

Meltdown & Spectre

Klaus Kusche

Frühjahr 2018

Inhalt

- **Wiederholung Rechnerarchitektur:
Architektur moderner Prozessoren
(Core, MMU, Cache, Speicher)**
- **Meltdown & Spectre:**
 - **Geschichte**
 - **Klassifizierung**
 - **Funktion und Wirkung**
 - **Betroffene Prozessoren**
 - **Gegenmaßnahmen**

Ziel

- Verständnis der Funktion der “Lücken”
- Verständnis der Funktion der Gegenmaßnahmen
- ... und ihrer Folgen, Umsetzung, ...

==> Einschätzung der Bedrohung

==> Entscheidung über Maßnahmen

(die z.B. in meinem Fall “nein” war!)

Voraussetzungen

- ***Rechner-Architektur***
Aufbau und Funktionsweise eines Prozessors
- ***Assembler*** oder zumindest
Implementierung von Programmen
bzw. Programmiersprachen
- Grundzüge ***Betriebssystem***
(Speicherverwaltung,
Unterschied normale Programme / Kernel)

Speicher (RAM)

- Folge von **Bytes**
(Speicherplätze für Werte 0 ... 255)
- Jedes Byte hat eine **Adresse**
(“*Hausnummer*”, fortlaufend 0 ... ein paar Mrd.)
- Physikalisch eigene Chips,
logisch weit weg von der CPU:

Core $\langle == \rangle$ MMU $\langle == \rangle$ Cache $\langle == \rangle$ RAM

- Viel **langsamer** als die CPU:
1 Zugriff ... ~200 Takte = **mehrere 100 Befehle!**

Cache (1)

- **Zwischenspeicher** in der CPU für RAM-Daten
- Ziel: Zugriff **schneller** machen:

	Zugriffszeit (Takte)
Register in der CPU:	0
L1 Cache:	2-6
L2 Cache:	10-25
L3 Cache:	25-75
RAM:	120-1000

Cache (2)

- Viel *kleiner* als das RAM
(16 KB ... 16 MB im Vergleich zu 4 GB ... ???)

==> Enthält nur einen kleinen Teil der RAM-Daten

... und zwar die

zuletzt benutzten!

==> Neue Zugriffe “*verdrängen*”

Daten älterer Zugriffe aus dem Cache

Verwaltung: In *Cachelines*

= Blöcken zu 32...128 aufeinanderfolgenden Bytes

Cache (3)

Erster Zugriff auf eine Speicherstelle

- holt die Daten ***aus dem RAM***
- kopiert betreffende Cacheline in den Cache

Weitere Zugriffe kurz darauf auf **dieselben Daten**

- holen die Daten ***aus dem Cache***

Für Anwendungen logisch ***völlig unsichtbar***:

Ein Programm kann **logisch nicht feststellen**

- ob & wie viel Cache vorhanden ist
- **ob bestimmte Daten im Cache sind** oder nicht

Cache (4)

Aber:

Die CPU hat takt-genaue Timer
(Auflösung ~0,25 ns)

==> *Im Cache / nicht im Cache* ist durch

Messung der Zugriffszeit

vom Programm aus feststellbar:

< 10 ns im Vergleich zu > 50 ns

Cache (5)

Der Cache ist für

- alle Programme
- alle VM's
- das Betriebssystem

gemeinsam!

*==> Jeder "spürt" die Auswirkungen
der Cache-Aktivitäten
aller anderen!*

(als Beschleunigung / Verlangsamung seiner Zugriffe)

Core (1)

Das “*klassische*” Prozessor-Modell:

- Befehl holen & dekodieren
- Operanden holen
- Operation ausführen
- Ergebnis speichern
- nächsten Befehl holen ...

... wäre bei gleicher Halbleiter-Technologie

einige 100 mal langsamer

als aktuelle Prozessoren wirklich sind!

Core (2)

Tricks:

- *Pipelining* ==> hier nicht entscheidend
- 4 ... 16 Rechenwerke (Funktionseinheiten) pro Core ("*superskalare*" CPU)
z.B. Intel: 4 * int, 2 * float, 3 * load/store, ...
==> Ziel:

*Möglichst viele Rechenwerke
möglichst ständig
(sinnvoll) "beschäftigen"!*

Core (3)

- *Hyperthreading*: 1 Core verarbeitet Befehle von 2 Programmen gleichzeitig / abwechselnd
- Mehrere Befehle *parallel* ausführen
- Unabhängige Befehle umordnen / vorziehen (*out-of-order*-Architektur)
- Befehle “*auf Verdacht*” / “*auf Vorrat*” ausführen (*speculative* Execution)

Bei modernen Prozessoren:

Bis zu 200 Befehle gleichzeitig “in-flight”!

Core (4)

“*Befehl ausführen*” und
“*Ergebnis für das Programm sichtbar machen*”

sind **zwei getrennte Schritte**

(ev. zeitlich viele Takte auseinander!!!)

Bei “spekulativer” Ausführung,
wenn sich später herausstellt
“***Spekulation war falsch***”:

*Ergebnis / Wirkung des Befehls
wird **nie sichtbar gemacht!***

Core (5)

Spekulative Ausführung

... ist logisch völlig transparent / unsichtbar,
d.h. vom Programm aus nicht feststellbar
(Terminologie Spectre: “*Transiente Befehle*”),

d.h. die Hersteller machen nichts “falsch”,

... aber eine fehlgeschlagene spekulative Ausführung
hinterlässt durch ihre Speicher-Lese-Zugriffe
trotzdem Spuren und Nebenwirkungen:

- Daten werden ev. **in den Cache geladen**
- MMU-Zustand ändert sich ev.

Core (6)

Sprungbefehl mit fixem Ziel: *0 Takte*

Sprungbefehl mit variablem Ziel: *10-25 Takte*

Was hat ein *variables Ziel*?

- Bedingte Sprünge: **if**, Schleifenbedingungen
- Indirekte Sprünge: **switch**, Interpreter & JIT's
- Indirekte Call-Befehle:
Methodenaufrufe (von virtuellen Methoden)
Aufrufe von Code in Shared Libraries

Core (7)

Trick: Sprungvorhersage = “***Branch Prediction***”

- Versucht, das Sprungziel vorherzusagen,
und beginnt am vorhergesagten Ziel
mit spekulativer Ausführung
- Wenn sich Vorhersage als richtig erweist:
 - Ergebnis der spekulativen Befehle wird sichtbar
 - Sprungdauer 0 Takte
- Wenn sich Vorhersage als falsch erweist:
 - Spekulativ ausgef. Befehle werden verworfen
 - Richtiger Sprung wird nachgeholt: 12-25 Takte

Core (8)

Wie funktioniert die Sprungvorhersage?

U.a. großer **“Sprungziel-Cache”**

(**“BPB”** = *“branch prediction buffer”*):

*“Wohin sprang derselbe Sprungbefehl
bei seiner vorigen Ausführung?”*

==> Die Sprungvorhersage ist

“trainierbar”!

*Bei genauer Kenntnis der internen BPB-Organisation kann man
sogar fremde statt demselben Sprungbefehl trainieren!*

Core (9)

Der BPB ist pro Core,
nicht programmbezogen!

==> Ein Programm kann Sprünge/Calls

- seines "Hyperthreading-Zwillings",
- des Kernels
- eines nachfolgenden Programmes
- theoretisch sogar des VM-Hosts
und eines anderen VM-Guests

beeinflussen!

Core (10)

Der *Microcode*

- steuert die Dekodierung und Abarbeitung von komplexen Maschinenbefehlen in der CPU
- ist ein “Hardware-Programm”
(entfernte Ähnlichkeit mit dem Programm eines FPGA)
- wird früh bei jedem Booten in die CPU geladen bzw. aktualisiert (vom BIOS oder vom OS aus)

Core (11)

Microcode - Updates:

- Normalerweise: ***BIOS-Update***
- Bei Linux und macOS:
Auch via ***Betriebssystem-Update*** möglich
(statt BIOS-Update)
- Microsoft:
 - Könnte die Updates mit Windows installieren, machte es aber bis April gar nicht und jetzt nur teilweise (rechtl. Gründe?)
 - Ausnahme: Surface: *Update von Microsoft*

Core (12)

Barrier-Befehle

serialisieren bzw. verzögern

die parallele / spekulative / ungeordnete
Ausführung von Befehlen

==> *Das Programm wird langsamer!*

Beispiel x86: **LFENCE** (fence = Zaun, Gatter, Hindernis)

*“Beginne alle nachfolgenden Befehle erst dann,
wenn alle Lese-Befehle vor dem **LFENCE**
ihre Daten fertig gelesen haben!”*

MMU (1)

“Memory Management Unit”

2 Haupt-Aufgaben:

- **Zuordnung** *“virtuelle Adressen”* (im Programm)
=> *“reale Adressen”* (im RAM),

Dabei:

Erkennung *ungültiger* virtueller Adressen

- **Prüfung** der *Zugriffsrechte*:

“Darf das Programm diese Speicherstelle lesen / schreiben/ als Befehl ausführen?”

MMU (2)

Steuerung der MMU durch

“Page tables”

- Werden vom Betriebssystem verwaltet
- Sind groß (viele MB)
- Liegen im normalen RAM,
mit einem eigenen Cache in der MMU:

TLB = “Translation Lookaside Buffer”

MMU (3)

==> Adress-Zuordnung & Rechte-Prüfung kann

- ***schnell*** sein (wenige Takte),
wenn die Pagetable-Daten *schon im TLB* liegen
- ***sehr lange dauern*** (*hunderte* Takte!)
wenn die Pagetable-Daten *zuerst aus dem RAM*
in den TLB geladen werden müssen

==> Eigentlicher *Daten-Zugriff* muss *warten*

==> Erkennung “*Adr ungültig*” oder “*Rechte fehlen*”
kann lange *dauern!*

MMU (4)

Nur bei Intel-CPU's (?!):

Adress-Zuordnung und Rechte-Prüfung sind
technisch & zeitlich getrennt!

Folge:

- Daten werden gleich nach der Adress-Zuordnung geladen und auch (spekulativ!) verwendet
- Abbruch wegen fehlender Rechte kommt (unter bestimmten Umständen) erst später!

MMU (5)

“Klassische” Pagetable-Organisation (früher):

- Eine Pagetable pro Prozess (Programm)
nur für die Speicherbereiche dieses Programms
- Plus eine eigene Pagetable für den Kernel

=> Jedes Programm “sieht”
nur seine eigenen Speicherbereiche

=> Zugriffe auf Betriebssystem-Daten
oder fremde Daten:

“Ungültige Adresse”

MMU (6)

Das “Klassische Modell” erfordert

Zwei Pagetable-Umschaltungen

(vom Programm zum Kernel und zurück)

pro Syscall, pro Interrupt, ...

Problem:

Pagetable-Umschaltungen kosten Zeit!!!

Einige hundert Takte... (bei neuen CPU's mit “PCID-Feature”),
...viele tausend Takte! (bei alten CPU's, kompletter TLB-Flush!)

MMU (7)

- 1.) Zusätzliches Speicher-Rechte-Bit, HW-geprüft:
“Zugriff nur erlaubt, wenn die CPU gerade im Kernel-Modus arbeitet”
- 2.) In jeder Prozess-Pagetable wird
der **gesamte Betriebssystem-Speicherbereich**
(und damit das gesamte RAM,
d.h. die Daten aller Programme!)
 - als “gültig” eingetragen
 - aber als “Kernel only” markiert.

MMU (8)

==> **Schneller!!!**

(weil keine Pagetable-Umschaltung
beim Wechsel in den Kernel mehr nötig ist!)

==> “Im Prinzip” gleicher Effekt ...

... *aber im Detail nicht!*

==> Jedes Programm “sieht” jetzt alles!

==> Zugriffe auf Betriebssystem-Daten
oder fremde Daten:

“Gültige Adresse, aber Zugriff verboten”

Fragen?

Wer & Wann?

- 2 Teams unabhängig voneinander:
 - Internationales Forscher-Team
(AT / US / AU, 4 Uni's und 3 Firmen)
(hatten schon "Rowhammer" gefunden)
 - Google
- Entdeckung Frühjahr 2017
- Seit Sommer 2017 HW/SW-Herstellern bekannt
- Seit Jänner 2018 veröffentlicht
in wissenschaftlichen Publikationen & CVE's

Klassifizierung (1)

- ***HW-basiert***, SW ist “unschuldig”
 - ==> Auf allen Betriebssystemen
(Windows, Linux, macOS / iOS, Android, ...)
 - ==> Trifft potentiell alle Programmier-Sprachen
(meist C/C++, sogar in JavaScript ausgenutzt!)
und alle Anwendungen
 - ==> Virenschutz usw. hilft nicht
- Aber:** Nicht auf “kleinen” Prozessoren
(Mikrocontroller wie Atmel, embedded ARM usw.)

Klassifizierung (2)

- **Lokal:**

Exploit-Code muss am betroffenen Rechner laufen
(keine direkte Ausnutzung über das Netz)

- Entweder:

Angreifer hat lokalen Rechner-Zugang

- Oder:

JavaScript o.ä. via Webseiten, in Mails, ...

==> *indirekt doch "remote"*,

aber JavaScript wurde sofort gefixt!

(was ist mit Makros in Dokumenten???)

Klassifizierung (3a)

- *Read-Only*, d.h. “*Datenklau*”

Hier: *Beliebiger / aller Daten aus dem RAM*

=> Nicht direkt von der Platte,
nicht aus Datenbanken,
nicht von anderen Rechnern

=> Kein “Denial of Service”
(kein Absturz, keine Daten-Löschung, ...)

=> Keine Daten- oder Code-Veränderung,
keine Änderung des Programm-Verhaltens

Klassifizierung (3b)

- ==> Auf Servern, Cloud-Systemen, ...:
Daten anderer Benutzer / Programme / Kunden
- ==> Passwörter, Zertifikate, Krypto-Schlüssel, ...
Webseiten-Inhalte & -Eingaben (z.B. TAN's)
- ==> Kernel-Daten, Code anderer Programme, ...
Hilft u.a., andere Sicherheitslücken zu finden
(Programm-Version, was liegt wo im RAM, ...)
- ==> Bei interpretierten Sprachen/Script-Sprachen
(Java, JavaScript, Python, ...):
Zugriff auf interne Daten des Interpreters, ...

Klassifizierung (4)

- **“Seitenkanal-Attacke”:**

Es liegt kein Fehler vor

- in der spezifizierten logischen Funktionalität

- oder in deren Implementierung

Die Hardware-Logik “arbeitet wie definiert”:

Es gibt kein direktes unberechtigtes Auslesen

==> Das Abgreifen der Daten erfolgt nicht direkt, sondern über die Messung physikalischer Effekte:
Stromverbrauch, Funkabstrahlung, Zeit, ...

Hier: **Zeit**

Grundidee (1)

// Vorbereitung:

*// Daten des Arrays **spy** aus dem Cache werfen*

*// **Eigentlicher Angriffscode:***

*// (spekulativ ausgeführt, **x** und **y** werden verworfen!)*

... // je nach Angriff, z.B. Sprungbefehl

***x** = ... // gewünschtes Byte aus dem Speicher laden*

y** = **spy**[**x** * **n**]; // mit **n** \geq **64

*// lädt genau eine Cacheline (64 Bytes) von **spy** in den Cache*

*// trifft für jeden Wert von **x** eine andere Cacheline*

*// (wegen **n** \geq **64** liegen die Elemente **spy**[**x** * **n**]*

// mindestens 64 Bytes auseinander!)

Grundidee (2)

// Auswertung / Messung:

```
for (i = 0; i < 256; ++i) {  
    ... // Zeitmessung Start  
    y += spy[i * n];  
    ... // Zeitmessung Ende:  
        // Durchlauf mit i gleich x ist schnell,  
        // weil das Element spy[i * n] schon im Cache liegt  
        // alle anderen Durchläufe sind langsam!  
        // ==> Wert des ursprünglich angegriffenen x ermittelbar!  
}
```

Meltdown: Idee

“Rogue Data Cache Load”

Nutzt das Verhalten der *MMU*
bei *gültigen, aber unzulässigen Zugriffen:*

Page Fault

(Programm-Abbruch)

kommt manchmal zu spät:

Erst, nachdem die *“verbotenen” Daten*
schon (spekulativ, unsichtbar)
gelesen & *verwendet wurden!*

Meltdown: Details

```
// Caches und TLB präparieren,  
// damit das Timing passt:  
// Zu *evil_pointer* gehörende Pagetable-Daten  
// dürfen nicht im TLB sein  
// (das ist der Normalfall, weil erstmalig zugegriffen wird)  
  
x = *evil_pointer; // Löst Signal bzw. Exception aus  
y = spy[...x...];  
    // wird zwar nie sichtbar,  
    // aber ev. trotzdem noch ausgeführt!  
    // (weil die MMU-Rechte-Prüfung lange dauert)  
  
// Exception abfangen, Cache auswerten
```

Meltdown: Effekt

- Liest alle Daten mit gültiger virtueller Adresse (auch ohne Rechte, z.B im Kernel-Adressbereich), und damit beliebige / ***alle RAM-Inhalte***
= Alle Kernel-Daten (incl. Filesystem-Caches),
alle Daten aller Applikationen, ...
- Relativ schnell: > 500 KB/sec
=> in ein paar Stunden
sind einige GB RAM lesbar

Inzwischen aktiv genutzt!

Meltdown: Opfer

- Alle Intel-Core-CPU's seit mindestens 2008
- Intel Atom nicht
- AMD nicht
- Nur modernster ARM-Core A75, andere nicht
- Sonst bisher nichts bekannt

- Bei ARM A15, A57 und A72:
Ähnlicher Bug erlaubt normalen Programmen
Zugriff auf Kernel-Register

Meltdown: Fix / Kurzfristig (1)

- Kein Mikrocode-Fix möglich
- Kein Fix auf Anwendungsebene nötig
- Im Betriebssystem:

***Pagetables
von normalen Prozessen
und Kernel wieder trennen!***

=> Komplexer Umbau!

(in Win 7 und Win Server 2008 desaströs schiefgelaufen!!!)

Meltdown: Fix / Kurzfristig (2)

Status:

- ***Intel & ARM:***

Windows, praktisch alle Linux, aktuelles macOS, aktuelles original Google Android & iOS:

Endgültiger Fix schon ausgeliefert!

(Windows-Fix legte fast alle Virens Scanner & manche Treiber lahm, machte Probleme auf alten AMD-Proz., ließ auf Win 7 und Win 2008 bis April eine gewaltige Lücke...)

Meltdown: Fix / Kurzfristig (3)

- 0-10 % Perf.-Verlust bei normalen Programmen
- max. 30 % Perf.-Verlust bei Servern (DB, Web), Spielen, I/O-lastigem Code
(manche Kunden beobachten auf ihren Cloud-Servern bis zu 50 % ?!)
- \geq 50 % Einbruch bei Benchmarks schneller SSD's und anderen synthetischen Tests
(75 % bei macOS ???)
- Abhängig vom Prozessor-Alter:
Verlust bei alten Intel-Prozessoren ohne PCID-Feature deutlich höher!

Meltdown: Fix / Langfristig

Änderung des MMU-Verhaltens wäre sinnvoll
und müsste fast “verlustfrei” möglich sein!

==> “Verbotene” Daten gar nicht laden,
Adresse & Rechte gleichzeitig auswerten!

Aber:

Umfangreiche HW-Architektur-Änderung,
erst in zukünftigen Generationen,
nicht via Microcode möglich!

Spectre I: Idee

“Bounds Check Bypass”

Nutzt die *Sprungvorhersage*
bei *bedingten Sprüngen*:

Trainiert die Sprungvorhersage
bei einem *normalen Array-Indexgrenzen-if*,
macht dann einen Zugriff mit “*bösem*” Index

*==> Prozessor macht spekulativ
den “Index ist ok”-if-Zweig
und greift auf Daten
weit außerhalb des Arrays zu!*

Spectre I: Detail

Suche irgendwo im bestehenden Programmcode (incl. Libraries) Code mit folgender Struktur, rufe ihn mit ein paar richtigen value_from_user und dann mit einem falschen auf!

```
if (value_from_user < arr1_size) {  
    // if dauert lang, wenn arr1_size nicht im Cache ist!  
    // ==> auch für falsches value_from_user  
    // werden die folgenden Befehle spekulativ ausgeführt!  
    x = arr1[value_from_user];  
        // greift für falsches value_from_user  
        // auf jeden beliebigen Speicher zu!  
    y = arr2[...x...];  
    ...  
}
```

Spectre I: Effekt & Opfer

Liest alle RAM-Daten des eigenen Programmes (interessant z.B. bei Interpretern, Browsern)...

... auf praktisch allen “großen, modernen CPU’s”:

- Alle Intel Core und AMD seit “ewig” (ev. 1995?)
- Intel Atom seit 2014
- Fast alle “Handy- und Server-ARM-Prozessoren” (einzige Ausnahme: Cortex A53 Core, z.B. Raspi)
- IBM Power, Sun/Oracle/Fujitsu Sparc, MIPS, ... (zSeries ???)

Spectre allg.: Fix (1)

Verhalten ist implizite, gewünschte Konsequenz moderner Prozessorarchitektur und “korrekt”:

Code verhält sich “wie geplant / gewünscht”!

==> Es gibt **keinen zentralen “Fix”**,
der die Ausführung bestehenden Codes
(ohne Änderungen am Code)
ohne essentielle Verluste “repariert”!

= ohne dass das Programm zu langsam wird!

*“Branch Prediction” komplett ausschalten
wäre **viel** (> Faktor 5!) zu langsam!*

Spectre allg.: Fix (2)

Nur “Workarounds” möglich:

Bestehender Code muss

durch anderen Code ersetzt

werden, der dasselbe Ergebnis erzielt,
aber

- mit anderen “internen Abläufen”
(z.B. ohne die Branch Prediction zu benutzen)
- oder mit “expliziten Bremsen”
(nur dort wo es wirklich nötig ist).

Spectre allg.: Fix (3)

1.) Falls nötig:

Per Microcode-Update

- Neue Maschinenbefehle einführen
(bei Intel 3 Stück für Spectre II,
bei ARM einer für Spectre I)
z.B. “Lösche / ignoriere Sprungvorhersage”
oder neue Barrier-Befehle
- Ausführung bestehender Maschinenbefehle
an neue Befehle anpassen
z.B. neue Status-Bits berücksichtigen

Spectre allg.: Fix (4)

2.) Alle Compiler (und JIT's) umbauen, sodass sie für **if**'s, virtuelle Methoden-Aufrufe, **switch**, ...

anderen Code erzeugen

(ev. mit den "neuen" Instruktionen)

3.) Gesamten Code neu kompilieren!!!

- Betriebssystem
- Alle Anwendungen
- Alle Libraries
- ...

(bestehenden Assembler-Code händisch durchgehen & fixen!)

Spectre allg.: Fix (5)

Stand Ende Jänner:

- Microcode-Fixes sind teilweise verfügbar, aber buggy (alle wieder zurückgezogen)!
- Modifizierte Compiler (gcc, llvm, msvc) sind im Beta-Test
- Betr. JIT's gibt es keine Info
- Frisch compilierte Anwendungen sind noch nicht in Sicht...

==> Das wird eine Sache von Monaten / Jahren!

Spectre allg.: Fix (6)

Aktueller Stand:

- Korrigierte Microcode-Fixes werden schön langsam freigegeben, aber noch nicht für alle Prozessoren...
- Modifizierte Compiler sind ausgeliefert, aber die Änderungen sind teilweise unvollständig
- Betriebssysteme sind mit den neuen Compilern frisch compiliert, Anwendungen meist nicht
- Sonst kaum offizielle Info's...

Spectre I: Fix (1)

Intel & AMD:

Theoretisch, für 100 % Schutz:

Zwischen jedem bed. Sprung & Speicherzugriff:

LFENCE

Effekt: Speicherzugriffs-Befehl muss warten,
bis die **if**-Bedingung fertig ausgewertet ist.

LFENCE gibt's schon länger

=> kein Microcode-Update nötig!

Spectre I: Fix (2)

Performance-Auswirkungen:

Ja (weniger spekulative “Vorausarbeit”),
katastrophal: > **60 % langsamer**

==> Nicht bei jedem bedingten Sprung einbauen!

==> Compiler muss überlegen:

“Wo brauche ich das *LFENCE*, wo nicht?”

==> Abwägung Performance / Sicherheit!

1. Versuch von MSVC:

Performance-Verlust gering, aber viele Lücken blieben offen!

Spectre I: Fix (3)

ARM:

2 verschiedene Fixes, je nach Core-Typ:

- **Neuer Barrier-Befehl CSDB** (ähnlich Intel)

==> Microcode-Update nötig!

- Conditional-Move-Befehl

statt normalem Move-Befehl verwenden

(würde auch bei x86/AMD wirken, aber nicht bei Intel?)

==> Für universellen Code (für alle ARM-Cores)
beides gleichzeitig einbauen?!

Spectre II: Idee

“Branch Target Injection”

Nutzt die *Sprungvorhersage*
bei indirekten Sprüngen und Calls:

Trainiert die Sprungvorhersage
eines fremden Sprungbefehls
(oder des eigenen Befehls bei Hyperthreading),
um ein beliebiges Codestück anzuspringen,
das mit dem if/Call gar nichts zu tun hat!

==> Ein anderes Programm (bzw. der zweite Thread)
führt später den “falschen” Code spekulativ aus!

Spectre II: Details

Man muss im angegriffenen Programm wieder ein vorhandenes Codestück suchen, das abhängig von einem bestimmten Wert auf irgendwelche anderen Daten zugreift (und dabei Spuren im Cache hinterlässt).

Die Adresse dieses Codestückes
trainiert man in die Sprungvorhersage...

Extremfall:

Google hat den BPF-Interpreter im Linux-Kernel dazu gebracht, beliebigen BPF-Code (vom User geladen & bösartig) auszuführen. Der BPF-Code hat dann Zieldaten gelesen und den Cache befüllt.

Spectre II: Effekt & Opfer

Liest *alle RAM-Daten* eines fremden Programmes oder des Kernels.

- Alle Intel Core seit ??? (sehr modellabhängig!)
- AMD: “Theoretisch ja, aber ...” (trotzdem Patch?)
- Intel Atom nein?
- Fast alle “Handy- und Server-ARM-Prozessoren” (einzige Ausnahme: Cortex A53 Core, z.B. Raspi)
- Andere: unbekannt

Ausnutzung erfordert komplexe Analyse pro CPU!

Spectre II: Fix

- Viel Chaos und Verwirrung
 - Noch weit entfernt von endgültiger Lösung
 - Intel:
Bisher keine Lösung mit Barrier-Befehlen?
(ARM: ???)
- ==> Zwei (oder mehr) unabhängige Lösungsansätze,
grundverschieden!
- ==> Derzeit meist nur vorläufige / teilweise Patches,
Probleme mit Nebenwirkungen

Spectre II: Fix Variante 1 (1)

*“Retpolines” =
“Return trampolines”*

Statt indirektem Sprung/Call-Befehl: (stark vereinfacht!)

- Zieladresse in ein Register legen
- Fixe Hilfs-Funktion aufrufen
- Hilfsfunktion ersetzt Return-Adresse am Stack durch Zieladresse
- ... und macht Return

==> Return springt zum ursprünglichen Ziel!

Spectre II: Fix Variante 1 (2)

- Nutzt “Return Prediction” statt Branch Prediction
- Kommt ohne neue Maschinenbefehle aus
=> Kein Mikrocode-Update nötig!
- Braucht dafür neue Compiler-Versionen

Aber ***kostet Performance***:

- Mehr Befehle
- “Return Prediction” sagt immer falsch voraus:
Der Code enthält dort eine zusätzliche Warteschleife,
die angesprungen und spekulativ ausgeführt wird,
bis der Spekulations-Fehler erkannt
und das echte Ziel angesprungen wird

Spectre II: Fix Variante 1 (3)

- Auf älteren Intel-Prozessoren:
Funktioniert, mittlere Performance-Verluste
 - Auf neueren Intel-Prozessoren:
 - Relativ hohe Performance-Verluste!
 - Ab Skylake: Auch Return ist Spectre-anfällig!
(nutzt ab 16 Returns nacheinander die Branch Prediction
==> Linux hat ziemlich komplexen Workaround)
- ==> Auf Intel trotzdem *neuer Microcode* nötig?!**

Spectre II: Fix Variante 1 (4)

- Schützt zu wenig
in bestimmten Fällen von VM's
=> In diesen Fällen ist Fix Variante 2
zusätzlich nötig!
... oder bei AMD Ryzen:

Nutzung der
Hauptspeicher-Verschlüsselung

=> Jede VM bekommt einen
anderen Verschlüsselungs-Key

Spectre II: Fix Variante 1 (5)

- Inkompatibel mit
neuem Sicherheitsmechanismus “CET”
in nächster Prozessor-Generation von Intel:

*CET klassifiziert Retpoline
als bösartigen Angriffscode!*

(weil Überschreiben der Return-Adresse am Stack
bisher nur in böartigem Code vorkam)

Spectre II: Fix Variante 2 (1)

*Zusätzliche, neu eingeführte Befehle
zur Kontrolle / Abschaltung / Leerung
der Sprungvorhersage*

Bei Intel:

- Indirect Branch Prediction Barrier **IBPB**
- Indirect Branch Restricted Speculation **IBRS**
- Single Thread Indirect Branch Predictor **STIBP**

Spectre II: Fix Variante 2 (2)

Wo sind die Befehle nötig?

- Auf jeden Fall im Betriebssystem
(bei jedem Syscall, bei jedem Prozesswechsel):

Gegen Angriffe auf das Betriebssystem
und auf fremde Programme (nur teilweise!)

- Auch in den Anwendungen
bei jedem einzelnen
indirekten Call- oder Sprung-Befehl???

Für Schutz gegen alle Varianten?

Spectre II: Fix Variante 2 (3)

- Microcode-Update Voraussetzung, sonst Absturz wegen “*unbekanntem Befehl*”
 - Erster Microcode-Update wurde zurückgezogen, machte Prozessoren instabil, seit Mitte März ok
 - Auf neueren Intel-Prozessoren:
Mittlere bis hohe Performance-Verluste
(angeblich > 4000 Takte bei Prozesswechsel?)
 - Auf älteren Intel-Prozessoren:
Sehr hohe Performance-Verluste!
- ==> Fast immer & überall langsamer als Variante 1!

Spectre II: Fix Variante 2 (4)

Auf AMD:

- Nachrüstbarkeit der neuen Befehle via μ Code:
 - Ryzen: Ja (heimlich ausgeliefert?!))
 - Ab Bulldozer: Teilw. ausgeliefert
(erst seit Mitte April verfügbar)
 - Auf älteren Prozessoren: Ziemlich sicher nicht!
- Bisher keine Aussagen über Performance!
(dafür läuft Fix Variante 1
auf AMD Ryzen sehr performant!)

Spectre II: Fix Variante 3

Auf AMD hilft auch ein einfaches **LFENCE**:

MOV *reg, zieladresse*

LFENCE

JMP *reg // oder CALL reg*

==> Verhindert das Problem nicht grundsätzlich,
aber macht den Spekulations-Zeitraum so kurz,
dass sich keine Ausnutzung mehr ausgeht.

Aber:

- Geschwindigkeitsauswirkungen unklar

- Keine Aussagen von Intel

 betr. Wirksamkeit auf Intel-Prozessoren

Spectre II: Welcher Fix?

- Microsoft / *Windows* setzt anscheinend ***ausschließlich auf Variante 2***
- *Linux* und Google / *Android* implementieren ***fast nur Variante 1*** (Variante 2 nur für VM's usw.?)
(ganz böse Mail von Linus Torvalds...)
- *AMD* möchte ***nur Variante 1***
(wegen wesentlich besserer Performance)
- *Intel* möchte ***ausschließlich Variante 2***
(wegen Inkompatibilität von Variante 1 mit CET & weil Variante 1 ab Skylake nur teilweise hilft)

Spectre II: Andere CPU's

- Auf **ARM**:

Anscheinend auch Variante 2:

Neuer Befehl via μ Code-Update:

BTB leeren

- Auf **IBM zSeries**:

Retpoline-ähnlicher Ansatz

Logistisches Problem (1)

Eigentlich bräuchte man

- verschiedene Betriebssystem-Versionen
- verschiedene Anwendungs-.exe's und -.dll's
 - für jede einzelne Intel-/ARM-Prozessor-Generation
 - für AMD-Prozessoren
 - für Systeme mit / ohne BIOS-Update

=> Problem

- bei “universellen” .exe's
- bei “universellen” VM's
- beim “Übersiedeln” von VM's
 - auf einen Server mit anderen Prozessoren

Logistisches Problem (2)

- ==> Logistisch nicht beherrschbar???
- ==> Prozessoren, Compiler und BIOS müssen wohl beides unterstützen?
- ==> “Auf Verdacht” doppelt Behebungen einbauen?!
- ==> *Noch mehr Performance-Problem!*

“Saubere” Lösung?

Bisher keine erkennbare Diskussion?

Aber:

Ankündigung von Intel, bis Jahresende Cores mit ***Hardware-Architektur-Fix*** zu liefern!

Meine Ideen:

- Branch Prediction Prozess-spezifisch machen
- Branch Prediction so umbauen, dass sich Training nicht auf andere Sprünge auswirken kann

Anwendungs-Änderungen

Z.B. in Chrome, Firefox, ...:

- JavaScript-Timerfunktionen
viel ungenauer gemacht (wegen Cache-Messung)!
- Sprach-Feature “*SharedArrayBuffer*” entfernt

==> Kein Performance-Verlust, aber Feature-Verlust!

- Ein separater Prozess pro Webseite / pro Tab
(aus Sicherheitsgründen allgemein empfohlen!)

==> Mehr Speicherbedarf, etwas Laufzeit-Overhead

Spectre allg.: Weitere Zugriffe (1)

“*Draufgaben*” von Google für Spectre I / II auf Linux:

- Spectre I & moderner Intel Core:
User-Programm liest Kernel-RAM
- Spectre I & AMD:
Ebenso, aber nur wenn BPF-JIT im Kernel aktiv
(BPF-JIT: Firewall-Regelcode-Interpreter)
- Spectre II & moderner Intel Core:
Root-Programm im VM Guest
liest Kernel-RAM des VM Host (!)

... und da kommt sicher noch mehr!

Spectre allg.: Weitere Zugriffe (2)

*Spectre liest auch Daten
aus **Intel-SGX-Enklaven!***

Intel SGX = “*Software Guard Extensions*”:

Schützt bestimmte Speicherbereiche (“Enklaven”) für *hochsensible Daten und Programme* vor Zugriffen durch Betriebssystem, Debugger, ...

Anwendungen von SGX:

- ***Keystores*** und Krypto-Algorithmen
- ***DRM***-Implementierungen

Literatur (1)

- Publikationen TU Graz:

<https://meltdownattack.com/>

<https://meltdownattack.com/meltdown.pdf>

<https://spectreattack.com/spectre.pdf>

- Publikation Google:

<https://googleprojectzero.blogspot.de/2018/01/reading-privileged-memory-with-side.html>

- CVE's:

CVE-2017-5753, CVE-2017-5715, CVE-2017-5754

Literatur (2)

- Intel:

<https://newsroom.intel.com/wp-content/uploads/sites/11/2018/01/Intel-Analysis-of-Speculative-Execution-Side-Channels.pdf>

- ARM:

<https://developer.arm.com/support/security-update/download-the-whitepaper>

<https://developer.arm.com/support/security-update>

Literatur (3)

- AMD:

<https://www.amd.com/en/corporate/speculative-execution>

<http://developer.amd.com/wordpress/media/2013/12/Managing-Speculation-on-AMD-Processors.pdf>

Literatur (4)

- “Retpolines”:

<https://support.google.com/faqs/answer/7625886>

- Intel zu Retpolines:

<https://software.intel.com/sites/default/files/managed/1d/46/Retpoline-A-Branch-Target-Injection-Mitigation.pdf>

“The end”

Fragen?