From bottom to top: Exploiting hardware side channels in web browsers

Clémentine Maurice, Graz University of Technology July 4, 2017—RMLL, Saint-Étienne, France



+ Secure Systems team: Daniel Gruss, Michael Schwarz, Peter Pessl

• safe software infrastructure does not mean safe execution

- safe software infrastructure does not mean safe execution
- information leaks because of the underlying hardware

- · safe software infrastructure does not mean safe execution
- information leaks because of the underlying hardware
- these vulnerabilities can also be exploited at a high level

- · safe software infrastructure does not mean safe execution
- information leaks because of the underlying hardware
- these vulnerabilities can also be exploited at a high level
- like a web browser

- · safe software infrastructure does not mean safe execution
- information leaks because of the underlying hardware
- these vulnerabilities can also be exploited at a high level
- like a web browser
- because JavaScript is nothing more than code executing on your machine :)

1. What are micro-architectural side channels?

- 1. What are micro-architectural side channels?
- 2. How can I use DRAM to create a covert channel?



- 1. What are micro-architectural side channels?
- 2. How can I use DRAM to create a covert channel?
- 3. How can I do that in JavaScript?!





- no "bug" in the sense of a mistake \rightarrow lots of performance optimizations

- no "bug" in the sense of a mistake \rightarrow lots of performance optimizations
- via power consumption, electromagnetic leaks

Sources of leakage



- no "bug" in the sense of a mistake \rightarrow lots of performance optimizations
- via power consumption, electromagnetic leaks
 - $\rightarrow\,$ targeted attacks, physical access

- no "bug" in the sense of a mistake \rightarrow lots of performance optimizations
- via power consumption, electromagnetic leaks
 - $ightarrow\,$ targeted attacks, physical access
- via shared hardware and microarchitecture

- no "bug" in the sense of a mistake \rightarrow lots of performance optimizations
- via power consumption, electromagnetic leaks
 - $ightarrow\,$ targeted attacks, physical access
- · via shared hardware and microarchitecture
 - $\rightarrow \ \text{remote attacks}$



DRAM and side channels

DRAM organization













• DRAM internally is only capable of reading entire rows

- DRAM internally is only capable of reading entire rows
- capacitors in cells discharge when you "read the bits"
- · buffer the bits when reading them from the cells
- write the bits back to the cells when you're done

- DRAM internally is only capable of reading entire rows
- capacitors in cells discharge when you "read the bits"
- · buffer the bits when reading them from the cells
- write the bits back to the cells when you're done
- $\rightarrow \mbox{ row buffer}$



DRAM bank



CPU wants to access row 1



CPU wants to access row 1

ightarrow row 1 activated







DRAM bank

CPU wants to access row 1

ightarrow row 1 activated

ightarrow row 1 copied to row buffer



DRAM bank



CPU wants to access row 2



CPU wants to access row 2

ightarrow row 2 activated





CPU wants to access row 2

ightarrow row 2 activated

ightarrow row 2 copied to row buffer



DRAM bank

CPU wants to access row 2

- ightarrow row 2 activated
- ightarrow row 2 copied to row buffer



DRAM bank



CPU wants to access row 2 \rightarrow row 2 activated \rightarrow row 2 copied to row buffer

 \rightarrow slow (row conflict)


DRAM bank



CPU wants to access row 2-again



DRAM bank



CPU wants to access row 2—again \rightarrow row 2 already in row buffer



DRAM bank

CPU wants to access row 2—again \rightarrow row 2 already in row buffer



DRAM bank



CPU wants to access row 2—again \rightarrow row 2 already in row buffer \rightarrow fast (row hit)



DRAM bank



row buffer = cache





row buffers are caches

- row buffers are caches
- we can observe timing differences

- row buffers are caches
- we can observe timing differences
- how to exploit these timing differences?

- row buffers are caches
- we can observe timing differences
- how to exploit these timing differences?
- target addresses in the same channel, rank and bank

- row buffers are caches
- we can observe timing differences
- how to exploit these timing differences?
- target addresses in the same channel, rank and bank
- but DRAM mapping functions are undocumented

- row buffers are caches
- we can observe timing differences
- how to exploit these timing differences?
- target addresses in the same channel, rank and bank
- but DRAM mapping functions are undocumented
- \rightarrow we reverse-engineered them! O https://github.com/IAIK/drama

P. Pessl et al. "DRAMA: Exploiting DRAM Addressing for Cross-CPU Attacks". In: USENIX Security Symposium. 2016

• infer behavior from memory accesses similarly to cache attacks

- infer behavior from memory accesses similarly to cache attacks
- works across VMs, across cores, across CPUs

- infer behavior from memory accesses similarly to cache attacks
- works across VMs, across cores, across CPUs
- · covert channels and side-channel attacks

- infer behavior from memory accesses similarly to cache attacks
- works across VMs, across cores, across CPUs
- · covert channels and side-channel attacks
- covert channel: two processes communicating with each other
 - not allowed to do so, e.g., across VMs

- infer behavior from memory accesses similarly to cache attacks
- works across VMs, across cores, across CPUs
- · covert channels and side-channel attacks
- covert channel: two processes communicating with each other
 - not allowed to do so, e.g., across VMs
- side-channel attack: one malicious process spies on benign processes
 - e.g., spies on keystrokes

Core" 17

DRAM bank		
00000000	00000000	
00000000	00000000	
00000000	00000000	
00000000	00000000	
00000000	00000000	
row buffer		

sender and receiver agree on one bank receiver continuously accesses a row *i*



sender and receiver agree on one bank receiver continuously accesses a row *i*



DRAM bank

00000000	00000000
00000000	00000000
00000000	00000000
00000000	00000000
00000000	00000000
00000000	000000000

sender and receiver agree on one bank receiver continuously accesses a row *i*

case #1: sender transmits 1





sender and receiver agree on one bank receiver continuously accesses a row *i*

case #1: sender transmits 1

sender accesses row $j \neq i$

сору



DRAM bank

00000000	00000000
00000000	00000000
00000000	00000000
00000000	00000000
00000000	00000000
00000000	00000000

sender and receiver agree on one bank receiver continuously accesses a row *i*

case #1: sender transmits 1

sender accesses row $j \neq i$





sender and receiver agree on one bank receiver continuously accesses a row *i*

case #1: sender transmits 1

sender accesses row $j \neq i$

next receiver access \rightarrow copy row buffer

сору

Core" 17



DRAM bank 00000000 00000000 00000000 00000000 00000000 00000000 ... 00000000 00000000 ... 00000000 00000000

sender and receiver agree on one bank receiver continuously accesses a row *i*

case #1: sender transmits 1

sender accesses row $j \neq i$

next receiver access \rightarrow copy row buffer

 $\rightarrow \textbf{slow}$



DRAM bank

00000000	00000000
00000000	00000000
00000000	00000000
00000000	00000000
00000000	00000000
00000000	000000000

sender and receiver agree on one bank receiver continuously accesses a row *i*

case #2: sender transmits o



DRAM bank

00000000	00000000
00000000	00000000
00000000	00000000
00000000	00000000
00000000	00000000
00000000	000000000

sender and receiver agree on one bank receiver continuously accesses a row *i*

case #2: sender transmits o

sender does nothing





sender and receiver agree on one bank receiver continuously accesses a row *i*

case #2: sender transmits o

sender does nothing

next receiver access \rightarrow already in buffer

Core" 17



DRAM bank

sender and receiver agree on one bank receiver continuously accesses a row *i*

case #2: sender transmits o

sender does nothing next receiver access \rightarrow already in buffer

ightarrow fast

Two applications can covertly communicate with each other But can we use that for spying?





spy and victim share a row *i*



spy and victim share a row *i*

case #1

spy accesses row $j \neq i$, copy to row buffer



spy and victim share a row *i*

case #1

spy accesses row $j \neq i$, copy to row buffer **victim accesses** row *i*, copy to row buffer



spy and victim share a row *i*

case #1

spy accesses row $j \neq i$, copy to row buffer victim accesses row i, copy to row buffer spy accesses row i, no copy





DRAM bank 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000000 0000000 0000000

spy and victim share a row *i*

case #1

spy accesses row $j \neq i$, copy to row buffer **victim accesses** row i, copy to row buffer spy accesses row i, no copy

ightarrow fast



spy and victim share a row *i*

case #2

spy accesses row $j \neq i$, copy to row buffer
DRAMA side-channel attacks





spy and victim share a row *i*

case #2

spy accesses row $j \neq i$, copy to row buffer **no victim access** on row *i*

DRAMA side-channel attacks



spy and victim share a row *i*

case #2

spy accesses row $j \neq i$, copy to row buffer **no victim access** on row *i* spy accesses row *i*, copy to row buffer

DRAMA side-channel attacks





DRAM bank 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000000 0000000 0000000

spy and victim share a row *i*

case #2

spy accesses row $j \neq i$, copy to row buffer **no victim access** on row *i* spy accesses row *i*, copy to row buffer

 $\rightarrow \text{slow}$

Spying on keystrokes on the Firefox URL bar

- side-channel: template attack
 - allocate a large fraction of memory to be in a row with the victim
 - profile memory and record row-hit ratio for each address



I'm sure we'll need to write a lot of C code At least we're safe with JavaScript!

Member Rowhammer.js?



DRAM covert channels in JavaScript?

• JavaScript is code executed in a sandbox

- JavaScript is code executed in a sandbox
- can't do anything nasty since it is in a sandbox, right?

- JavaScript is code executed in a sandbox
- can't do anything nasty since it is in a sandbox, right?
- except side channels are only doing benign operations

- JavaScript is code executed in a sandbox
- can't do anything nasty since it is in a sandbox, right?
- except side channels are only doing benign operations
 - 1. accessing their own memory

- JavaScript is code executed in a sandbox
- can't do anything nasty since it is in a sandbox, right?
- except side channels are only doing benign operations
 - 1. accessing their own memory
 - 2. measuring time

Challenges with JavaScript







1. No knowledge about physical addresses

2. No instruction to flush the cache

3. No high-resolution timers

- OS optimization: use Transparent Huge Pages (THP, 2MB pages)
- last 21 bits (2MB) of physical address
- = last 21 bits (2MB) of virtual address

- OS optimization: use Transparent Huge Pages (THP, 2MB pages)
- last 21 bits (2MB) of physical address
- = last 21 bits (2MB) of virtual address
- $\rightarrow\,$ which JS array indices?

#1. Obtaining the beginning of a THP



- physical pages for these THPs are mapped on-demand
- ightarrow page fault when an allocated THP is accessed for the first time

D. Gruss et al. "Practical Memory Deduplication Attacks in Sandboxed JavaScript". In: ESORICS'15. 2015.

- we now know the last 21 bits of physical addresses
- enough for most systems, e.g., Sandy Bridge with DDR3





- measure DRAM timing
- only non-cached accesses reach DRAM
- no clflush instruction
- $\rightarrow~{\rm evict}~{\rm data}~{\rm with}~{\rm other}~{\rm memory}~{\rm accesses}$



D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.



D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.



D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.



D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.



D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.



D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.



D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.



D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.



D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.



• it's a bit more complicated than that: replacement policy is not LRU

D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.



- it's a bit more complicated than that: replacement policy is not LRU
- but we already solved this problem before :)

D. Gruss et al. "Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript". In: DIMVA'16. 2016.

• measure small timing differences: need a high-resolution timer

- measure small timing differences: need a high-resolution timer
- native: rdtsc, timestamp in CPU cycles

- measure small timing differences: need a high-resolution timer
- native: rdtsc, timestamp in CPU cycles
- JavaScript: performance.now() has the highest resolution

- measure small timing differences: need a high-resolution timer
- native: rdtsc, timestamp in CPU cycles
- JavaScript: performance.now() has the highest resolution

performance.now()

[...] represent times as floating-point numbers with up to microsecond precision. — Mozilla Developer Network

High-resolution timers in JavaScript

• before September 2015: performance.now() had a nanosecond resolution

Y. Oren et al. "The Spy in the Sandbox: Practical Cache Attacks in JavaScript and their Implications". In: CCS'15. 2015. https://www.mozilla.org/en-US/security/advisories/mfsa2015-114/

- before September 2015: performance.now() had a nanosecond resolution
- Oren et al. demonstrated cache side-channel attacks in JavaScript

Y. Oren et al. "The Spy in the Sandbox: Practical Cache Attacks in JavaScript and their Implications". In: CCS'15. 2015. https://www.mozilla.org/en-US/security/advisories/mfsa2015-114/
- before September 2015: performance.now() had a nanosecond resolution
- Oren et al. demonstrated cache side-channel attacks in JavaScript
- "fixed" in Firefox 41: rounding to 5 µs

Y. Oren et al. "The Spy in the Sandbox: Practical Cache Attacks in JavaScript and their Implications". In: CCS'15. 2015. https://www.mozilla.org/en-US/security/advisories/mfsa2015-114/

Firefox < 41 (1 ns) $1 \cdot 10^{-3}$









Firefox < 41 (1 ns) $1 \cdot 10^{-3}$



D. Kohlbrenner et al. "Trusted Browsers for Uncertain Times". In: USENIX Security Symposium. 2016

microsecond resolution is not enough

M. Schwarz et al. "Fantastic Timers and Where to Find Them: High-Resolution Microarchitectural Attacks in JavaScript". In: FC'17. 2017.

- microsecond resolution is not enough
- two approaches

M. Schwarz et al. "Fantastic Timers and Where to Find Them: High-Resolution Microarchitectural Attacks in JavaScript". In: FC'17. 2017.

- microsecond resolution is not enough
- two approaches
 - 1. recover a higher resolution from the available timer

M. Schwarz et al. "Fantastic Timers and Where to Find Them: High-Resolution Microarchitectural Attacks in JavaScript". In: FC'17. 2017.

- microsecond resolution is not enough
- two approaches
 - 1. recover a higher resolution from the available timer
 - 2. build our own high-resolution timer

M. Schwarz et al. "Fantastic Timers and Where to Find Them: High-Resolution Microarchitectural Attacks in JavaScript". In: FC'17. 2017.





• measure how often we can increment a variable between two timer ticks



• to measure with high resolution



- to measure with high resolution
 - start measurement at clock edge



- to measure with high resolution
 - start measurement at clock edge
 - increment a variable until next clock edge



- to measure with high resolution
 - start measurement at clock edge
 - increment a variable until next clock edge
- Firefox/Chrome: 500 ns, Tor: 15 μs

·	
•	
•	
•	





ightarrow padding so the slow function crosses one more clock edge than the fast one

Recovering resolution: Edge thresholding



Recovering resolution: Edge thresholding



nanosecond resolution

Recovering resolution: Edge thresholding



- nanosecond resolution
- Firefox/Tor: 2 ns, Edge: 10 ns, Chrome: 15 ns

• goal: counter that does not block main thread

- goal: counter that does not block main thread
- baseline setTimeout: 4 ms (except Edge: 2 ms)

- goal: counter that does not block main thread
- baseline setTimeout: 4 ms (except Edge: 2 ms)
- CSS animation \rightarrow increase width of element as fast as possible

- goal: counter that does not block main thread
- baseline setTimeout: 4 ms (except Edge: 2 ms)
- CSS animation \rightarrow increase width of element as fast as possible
- timestamp = width of element

- goal: counter that does not block main thread
- baseline setTimeout: 4 ms (except Edge: 2 ms)
- + CSS animation \rightarrow increase width of element as fast as possible
- timestamp = width of element
- but animation limited to 60 fps \rightarrow 16 ms resolution

• JavaScript can spawn new threads called web worker

- JavaScript can spawn new threads called web worker
- web worker communicate using message passing

- JavaScript can spawn new threads called web worker
- web worker communicate using message passing
- let worker count and request timestamp in main thread

- JavaScript can spawn new threads called web worker
- web worker communicate using message passing
- let worker count and request timestamp in main thread
- possibilities: postMessage, MessageChannel Or BroadcastChannel

- JavaScript can spawn new threads called web worker
- web worker communicate using message passing
- let worker count and request timestamp in main thread
- possibilities: postMessage, MessageChannel Or BroadcastChannel
- microsecond resolution (even on Tor and Fuzzyfox)
• experimental feature to share data: SharedArrayBuffer

- experimental feature to share data: SharedArrayBuffer
- web worker can simultaneously read/write data

- experimental feature to share data: SharedArrayBuffer
- web worker can simultaneously read/write data
- no message passing overhead

- experimental feature to share data: SharedArrayBuffer
- web worker can simultaneously read/write data
- no message passing overhead
- one dedicated worker for incrementing the shared variable

- experimental feature to share data: SharedArrayBuffer
- web worker can simultaneously read/write data
- no message passing overhead
- one dedicated worker for incrementing the shared variable
- Firefox/Fuzzyfox: 2 ns, Chrome: 15 ns

Building a timer: Is it good enough?



Building a timer: Is it good enough?



 \rightarrow we can distinguish cache hits from cache misses (only \approx 150 cycles difference)!

Take-away



- idea is not new: Wray (1992)
- we also exploited it in other contexts
 - on ARM
 - inside an SGX enclave

J. C. Wray. "An analysis of covert timing channels". In: Journal of Computer Security 1.3-4 (1992), pp. 219–232.

M. Lipp et al. "ARMageddon: Cache Attacks on Mobile Devices". In: USENIX Security Symposium. 2016.

M. Schwarz et al. "Malware Guard Extension: Using SGX to Conceal Cache Attacks". In: DIMVA'17. 2017.

DRAM covert channels in JavaScript!

• sender: native application in a VM

- sender: native application in a VM
- receiver: JavaScript in a web page on the host

- sender: native application in a $\ensuremath{\mathsf{VM}}$
- receiver: JavaScript in a web page on the host
- sender and receiver select the same bank

- sender: native application in a VM
- receiver: JavaScript in a web page on the host
- sender and receiver select the same bank
- sender and receiver select a different row inside this bank

- sender: native application in a $\ensuremath{\mathsf{VM}}$
- receiver: JavaScript in a web page on the host
- sender and receiver select the same bank
- sender and receiver select a different row inside this bank
- sender transmits o by doing nothing and 1 by causing row conflict

- sender: native application in a VM
- receiver: JavaScript in a web page on the host
- sender and receiver select the same bank
- sender and receiver select a different row inside this bank
- sender transmits o by doing nothing and 1 by causing row conflict
- receiver measures access time for its row: fast \rightarrow 0, slow \rightarrow 1



• communication based on 11-bit packets, with 5-bit of data



- communication based on 11-bit packets, with 5-bit of data
- packet starts with a 2-bit preamble



- communication based on 11-bit packets, with 5-bit of data
- packet starts with a 2-bit preamble
- data integrity checked by an error-detection code



- communication based on 11-bit packets, with 5-bit of data
- packet starts with a 2-bit preamble
- data integrity checked by an error-detection code
- sequence bit indicates whether it is a retransmission or a new packet

transmission of approximately 11 bits/s

- transmission of approximately 11 bits/s
- can be improved using

- transmission of approximately 11 bits/s
- can be improved using
 - fewer re-transmits

- transmission of approximately 11 bits/s
- can be improved using
 - fewer re-transmits
 - error correction

- transmission of approximately 11 bits/s
- can be improved using
 - fewer re-transmits
 - error correction
 - + multithreading \rightarrow multiple banks in parallel

- transmission of approximately 11 bits/s
- can be improved using
 - fewer re-transmits
 - error correction
 - + multithreading \rightarrow multiple banks in parallel
- native code: 596 kbit/s cross CPU and cross VM



Conclusion

• information leaks because of the underlying hardware

- information leaks because of the underlying hardware
- vulnerabilities exploitable at the browser level

- information leaks because of the underlying hardware
- vulnerabilities exploitable at the browser level
- running arbitrary JavaScript allows building high-resolution timers

- information leaks because of the underlying hardware
- vulnerabilities exploitable at the browser level
- running arbitrary JavaScript allows building high-resolution timers
- hard to mitigate without reducing functionality

Thank you!

Contact

- ✓ clementine@cmaurice.fr
- ✓ @BloodyTangerine

From bottom to top: Exploiting hardware side channels in web browsers

Clémentine Maurice, Graz University of Technology July 4, 2017—RMLL, Saint-Étienne, France