

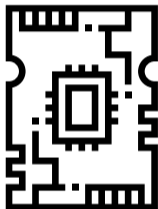


Page Cache Attacks

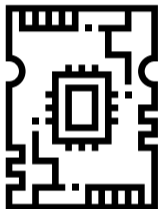
Microarchitectural Attacks on Flawless Hardware

Daniel Gruss, Trishita Tiwari, Michael Schwarz, Erik Kraft

- Modern CPUs contain multiple **microarchitectural elements**



- Modern CPUs contain multiple **microarchitectural elements**

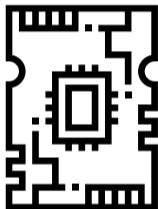


Caches and buffers



Predictors





- Modern CPUs contain multiple **microarchitectural elements**



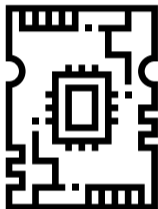
Caches and buffers



Predictors



- **Transparent** for the programmer



- Modern CPUs contain multiple **microarchitectural elements**



Caches and buffers

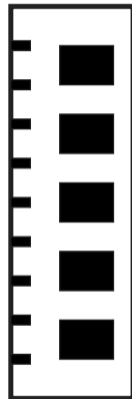
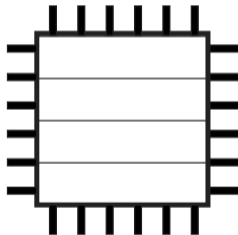


Predictors



- **Transparent** for the programmer
- Timing optimizations → side-channel leakage

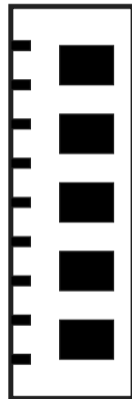
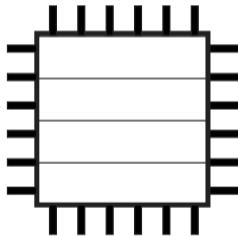
```
printf("%d", i);  
printf("%d", i);
```

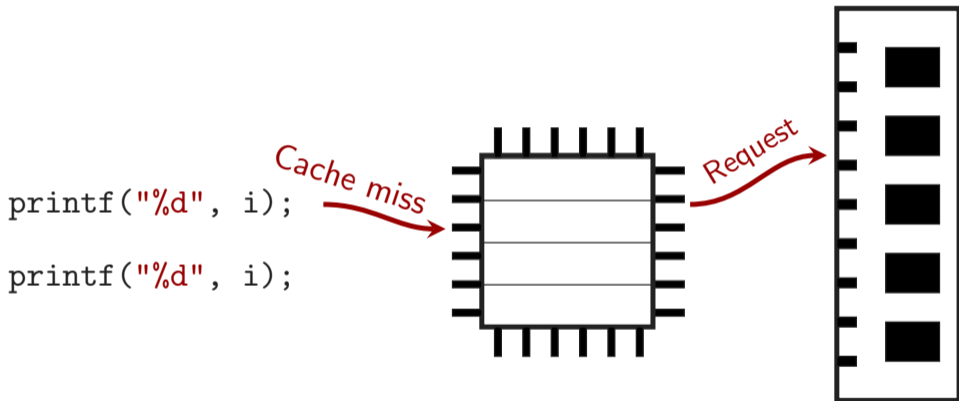


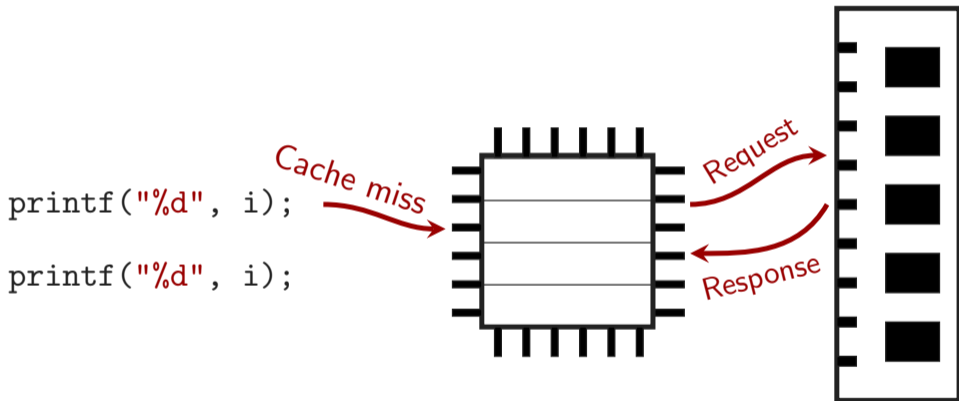
```
printf("%d", i);
```

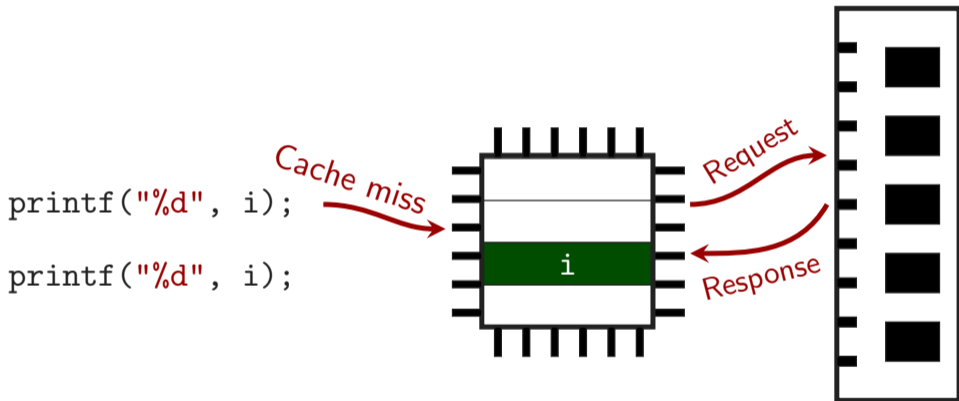
```
printf("%d", i);
```

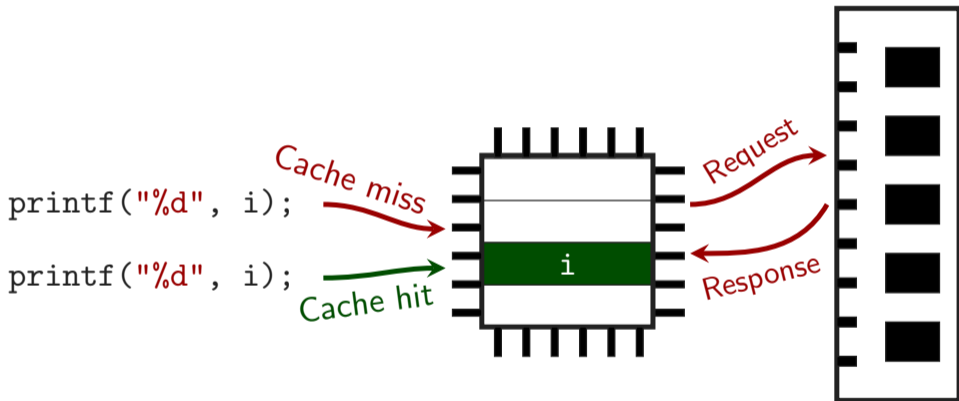
Cache miss

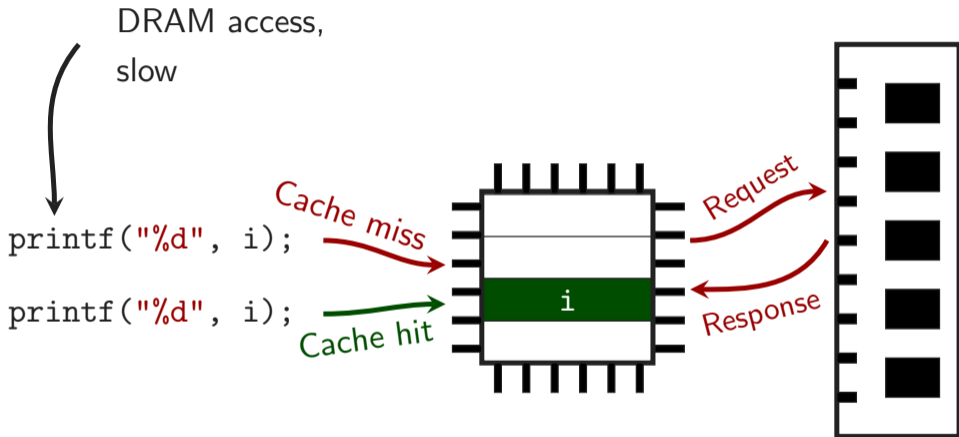


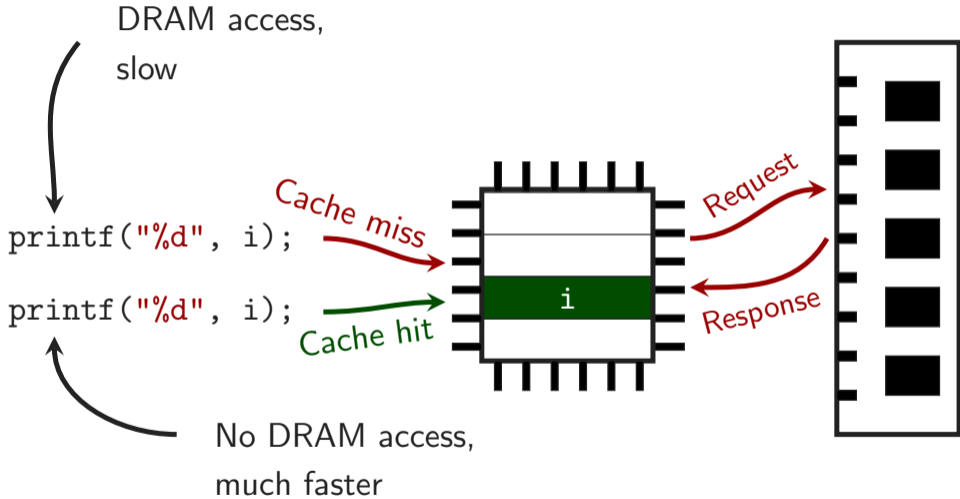


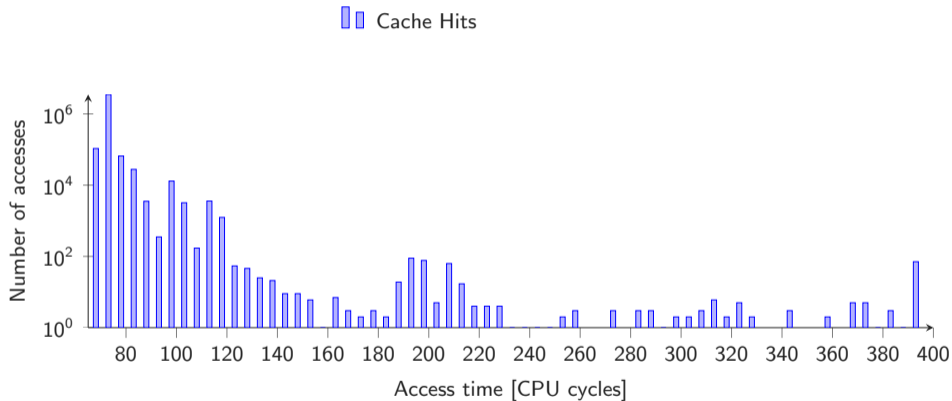


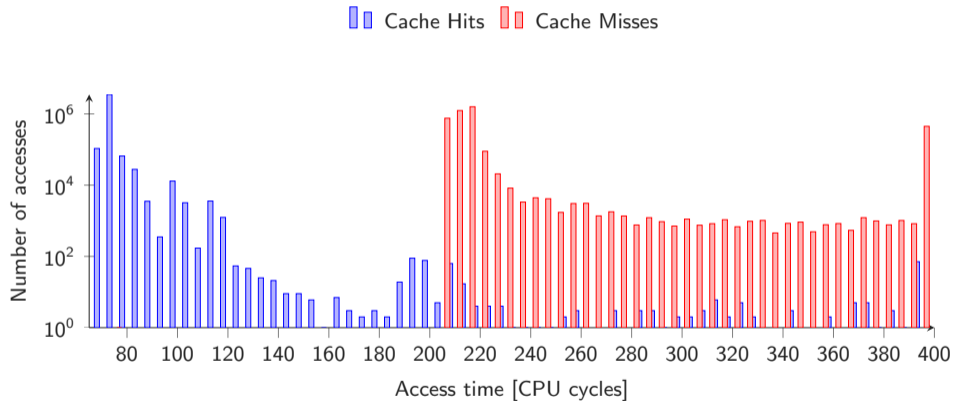


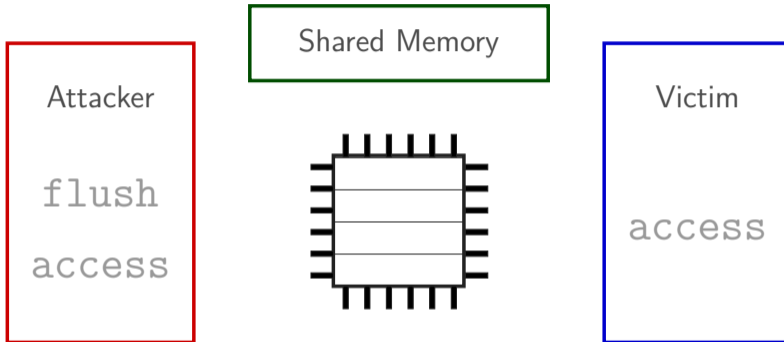


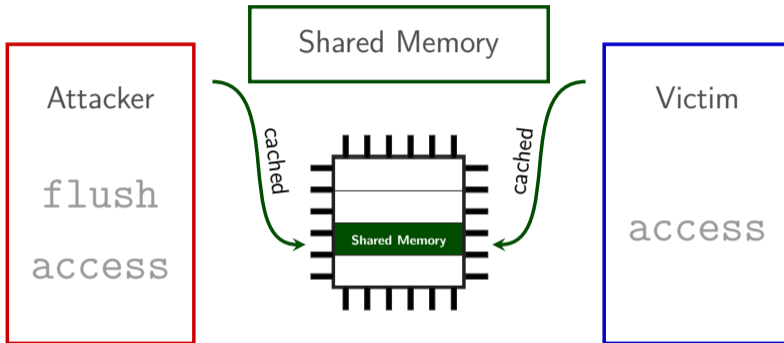


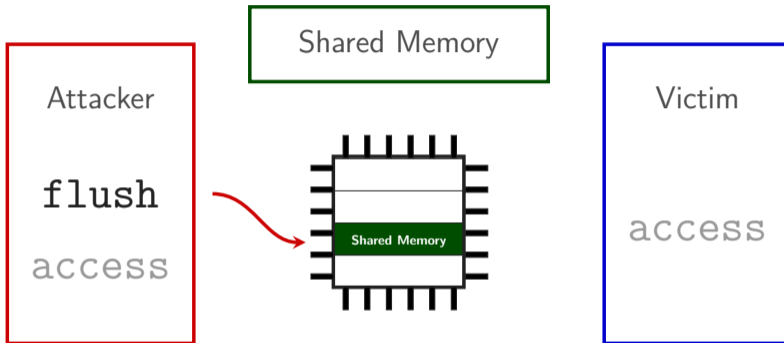


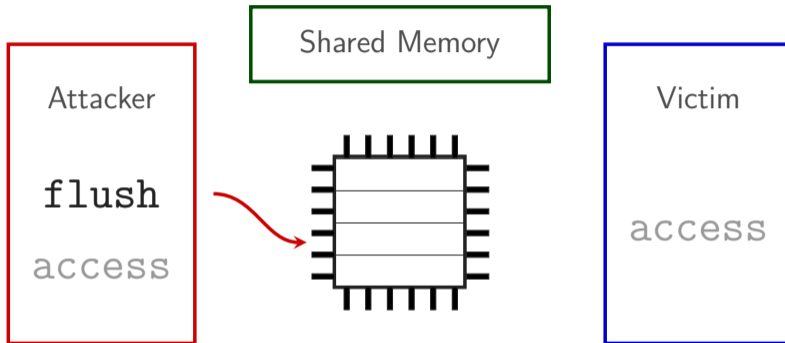


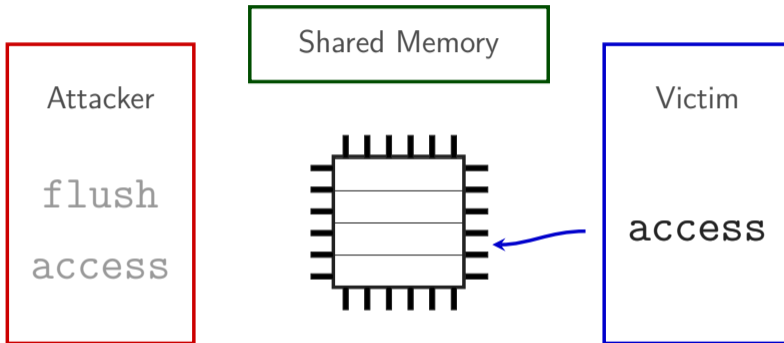


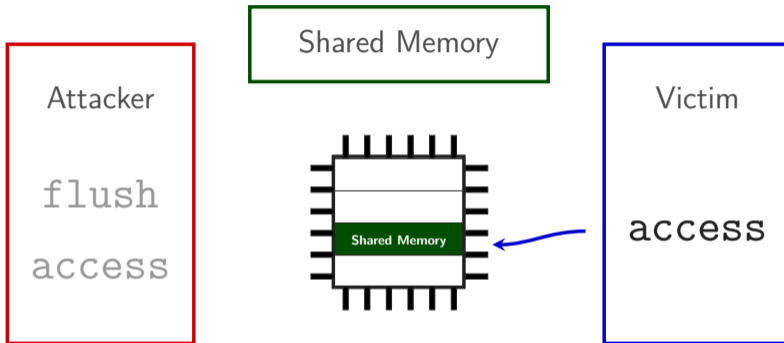


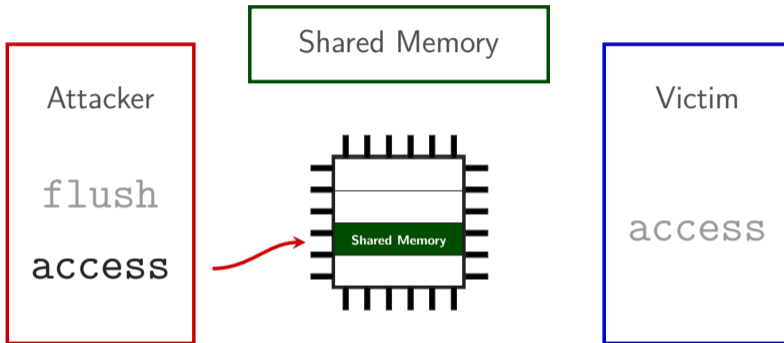


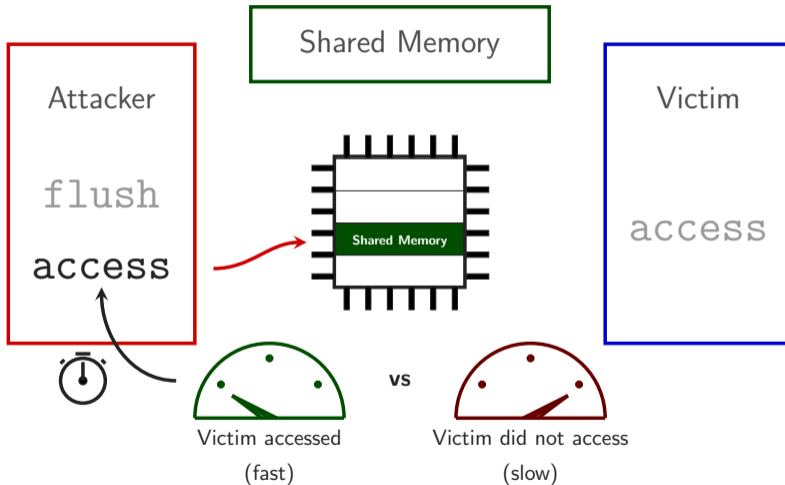


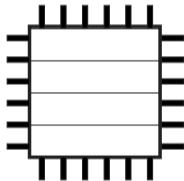


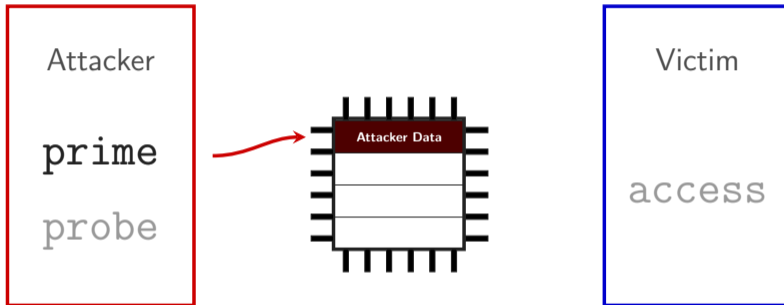


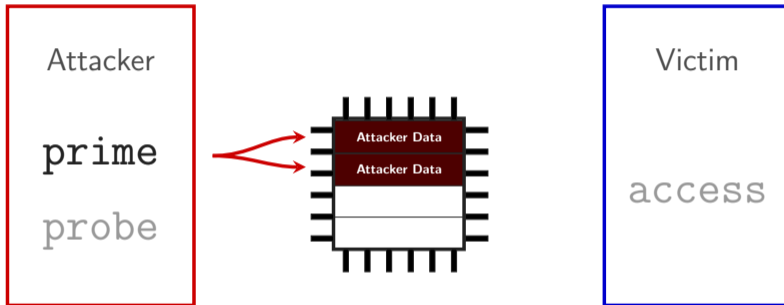


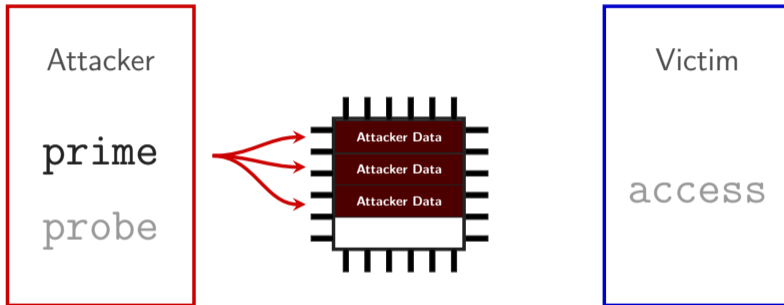


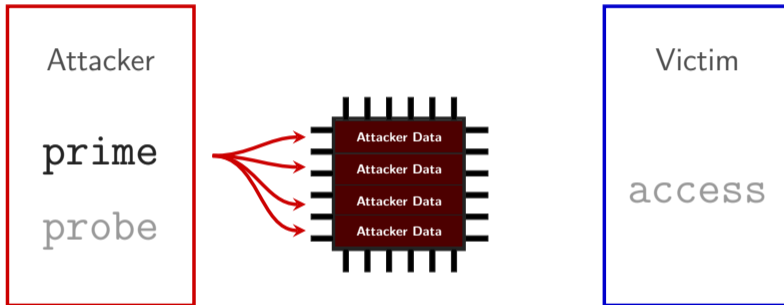


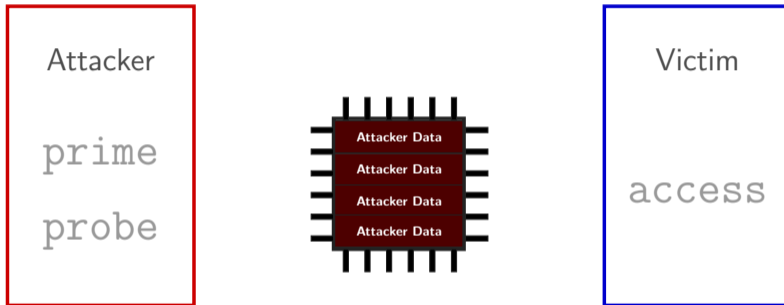


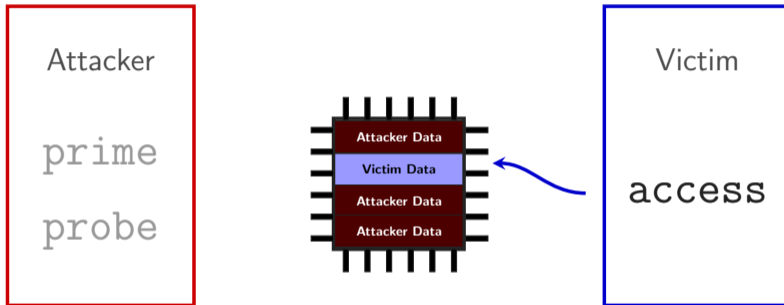


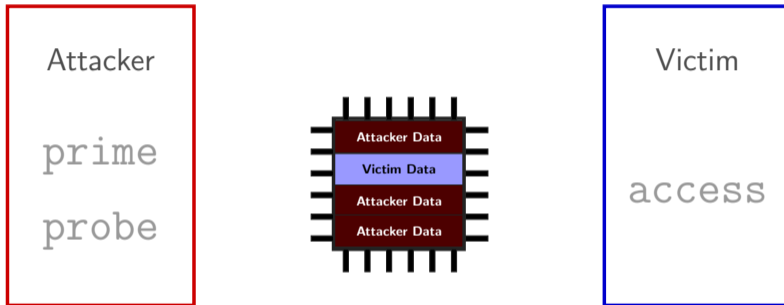


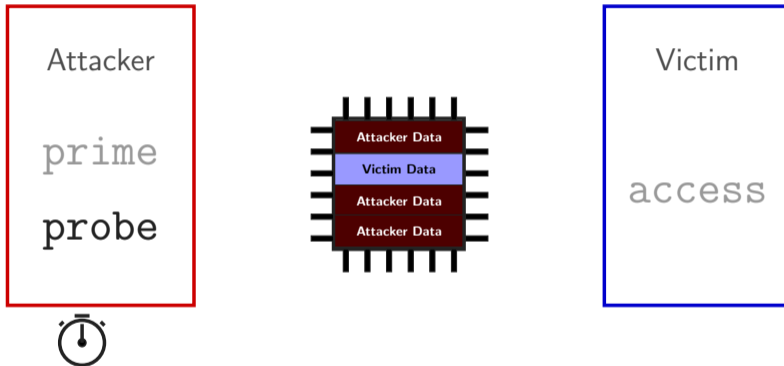


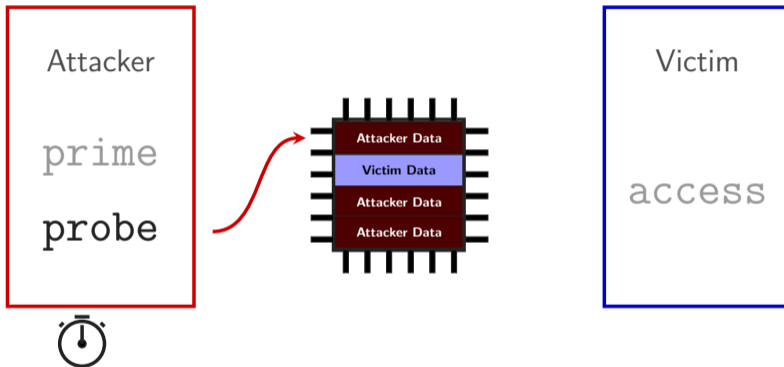


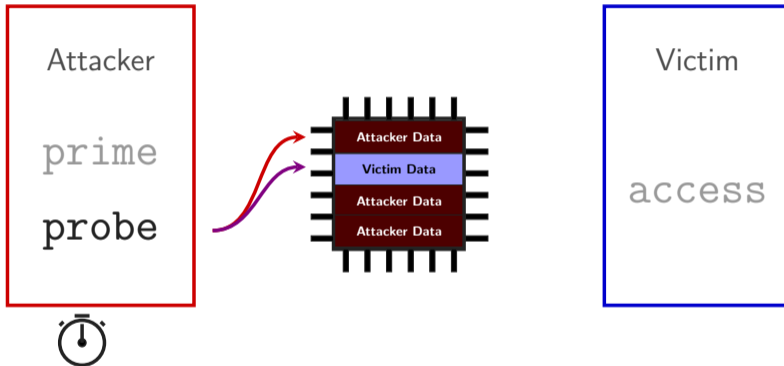


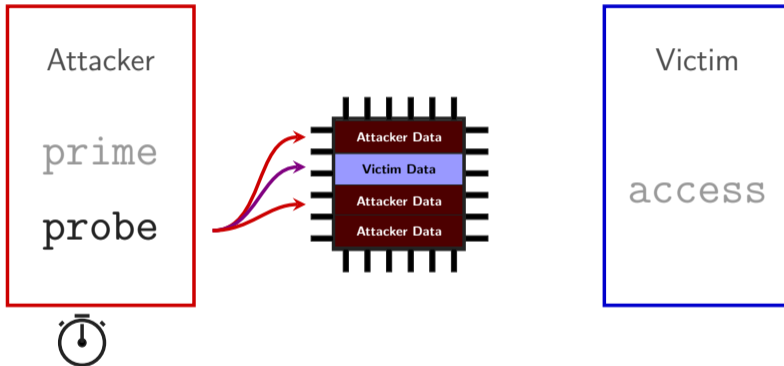


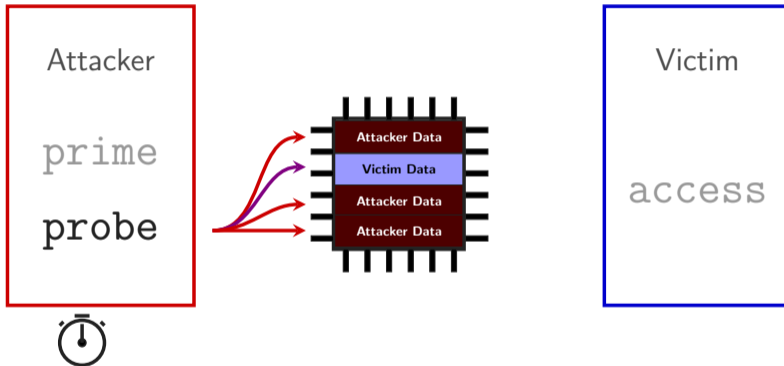


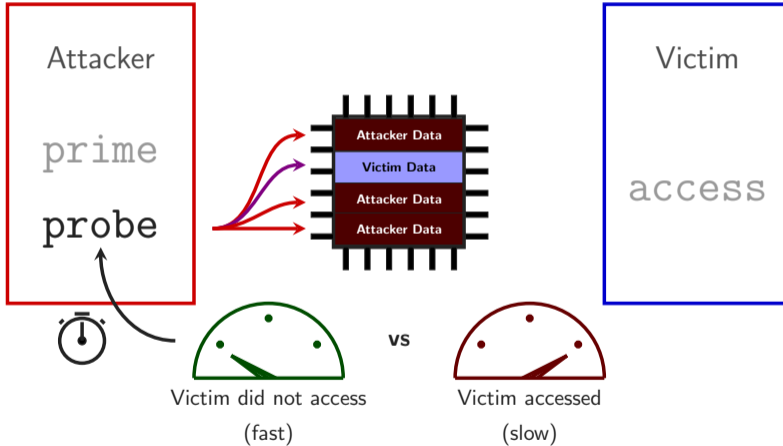


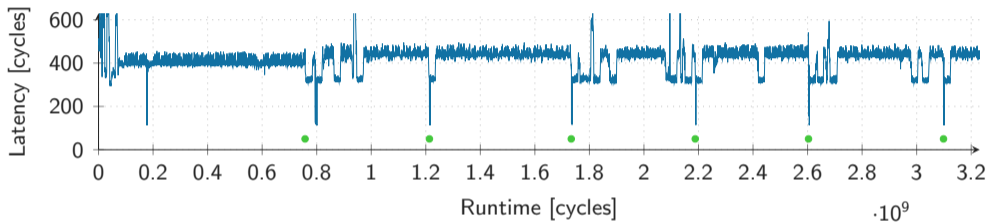




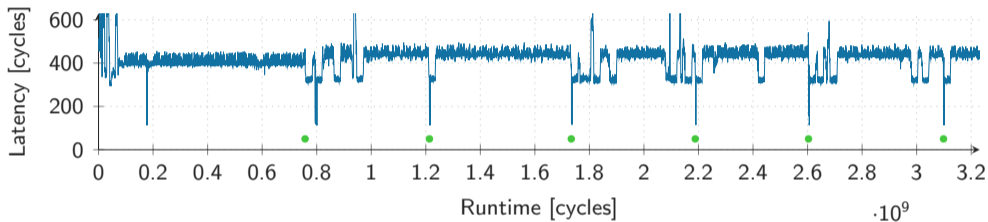




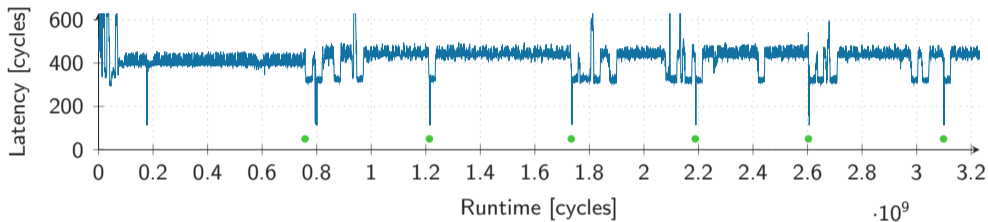




- Key presses trigger code execution in shared library (e.g., `libgdk`)



- Key presses trigger code execution in shared library (e.g., `libgdk`)
- Flush+Reload does not reveal actual key, only **time difference** between keys



- Key presses trigger code execution in shared library (e.g., `libgdk`)
- Flush+Reload does not reveal actual key, only **time difference** between keys
- → Recover text with machine learning





- Deeply rooted in hardware \rightarrow no real fixes



- Deeply rooted in hardware \rightarrow no real fixes
- More isolation \rightarrow make exploitation harder



- Deeply rooted in hardware \rightarrow no real fixes
- More isolation \rightarrow make exploitation harder
- Attacks on design difficult to fix



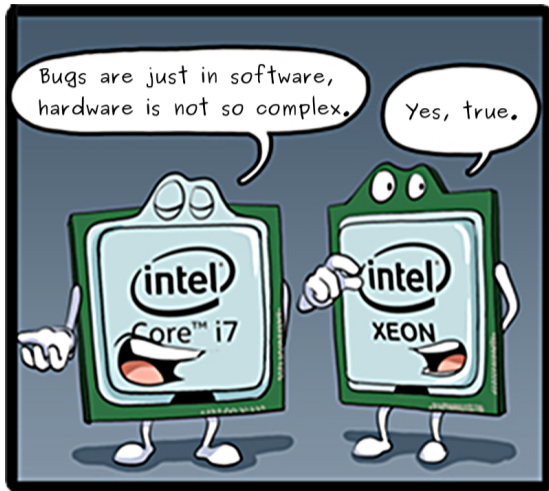
- Deeply rooted in hardware \rightarrow no real fixes
- More isolation \rightarrow make exploitation harder
- Attacks on design difficult to fix
 - Caches \rightarrow we want timing differences



- Deeply rooted in hardware \rightarrow no real fixes
- More isolation \rightarrow make exploitation harder
- Attacks on design difficult to fix
 - Caches \rightarrow we want timing differences
 - Prediction \rightarrow we don't want stalls



- Deeply rooted in hardware \rightarrow no real fixes
- More isolation \rightarrow make exploitation harder
- Attacks on design difficult to fix
 - Caches \rightarrow we want timing differences
 - Prediction \rightarrow we don't want stalls
- So far: fixing symptoms



Original image from commitstrip.com



Thought experiment: what if there were no hardware side channels?

OS

CPU

OS

CPU

μ-arch. Interface



Non-standard Syscall

x86

ISA Extension

OS

CPU

μ -arch. Interface



Non-standard Syscall

x86

ISA Extension

μ -Architecture



Software Caches



Hardware Caches

OS

CPU

μ -arch. Interface



Non-standard Syscall

x86

ISA Extension

μ -Architecture



Software Caches



Hardware Caches



Software Prefetcher



Hardware Prefetcher

OS

CPU

μ -arch. Interface



Non-standard Syscall

x86

ISA Extension

μ -Architecture



Software Caches



Hardware Caches



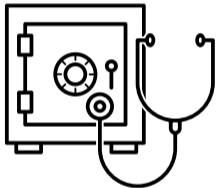
Software Prefetcher

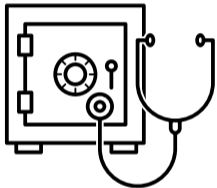


Hardware Prefetcher

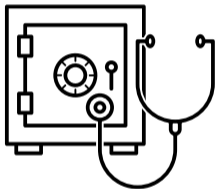




- Hardware → Software?

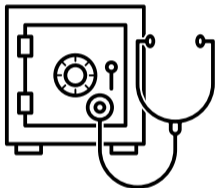






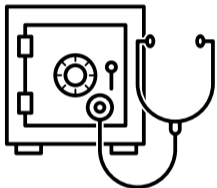
- Hardware → Software?
- Hardware-Agnostic Side Channel through the OS Page Cache





- Hardware → Software?
- Hardware-Agnostic Side Channel through the OS Page Cache
- Temporal resolution:
 -  $2 \mu\text{s}$ (≤ 6.7 measurements per second)
 -  466 ns (≤ 223 measurements per second)



- Hardware → Software?
- Hardware-Agnostic Side Channel through the OS Page Cache
- Temporal resolution:
 -  2 μ s (≤ 6.7 measurements per second)
 -  466 ns (≤ 223 measurements per second)
- Spatial resolution of 4 KiB



- Hardware → Software?
- Hardware-Agnostic Side Channel through the OS Page Cache
- Temporal resolution:
 -  $2 \mu\text{s}$ (≤ 6.7 measurements per second)
 -  466 ns (≤ 223 measurements per second)
- Spatial resolution of 4 KiB
- Various attacks: PHP RNG, UI-Redress, Windows ASLR, Keystroke Timings, Covert channels (local + remote)



- Managed by **operating system**



- Managed by **operating system**
- Buffers **file pages** in RAM for **faster accesses**



- Managed by **operating system**
- Buffers **file pages** in RAM for **faster accesses**
- Ideally all file pages in page cache

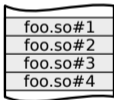


- Managed by **operating system**
- Buffers **file pages** in RAM for **faster accesses**
- Ideally all file pages in page cache
- Implemented by all major operating systems





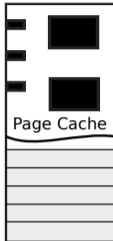
Victim



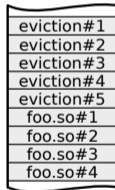
Address Space



Disk



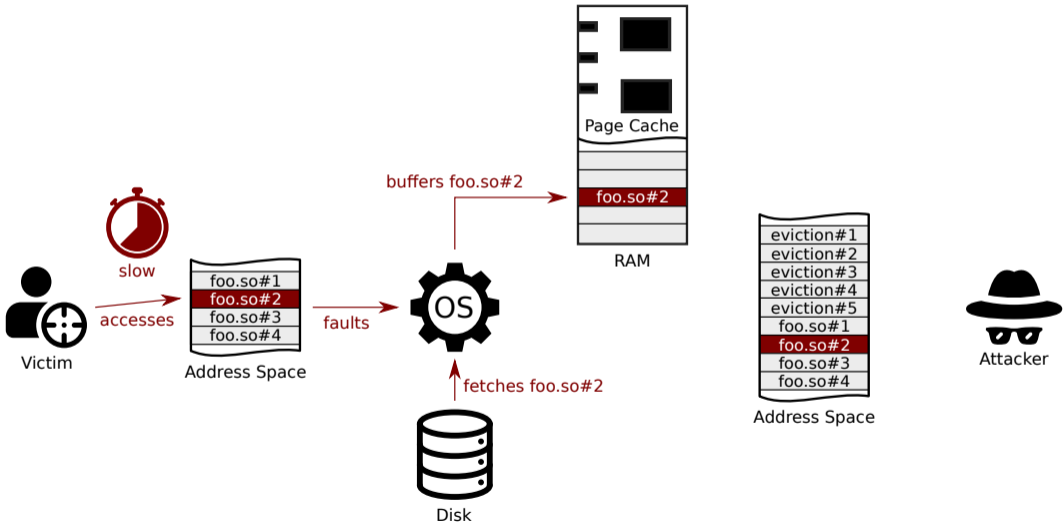
RAM



Address Space



Attacker





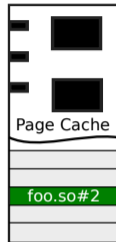
Victim



Address Space



OS



RAM



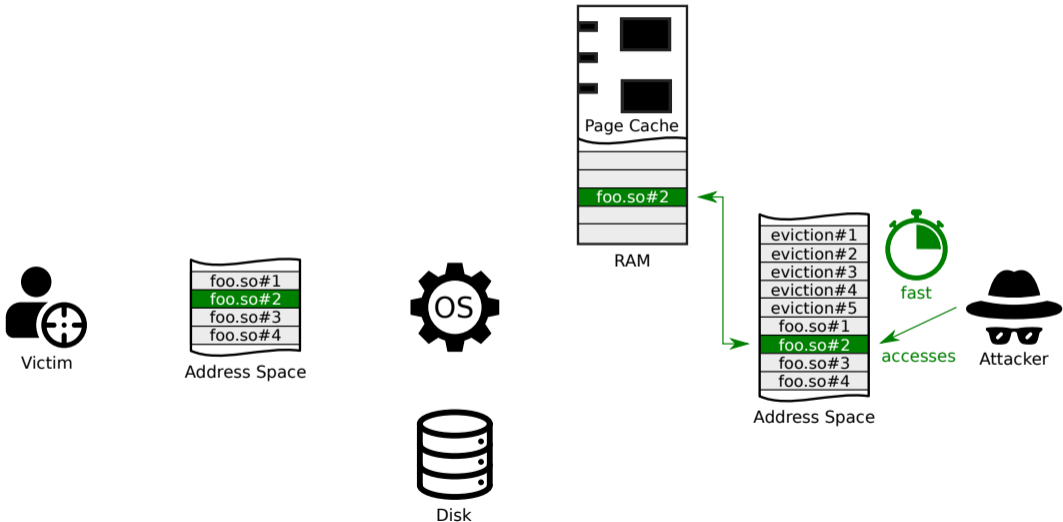
Address Space

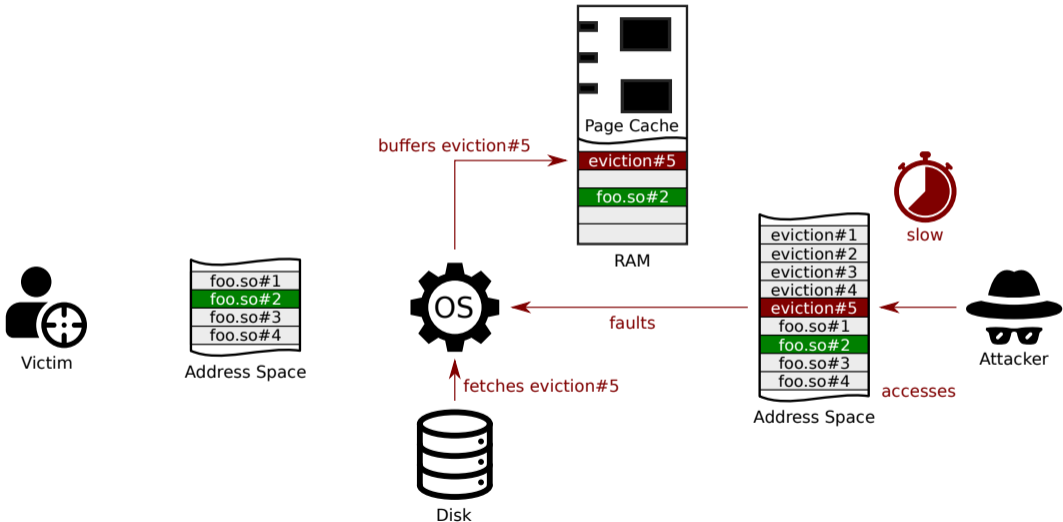


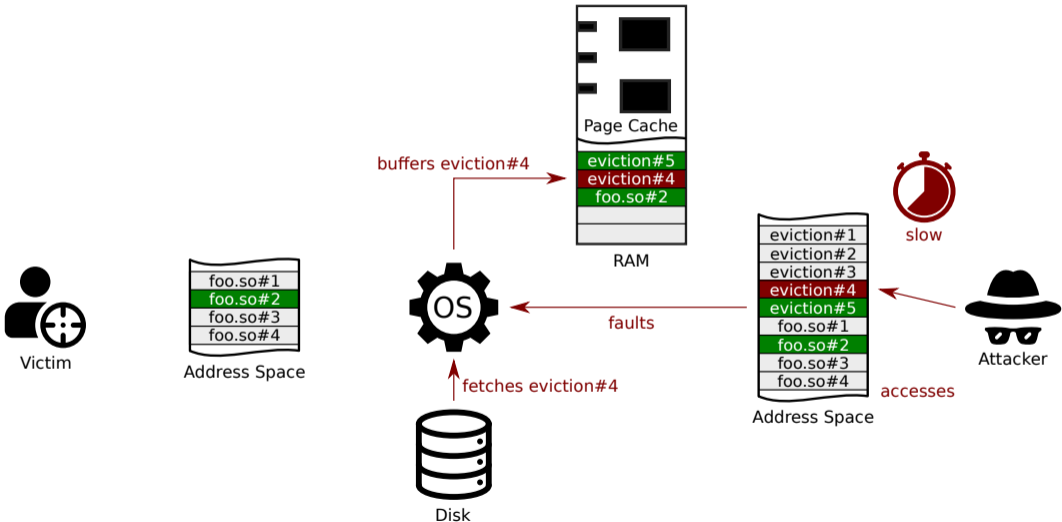
Attacker

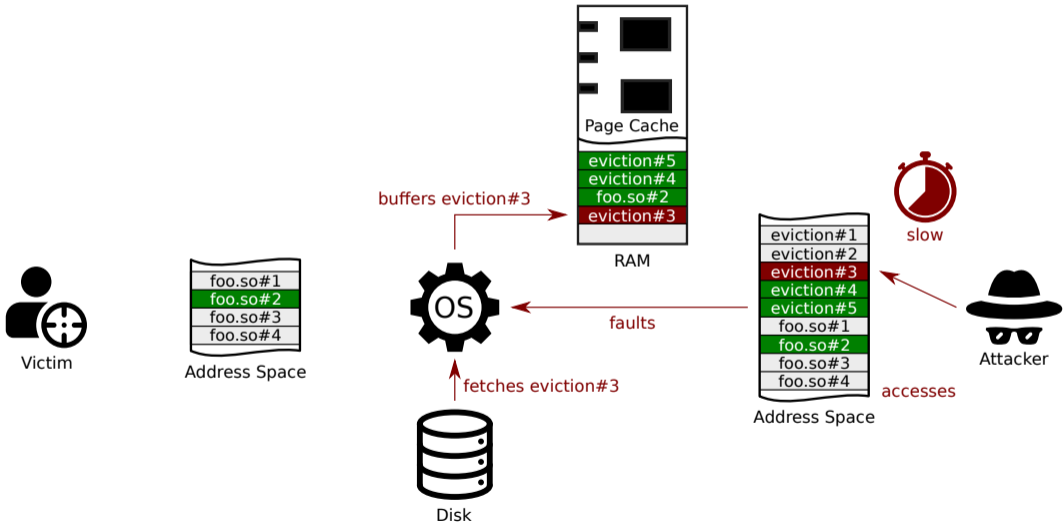


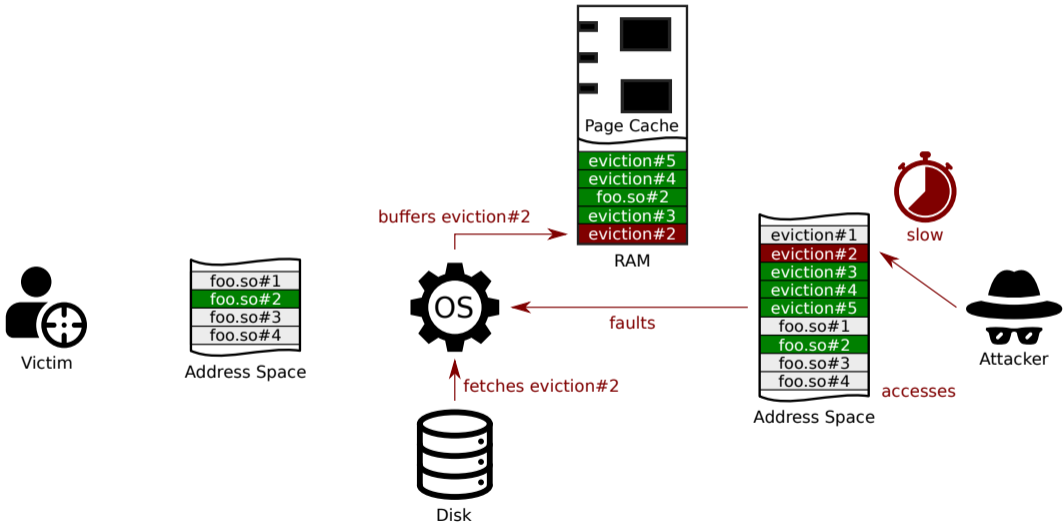
Disk

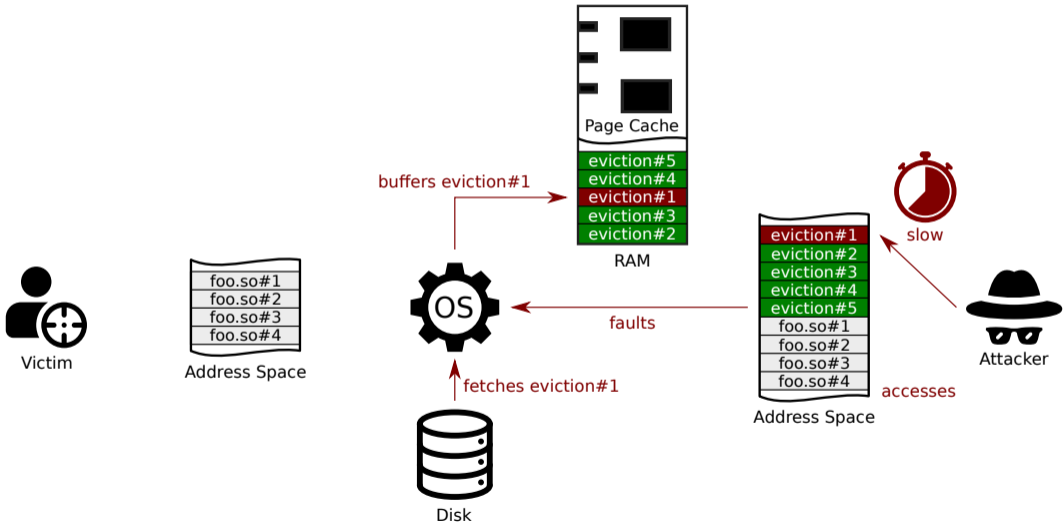






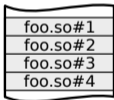








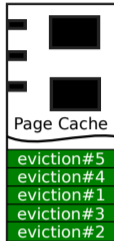
Victim



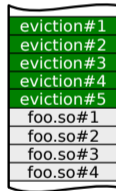
Address Space



OS



RAM



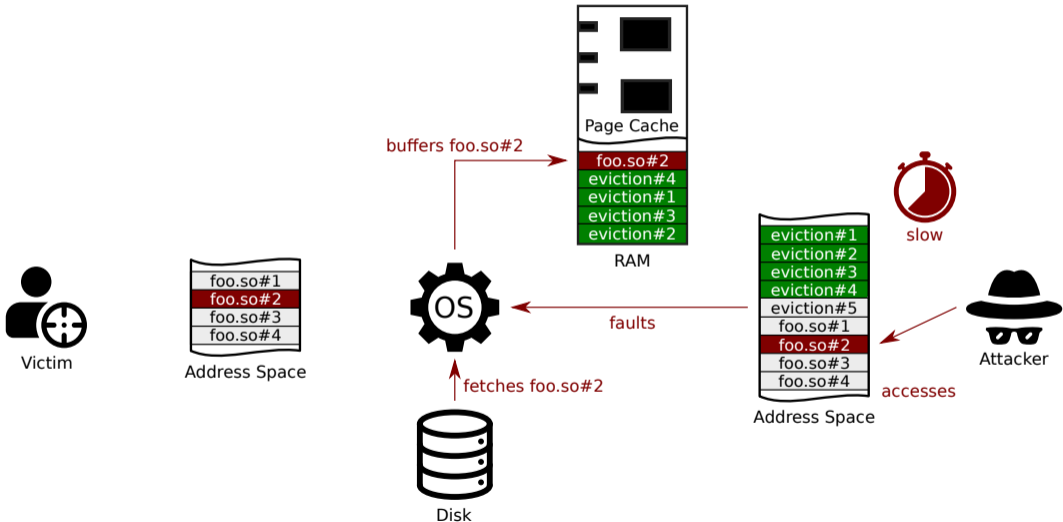
Address Space



Attacker



Disk





I AM

WATCHING YOU

First idea:



First idea:

- Measure page **access time**



First idea:

- Measure page **access time**
- Buffers pages in page cache → **destructive**



First idea:



- Measure page **access time**
- Buffers pages in page cache → **destructive**
- Eviction always necessary → **lower** average resolution



First idea:

- Measure page **access time**
- Buffers pages in page cache → **destructive**
- Eviction always necessary → **lower** average resolution

Better:



First idea:

- Measure page **access time**
- Buffers pages in page cache → **destructive**
- Eviction always necessary → **lower** average resolution

Better:

- Use **APIs** provided by the operating system



First idea:

- Measure page **access time**
- Buffers pages in page cache → **destructive**
- Eviction always necessary → **lower** average resolution

Better:

- Use **APIs** provided by the operating system
 - `mincore` 



First idea:

- Measure page **access time**
- Buffers pages in page cache → **destructive**
- Eviction always necessary → **lower** average resolution

Better:



- Use **APIs** provided by the operating system
 - mincore 
 - QueryWorkingSetEx 

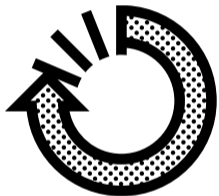


First idea:

- Measure page **access time**
- Buffers pages in page cache → **destructive**
- Eviction always necessary → **lower** average resolution

Better:

- Use **APIs** provided by the operating system
 - mincore 
 - QueryWorkingSetEx 
- Non-destructive → **higher** average resolution



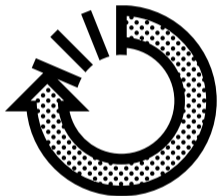
- Necessary for detecting **multiple** accesses



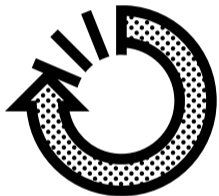
- Necessary for detecting **multiple** accesses
- **Bottleneck** of side channel



- Necessary for detecting **multiple** accesses
- **Bottleneck** of side channel
- Ideal strategy depends on memory management implementation

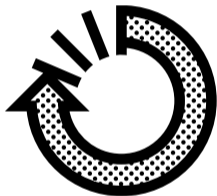




- Necessary for detecting **multiple** accesses
- **Bottleneck** of side channel
- Ideal strategy depends on memory management implementation
 - Differences in **page replacement**

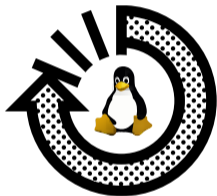


- Necessary for detecting **multiple** accesses
- **Bottleneck** of side channel
- Ideal strategy depends on memory management implementation
 - Differences in **page replacement**
 - Global CLOCK-Pro like algorithm

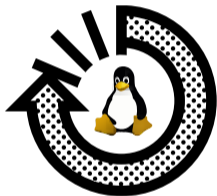




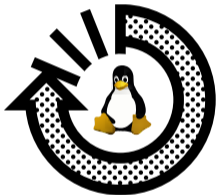
- Necessary for detecting **multiple** accesses
- **Bottleneck** of side channel
- Ideal strategy depends on memory management implementation
 - Differences in **page replacement**
 - Global CLOCK-Pro like algorithm 
 - Per-process working sets with Aging algorithm 



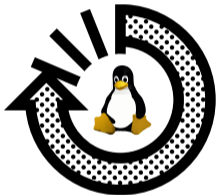
- Access pages until target page is replaced



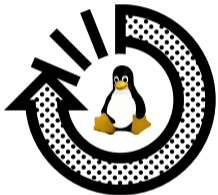
- Access pages until target page is replaced
- Basic eviction set: Large **memory-mapped file**



- Access pages until target page is replaced
- Basic eviction set: Large **memory-mapped file**
- -01: Add pages **already** in page cache

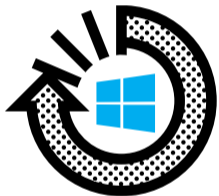


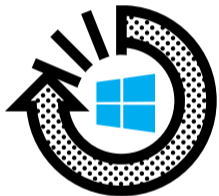
- Access pages until target page is replaced
- Basic eviction set: Large **memory-mapped file**
- -01: Add pages **already** in page cache
- -02: Fill memory with **anonymous dirty pages**



- Access pages until target page is replaced
- Basic eviction set: Large **memory-mapped file**
- -01: Add pages **already** in page cache
- -02: Fill memory with **anonymous dirty pages**
- Average run time down to **149 ms** depending on optimisations

- Page cache eviction \leftrightarrow target page drops out of **all** working sets

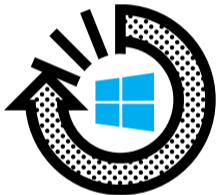




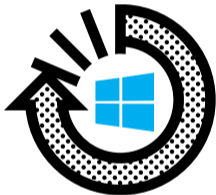
- Page cache eviction \leftrightarrow target page drops out of **all** working sets
 - Previous approach **slow...**



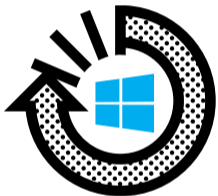
- Page cache eviction \leftrightarrow target page drops out of **all** working sets
 - Previous approach **slow...**
- Optimizations:
 - Increase ws size + memory pressure \rightarrow **self-eviction** (<2s)



- Page cache eviction \leftrightarrow target page drops out of **all** working sets
 - Previous approach **slow...**
- Optimizations:
 - Increase ws size + memory pressure \rightarrow **self-eviction** (<2s)
 - Evicting any page in **other** processes
 - \rightarrow `SetProcessWorkingSetSize` (4.48 ms)

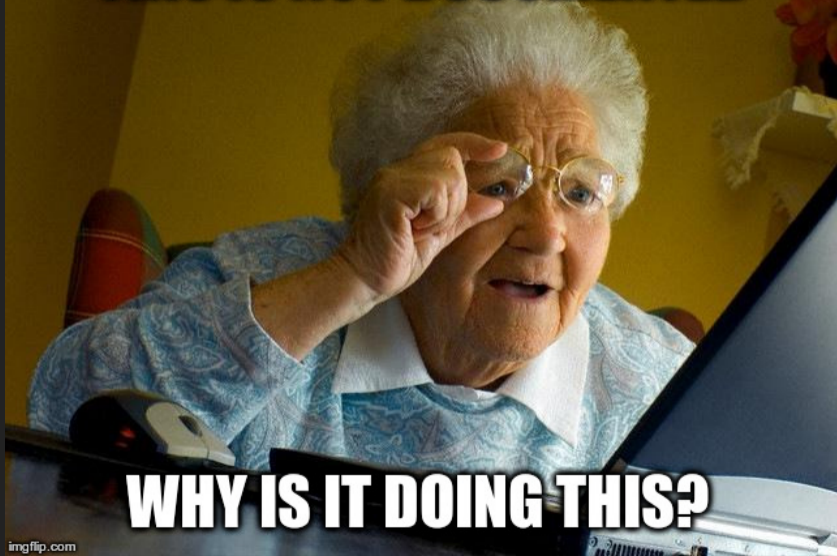


- Page cache eviction \leftrightarrow target page drops out of **all** working sets
 - Previous approach **slow...**
- Optimizations:
 - Increase ws size + memory pressure \rightarrow **self-eviction** (<2s)
 - Evicting any page in **other** processes
 - \rightarrow `SetProcessWorkingSetSize` (4.48 ms)
 - for processes with same integrity level as attacker



- Page cache eviction ↔ target page drops out of **all** working sets
 - Previous approach **slow...**
- Optimizations:
 - Increase ws size + memory pressure → **self-eviction** (<2s)
 - Evicting any page in **other** processes
 - SetProcessWorkingSetSize (4.48 ms)
 - for processes with same integrity level as attacker
 - Evicting pages you have in your **own** working set
 - VirtualUnlock (17.69 μ s)

THIS IS NOT DOCUMENTED



WHY IS IT DOING THIS?

- **Shared file** as information carrier





- **Shared file** as information carrier
- File page **presence** in page cache \leftrightarrow message bits



- **Shared file** as information carrier
- File page **presence** in page cache \leftrightarrow message bits
- Additional file pages for transmission control



- **Shared file** as information carrier
- File page **presence** in page cache \leftrightarrow message bits
- Additional file pages for transmission control

Different implementation approaches:



- **Shared file** as information carrier
- File page **presence** in page cache \leftrightarrow message bits
- Additional file pages for transmission control

Different implementation approaches:

OS	Eviction	Observation	Speed
Linux	like side channel	mincore	20.20 kB/s
	madvise	mincore	81.16 kB/s
	posix_fadvise		
Windows	process WS VirtualUnlock	QueryWorkingSetEx (ShareCount)	100.11 kB/s



- **Shared file** as information carrier
- File page **presence** in page cache \leftrightarrow message bits
- Additional file pages for transmission control

Different implementation approaches:

OS	Eviction	Observation	Speed
Linux	like side channel	mincore	20.20 kB/s
	madvise	mincore	81.16 kB/s
	posix_fadvise		
Windows	process WS VirtualUnlock	QueryWorkingSetEx (ShareCount)	100.11 kB/s

- **Low** bit error rate for all approaches (down to 0.000 03 %)

- Targets **seeding** of PHP PRNG





- Targets **seeding** of PHP PRNG
- **microtime** used as seed by some frameworks



- Targets **seeding** of PHP PRNG
- **microtime** used as seed by some frameworks
 - Returns UNIX timestamp in microseconds



- Targets **seeding** of PHP PRNG
- **microtime** used as seed by some frameworks
 - Returns UNIX timestamp in microseconds
- Later PRNG used for cryptographic operations :(



- Targets **seeding** of PHP PRNG
- **microtime** used as seed by some frameworks
 - Returns UNIX timestamp in microseconds
- Later PRNG used for cryptographic operations :(
- Side channel used to **detect** microtime call



- Targets **seeding** of PHP PRNG
- **microtime** used as seed by some frameworks
 - Returns UNIX timestamp in microseconds
- Later PRNG used for cryptographic operations :(
- Side channel used to **detect** microtime call
 - Seed recoverable



- Targets **seeding** of PHP PRNG
- **microtime** used as seed by some frameworks
 - Returns UNIX timestamp in microseconds
- Later PRNG used for cryptographic operations :(
- Side channel used to **detect** microtime call
 - Seed recoverable
- `zif_microtime` on page 781 of `php-fpm` executable



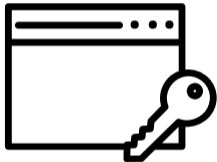
- Targets **seeding** of PHP PRNG
- **microtime** used as seed by some frameworks
 - Returns UNIX timestamp in microseconds
- Later PRNG used for cryptographic operations :(
- Side channel used to **detect** microtime call
 - Seed recoverable
- `zif_microtime` on page 781 of `php-fpm` executable
 - PHP 7.3.5, depends on build environment/configuration

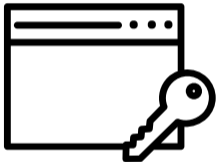


- Targets **seeding** of PHP PRNG
- **microtime** used as seed by some frameworks
 - Returns UNIX timestamp in microseconds
- Later PRNG used for cryptographic operations :(
- Side channel used to **detect** microtime call
 - Seed recoverable
- `zif_microtime` on page 781 of `php-fpm` executable
 - PHP 7.3.5, depends on build environment/configuration
- Average detection accuracy: **± 1 ms**

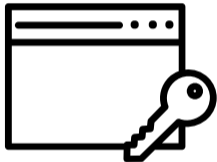


- Detect opening of interesting **window**

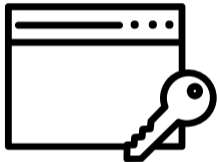




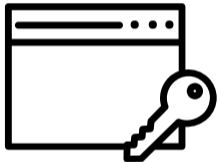
- Detect opening of interesting **window**
 - e.g. authentication windows



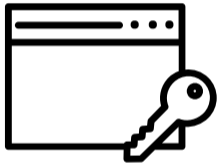
- Detect opening of interesting **window**
 - e.g. authentication windows
- Overlay original window with **fake**



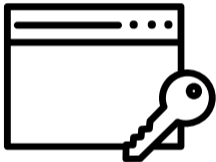
- Detect opening of interesting **window**
 - e.g. authentication windows
- Overlay original window with **fake**
- Side channel used as a **trigger**



- Detect opening of interesting **window**
 - e.g. authentication windows
- Overlay original window with **fake**
- Side channel used as a **trigger**
- Provides **low latency** → hardly noticeable



- Detect opening of interesting **window**
 - e.g. authentication windows
- Overlay original window with **fake**
- Side channel used as a **trigger**
- Provides **low latency** → hardly noticeable
- Tested with root authentication window on Ubuntu 16.04



- Detect opening of interesting **window**
 - e.g. authentication windows
- Overlay original window with **fake**
- Side channel used as a **trigger**
- Provides **low latency** → hardly noticeable
- Tested with root authentication window on Ubuntu 16.04
 - Page 6 of binary `polkit-gnome-authentication-agent-1`





- Identified as **CVE-2019-5489**



- Identified as **CVE-2019-5489**
- Linux and Windows deployed countermeasures



- Identified as **CVE-2019-5489**
- Linux and Windows deployed countermeasures



PATCHES

PATCHES EVERYWHERE



- Higher **privilege** for QueryWorkingSetEx on **other** processes



- Higher **privilege** for QueryWorkingSetEx on **other** processes
 - No direct querying of working set state



- Higher **privilege** for QueryWorkingSetEx on **other** processes
 - No direct querying of working set state
- ShareCount **hidden** for unprivileged users



- Higher **privilege** for QueryWorkingSetEx on **other** processes
 - No direct querying of working set state
- ShareCount **hidden** for unprivileged users
 - No indirect querying of working set state



- Higher **privilege** for QueryWorkingSetEx on **other** processes
 - No direct querying of working set state
- ShareCount **hidden** for unprivileged users
 - No indirect querying of working set state
- No non-destructive probing of higher-integrity processes



- Higher **privilege** for QueryWorkingSetEx on **other** processes
 - No direct querying of working set state
- ShareCount **hidden** for unprivileged users
 - No indirect querying of working set state
- No non-destructive probing of higher-integrity processes
→ **weaker** attack



- Higher **privilege** for QueryWorkingSetEx on **other** processes
 - No direct querying of working set state
- ShareCount **hidden** for unprivileged users
 - No indirect querying of working set state
- No non-destructive probing of higher-integrity processes
→ **weaker** attack
 - If QueryWorkingSetEx only possible leakage source



- Higher **privilege** for QueryWorkingSetEx on **other** processes
 - No direct querying of working set state
- ShareCount **hidden** for unprivileged users
 - No indirect querying of working set state
- No non-destructive probing of higher-integrity processes
 - **weaker** attack
 - If QueryWorkingSetEx only possible leakage source
 - Page-cache eviction already harder than on Linux



- `mincore` only reveals information for **writable** file pages



- `mincore` only reveals information for **writeable** file pages
 - Read-only files excluded → **shared** libraries, executables



- `mincore` only reveals information for **writable** file pages
 - Read-only files excluded → **shared** libraries, executables
 - Merged with the release of the 5.1.4 kernel



- `mincore` only reveals information for **writeable** file pages
 - Read-only files excluded → **shared** libraries, executables
 - Merged with the release of the 5.1.4 kernel
- Non-destructive probing no longer possible?



- `mincore` only reveals information for **writeable** file pages
 - Read-only files excluded → **shared** libraries, executables
 - Merged with the release of the 5.1.4 kernel
- Non-destructive probing no longer possible?
 - No, **`preadv2`** with `RWF_NOWAIT` leaks same information



- `mincore` only reveals information for **writeable** file pages
 - Read-only files excluded → **shared** libraries, executables
 - Merged with the release of the 5.1.4 kernel
- Non-destructive probing no longer possible?
 - No, **`preadv2`** with `RWF_NOWAIT` leaks same information
 - Not fixed yet



- We **want** the performance optimizations



- We **want** the performance optimizations
- Many side-channel attacks exploit **intended behavior**



- We **want** the performance optimizations
- Many side-channel attacks exploit **intended behavior**
- Often a **trade-off** between security and performance



- We **want** the performance optimizations
- Many side-channel attacks exploit **intended behavior**
- Often a **trade-off** between security and performance
- Every optimization is potentially a side channel



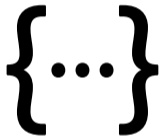
- We won't get rid of side channels



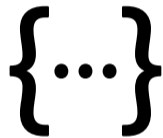
- We won't get rid of side channels
- More optimizations → more side channels



- We won't get rid of side channels
- More optimizations → more side channels
- More attacks on the “OS microarchitecture”



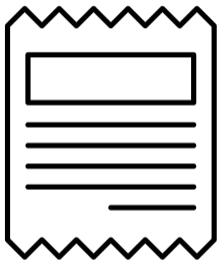
- **Abstraction** leads to side channels



- **Abstraction** leads to side channels
- Software-cache attacks are similar to hardware-cache attacks



- **Abstraction** leads to side channels
- Software-cache attacks are similar to hardware-cache attacks
- Finding **countermeasures** is **difficult**



We want to thank James Forshaw for helpful discussions on COM use cases and Simon Gunacker for early explorative work on this topic. Daniel Gruss and Michael Schwarz were supported by a generous gift from ARM and also by a generous gift from Intel. Ari Trachtenberg and Trishita Tiwari were supported, in part, by the National Science Foundation under Grant No. CCF-1563753 and Boston University's Distinguished Summer Research Fellowship, Undergraduate Research Opportunities Program, and the department of Electrical and Computer Engineering. Any opinions, findings, and conclusions or recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of the funding parties.



Page Cache Attacks

Microarchitectural Attacks on Flawless Hardware

Daniel Gruss, Trishita Tiwari, Michael Schwarz, Erik Kraft