation in chapter 6 and draw our conclusions about our work in chapter 7.



# Background

This chapter gives the background information necessary to understand our PAGE SACRIFICE mitigation strategy against Rowhammer based attacks. First, we briefly revisit the fundamental operating system principles, as they can be found in any undergraduate course book such as Operating Systems: Internals and Design Principles by William Stallings [6]. Secondly, we discuss the Rowhammer problem and thirdly, we explain techniques necessary to flip bits using Rowhammer: cache flushing and cache eviction strategies.

## 2.1 Hardware Architecture

In today's computers the *memory bus* connects the central processing unit (CPU) with the *Dual In-line Memory Modules* (DIMM).

Kim et al. [4] and Brasser et al. [7] summarize the low-level structure of DIMMs as following. Figure 2.1 gives a high-level view about the logical structure of the interaction of CPU and memory modules, as well as the structure of memory modules. DIMMs follow a hierarchic structure, they consist of one or two *ranks*, which correspond to its front and backside. A rank consists of multiple interconnected *chips* containing one or more memory *banks*. Memory banks consist of horizontal *rows* and vertical *columns* of memory *cells*. The *wordline* connects the cells horizontally and the *bitline* connects them vertically. Memory cells are capacitors whose charged or uncharged state represents a binary data value (zero or one). The *access-transistor* manages the read and write access to the word lines and bit lines. Figure 2.2 depicts the structure of a memory bank.

A row of cells has typically the length of one *word*, the unit in which the CPU processes data. Today, the word size is typically 32 bit, 64 bit and for some architectures even 128 bit.

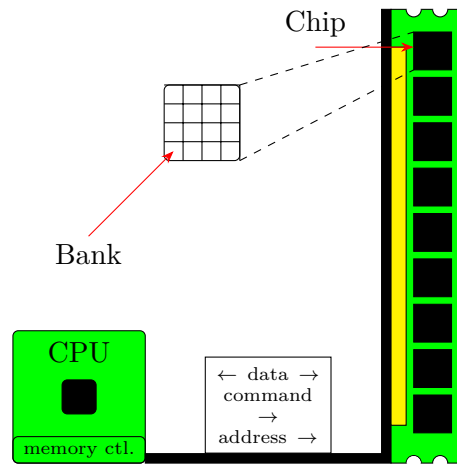


Figure 2.1: The CPU is physically connected to the primary memory via the memory bus.

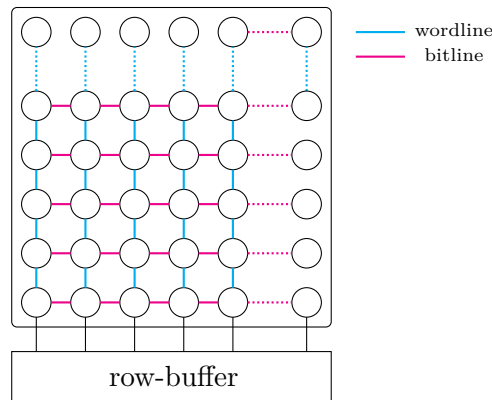


Figure 2.2: Structure of a DRAM bank.

### Memory Access

When the memory controller of the CPU issues a command to read a word, the reading operation triggers higher voltage on the word line, which in turn enables all access-transistors, connects the cells to the respective bit line and transfers the charge to the *row-buffer*. This is commonly referred to as *opening* a memory row. The row-buffer serves as interface to the memory bank, it carries out any operation on the data stored in a bank. Since opening rows is a destructive operation, the row-buffer immediately restores the data in the row, to avoid data loss. Before another row can be accessed, the currently open row needs to be *closed*, i.e. the voltage of the word line is lowered again and the row-buffer is cleared.

## 2.2 CPU Cache and Main Memory Structure

Modern CPUs perform operations on integrated *registers*. Hence, data has to be loaded into registers prior to any data operation. As CPUs perform operations on registers faster than they can load blocks from primary memory, it is sensible to store frequently used data blocks in faster but smaller *cache* memory, which is physically closer to the CPU. Commonly, CPUs employ several *cache levels*; If a block gets less frequently used, the CPU moves it to a larger but slower cache level; Typically CPUs use three cache levels, some cheaper CPUs use only two. Usually each core has its own first and second level cache, the third level is shared among all CPU cores. If a block is not accessed for a certain amount of time finally, the CPU removes it from cache. In case the data in it was modified, it is written back to main memory, otherwise it is simply removed as the data is still available in main memory. Figure 2.3 illustrates three-level caching.

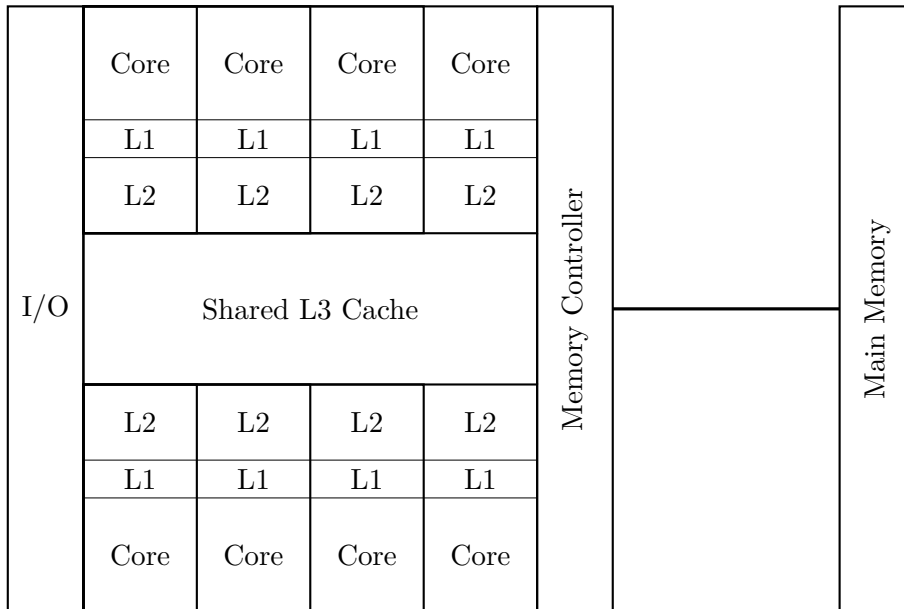


Figure 2.3: Schematic representation of an Intel Core i7-5960X (cf. Operating Systems: Internals and Design Principles [6]).

The frequency the CPU loads a certain block to its registers is often described with a metaphor: frequently used blocks are called *hot* blocks and less frequently are called *cold* blocks. The *principal of locality*, states that in practice memory blocks tend to reference physically close memory blocks [6]. Hence, data currently processed by the CPU is likely to reference data that can be loaded quickly from cache.

Since CPU caches have rather limited storage capacities compared to the amount of main memory, the processor removes less frequently used blocks (i.e. *cold* blocks) by more frequently ones (*warm* blocks). The replacement strategies vary from CPU model to CPU model — and manufacturers generally do not publish them.

## 2.3 Memory Management

Multitasking operating systems subdivide and manage main memory to accommodate multiple processes. This fundamental task of modern operating systems is called *memory management* [6]. In the following we discuss how memory management eases application development by providing *virtual memory* and how memory is subdivided for this purpose using the *paging* technique.

### 2.3.1 Virtual Memory

Modern memory management provides *virtual memory*, such that each process can virtually address the entire main memory. The operating system instructs the CPU to translate virtual addresses used by applications to physical addresses. This abstraction dramatically eases application and compiler development, as software developers do not need to take care about address translation and memory management. If a program requires more memory than the operating system can provide, it can free memory blocks that it can easily reload (e.g. from secondary memory). Additionally, less frequently used memory blocks can be temporarily *swapped-out*, i.e removed from memory and stored on secondary storage. This means that only the parts of an executable files that are currently processed have to be present in main memory; The same applies to files stored on secondary storage. The actual value of virtual addresses depends on the operating system and the processor architecture, however virtual addresses can be imagined to be enumerated from 0 to  $n$ , where  $n$  is the number of blocks the memory is subdivided into. Even if memory addresses are continuous from an application's perspective, the memory blocks may be scattered across the whole memory. If virtually continuous blocks are physically stored in two or more chunks, they are called *fragmented*. The more logically continuous blocks are fragmented, the more address translation overhead is necessary to translate virtual addresses to physical addresses. Therefore, operating systems usually employ techniques to minimize fragmentation (see section 4.1). To be able to remove parts of continuous blocks from main memory, large blocks have to be subdivided into smaller blocks. This can be achieved using the memory *segmentation* technique or the *paging* technique [6]. As paging is the de facto standard for implementing virtual memory we do not discuss segmentation any further.

## 2.4 Paging

The prevailing method to implement virtual memory is *paging*. It divides physical memory into fixed-size chunks referred to as *frames* or *page frames*. A page frame typically comprises of the data stored in one *virtual page*. Note that were clear from the context the literature often uses the term *page* for virtual pages and page frames. To map addresses of virtual page to page frame addresses it is necessary to *translate* the addresses. The information which virtual address maps to which physical address is stored in a multi-level hierarchy of *page tables* (see section 2.4.1). The processor

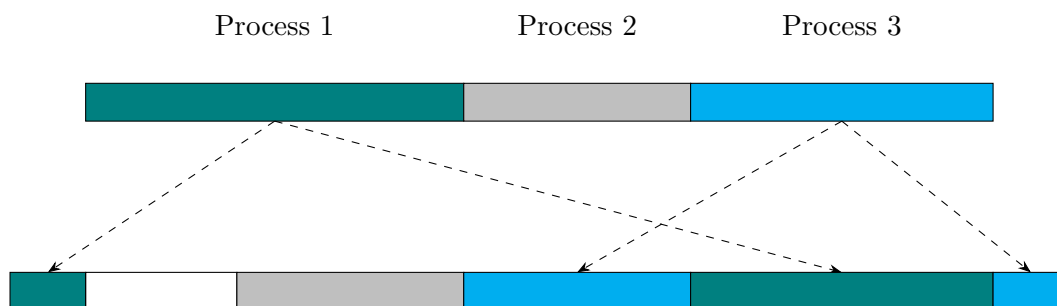


Figure 2.4: A possible mapping of virtual addresses to physical ones. Process 1 and process 3 are fragmented.

architecture defines their hierarchy as well as their structure; Modern processors typically have a component for transparently translating virtual to physical addresses <sup>1</sup>. On most operating systems there is an own page table hierarchy for each process.

Even though the address translation process is done transparently by hardware, it is the operating system's responsibility to ensure that the lower level page tables are present in memory. When a requested page is not in main memory a *page fault* occurs: The processor raises an interrupt and executes the *page fault handler* routine of the operating system. This routine loads the requested page and updates the translation tables accordingly. A page fault is — contrary to what the name suggests — not an error but a normal operation. Typically, data is not loaded into memory until it gets used; When the processor requests data not yet present in memory it raises a page fault to load it.

When data in a page gets modified the page is marked *dirty* until the data is updated on backing storage (e.g. the hard disk). For example, this happens when a user modifies a file and stores the modifications; The file is first updated in main memory and the pages containing the modified data are marked dirty until they get updated on secondary memory.

Because virtual-to-physical address translation is such a common operation, most modern processors cache virtual-to-physical address translations in the *translation lookaside buffer* (TLB). Because memory references tend to be local (cf. section 2.2) using a cache for address translations is highly effective.

### 2.4.1 Page Tables Hierarchies

A simple page table hierarchy is a virtual address consisting of a pointer to a *page table entry* (PTE) and an offset. Physical addresses can be calculated by adding the address in the PTE and the offset. Additionally to address bits, page table entries usually contain

<sup>1</sup>The part of the CPU responsible for the memory translation is often referred to as *memory management unit* (MMU). On modern CPUs the MMU is incorporated into the CPU.

extra bits indicating if the page frame is currently loaded in main memory and if it is marked dirty (i.e. modified but not yet updated on the backing storage). The page table must be permanently available in main memory. Figure 2.5 depicts a simple paging hierarchy.

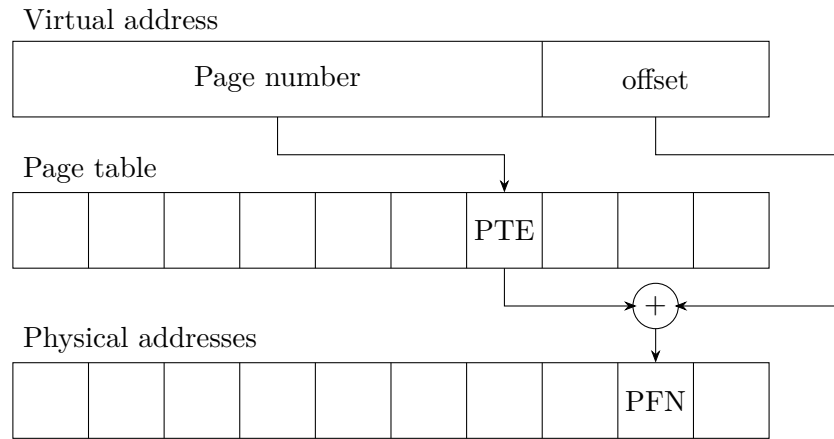


Figure 2.5: Simple paging, the physical address is calculated by adding the page table entry and the offset.

Current processors and operating systems extend the abstraction of virtual addresses, by employing multiple levels of page tables. Multiple levels of page table make it possible to store page tables themselves in virtual memory, consequently not all page tables have to be present in main memory, because page frames holding page tables can be temporarily moved to secondary storage, like any other page. However, the highest level page table must always be present in main memory. Typically, operating systems lazily generate parts of page tables when needed. Lazy generation of page tables is an important precondition for many Rowhammer attacks (see section 4.2).

Current x86\_64 processors support four page levels [8]. A virtual address consists of 48 bits, of which the first nine bits (39-47) point to the top level page table, the next nine (30-38) point to the third level, bits 21-29 to the second level, bits 12-20 point to the page table entry and finally the least significant bits (0-11) are the offset. In Linux the top most page table level is termed *page global directory (PGD)*, the fourth level *page upper directory (PUD)*, the third *page middle directory (PMD)* and the second *page table entry (PTE)*. Figure 2.6 depicts the four-level page hierarchy of the x86\_64 architecture using the Linux kernel parlance. Other operating systems such as FreeBSD use different terms [9].

The x86 architecture supports only three page table levels, to support this the Linux kernel skips the PUD and the PMD level, when running on this architecture. The following table summarizes the page tables the Linux kernel uses on the x86 and the x86\_64 architecture. For the sake of completeness it is important to mention that this summary does not discuss page address extension (PAE), a feature of some x86 processors,

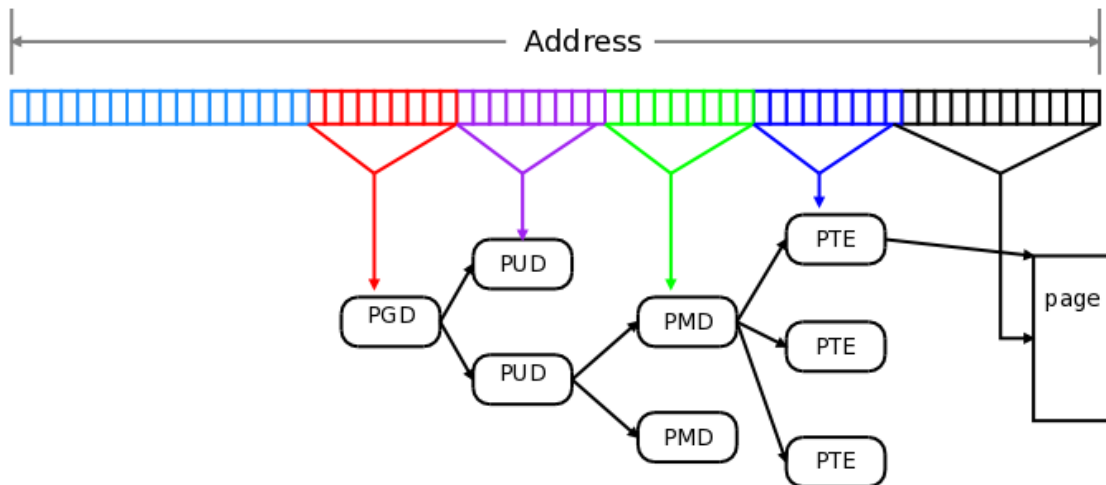


Figure 2.6: Four-level page table hierarchy of Linux for a x86\_64 processor. Graphic by Jonathan Corbet <https://lwn.net/Articles/717293/>.

allowing to use more than 32 bit (the processors word length) for addressing memory. Using more than 32 bits makes it possible to address more than four gigabyte of memory.

		Bits used				
Level		PGD	PUD	PMD	PTE	page
Architecture	x86_64	39-47	30-38	21-29	12-20	0-11
	x86	22-31			12-21	0-11

Table 2.1: Bits used for four-level paging by the Linux kernel

During a context switch the operating system stores the address of the top most page table associated with a process to a special register, so the CPU can transparently translate addresses; On the x86\_64 architecture this is the `cr3` register.

In 2017 Intel announced CPUs supporting five-level paging, using 52 address bits [10]. Even though no hardware is available yet, the Linux kernel already implements five page table levels [11]. To support current processors with four-level paging, the top most level is simply a pointer to the fourth level.

## 2.5 Huge Pages

At hardware level only page frames of fixed size are addressable. On the x86 and the ARM architecture this is typically 4 KB. However, recent CPUs support larger pages (known as *huge pages*), which can be used to allocate large amounts of consecutive memory (e.g. up to 1 GB on ARM). User land applications can request huge pages from the kernel using the `mmap()` or the shared memory systems calls `shmget` and `shmat`. Because

huge pages can be addressed with a single address they reduce memory management overhead.

Many Rowhammer based attacks like Phys Feng Shui (see section 4.3) aim to flip bits at vulnerable memory locations, because this usually requires the smallest addressable page size, huge pages are not discussed further in this thesis.

### 2.6 Page Cache

As of Linux 2.4 there is a unified cache for all block device I/O. Caching I/O to block devices is beneficial for two reasons: first, pages can be ordered, so the operating system can write them back to secondary memory more efficiently; second, because memory references tend to cluster (cf. section 2.2), it is very likely that pages can be served from the page cache when requested again. Other operating systems provide similar features [6].

Because data in the page cache mostly consists of pages, that can be reloaded from secondary storage at any time, operating systems typically aim to reserve as much unused memory as reasonable for this cache. The Linux kernel does not count data cached in the page cache to the amount of used memory. It is important to mention that this behavior is Linux specific, as other operating systems do count their I/O caches to their memory usage statistics. For instance, FreeBSD usually displays much less free memory than Linux, because I/O caches are counted as used memory [12].

### 2.7 Kernel Space And User Space

Most modern processors support at least two modes of operation: the *user mode* and the *kernel mode*. In user mode the processor executes normal applications; in *kernel mode* operating system code is executed. Code running in kernel mode has full access to all memory locations and is privileged to execute all machine instructions. User land applications can use functionality only available in kernel mode (such as accessing files, communicating over the network or allocating memory) by using interfaces provided by the operating system for this purpose. These interfaces are called *system calls*. When system calls are executed the operating system is responsible for setting and clearing a control bit, indicating the mode in which the processors shall operate. If user mode applications try to directly execute instructions which are not allowed in user mode, an interrupt is raised and the operating system typically ends the process [6].

### 2.8 Rowhammer

At hardware level dynamic random-access memory (DRAM) is organized in two dimensional arrays of memory cells. A memory cell stores one bit of information in a capacitor, whose charge state represents a binary value. Memory banks group memory cells in rows



and columns, individual rows and columns can be accessed via an access transistor (cf. section 2.1).

Because capacitors lose their charge, they have to be periodically refreshed at least every 64 milliseconds according to the DDR 3 SDRAM specification [13]; this refreshing operation is done systematically, row-by-row.

Memory cells are placed so closely next to each other, that accessing memory cells may cause disturbance errors in neighboring cells. Such disturbance errors normally do not occur under normal circumstances and not every memory module is vulnerable to them. Vulnerability to disturbance errors depends on various properties, such as manufacturer, manufacturing date and the refresh rate of the memory module.

Often disturbance errors can cause interference among memory cells, in a way that a bit changes its value; this change is called a *bit flip*. Even though not every memory module is vulnerable, tests showed that about 85% of DDR3 memory modules are vulnerable [4]. It is important to mention that newer DDR4 memory modules are also vulnerable to bit flips [5].

Kim et al. [4] demonstrated that bit flips in memory cells may be triggered by performing a large number of consecutive reading operations on physically neighboring memory rows. This process is known as *Row hammer* or *Rowhammer*. The literature also refers to it as *rowhammering*, *row hammering* or *hammering* a memory cell. These terms were introduced by a patent, assigned to Intel [14] and are now common terms.

As not every memory cell is vulnerable to row hammering, attackers first have to search for vulnerable memory cells. Once a vulnerable cell has been found, it is very likely, that the bit flip can be repeated later. Repeatability is a necessary precondition for row hammering.

Even memory error correction-codes (ECC) do not reliably protect from this issue. Kim et al. [4] state:

While most words have just a single victim, there are also some words with multiple victims. This has an important consequence for error-correction codes (ECC). For example, SECDED (single error-correction, double error-detection) can correct only a single-bit error within a 64-bit word. If a word contains two victims, however, SECDED cannot correct the resulting double-bit error. And for three or more victims, SECDED cannot even detect the multi-bit error, leading to silent data corruption. Therefore, we conclude that SECDED is not failsafe against disturbance errors.

To be able to hammer memory rows it is necessary to bypass all CPU caches, because otherwise the data stored in the target cells would be loaded from these caches as it is frequently accessed. Several techniques for bypassing CPU caches exist, we discuss them in section 2.9.

Hammering one row to introduce bit flips in a neighboring row is known as *single-sided row hammering*. The chances to cause bit flips increase if not only one neighboring memory row is hammered, but if both neighboring cells are hammered. Hammering from two sides is a second variant of Rowhammer known as *double-sided Rowhammer*. A third variant for some memory modules that allow to keep the connection to a memory row open is *one-location hammering* [5]. We introduce these three variants in the following.

The ability to cause reproducible bit flips itself, is not sufficient for performing a serious attack (expect attacks aiming to corrupt data). Chances to cause a bit flip in critical data are relatively low, especially as many computer systems implement some sort of checksums for critical data. Attackers usually scan the whole available memory until they find one or more bit flips, with suitable alignment for their need. Then they trick the operating system to place critical data at a memory location vulnerable to a bit flip. Several techniques for doing that exist. We discuss them in section 4.2.

### 2.8.1 Single-sided Row Hammering

Flipping bits in memory cells by row hammering one neighboring memory cell is known as single-sided row hammering. The technique is more efficient if multiple memory locations get hammered simultaneously. Dullien and Seaborn successfully induced bit flips in neighboring cells by hammering four or eight locations at once [15].

Listing 2.1 shows x86 assembly code [4][15] for hammering two memory locations  $X$  and  $Y$ . To successfully perform bit flips the addresses  $X$  and  $Y$  must point to different memory rows in the same memory bank. The first line defines a label `code1a`: it marks a location to which the code can later return to; the label therefore marks the begin of a loop. Line 2 reads the content from address  $X$  to the `eax` register. Similarly, line 3 reads from address  $Y$  to the `ebx` register. Line 4 and 5 force the CPU to flush the data stored in  $X$  and  $Y$  from CPU caches. The last line jumps back to the label `code1a` in line 1 and the code is executed again beginning with line 2.

Because the CPU cache is emptied by the `clflush` instruction in every iteration of the loop, the data stored in  $X$  and  $Y$  is always read from the RAM (i.e.  $X$  and  $Y$  are hammered), possibly more often than the refresh cycle of a capacitor storing a single bit.

An advantage of single-sided row hammering is that no knowledge about the underlying memory architecture is required, in order to successfully flip bits in neighboring memory rows, because each memory row has at least one neighbor.

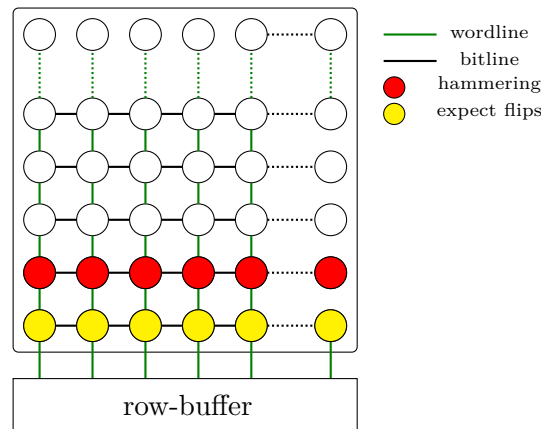


Figure 2.7: Single-sided Rowhammer is the most simple row hammering technique. Its advantage is that no knowledge about the inner memory structure is required.

```

1 code1a:
2   mov (X), %eax // Read from address X
3   mov (Y), %ebx // Read from address Y
4   clflush (X)   // Flush cache for address X
5   clflush (Y)   // Flush cache for address Y
6   jmp code1a

```

Listing 2.1: The x86 assembly used to perform single-sided row hammering [4][15].

### 2.8.2 Double-sided Row Hammering

Double-sided hammering is the Rowhammer technique with the highest rate of successful bit-flips. Because, both neighboring memory rows are hammered the changes of successfully flipping a bit in a victim row increase. Hammering from both neighboring cells requires precise knowledge about the physical geometry of the hammered memory module, as attackers must be able to allocate exactly both neighboring memory rows of the victim row.

Seaborn and Dullien found out that hammering 256 KB above and below victim memory rows is highly effective in order to induce bit flips [15]. However, they tested several memory modules of only one vendor and suspected that this value differs from manufacturer to manufacturer. Double-sided row hammering is the hammering technique with the highest chance of inducing bit flips. Figure 2.8 depicts a double-sided row hammer attack.

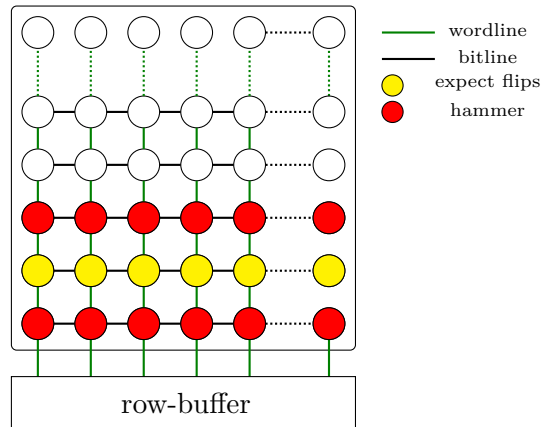


Figure 2.8: Double-sided row hammering induces bit flips by performing reading operations from both physically neighboring memory rows. This technique induces more bit flips than single-sided row hammering and One-location Hammering.

### 2.8.3 One-location Row Hammering

One-location row hammering makes use of performance optimizations of newer memory controllers, which close the connection between row and a row buffer earlier than necessary. Because the connection gets closed too early, repeatedly reading from a memory location causes a memory row to be repeatedly opened and closed. Frequently, re-opening memory rows influences the charge of memory cells sufficiently to cause bit flips [5]. As a consequence, the one-location row hammering technique is not based on inducing interference among memory rows, like single sided- or double sided row hammering, but on effects caused by repeatedly reopening a memory row. As it is the case for single sided- and double-sided row hammering bit flips are reproducible once they are found. Gruss et al. successfully performed one-location row hammering on DDR3 and DDR4 memory modules [5]. Similar to single-sided row hammering this, one-location row hammering does not require knowledge about the underlying memory geometry; to cause bit flips it is sufficient to hammer a random location. Figure 2.9 shows where bit flips can be expected when performing one-location hammering.

An advantage of one-location hammering is that Rowhammer mitigation techniques based on observing memory patterns (see section 3.5) are not able to detect hammering [5]. One-location hammering is the least effective row hammering method, as it produces less bit flips than the other two methods.

### 2.8.4 Bit flip Attacks Before Rowhammer

The potential exploitation of bit flips was known long before Kim et al. [4] discovered the Rowhammer bug in 2014. In 2003 an attack was presented allowing a specially crafted

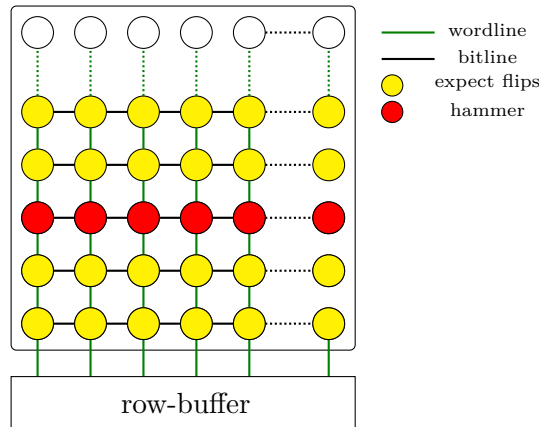


Figure 2.9: One-location hammering can induce bit flips even in non-neighboring memory rows — which makes it a powerful row hammering technique [5]. In contrast to single-sided and double-sided row hammering bit flips are not caused by row interference, but by effects caused from repeatedly re-opening a memory location.

Java application to take full control over the Java Virtual Machine <sup>2</sup> or Microsoft’s .NET <sup>3</sup> virtual machine in many cases, once a flip occurs [16]. Even though bit flips are rare, attackers can exploit them if the circumstances allow them to run their attack for a long period of time. The authors state that in 1996 one bit flip per month is to be expected on a end-user PC. 70% of the flips, occurring in rare events such as hardware defect or influence of cosmic rays are suitable for their attack.

## 2.9 CPU Cache Flushing and Eviction Strategies

Modern processor architectures cache frequently accessed data in smaller but faster cache memory (see section 2.2). Because Rowhammer attacks need to repeatedly perform reading operations directly from memory, attackers need to circumvent processor caches in order to cause bit flips in memory. Circumvention can either be done by *flushing* or *evicting* data from the cache. Flushing means to force removal of the data block from cache by using CPU instructions. Eviction means loading other data into the cache until the CPU replaces the target data block in course of normal operation. Both methods make it necessary, that the CPU loads the data again from main memory if it is subsequently accessed. Reloading data from memory triggers a read operation on memory the memory module. Repeatedly evicting and reloading data may induce bit flips as explained in section 2.8. Some operating systems provide direct uncached memory access, in such cases the cache does not need to be circumvented for flipping bits in memory. In the following we discuss various techniques for flushing and evicting data from CPU caches.

<sup>2</sup><https://www.java.com>

<sup>3</sup><https://www.microsoft.com/net>

### 2.9.1 Flushing CPU Caches using CPU Instructions

The easiest way to circumvent CPU caches is to use special instructions that allow to remove data from cache. x86 processors supporting the SSE2<sup>4</sup> instruction set extension (i.e. virtually any modern x86 compatible CPU) provide the `clflush` instruction, allowing user-level code to flush variables from CPU cache. Repeatedly reading from the same memory location and flushing the cache afterwards can cause a bit flip (cf. section 2.1).

### 2.9.2 CPU Cache Eviction using Memory Access Patterns

On platforms without instructions to flush CPU caches, special methods to circumvent CPU caches are necessary to cause repeated read operations directly from memory. Cache eviction strategies based on memory access patterns aim to find sequences of read operations that replace cached data. For row hammering it is necessary to subsequently read data in the target memory cell directly from the memory module. This can be achieved by repeatedly evicting and reloading the same data [18] [19].

Even though generic eviction strategies exist (see section 2.9.5), efficient eviction strategies require precise knowledge about the replacement algorithm the CPU on the victim system uses. Replacement algorithms usually replace less frequently used slots, however they do not exactly employ a least recently used strategy [18]. CPU manufacturers usually do not publish documentation on the cache replacement strategies they employ. Nevertheless, it is possible to observe the CPU cache's timing behavior and draw conclusions about the employed eviction strategy. Based on these observations an *eviction set* can be created. An eviction set is a set of data at that, when loaded in into CPU caches, from the correct locations and in the correct sequence, replaces specific data in the CPU cache. The sequence in which an eviction set has to be loaded to remove certain data from cache is called *access pattern*. The combination of eviction set and access pattern is called *eviction strategy*.

Once attackers have crafted a suitable eviction strategy they can remove data of their choice from CPU cache efficiently and reliably. Repeatedly evicting and reloading data from memory causes repeated reading operations from main memory. A loop of evicting and reloading is fast enough for row hammering [5].

Gruss et al. defined four types of eviction strategies [19]:

- *Static eviction set and static access pattern*: Eviction set and access pattern are pre-computed prior to performing the actual attack. Pre-computing requires precise knowledge about the microarchitecture used on the target system, however no searching for necessary parameters on the victim machine is necessary. This approach is similar to the eviction presented by Awke et al. (see section 2.9.3).

---

<sup>4</sup>SSE 2 stands for Streaming SIMD Extensions 2, it is an extension to the standard x86 instruction set. [17]

- *Dynamic eviction set and static access pattern:* Based on knowledge about the replacement strategy of the victim system attackers pre-compute access patterns and generate eviction sets at run time. This method is very efficient for attacks as pre-computation of access patterns, needs to be done only once for every target architecture. Various algorithms exist for finding static access patterns [20] [21] [19].
- *Dynamic eviction set and dynamic access pattern:* This eviction strategy is essentially brute-force search for eviction set and access pattern at run time. It comes at the cost of long run time but has the advantage, that attackers do not need to know the microarchitecture of the target system. Consequently, this strategy allows wide ranging, fully automated attacks.
- *Static eviction set and dynamic access pattern:* Using a pre-defined evicting set and a randomly generated access pattern. According to the authors this is similar to randomly testing static access patterns.

### 2.9.3 Cache Eviction based on Knowledge about the target Microarchitecture

Aweke et al. [18] created an eviction strategy for Intel Sandy Bridge processors, by loading arbitrary addresses and subsequently searching for other addresses the CPU maps to the same cache slot. Reverse engineered documentation about the Intel Sandy Bridge microarchitecture provided them with information about the exact replacement strategies and enabled them to craft effective precisely access patterns. The exact functioning requires a detailed discussion of the microarchitecture and is therefore out of scope of this thesis. However, the basic idea is similar to the generic eviction strategy presented by Gruss et al. (see section 2.9.5).

### 2.9.4 Direct Memory Access on Android Devices

Most Google Android <sup>5</sup> devices such as smartphones use the ARM architecture. ARM does not provide instructions to flush data — like `clflush` (cf. section 2.9.1) — from CPU caches. To evict data from CPU caches Van der Veen et al. constructed a cache eviction strategy for ARMv7 and ARMv8 CPUs, similar to the method presented by Aweke et al. [18] (see section 2.9.3), however it turned out that this technique is too slow for practical use [22]. Consequently, page cache eviction is not a reliable method for attacking Android devices. However, gaining direct uncached memory access on Android is possible: With Android 4.0 ICS (Ice Cream Sandwich) Google introduced the ION memory manager [23], a memory manager for the Android kernel (which is essentially a modified Linux kernel), allowing unprivileged applications directly access memory via the `/dev/ion` device. Consequently, ION can be used to hammer memory cells [22].

---

<sup>5</sup><https://www.android.com/>

### 2.9.5 Generic Cache Eviction Strategies

Gruss et al. developed a platform independent, generic method for finding eviction strategies by observing memory access patterns [19]. They describe their eviction strategy as follows:

The goal is to find a eviction strategy that matches the undocumented replacement policy the closest, including the possibility of policy switches.

Essentially, the method is a brute-force like-search for cache eviction strategies. Searching is done by observing timing differences when iterating over arrays of suitable size. Finding an optimal eviction strategy is not feasible in reasonable time, therefore the authors reduced the search space using heuristics. As this strategy allows searching for eviction strategies independently from the execution environment, it enables attackers to evict data from cache slots from high-level run time environments such as JavaScript [19]. The authors presented two ways to determine eviction sets and access patterns suitable to hammer memory rows from JavaScript run time environments of web browsers.



## State of the Art

This chapter summarizes defense strategies for attacks based on Rowhammer-induced bit flips. First, we introduce B-CATT [24] and G-CATT [25] two methods aiming to prevent attackers from accessing row hammerable memory locations. In section 3.3 we explain Kernel Page-Table Isolation [26] (formerly known as *KAISER*), a method to separate kernel page tables from user land page tables. This separation makes it impossible for attackers to row hammer kernel memory. Furthermore, section 3.5 introduces *ANVIL*, a strategy to detect Rowhammer attacks and stop the attacker’s process. Section 3.4 presents *MASCAT* an approach detecting patterns of Rowhammer attacks in executable binaries. Finally, in section 3.5.2, we provide a brief overview about proposed hardware-based defenses.

### 3.1 B-CATT

The principal idea of *B-CATT*<sup>1</sup> is fairly simple, before the boot loader loads the operating system, it scans the entire memory for locations vulnerable to bit flips, that means it row hammers the whole memory. Brasser et al. implemented B-CATT by adapting the well-known GRUB2<sup>2</sup> boot loader. During the boot process the firmware (typically BIOS or UEFI<sup>3</sup> on newer machines) reports the areas of available memory to the boot loader, which reports it to the operating system. Before GRUB reports the operating system about the available memory blocks, B-CATT scans the entire memory for bit flips and reports only non-vulnerable memory blocks to the operating system. Therefore, this defense is fully transparent to the operating system — hence no modification of operating systems is required. Being an operating system independent defense is the

---

<sup>1</sup>Brasser et al. explain the name as following: “The name B-CATT is composed of two parts: B refers to our **B**ootloader based solution, CATT abbreviates **CA**n’t **T**ouch **T**his.”

<sup>2</sup><https://www.gnu.org/software/grub/>

<sup>3</sup><https://www.uefi.org/>

biggest advantage of B-CATT. The authors reported only a minor memory overhead of less than one megabyte on their testing systems. They did extensive benchmarks using various benchmarks suites and could not detect any measurable impact on the systems performance.

B-CATT is considered an effective, but impractical defense for Row hammer attacks. Kim et al. [4] observed that 95% of the main memory would be blocked by B-CATT in realistic scenarios, Gruss et al. verified this result [5]. Moreover, B-CATT requires row hammering the entire memory each time the system boots.

## 3.2 G-CATT

G-CATT (*Generic CAN't Touch This*) is a further development of B-CATT (see section 3.1) by the same authors [25]. The defense is like its predecessor, a low-level defense against Rowhammer attacks; However unlike its predecessor, G-CATT does not require row hammering at boot time and is implemented in the operating systems, physical memory allocator rather than the boot loader. The defense does not prevent row hammering, but aims to limit the negative consequences of Rowhammer based attacks.

G-CATT divides the physical memory in two areas: A kernel space area and user space area. The two areas are physically kept separate during the whole run time. Hence, memory storing kernel data structures (such as page tables) can only be allocated in the kernel area. As attackers can only allocate user space memory they row hammer attacks such as Phys Feng Shui (see section 4.3 for an explanation of the attack) are mitigated.

Brasser et al. implemented a prototype of G-CATT for Linux, due to the generic design of the Linux kernel, the prototype is architecture independent, its authors successfully mitigated Rowhammer based attacks on ARM and x86\_64 systems. To study the performance impact the authors did a large number of benchmarks and concluded that G-CATT has no significant impact on the system's computational power or memory usage. Even though the defense has no performance impact, it can only mitigate Row hammer attacks targeting the kernel memory. Gruss et al. presented a Rowhammer based privilege escalation attack using only user space programs [5]. Furthermore, in early 2018 Cheng et al. [27] exploited the Rowhammer bug in the presence of G-CATT using special purpose memory that is shared between user space and kernel space (e.g. video buffers).

## 3.3 Kernel Page-Table Isolation

The idea of Kernel Page-Table Isolation (KPTI) <sup>4</sup> is similar to G-CATT (see section 3.2): kernel space and user space memory are physically separated. The difference is that

---

<sup>4</sup>KPTI is also known under the name KAISER (Kernel Address Isolation to have Side-channels Efficiently Removed), however the Linux kernel developers decided to name their implementation KPTI [28].

KPTI does not physically separate kernel memory and user space memory; it rather uses separate page tables for kernel space and user space. This separation makes hammering kernel memory impossible, as page tables of user space applications are stored in separate physical locations; hammering attacks requires them to be physically neighboring (see section 2.8). Before KPTI was available, operating systems used the same top-level page table for the kernel and user space applications. Using only one top-level page table is problematic; there are at least three published attacks not related to the Rowhammer vulnerability based on this behavior [29][26].

The basic principle of KPTI is quite simple, each time the context is switched (i.e. the processor changes the process it currently executes) from a user space process to a kernel process (or vice versa), the kernel changes a pointer to the top-level page table to another independent page table hierarchy. As, the x86 architecture requires some parts of the kernel (e.g. interrupt handler routines, system calls, etc.) to be mapped in the user space addresses. Therefore, KPTI minimizes these parts to an absolute minimum necessary for compatibility. The Linux kernel developers included KPTI in the kernel [28] and backported it to the long-term support releases.

The authors of KAISER reported a performance overhead of only 0.28% [26].

### 3.4 Static Binary Analysis

*MASCAT* [30] is a tool for binary analysis allowing to scan software binaries for unusual instruction patterns. As some instruction sequences are typical for low-level attacks, harmless code can often be distinguished from attack code. The tool is capable of detecting a wide range of attack patterns, including side-channel attacks and Rowhammer attacks.

Single-sided and double-sided row hammering (see section 2.8) typically uses many cache flush instructions (e.g. `clflush` on the x86 architecture) in loops. Another indicator for detecting potential Rowhammer attacks are instructions triggering direct, uncached memory access in loops (e.g. `movnti` and `movntdq`). *MASCAT* does not only rely on these patterns to detect row hammering: it scans for unusual instruction patterns containing timing instructions such as `rdtscp` and `rdtsc`. These instructions are often used by timing side-channel attacks, which may be used to craft eviction strategies necessary for row hammering (see section 2.9).

Gruss et al. [5] demonstrated that static binary analysis can be circumvented by hiding attack code in secure memory areas, such as SGX enclaves<sup>5</sup>. SGX enclaves allow user space applications to protect memory areas by encrypting them [31]. The CPU encrypts data and code completely independent of the operating system, therefore, even code running with maximum operating system privileges or in kernel mode cannot an enclave's content. Currently, SGX enclaves are supported exclusively by recent Intel CPUs [31], however some AMD processors provide similar features [32].

<sup>5</sup><https://software.intel.com/en-us/sgx>

As MASCAT supports only scanning for instruction patterns and does not perform semantic analysis or code execution, it cannot guarantee that a certain code pattern actually means an attack. Hence, it cannot detect all kinds of attacks; skilled attackers knowing the target microarchitecture likely can circumvent any static analysis tool; the authors of MASCAT state:

Note that the goal of our design is not to detect all possible microarchitectural attack designs (indeed, if an expert attacker knows the approach taken by our tool, he can always find a way to bypass it) but rather to give a good framework on the detection of existing attacks.

Consequently, static analysis is a valuable tool, but no ultimate solution for Rowhammer based attacks.

### 3.5 Rowhammer Defenses based on CPU Performance Counters

ANVIL is a Rowhammer mitigation technique aiming to detect and prevent single-sided and double-sided Rowhammer attacks [18]. A kernel module monitors the behavior of all operating system processes; if the module detects potential row hammering it thwarts the attack by refreshing hammerable memory cells in time.

The kernel module counts the number of last-level cache misses, if the number exceeds a predefined limit, it instructs the processor to monitor which virtual memory addresses contain the code responsible for the increased amount of cache misses. Counting cache misses and evaluating the virtual memory addresses requires special CPU monitoring features. Even though the proof of concept implementation of ANVIL relies on Intel performance counters <sup>6</sup>, the authors point out that their defense can be implemented on any processor architecture providing such features; they state that this includes some AMD processors.

ANVIL makes use of the fact Intel performance monitor counters do not only report the number of cache misses, but also if the data was loaded from CPU cache or from main memory.. Consequently, the kernel modules can distinguish if a page was loaded from CPU cache or from main memory. If a page is repeatedly loaded from memory, even though it is frequently used, ANVIL suspects that the page is used for row hammering (cf. section 2.8). As the mitigation is implemented as kernel modules it always runs in kernel mode and can therefore access the entire memory. Therefore, it can access the data structures of the kernel's processes management (`struct task_struct`), and figure out which process is responsible for the observed suspicious behavior.

ANVIL refreshes the suspected victim cell performing a reading operation on it. To do so it relies on information about how the RAM maps physical memory addresses to memory

---

<sup>6</sup><https://software.intel.com/en-us/articles/intel-performance-counter-monitor>

cells. In general RAM manufacturers do not publish documentation on this, therefore the authors of ANVIL relied on reverse engineered mapping information.

The defense has generally little overhead and no significant impact on normal applications, the authors report a peak overhead (measured by SPEC CPU2006), of 3.18% and 1.7% on average.

### 3.5.1 Other Rowhammer defense approaches using Performance Counters

The Linux Kernel developers discussed if a Rowhammer defense based on performance counters should be included in the kernel [33]

Run time detection of single-sided and double-sided row hammering using performance counters is also supported by intrusion detection systems such as HexPADS [34] and CloudRadar [35].

### 3.5.2 Hardware-based Defenses

Various researches proposed based defenses against Rowhammer based attacks. Most of them provide reliable protection, however the enormous number of vulnerable devices employed globally makes hardware replacement impossible.

- In the first paper describing the Rowhammer bug, Kim et al. [4] suggested *PARA*, a state-less, probabilistic Rowhammer defense for low-level memory controllers; Each time a row is opened or closed, one neighboring cell is also opened and therefore refreshed. During single-sided and double-sided Rowhammer attacks the neighboring cells are opened and closed repeatedly, which finally can cause disturbances eventually leading to bit flips (cf. section 2.8). *PARA* randomly refreshes neighboring memory rows, each time a row is opened. Because frequent accesses cause many random refreshes, the authors conclude that it is statistically guaranteed, that row hammering is prevented, because the neighboring cells of victim cells will get refreshed in time. To the best of our knowledge, *PARA* was not adapted by the industry.
- *ARMOR (A Run-time Memory Hot-Row DetectOR)* [36] is another hardware solution eliminating the possibility to hammer memory rows. The authors suggest to add an additional cache for memory cells to memory controllers. Consequently, hammering neighboring cells is prevented, because the data stored in frequently accessed memory cells is read from the hardware-cache, rather than from memory modules. Because, there is no software control about the hardware cache attackers cannot circumvent it.
- The LPDDR4 specification [37] defines two strategies against Rowhammer attacks.

- *Target Row Refresh (TRR)* which introduces a counter for each memory row. Each time a memory row is accessed the counter of the neighboring cells is incremented; Once the number of accesses exceeds a threshold the memory modules refreshes neighboring memory cells.
- *Maximum Activation Count (MAC)* sets an upper limit how often a cell can be accesses before it has to be refreshed.

Both methods prevent row hammering, as memory cells are refreshed before memory disturbances suitable for causing bit flips can occur.

#### 3.5.3 Rowhammer related Firmware Updates

Several hardware manufacturers (e.g. HP [38], Lenovo [39] and Apple [40]) provided firmware updates doubling the refresh rates of memory cells. The DDR3 specification requires memory modules to refresh memory rows at least every 64 milliseconds [13]. Doubling the refresh rates (i.e. using refresh intervals of 32 milliseconds) makes bit flips less likely — but it does not make them impossible. Aweke et al. [18] were able to perform double-sided row hammering, even with memory modules configured to refresh memory rows every 15ms (i.e. more than four times of the required refresh rate). Furthermore, the authors state that increasing the refresh interval comes at the cost of increased power usage and reduced data throughput. For these reasons firmware updates are no effective solution for the Rowhammer vulnerability. Refresh rates of 8ms are necessary to reliably prevent row hammering [33].

# CHAPTER 4

## Design

Attackers can cause serious malfunctioning of computing systems by inducing bit flips at random memory locations; if they are able precisely control in which data they flip bits, they can craft much more powerful attacks.

This chapter discusses the most important techniques — publicly known by the time of writing this thesis — to exploit the Rowhammer vulnerability. First, we discuss the buddy system for physical memory management, as it is essential for understanding privilege escalation attacks as well as our proposed solution (see chapter 5). Then we discuss *Flip Feng Shui* and *Phys Feng Shui*, two attacks based on the deterministic nature of the buddy system. Furthermore, we discuss Linux kernel internals relevant for the prototype of PAGE SACRIFICE, that we implemented (see chapter 5). Finally, we provide a discussion of software diversity, as it is the fundamental concept used for our Rowhammer defense.

### 4.1 The Buddy System

A fundamental task of modern operating systems is allocating physical memory for processes. Typically, operating systems support physical memory allocations in various different sizes. The *buddy system* is an algorithm for efficiently managing allocations of different sizes. Essentially, the algorithm splits memory blocks into halves until it finds a block of optimal size; similarly, it merges previously split blocks to larger blocks when they get freed. When blocks larger than the maximum block size are requested, several (ideally consecutive) blocks are allocated. Because block sizes have to be permanently multiplied divided by 2, their size is usually a power of two.

Empirical experiments by Donald Knuth showed that the buddy system works very efficiently, he also states that it provides good memory balance — as during his experiments no memory overflows occurred until 95% of the memory was reserved [41].

A typical size for the smallest addressable memory block that can be addressed by hardware (i.e. a page frame) is 4096 (4 KB), older architectures may use different values [6]. Some architectures even support page sizes of different size (cf. section 2.5). During memory initialization the memory initialization routine of the operating system divides the whole usable memory into blocks of size  $2^n \times page\_size$ . On Linux  $n$  currently defaults to 11 on x86, which results in memory blocks of size  $2^{11} \times 4096$  bytes (8 MB), for  $page\_size = 4096$ . When the memory allocation subsystem processes a memory request, smaller than the maximum block size, it splits memory blocks in two parts, until further division would result in a memory block too small to fulfill the memory request. The remaining half of a split block is marked as *buddy* of the block that is used to fulfill the request. When a block is freed and its buddy is also free they are *merged* (or *coalesced*) to a block of double size. After merging, the resulting block has the size of both buddies (i.e.  $buddy\_size \times 2$ ). If the buddy of the resulting free block is also free, they are also merged; Merging continues until one buddy is not free or the maximal block size is reached. If no free block of a certain size is available the next largest block is split.

Figure 4.1 depicts the functioning of the buddy system, using a maximum block size of 2048 KB =  $2^{11}$  in nine steps.

- (1) Initially a free memory block of maximum size is available.
- (2) The buddy system receives a request for a 512 block. It splits the 2048 block into two 1024 blocks; As a 1024 block can accommodate two 512 blocks it is split. Without the last split operation one 512 block would remain unused.
- (3) The system gets a request for a 256 block, as the smallest available block is the buddy of the 512 allocated in step 2, this block gets split into two 256 blocks.
- (4) A process requested a 128 block, consequently the free 256 block is split into two 128 buddies.
- (5) The 256 block is freed. As the 128 block is still in use, no merging takes place yet.
- (6) Another 128 block gets allocated. As a suitable block is already available the buddy system immediately reserves it to fulfill the request.
- (7) Both 128 blocks are freed. This triggers a several merge operations.
- (8) Because both 128 blocks are free, they are merged to one 256 block.
- (9) Since the buddy of the block resulting from the merging operation in step 8 is also free, they are merged to a 512 block.

Due to its high efficiency the buddy system is used in many general purpose operating systems [42][43].



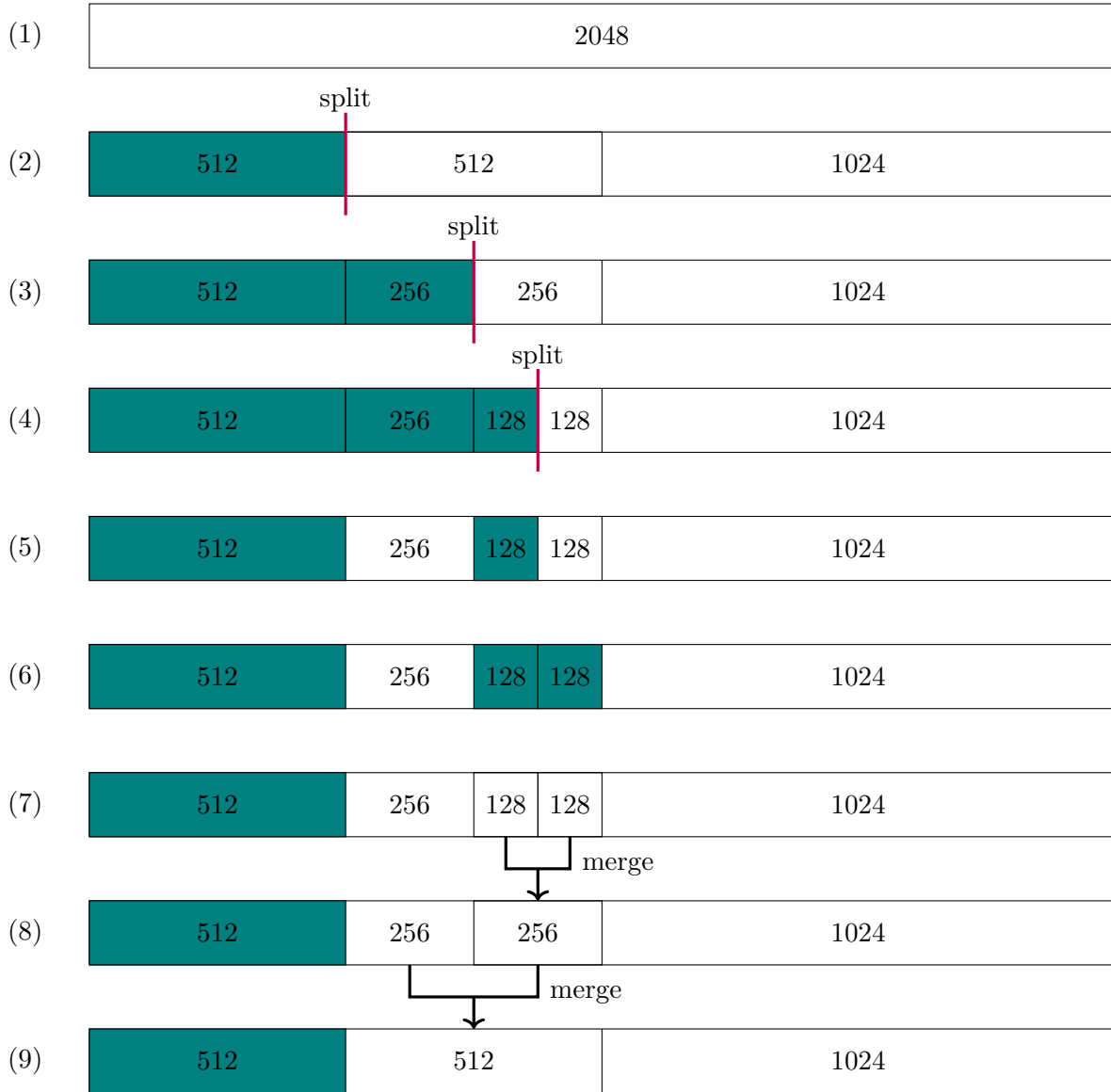


Figure 4.1: Example of the buddy allocation principle.

## 4.2 Rowhammer-based Attacks

This section gives an overview about state-of-the-art attacks based on Rowhammer-induced bit flips.

On a high-level all attacks based on row hammering can be divided into three steps.

- *Memory templating*: The attackers scan the entire memory for blocks vulnerable to

row hammering.

- *Memory massaging*: The attackers force that the operating system stores critical data in a location where they can induce bit flips.
- *Exploitation*: The neighboring cells of the victim cell (which are always in the control of the attacker) are hammered and a bit flip in the critical data is induced).

#### 4.2.1 Escaping the Sandbox of Google Chrome using Rowhammer

In 2015 Seaborn et al. published the first privilege escalation exploit based on the Rowhammer vulnerability. The exploit allows malicious web site operators to gain root privileges by escaping from the Native Client (NaCl)<sup>1</sup> sandboxing environment of the Google Chrome browser<sup>2</sup>. NaCl allows execution of native machine code instructions directly from the web browser. Using native machine code instructions allows attackers to use the `clflush` instruction for row hammering (see section 2.8 for a detailed explanation).

Once attackers found a memory location vulnerable to row hammering, a second step is necessary to flip a bit at a critical location. To do so they repeatedly mapped the same file to memory using the `mmap` system call. This repeated mapping causes many second-level page table allocations (cf. section 2.4), therefore the whole memory is *sprayed* with second level page tables. This makes it likely that a page table gets placed at a memory location vulnerable to a bit flip. Flipping the right bit in a second level page table (i.e. the page middle directory (PMD) on Linux), changes the pointer of the page table entry (PTE, see section 2.4) such that it points to a memory area of the attackers process, figure 4.2 depicts this. Once the PMD points to memory under the control of the attackers they can set the memory addresses the PTE points to. Consequently, they can scan the whole memory for critical data. Seaborn et al. suggest to overwrite parts of executable binaries which can be executed by unprivileged users, but require with root privileges for special purposes<sup>3</sup> Those binaries can then be executed by the attackers and their code gets executed with root privileges.

To prevent row hammering from the NaCl environment Google disallowed the use of the `clflush` instruction [45]. However, Gruss et al. showed that memory attacks from web browsers are possible, even without using the `clflush` instruction [19] (cf. section 2.9.5).

#### 4.2.2 Flip Feng Shui

*Flip Feng Shui (FFS)* is an attack technique against virtualization servers hosting several virtual machines on the same physical host [46]. Attackers can force data to physical

---

<sup>1</sup><https://developer.chrome.com/native-client>

<sup>2</sup><https://www.google.com/chrome/>

<sup>3</sup>On Linux and UNIX operating systems, these are applications with the `setuid` bit set. `setuid` stands for set-user-ID and allows the executable to run with privileges different from the user who started it [44]. For example, this is necessary when a normal user wants to change the login password — as the password file is only accessible by the root user.

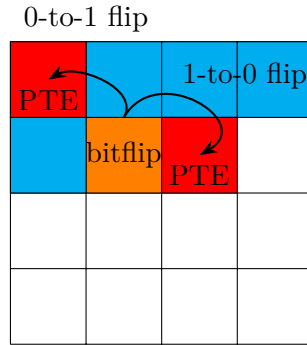


Figure 4.2: No matter in which direction a bit in the page table flips, it always points to attacker controlled memory.

locations where they can induce bit flips. The attack exploits the fact that virtualization software typically stores identical memory blocks only once — even if multiple virtual machines use it independently. Typically, virtualization software *deduplicates* memory: it regularly checks for identical memory blocks, frees all but one of them and updates all reference pointers to the removed blocks, such that they point to the remaining block containing the identical data. This is a common situation, since operating system kernels, application software and libraries need to be mapped in the memory of each virtual machine. The more identical operating systems instances run on one physical host the more blocks are deduplicated. Naturally, attackers can use various techniques (such as operating system fingerprinting, or header information of services, etc.) to learn which operating systems and applications run on virtual machines on the same physical host, and setup an identical or similar configuration. Once the attackers have done so, they can assume that the virtualization software will eventually deduplicate their memory blocks. Not attackers can craft more blocks containing data that is known to be in the victim machine, and induce bit flips after deduplication. Because the memory block is only stored once, the bit also flips on the victim machine.

Razavi et al. [46], presented a proof of concept implementation for manipulating SSH keys and domain names. When the attackers know about public SSH keys, a single bit flip makes calculation of the matching private key feasible in many cases [46]. If calculating a private key is successful attacks can login into the server. Bit flips in domain names are critical, can be used to change the domains of software repositories; For example, `ubuntv.com` differs from `ubuntu.com` by just one bit. Therefore, attackers can compromise the system by registering a domain and setting up a software repository.

### 4.3 Phys Feng Shui

Phys Feng Shui [22] is a method for placing second-level page tables at a vulnerable memory location, similar to Seaborn et al. (see section 4.2.1). The difference is that

attackers make use of the buddy allocator (cf. section 4.1) to force allocation of a second-level page table at a location vulnerable to bit flips.

Van der Veen et al. implemented a proof of concept exploit for Android devices [22]. Once their exploit has control over a PTE, it scans the entire memory for access control information (`struct cred` on Linux). As the attackers can read and write to the page (cf. section 4.2.1) they can simply increase their privileges (i.e. become root). The CPU's and operating system's access control functions cannot prevent the attacker's process from raising its own privileges, because it just changes a page in memory reserved for it, which is a perfectly valid and common operation. The whole process of Phys Feng Shui is illustrated in figure 4.3.

- (1) Exhaust the whole memory by allocating all available 512 KB blocks. Probe each allocated block for vulnerable bit flips with suitable alignment.
- (2) Allocate all remaining blocks (i.e. those smaller than 512 KB), by allocating 64 KB blocks. This step ensures, that the next 64 KB allocation will be done in the vulnerable block.
- (3) Free the vulnerable 512 KB block, so it is the only free memory block of the system.
- (4) Allocate eight 64 KB blocks. It is guaranteed that they will be mapped to the vulnerable 512 KB block.
- (5) To avoid running out of memory, immediately free all 512 KB blocks. The buddy system merges them to 1024 KB blocks. Then free the 64 KB block (containing the bit flip), so it can be allocated in the next steps.
- (6) Allocate 4 KB (page size) blocks until the vulnerable 64 KB block is used. The attackers can easily determine if the vulnerable block is used by reading internal kernel memory information from `/proc/zoneinfo` and `/proc/pagetypeinfo`, both are world readable.
- (7) Allocate 4 KB until the 4 KB block, containing the vulnerable bit flip is allocated.
- (8) Trigger a page table allocation by using the `mmap` system call with the `MAP_FIXED` flag.

#### 4.3.1 Memory Waylaying

Memory waylaying is a reliable alternative method for placing attacker-chosen memory blocks at vulnerable memory locations [5]. The method consists of two steps, the *Prefetch-based Prediction Oracle* telling attackers the location of data they want to row hammer, and repeated *Page Cache Evictions* to place the data at a location vulnerable to row hammering.

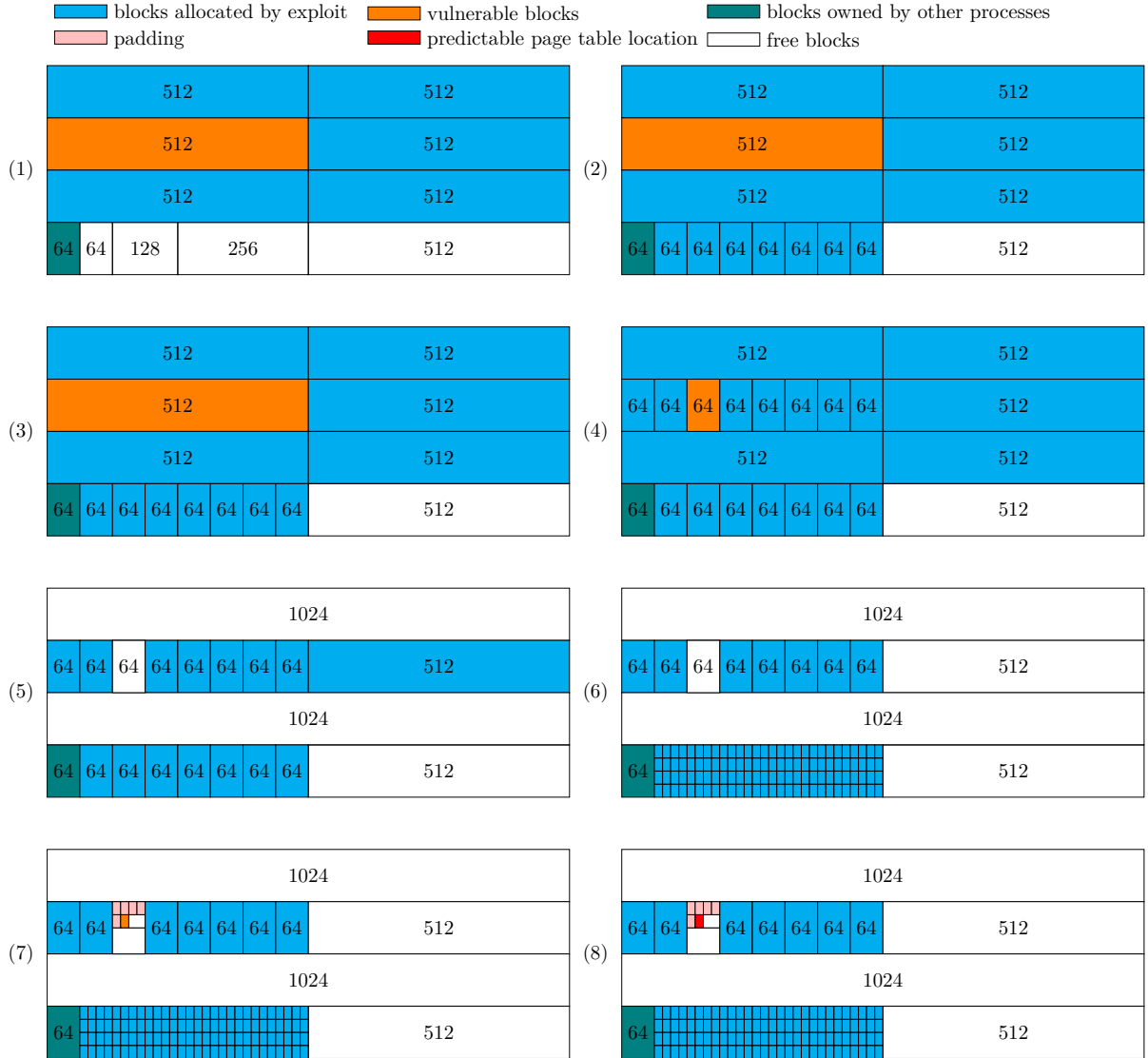


Figure 4.3: Example of Phys Feng Shui

### Prefetch-based Prediction Oracle

The first step is a timing side-channel attack telling attackers if two distinct virtual memory addresses map to the same physical address (i.e. page frame address). To be able to exploit this side-channel, attackers preliminary need to measure the time *prefetch* instructions need to load memory blocks from each CPU cache level and the main memory. Prefetch instructions allow programmers and compilers to load data into CPU caches in advance; Proper prefetching can drastically reduce cache misses and therefore increase execution speed. According to the authors Intel CPUs as well as ARM CPUs support prefetch instructions [5].

The side-channel attack works because most processors load prefetched data blocks to cache *without* without checking if the process is allowed to access them. Hence, the first step attackers need do is loading interesting data to the processor caches. For instance, if they want to find out the address of an executable file, it is sufficient to map it to memory (e.g. by using the `mmap` system call on Linux or UNIX) and flush or evict a block  $p$  containing the code of the executable binary from the CPU cache (cf. section 2.9). Then the attackers prefetch a random address  $\bar{p}$  and reload the previously flushed block. If the reload operation is fast enough (i.e corresponds with the load time for a cache level, learned in the first step) the attackers can conclude that  $p$  is loaded from cache with high probability. This is can only be the case if the page addresses by  $p$  is mapped by another application. The following list summarizes the three steps (cf. section [47]).

1. Flush an address  $p$  from page cache.
2. Prefetch an *inaccessible* address  $\bar{p}$ .
3. Reload  $p$ .

On a typical system 10 million such measurements can be done. Therefore, attackers are able to probe the entire memory, in other words then can iterative use all memory addresses for  $\bar{p}$  [47].

Gruss et al. call named step *prefetch-based prediction oracle*.

Instead of using the prefetch-based prediction oracle, attackers can learn if two memory addresses map to the same physical one by using a *Evict+Prefetch* side-channel attack. This attack is similar to the prefetch-based prediction oracle, the difference is that the attackers do not need to search for  $\bar{p}$ , instead they learn that  $p$  is used by a system call or library function [47].

The first of the attack requires more effort, as the blocks are evicted from cache, rather than flushed; This means the attackers fill the processor cache until they can assume the CPU has removed the target block from cache (see section 2.9 for details). In the second step the attackers perform steps to load a target address  $\bar{p}$  (e.g. by calling a library function or a system call). In the third step they prefetch  $p$  and measure the load-time differences. The following steps summarize the Evict+Prefetch method (cf. section [47]).

1. Evict a known address address  $p$ .
2. Execute function or system call, that accesses an inaccessible address  $\bar{p}$ .
3. prefetch  $p$ .

Initially, the prefetch-based prediction oracle and Evict+Prefetch side-channel attack were invented to bypass ASLR. However, as they can be used to tell attackers if two memory addresses are physically neighboring, they help attackers aiming to flip bits using Rowhammer [47] (see section 2.8).

### Page Cache Eviction

Evicting contents of the page cache (see section 2.6) is especially interesting for attackers intending to perform Rowhammer attacks, as the page cache stores binary executables, and evicting them may allow attackers to place them memory locations vulnerable to bit flips. Page cache eviction must not be confused with CPU cache eviction (section 2.9), which removes data from CPU caches, page cache eviction evicts data cached in primary memory that was previously loaded from secondary memory.

The *Page Cache Eviction* procedure [5] fills the entire page cache (see section 2.6) by mapping a very large file to memory and iterating over it. Therefore, it eventually evicts an executable interesting for attackers (e.g. the `sudo` binary). Because on Linux the page cache uses only memory which would remain unused otherwise (cf. section 2.6), this procedure is likely not to arise attention of system monitoring tools. For instance, the well-known `free` tool to show the amount of free memory available on the system does not count the page cache to the used memory.

Mapping very large files (several gigabyte) to memory and iterating over its content requires the kernel to remove other data from the page cache.

The `mincore` system call allows unprivileged to test if a page is currently cached in the page cache or not. This system call is also available on other UNIX-like operating systems such as FreeBSD and OpenBSD.

### Memory Waylaying

The prefetch-based prediction oracle (or rather the Evict+Prefetch side-channel) combined is called *Memory Waylaying*. Attackers need to repeatedly evict pages and check if a page in which they intend to flip a bit gets loaded to a location they can hammer. Because memory waylaying performs most operations on the page cache, it does not raise attention by increasing memory usage.

A disadvantage of the method is that it may take several hours or even days [5] until a page interesting for attackers gets loaded to a location where they can induce a bit flip.

### Memory Chasing

An alternative variant of memory waylaying is *memory chasing*. To enforce frequent relocation of pages containing the victim binary, memory chasing exploits the copy-on-write nature of the Linux kernel's process management. Gruss et al. describe the procedure as following [5].

1. Map the entire victim binary to memory using the `mmap` system call.
2. Invoke the `fork`<sup>4</sup> system call.
3. Overwrite parts of the binary in the child process. Due to the copy-on-write nature of the Linux kernel ensures that the binary gets copied to a new physical page.
4. Kill the parent process so the original unmodified pages are released.
5. Repeat this until the victim page is placed at a row hammerable vulnerable memory location. The location can be checked by using the prefetch-based prediction oracle described above.

Attackers now face the problem that with the procedure explained above the modified page (i.e. the page that child process wrote to) gets placed at the vulnerable location, therefore attackers need to trick the kernel to place the original (i.e. unmodified) page from the binary at this location. This can be done by evicting the page from the page cache using the page cache (e.g. using the strategy explained in section 4.3.1) and immediately mapping the original file again. Then the target binary is immediately released, the authors claim that this ensures that the same physical pages are used. This approach is considerably faster than memory waylaying, however it has significantly higher CPU usage and uses lots of `fork` system calls, hence it can be easier detected by intrusion detection software.

Gruss et al. implemented a proof of concept attack using memory waylaying and one-location row hammering to gain root access by flipping a bit in the `sudo` binary. The flip turns `jne` x86 assembly instruction turn into a `jne` instruction, or vice versa. Consequently another program branch is selected and the attacker can gain root access. In their analysis they figured out that the `sudo` binary has 29 locations where a bit flip can lead to a privilege escalation exploit [5].

This attack is an example for a user-space-to-user-space Rowhammer attack, it does not rely on tricking the kernel to place page tables (which belong to kernel memory space) to vulnerable locations.

---

<sup>4</sup>The `fork` system call creates a child process with an identical process image, which the kernel copies once the child process performs write operations (copy-on-write). The system call is available on most UNIX-like operating systems.



### 4.3.2 Page Types

The Linux kernel distinguishes different *mobility types* of pages, this distinction allows the kernel to group pages by their ability to be moved to other locations. This helps the kernel to store pages, logically belonging together, in physically close locations. Storing memory pages grouped is beneficial for caching, due to the principal of locality (see section 2.2). Furthermore, grouped storage helps to minimize fragmentation effects [42].

As of Linux 4.13 three different page types exist:

- `MIGRATE_UNMOVABLE` pages of core kernel components, their fixed location is essential for the operation of the kernel.
- `MIGRATE_RECLAIMABLE` pages containing data that can easily be reloaded. A typically example for pages of this type are pages containing data stored on secondary memory, which can be loaded again at any time.
- `MIGRATE_MOVABLE` pages without special position requirements, that can be moved to any physically memory location at any time. For example, pages containing data of user land applications are always movable. When pages are moved, the kernel updates the entry in the page tables referencing them. This happens completely transparently for user land applications.

### 4.3.3 Page Lists

To implement the buddy system (cf. section 4.1) the Linux kernel keeps lists of free pages for each combination, order and migrate type. For example, there is a list of movable pages of order  $n$ , one for reclaimable pages of order  $n$  and one for non-movable pages of order  $n$ , and so on.

Often systems with more than one processor have special memory areas, which are faster accessible for specific processors. These systems are called *Non-uniform memory access (NUMA)* systems [42]. When an application is executed on a specific processor, the kernel can move data and code to a memory area which is faster accessible for the processor.

Pages of order 0 are treated specially by the kernel, they are organized in page lists for each order. The kernel uses these lists for serving all memory allocations of order 0 (i.e. typically 4 KB). On non-NUMA systems (i.e. systems with only one CPU) the kernel uses only a single list.

The `rmqueue()` function (implemented in the file `mm/page_alloc.c`) is responsible for removing appropriate pages from page lists. Depending on how often the page was loaded into the CPU registers, `rmqueue()` takes the first or the last list element from the list. This is because the kernel guesses that frequently used pages (even if they are free), are still cached and consequently does not need to be reloaded to the TLB (see section 2.4).

Allocations work similarly for larger (physically continuous) allocations of more than one page. As the kernel removes one entry from the respective list, which represents the first page of the new allocation.

## 4.4 Software Diversity

The success of any Rowhammer attack depends on the attacker's ability to place a memory block at a location vulnerable to a bit flip (cf. section 2.8). If the memory placement strategy of the operating system places memory blocks in a manner, not predictable by user land applications, attackers cannot craft reliable strategies to trick operating systems into placing memory blocks at vulnerable locations. The main idea of this thesis is to apply the ideas from software diversity to the memory allocator of the Linux kernel, in order to make it as hard as possible for attackers to place memory blocks at vulnerable memory locations.

Software diversity is the idea to make the inner functioning of each copy of a software product as unique as possible, while the functionality remains equivalent. An example is a web browser, distributed over a website, each time an end user downloads a copy, the web server changes the control flow and variable locations in the web browser's binary, so each end user receives a different version of the binary. Therefore, an attacker cannot download the web browser and study its vulnerabilities and craft a large-scale automated attack (e.g. by exploiting a buffer overflow vulnerability), because the attack cannot predict where the memory area containing the buffer overflow is located. Naturally, a willing attacker who can access the binary version of the target system will — given enough time and persistence — eventually be successful, however the crafted exploit will most likely not work on a similar system using the same web browser.

In his 1992 paper, Cohen [48] explained the motivation behind software diversity:

The ultimate attack against any system begins with physical access, and proceeds to disassembly and reverse engineering of whatever programmed defenses are in place. Even with a cryptographic key provided by the user, an attacker can modify the mechanism to examine and exploit the key, given ample physical access. Eventually, the attacker can remove the defenses by finding decision points and altering them to yield altered decisions.

Without physical protection, nobody has ever found a defense against this attack, and it is unlikely that anyone ever will. The reason is that any protection scheme other than a physical one depends on the operation of a finite state machine, and ultimately, any finite state machine can be examined and modified at will, given enough time and effort. The best we can ever do is delay attack by increasing the complexity of making desired alterations.

Software diversity is also an effective method against side-channel attacks [49] and code reuse attacks [50].

# Design and Implementation of Page Sacrifice

In this chapter we describe how we implemented PAGE SACRIFICE, our proposed defense against attacks based on deterministic memory allocation. First, we give an overview about the core concepts, then we explain the functioning and effectiveness of our defense on the basis of *Phys Feng Shui* (explained in section 4.3). Finally, we discuss our prototype implementation for the Linux kernel in depth.

## 5.1 Overview

Essentially, PAGE SACRIFICE adds a random length padding both, before and after page frames and besides that it skips merging and splitting operations of the buddy system (see section 4.1). In detail padding and skipping operations work as described in the following.

- **Adding padding:** Before the memory allocator reserves a page frame, it does a random decision; Based on this decision a page frame is either reserved or skipped. Skipping means that the kernel ignores this page frame, to fulfill the request, it continues to either skip or reserve page frames until the decision is made to reserve a particular page frame or a maximum is reached. The kernel keeps track of which page frames were sacrificed while reserving each physical block. When the kernel frees page frames, it also frees all pages that were skipped while allocating it. In the following, we refer to skipped blocks are referred to as *sacrificed page frames* or simply *sacrificed pages*. Figure 5.1 illustrates a rather extreme case of PAGE SACRIFICE: A 128 KB block consisting of 32 4KB pages; Without PAGE SACRIFICE all pages are used, with PAGE SACRIFICE enabled only 11 are used and 21 are *sacrificed* pages, a sacrifice rate of ca. 66%

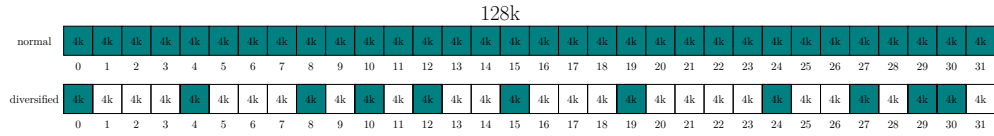


Figure 5.1: Witout PAGE SACRIFICE all blocks are deterministically allocated and used. With PAGE SACRIFICE not all pages are used, but their positioning is not predictable.

- **Skipping splitting and merging operations:** Whenever the buddy system splits a block into two buddies or merges them to a larger block, PAGE SACRIFICE makes a random decision if operation is skipped or not.
  - For **splitting operations** this means: instead of splitting blocks into two halves, the buddy system may keep its physical size and subsequently treats the block as if the splitting operation was done. This means that buddy system uses only one half of the block.
  - For **merging operations** this means: instead of merging previously split blocks the buddy systems may keep their current size. The buddy system subsequently continues to work like if one half of the block would not exist. In other words one half of the block remains sacrificed. As both buddies remain buddies they will eventually be merged later, when one of the buddies gets freed.

## 5.2 Preliminaries

To prove the functioning of our idea, we adapted the buddy system of the Linux kernel. The choice for Linux was natural as it the source code is available under a free software license.

The Linux kernel needs to store information about each page frame (cf. section 2.4), the data structure for this is `struct page`. It is important to understand, that the order and alignment of `struct page` is hardware-defined. Therefore it cannot be modified without adapting large parts of the kernel, however it is possible to append further data at its end. PAGE SACRIFICE makes use of that and extends `struct page` by a linked list used to keep track of page frame that were sacrificed during allocation of a page. This list is later used to free sacrificed pages, when the page gets freed.

The random decisions, if a splitting or merging operations shall be skipped and how many pages shall be sacrificed is done by calling the cryptographically secure kernel function `get_random_bytes()`. To avoid overhead by frequently calling this function, we buffer random bytes and refill the buffer, if not enough bits for the current random decision are available. We only use as many random bits as necessary for the current decision and remove the bits from the buffer by right shifting it. For example, if the kernel is configured to skip merging operations with probability of  $p = 0.5$ , only one bit is used.

### 5.2.1 Add Padding to Page Frames

The source code of the Linux kernel's physical memory allocation and the buddy system is implemented in several functions the `mm/page_alloc.c` file. The `rmqueue_pcplist` function is responsible for removing pages of order 0 from the CPU specific lists of free pages (cf. section 4.3.3). We modified the function to add random padding before and after the page that the functions returns. The following code listing shows the function including our modifications.

```

1 static struct page *rmqueue_pcplist
2 (
3 struct zone *preferred_zone, struct zone *zone,
4 unsigned int order, gfp_t gfp_flags, int migratetype) {
5 ...
6 local_irq_save(flags);
7 pcp = &this_cpu_ptr(zone->pageset)->pcp;
8 list = &pcp->lists[migratetype];
9
10 sacrifice_stats_total_alloc++;
11
12 sacrifice_pages(list, &(pcp->count), &preceding_sacrifice);
13
14 page = __rmqueue_pcplist(zone, migratetype, cold, pcp,
15 list);
16
17 sacrifice_pages(list, &(pcp->count), &succeeding_sacrifice);
18
19 add_sacrifice_to_page(page, &preceding_sacrifice,
20 &succeeding_sacrifice);
21 ...
22 local_irq_restore(flags);
23 return page;
24 }

```

Listing 5.1: The modified kernel function for allocation of physical page frames. Note that variable declarations are skipped for simplicity.

- Line 6-8: Unmodified kernel code, `local_irq_save` disables interrupts for the current CPU. The following other two lines get a pointer to the list of free pages, associated with the current CPU.
- Line 10: Updates statistics

- Line 12: Sacrifice page frames physically preceding the page that actually gets reserved for the request. The `sacrifice_pages` function does a random decision how many pages are sacrificed and removes them from the list of free pages.
- Line 14: Unmodified kernel code, remove the page used to fulfill the memory request.
- Line 17: Similar to Line 12, we sacrifice physically succeeding pages.
- Line 19: Associates the sacrificed pages with the page that actually gets used.
- Line 21: For clarity, we skip unmodified kernel code not relevant for explaining PAGE SACRIFICE.
- Line 22: Finally, interrupts are enabled again on the current CPU.

### 5.2.2 Implementation of Skipping Split Operations

The buddy system splits blocks into smaller blocks until it reaches a suitable block size (see section 4.1). Splitting is done by the `expand` function<sup>1</sup>

```
1 static inline void expand(  
2 struct zone *zone, struct page *page, int low, int high,  
3 struct free_area *area, int migratetype) {  
4 unsigned long size = 1 << high;  
5 while (high > low) {  
6     area--; high--; size >>= 1;  
7     ...  
8     sacrifice_stats_total_splits++;  
9     if (high == 0 && make_random_sacrifice_decision()) {  
10         sacrifice_stats_splits++;  
11         return;  
12     }  
13 list_add(&page[size].lru, &area->free_list[migratetype]);  
14 area->nr_free++;  
15 set_page_order(&page[size], high);  
16 } }
```

Listing 5.2: Skipping split operations happens in the `expand` function

The `expand` function uses a while loop to iterate to the lowest block order suitable for the allocation. If the variable `high` has value 0, the function splits two blocks of order 1 into two of order 0, we make a random decision if the operation shall be done or not. The rest of the function is kernel internal bookkeeping and not part of PAGE SACRIFICE.

---

<sup>1</sup>Note that the function name `expand` might be misleading, `expand` means to expand the amount of free blocks of a certain order — which is what the split operation does.

### 5.2.3 Implementation of Skipping Merge Operations

When the kernel frees a physical page, it calls the `__free_one_page` function and checks, if the page has a buddy or not, if not it immediately continues freeing without merging, by jumping to the `done_merging` label. PAGE SACRIFICE reuses this label, according to a random decision the buddy blocks get merge or not. It is important to mention, that skipping merging operations, does not mean sacrificed (i.e. unused) memory, since the blocks can still be used.

```

1 static inline void __free_one_page(struct page *page,
2 unsigned long pfn,
3 struct zone *zone, unsigned int order,
4 int migratetype)
5
6 ...
7 if (!page_is_buddy(page, buddy, order))
8     goto done_merging;
9
10
11 if (order == 0 && make_random_sacrifice_decision()) {
12     sacrifice_stats_merges++;
13     goto done_merging;
14 }
15 ...

```

Listing 5.3: PAGE SACRIFICE skips merge operations in the `__free_one_page` function

### 5.2.4 When page sacrifice is not suitable

Direct memory access (DMA) allows hardware components to directly access memory, without loading data to the CPU first. Hardware components using memory areas directly via DMA, are not aware of PAGE SACRIFICE and consequently expect contiguous pages in memory. Therefore PAGE SACRIFICE is not suitable for memory areas that are used for DMA.

Furthermore, no pages are sacrificed when the memory allocator is called with the `__GFP_NOFAIL` bit set. This flag indicates that the system is in a critical state and must allocate usable memory at all cost. For instance this happens when memory is scarce and the kernel needs memory for critical operations.

For performance reasons allocations that are made with the `GFP_ATOMIC` flag set are also not sacrificed. It is important to note that only the kernel can request these flags, so disabling memory in these cases, does not mean that attackers can use them to weaken our defense.

### 5.2.5 Configuration

Because the Linux kernel has a broad user base and sacrificing pages for higher security may not be an option for every use case, the feature and default parameters can be configured in the kernel configuration. Fine grained control at run time is also possible via the `sysctl` interface. Users can configure the following parameters:

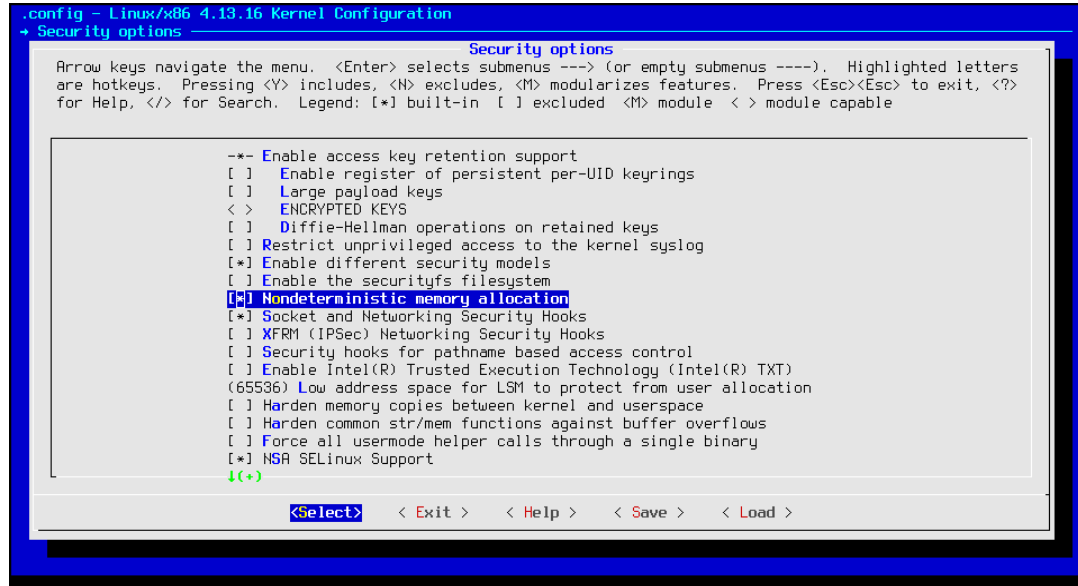


Figure 5.2: PAGE SACRIFICE can be enabled at compile time using standard Linux tools. Here the menuconfig tool shipped with the kernel is shown.

- `vm.page_sacrifice`: if set to a nonzero value PAGE SACRIFICE is enabled.
- `vm.page_sacrifice_max_at_once`: integer defining the maximum number of pages sacrificed during the allocation of one page. This defines the maximum size of a padding before and after an actually used page.
- `vm.page_sacrifice_prob_bits`: integer defining the number  $p$  of bits used for the random decision if a block shall be sacrificed or not. A bit is sacrificed if all  $p$  bits are set to zero. This means that each time a block is allocated it is sacrificed with probability of  $\frac{1}{2^p}$ . Accordingly, setting this parameter to 1 means: sacrifice with probability of  $\frac{1}{2}$ , to 2 means:  $\frac{1}{4}$  and so on.

### 5.2.6 Effectiveness of Page Sacrifice

When PAGE SACRIFICE is deployed exploits such as Phys Feng Shui (cf. section 4.3) cannot assume predictable locations of page frames. Figure 5.4 depicts an attempted Phys Feng Shui attack prevented by PAGE SACRIFICE. Additionally, determining if a



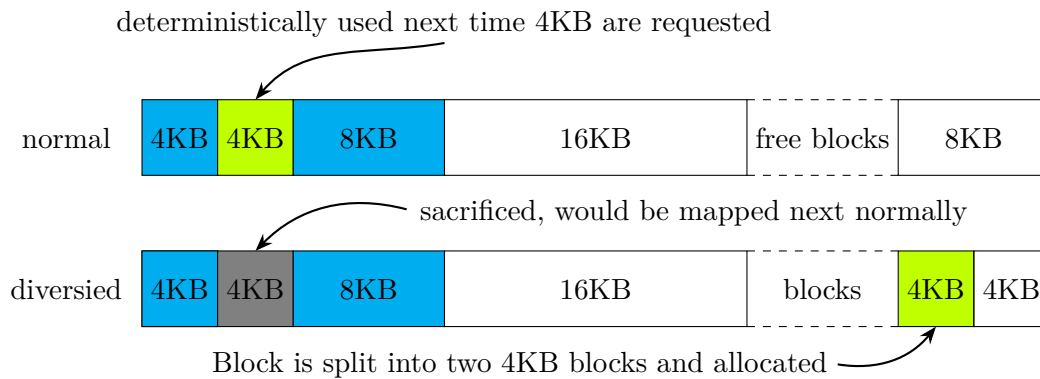


Figure 5.3: Often pages are stored at a physical location very different from where the original kernel would place it.

merging or split operation is performed becomes impossible. Another advantage of PAGE SACRIFICE is that pages are physically placed at locations far away from the usual. This is because the buddy systems does not know about the sacrificed pages, when a block is requested the buddy allocator assumes that pages are stored contiguously in memory. As this is not the case when using PAGE SACRIFICE pages may be stored at memory locations hard to predict. Figure 5.3 illustrates this effect (see also figure 4.3).

- (1) Like in figure 4.3, the whole memory is filled by allocating 512 KB blocks, here the split operation of the buddy system was skipped. So the physical 1024 KB block is seen as 512 KB block by the buddy system.
- (2) All remaining 64 KB blocks are filled. During these allocation operations two 64 KB blocks are sacrificed.
- (3) This step is equivalent to what would happen without PAGE SACRIFICE. The 512 KB block is freed.
- (4) Phys Feng Shui fills 64 KB blocks until the previously freed 512 KB block, containing the vulnerable page is filled.
- (5) This step is equivalent to what would happen, without PAGE SACRIFICE: The vulnerable 64 KB block is freed, and all remaining 512 KB blocks are freed, to avoid an out of memory situation.
- (6) Like without the PAGE SACRIFICE technique, 4 KB blocks are allocated, until the vulnerable 64 KB block is used. While allocating a large number of blocks of different size is sacrificed.
- (7) The vulnerable 64 KB block has been reached. The Phys Feng Shui inserts padding blocks until the vulnerable block is reached. Because, in this example two blocks

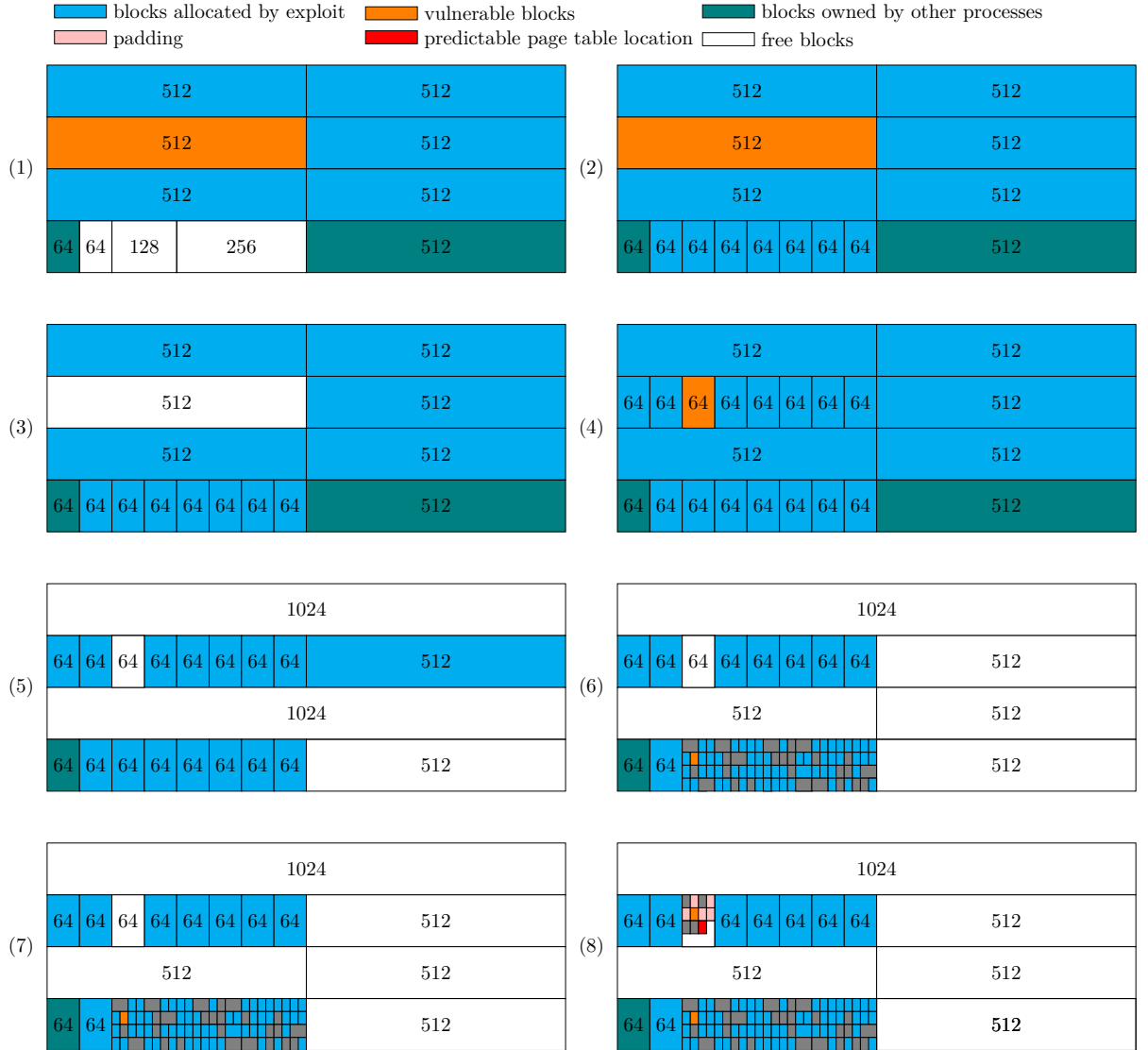


Figure 5.4: A Phys Feng Shui attack is repelled by PAGE SACRIFICE.

are sacrificed, two more padding blocks are inserted, physically right behind the vulnerable block.

- (8) The attacker forces a page table (PMD) allocation, in this case two pages are sacrificed and consequently the PMD is not stored at the block vulnerable to a bit flip.

In summary there are three reasons why Phys Feng Shui may fail. Each reason alone would prevent a successful privilege escalation exploit:

- **Unpredictable padding:** The number padding blocks required cannot be determined by attackers, because pages may be sacrificed when they are allocated.
- **Unpredictable page table location:** The page table allocation enforced by the attacker does not place the page table at a vulnerable physical block, because a random number of sacrificed pages is inserted.
- **Vulnerable block sacrificed** The block containing the vulnerable page is sacrificed and is therefore not used at all.



# Evaluation

This chapter evaluates and discusses the PAGE SACRIFICE technique we presented in the last chapter. We analyze the effects on memory usage and performance and evaluate the improved system security.

## 6.1 Memory Usage

It is in the nature of things that PAGE SACRIFICE increases memory usage. The rate of increase depends on the probability parameter  $p$ , which controls the probability for skipping buddy allocator operations and adding padding to pages. Table 6.1 shows memory usage for the Chromium and Mozilla Firefox web browsers with different parameters. To perform memory tests we set up a virtual machine using the Kernel Virtual Machine (KVM)<sup>1</sup>. We chose Debian <sup>2</sup> 9 as Linux distribution with a PAGE SACRIFICE-enabled kernel version 4.13, for performing our tests. The test environment ran only the minimum amount of processes, necessary to start our tests. Notably, the X Window System<sup>3</sup> as display server, the LightDM <sup>4</sup> display manager and the xterm <sup>5</sup> terminal emulator. After ten seconds our evaluation process killed the web browser's process and determined the amount of used memory and the total number of sacrificed pages. The total memory usage was determined by observing the total amount of physically available free memory, before and after the start of the application. It is important to note that the total memory usage is an approximate value, since the Linux does not keep track of the exact physical memory usage per process.

---

<sup>1</sup><https://www.linux-kvm.org>

<sup>2</sup><https://www.debian.org/>

<sup>3</sup><https://www.x.org/>

<sup>4</sup><https://freedesktop.org/wiki/Software/LightDM/>

<sup>5</sup><http://invisible-island.net/xterm/xterm.html>

Browser	$p$	Sacrificed pages	Mem.sacrificed	Total mem. usage	Overhead rate
Chromium	0.125	20771	81 MB	158 MB	51.27 %
	0.25	37706	147 MB	192 MB	76.56 %
	0.5	99472	389 MB	326 MB	119.33 %
	1	346016	1352 MB	875 MB	154.51 %
Firefox	0.125	12636	49 MB	245 MB	20.00 %
	0.25	34644	135 MB	308 MB	43.83 %
	0.5	94393	369 MB	478 MB	77.20 %
	1	263025	1027 MB	931 MB	110.31 %

Table 6.1: Memory overhead of web browsers

## 6.2 CPU Benchmarks

To quantify the impact of our solution on the computing power of the system, we ran the SPEC CPU2006<sup>6</sup> benchmark suite on a Intel Core i5-8350U processor on a notebook with 16 GB memory. The test system was running Debian<sup>7</sup> 9, with the GNU C compiler (GCC)<sup>8</sup> version 6.3.0 using Linux 4.13 with PAGE\_SACRIFICE enable as kernel. Figure 6.1 shows the performance degradation. A Detailed of our benchmarks can be found in appendix A.1.

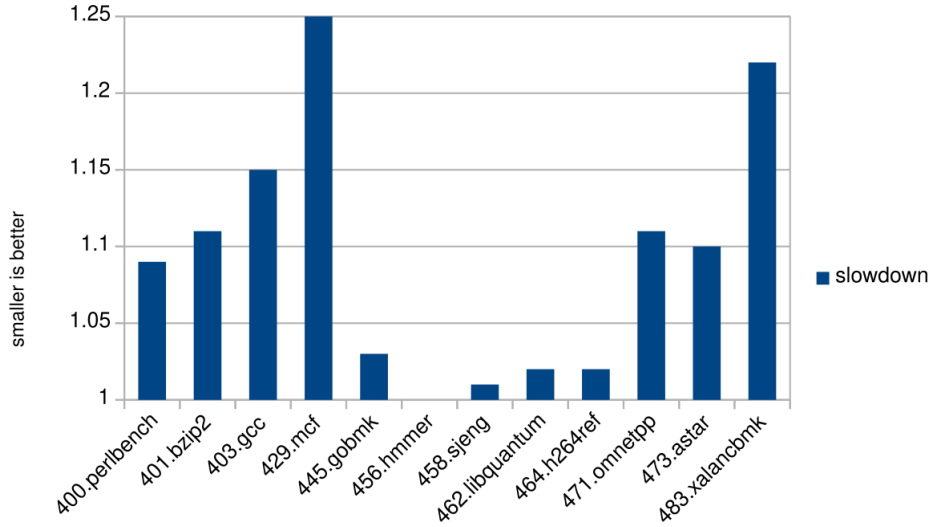


Figure 6.1: The performance degradation of PAGE\_SACRIFICE is negligible for many use cases.

<sup>6</sup><https://www.spec.org/cpu2006/>

<sup>7</sup><https://www.debian.org/>

<sup>8</sup><https://gcc.gnu.org/>

## 6.3 Security evaluation

A key idea of PAGE SACRIFICE is that processes cannot enforce to get the same page frame again. To evaluate the security of PAGE SACRIFICE, we examined how many attempts are required until Linux reserves the same physical page frame for a process, which constantly allocates and frees memory. The Linux kernel provides information on the allocated physical pages of a process in the `/proc/self/pagemap` file. This file is only accessible by the root user executing binaries with the `cap_sys_admin` captive set.

In our test environment, this allowed us to reliably access page mapping information, however in a real world attack scenario, attackers aim to escalate their privileges: so when attackers can already access the `/proc/self/pagemap` file, they most likely already have root access, hence there is not reason for further privilege escalation.

We wrote a test application simulating the Phys Feng Shui attack (see section 4.3). The tool probes how many attempts it takes until the Linux kernel assigns the same page frame, it does so by constantly allocating, filling and freeing heap memory. Because, the GNU libc library implements a cache for page allocations, we had to use a non-default version, with disabled cache. Table 6.2 shows the results of our experiments.

$p$	same page after n runs	run time
0.125	1	< 1 s
0.25	1	< 1 s
0.5	< 8	< 1 s
1	never	1.5 h

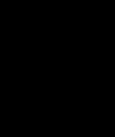
Table 6.2: Attempts needed to return the same physical page.

If at least one page is sacrificed the operating system never assigned the same page frame to the process. However, if we artificially increased the memory pressure, by running another process constantly allocating and freeing memory, we were able to force the operating system to assign the same page to the process simulating an attack. Due to lacking memory, the operating system kernel was forced to assign the same page frame again after less than 10000 iterations.

Hence, if attackers succeed to fill the entire memory and free a page frame vulnerable to a bit flip at the right moment, they will eventually be able to place a page table at a vulnerable location. Subsequently, they can hammer it to gain write-access to the whole memory, despite the PAGE SACRIFICE defense. However, in this case it is very likely that the attack fails because the attacking process runs of memory. Additionally, it is likely the attacker’s process abnormal behavior get detected by system monitoring services or system administrators.







# Conclusion

## 7.1 Contribution

In this thesis we discussed the fundamental techniques necessary for exploiting bit flips in main memory for privilege escalation attacks and discussed the underlying software and hardware principles.

We presented PAGE SACRIFICE (see chapter 5), a modification of the well-known buddy allocation mechanism, based on the principles of software diversity. Our defense dramatically increases the efforts for attackers trying to perform *Phys Feng Shui* attacks (see section 4.3).

We implemented our idea for the Linux kernel and tested it by running real world applications. Moreover, we evaluated the impact of our solution on computational power using the SPEC CPU2006 benchmark suite. To evaluate the effect on memory usage, we measured the memory usage of current web browsers, running on a Linux system equipped with PAGE SACRIFICE. While the impact on the computing power is rather moderate and might be acceptable for many use cases, we found out that the memory consumption of our solution presented in this thesis is extremely high.

## 7.2 Limitations

It turned out that sacrificing pages is not satisfying in terms of memory usage. Our experiments revealed that the memory pressure increases drastically for real world applications such as web browsers. However, we consider PAGE SACRIFICE as first part of an effective defense against Rowhammer attacks.

### 7.3 Future Work

We intend to develop PAGE SACRIFICE further to achieve an acceptable, constant memory overhead. This can either be done by keeping lists of pages allocated when memory is scarce or by constantly returning sacrificed memory in random order. Furthermore, we aim to approach the problems induced by Rowhammer and apply the ideas of software diversity to other core components of modern operating systems and compiler toolchains. Another motivation for our future research is that our work showed, that software diversity is a suitable, effective and promising method for invalidating assumptions of current and future side-channel attacks.

## SPEC CPU2006 Benchmarks

Benchmark	Run 1	Run 2	Run 3	avg.
400.perlbench	216 sec	217 sec	216 sec	216 sec
401.bzip2	370 sec	376 sec	375 sec	374 sec
403.gcc	237 sec	235 sec	236 sec	236 sec
429.mcf	314 sec	318 sec	359 sec	330 sec
445.gobmk	369 sec	377 sec	373 sec	373 sec
456.hmmer	285 sec	285 sec	285 sec	285 sec
458.sjeng	394 sec	406 sec	405 sec	402 sec
462.libquantum	249 sec	249 sec	245 sec	248 sec
464.h264ref	383 sec	384 sec	384 sec	384 sec
471.omnetpp	329 sec	339 sec	339 sec	336 sec
473.astar	349 sec	355 sec	361 sec	355 sec
483.xalancbmk	188 sec	195 sec	203 sec	195 sec

Table A.1: Detailed SPEC CPU2006 benchmarks with a PAGE SACRIFICE-enabled kernel.



# Bibliography

- [1] Moritz Lipp et al. “Meltdown: Reading Kernel Memory from User Space”. In: *27th USENIX Security Symposium (USENIX Security 18)*. 2018.
- [2] Paul Kocher et al. “Spectre Attacks: Exploiting Speculative Execution”. In: *40th IEEE Symposium on Security and Privacy (S&P’19)*. 2019.
- [3] US-CERT <https://www.us-cert.gov/ncas/alerts/TA18-141A>. *Side-Channel Vulnerability Variants 3a and 4*. May 2018. URL: <https://www.us-cert.gov/ncas/alerts/TA18-141A>.
- [4] Yoongu Kim et al. “Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors”. In: *Proceeding of the 41st Annual International Symposium on Computer Architecture*. ISCA ’14. Minneapolis, Minnesota, USA: IEEE Press, 2014, pp. 361–372. ISBN: 978-1-4799-4394-4. URL: <http://dl.acm.org/citation.cfm?id=2665671.2665726>.
- [5] Daniel Gruss et al. “Another Flip in the Wall of Rowhammer Defenses”. In: *2018 IEEE Symposium on Security and Privacy, SP 2018, Proceedings, 21-23 May 2018, San Francisco, California, USA*. 2018, pp. 245–261. DOI: 10.1109/SP.2018.00031. URL: <https://doi.org/10.1109/SP.2018.00031>.
- [6] William Stallings. *Operating Systems: Internals and Design Principles*. 9th. Upper Saddle River, NJ, USA: Prentice Hall Press, 2017. ISBN: 978-0134670959.
- [7] Ferdinand Brasser et al. “CAN’T Touch This: Software-only Mitigation Against Rowhammer Attacks Targeting Kernel Memory”. In: *Proceedings of the 26th USENIX Conference on Security Symposium*. SEC’17. Vancouver, BC, Canada: USENIX Association, 2017, pp. 117–130. ISBN: 978-1-931971-40-9. URL: <http://dl.acm.org/citation.cfm?id=3241189.3241200>.
- [8] *Four-level page tables merged*. Jan. 2005. URL: <https://lwn.net/Articles/117749/>.
- [9] *FreeBSD Architecture Handbook*. URL: <https://www.freebsd.org/doc/en/books/arch-handbook/vm-pagetable.html>.
- [10] Intel. *5-Level Paging and 5-Level EPT White Paper*. May 2017. URL: [https://software.intel.com/sites/default/files/managed/2b/80/5-level\\_paging\\_white\\_paper.pdf](https://software.intel.com/sites/default/files/managed/2b/80/5-level_paging_white_paper.pdf).

- [11] *Five-level page tables*. Mar. 2017. URL: <https://lwn.net/Articles/717293/>.
- [12] The FreeBSD Documentation Project. *Frequently Asked Questions for FreeBSD 10.X and 11.X*. URL: <https://www.freebsd.org/doc/en/books/faq/misc.html>.
- [13] Jedec Solid State Technology Association. *DDR3 SDRAM Standard*. July 2012. URL: <https://www.jedec.org/standards-documents/docs/jesd-79-3d>.
- [14] Kuljit S Bains et al. *Row hammer refresh command*. US Patent 9,236,110. 2016.
- [15] Mark Seaborn and Thomas Dullien. *Exploiting the DRAM rowhammer bug to gain kernel privileges*. Mar. 2015. URL: <https://googleprojectzero.blogspot.com/2015/03/exploiting-dram-rowhammer-bug-to-gain.html>.
- [16] S. Govindavajhala and A. W. Appel. “Using memory errors to attack a virtual machine”. In: *2003 Symposium on Security and Privacy, 2003*. May 2003, pp. 154–165. DOI: 10.1109/SECPRI.2003.1199334.
- [17] *Intel® Streaming SIMD Extensions Technology*. July 2017. URL: <https://www.intel.com/content/www/us/en/support/articles/000005779/processors.html>.
- [18] Zelalem Birhanu Aweke et al. “ANVIL: Software-Based Protection Against Next-Generation Rowhammer Attacks”. In: *Proceedings of the Twenty-First International Conference on Architectural Support for Programming Languages and Operating Systems*. ASPLOS ’16. Atlanta, Georgia, USA: ACM, 2016, pp. 743–755. ISBN: 978-1-4503-4091-5. DOI: 10.1145/2872362.2872390. URL: <http://doi.acm.org/10.1145/2872362.2872390>.
- [19] Daniel Gruss, Clémentine Maurice, and Stefan Mangard. “Rowhammer.Js: A Remote Software-Induced Fault Attack in JavaScript”. In: *Proceedings of the 13th International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment - Volume 9721*. DIMVA 2016. San Sebastián, Spain: Springer-Verlag New York, Inc., 2016, pp. 300–321. ISBN: 978-3-319-40666-4.
- [20] Yossef Oren et al. “The Spy in the Sandbox: Practical Cache Attacks in JavaScript and Their Implications”. In: *Proceedings of the 22Nd ACM SIGSAC Conference on Computer and Communications Security*. CCS ’15. Denver, Colorado, USA: ACM, 2015, pp. 1406–1418. ISBN: 978-1-4503-3832-5. DOI: 10.1145/2810103.2813708. URL: <http://doi.acm.org/10.1145/2810103.2813708>.
- [21] Fangfei Liu et al. “Last-Level Cache Side-Channel Attacks Are Practical”. In: *Proceedings of the 2015 IEEE Symposium on Security and Privacy*. SP ’15. Washington, DC, USA: IEEE Computer Society, 2015, pp. 605–622. ISBN: 978-1-4673-6949-7. DOI: 10.1109/SP.2015.43. URL: <https://doi.org/10.1109/SP.2015.43>.
- [22] Victor van der Veen et al. “Drammer: Deterministic Rowhammer Attacks on Mobile Platform”. In: *Proceedings of the 23rd Conference on Computer and Communications Security (CCS 2016)*. Oct. 2016.

- [23] Thomas M. Zeng. *The Android ION memory allocator*. Feb. 2012. URL: <https://lwn.net/Articles/480055/>.
- [24] Franz Ferdinand Brasser et al. “CAN’t Touch This: Practical and Generic Software-only Defenses Against Rowhammer Attacks”. In: *CoRR* abs/1611.08396 (2016). arXiv: 1611.08396. URL: <http://arxiv.org/abs/1611.08396>.
- [25] Franz Ferdinand Brasser et al. “CAN’t Touch This: Practical and Generic Software-only Defenses Against Rowhammer Attacks”. In: *CoRR* abs/1611.08396 (2016). arXiv: 1611.08396. URL: <http://arxiv.org/abs/1611.08396>.
- [26] Daniel Gruss et al. “KASLR is Dead: Long Live KASLR”. In: *Engineering Secure Software and Systems*. Ed. by Eric Bodden, Mathias Payer, and Elias Athanasopoulos. Cham: Springer International Publishing, 2017, pp. 161–176. ISBN: 978-3-319-62105-0.
- [27] Yueqiang Cheng, Zhi Zhang, and Surya Nepal. “Still Hammerable and Exploitable: on the Effectiveness of Software-only Physical Kernel Isolation”. In: *CoRR* abs/1802.07060 (2018). arXiv: 1802.07060. URL: <http://arxiv.org/abs/1802.07060>.
- [28] Jonathan Corbet. *The current state of kernel page-table isolation*. Dec. 2017. URL: <https://lwn.net/Articles/741878/>.
- [29] Dave Hansen. *KAISER: unmap most of the kernel from userspace page tables*. Oct. 2017. URL: <https://lwn.net/Articles/737940/>.
- [30] Gorka Irazoqui, Thomas Eisenbarth, and Berk Sunar. “MASCAT: Stopping Microarchitectural Attacks Before Execution”. In: *IACR Cryptology ePrint Archive 2016* (2016), p. 1196. URL: <http://eprint.iacr.org/2016/1196>.
- [31] Intel Corporation. *Intel® Software Guard Extensions (Intel® SGX)*. URL: <https://software.intel.com/en-us/sgx>.
- [32] David Kaplan, Jeremy Powell, and Tom Woller. *AMD Memory Encryption White Paper*. Mar. 2016. URL: [http://amd-dev.wpengine.netdna-cdn.com/wordpress/media/2013/12/AMD\\_Memory\\_Encryption\\_Whitepaper\\_v7-Public.pdf](http://amd-dev.wpengine.netdna-cdn.com/wordpress/media/2013/12/AMD_Memory_Encryption_Whitepaper_v7-Public.pdf).
- [33] Jonathan Corbet. *Defending against Rowhammer in the kernel*. Oct. 2016. URL: <https://lwn.net/Articles/704920/>.
- [34] Mathias Payer. “HexPADS: A Platform to Detect “Stealth” Attacks”. In: *Engineering Secure Software and Systems*. Ed. by Juan Caballero, Eric Bodden, and Elias Athanasopoulos. Cham: Springer International Publishing, 2016, pp. 138–154. ISBN: 978-3-319-30806-7.
- [35] Tianwei Zhang, Yinqian Zhang, and Ruby B. Lee. “CloudRadar: A Real-time Side-channel Attack Detection System in Clouds”. In: Sept. 2016.
- [36] M.Ghasempour, M. Lujan, and J.Garside. *ARMOR: A Run-time Memory Hot-Row Detector*. 2015. URL: <http://apt.cs.manchester.ac.uk/projects/ARMOR/RowHammer/armor.html>.

- [37] Jedec Solid State Technology Association. *Low Power Double Data Rate 4*. Mar. 2017. URL: <http://www.jedec.org/standards-documents/docs/jesd209-4b>.
- [38] Hewlett-Packard Inc. *HP Servers - "Rowhammer" Security Vulnerability*. Mar. 2015. URL: [https://support.hpe.com/hpsc/doc/public/display?docId=emr\\_na-c04639675](https://support.hpe.com/hpsc/doc/public/display?docId=emr_na-c04639675).
- [39] Lenovo Inc. *Row Hammer Privilege Escalation Lenovo Security Advisory: LEN-2015-009*. July 2016. URL: [https://support.lenovo.com/at/en/product\\_security/row\\_hammer](https://support.lenovo.com/at/en/product_security/row_hammer).
- [40] Apple Inc. *About the security content of Mac EFI Security Update 2015-001*. Jan. 2017. URL: <https://support.apple.com/en-us/HT204934>.
- [41] Donald E. Knuth. *The Art of Computer Programming, Volume 1 (3rd Ed.): Fundamental Algorithms*. Redwood City, CA, USA: Addison Wesley Longman Publishing Co., Inc., 1997. ISBN: 0-201-89683-4.
- [42] Wolfgang Mauerer. *Professional Linux Kernel Architecture*. Birmingham, UK, UK: Wrox Press Ltd., 2008. ISBN: 9780470343432.
- [43] Jason Evans. *A Scalable Concurrent malloc(3) Implementation for FreeBSD*. Apr. 2006. URL: <https://people.freebsd.org/~jasone/jemalloc/bsdcn2006/jemalloc.pdf>.
- [44] *setuid(2) Linux Programmer's Manual*. Sept. 2017.
- [45] JF Bastien. *Disallow the x86 "clflush" instruction due to DRAM "rowhammer" problem*. Mar. 2015. URL: <https://bugs.chromium.org/p/nativeclient/issues/detail?id=3944#c15>.
- [46] Kaveh Razavi et al. "Flip Feng Shui: Hammering a Needle in the Software Stack". In: *25th USENIX Security Symposium (USENIX Security 16)*. Austin, TX: USENIX Association, 2016, pp. 1–18. ISBN: 978-1-931971-32-4. URL: <https://www.usenix.org/conference/usenixsecurity16/technical-sessions/presentation/razavi>.
- [47] Daniel Gruss et al. "Prefetch Side-Channel Attacks: Bypassing SMAP and Kernel ASLR". In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, Vienna, Austria, October 24-28, 2016*. 2016, pp. 368–379. DOI: 10.1145/2976749.2978356. URL: <http://doi.acm.org/10.1145/2976749.2978356>.
- [48] Frederick B. Cohen. "Operating System Protection Through Program Evolution". In: *Comput. Secur.* 12.6 (Oct. 1993), pp. 565–584. ISSN: 0167-4048. DOI: 10.1016/0167-4048(93)90054-9. URL: [http://dx.doi.org/10.1016/0167-4048\(93\)90054-9](http://dx.doi.org/10.1016/0167-4048(93)90054-9).
- [49] Stephen Crane et al. "Thwarting Cache Side-Channel Attacks Through Dynamic Software Diversity." In: *NDSS*. 2015, pp. 8–11.



- [50] Andrei Homescu et al. “Profile-guided Automated Software Diversity”. In: *Proceedings of the 2013 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. CGO ’13. Washington, DC, USA: IEEE Computer Society, 2013, pp. 1–11. ISBN: 978-1-4673-5524-7. DOI: 10.1109/CGO.2013.6494997. URL: <http://dx.doi.org/10.1109/CGO.2013.6494997>.