# Modularity through Memory Virtualization

Prof. George Candea

*School of Computer & Communication Sciences*

# Objectives

- Understand memory virtualization

  - *the coolest form of modularization we have in operating systems*

- Understand more deeply the role played by naming in modularization

- Start developing a "system designer" mindset

# Outline

- Enforced modularity

- Page tables

  - *Directory for mapping memory names (VA) to memory locations (PA)*

- Caching

- Constants in systems design

# Enforced Modularity

# Modularity Recap

- ## Names or not?

  - *Memory address names a memory location*

  - *Pointer*

    - `void*` names a memory location
    - `int*` names the location of an integer

  - *Virtual address names a physical memory location*

- ## Enforced modularity

  - *module boundaries provided by a mechanism external to the modules*

# Enforced vs. Soft Modularity

- Module X interacts with module Y

  - *Encapsulation demands of X: module Y does not write to memory pages other than its own*

- Three options for enforced modularity?

  - *(a) Y is written carefully* ✘

  - *(b) The compiler used to compile Y guarantees it* ✔ *if X and Y share compiler*

  - *(c) Y's language runtime / interpreter guarantees it* ✔ *if X and Y share runtimes*

# Role of Trust in Enforcing Modularity

- Further modularization options for X and Y

  - *(d) Y runs in a separate process / address space*

  - *(e) Y runs in a separate VM*

  - *(f) Y runs on a separate physical machine*

  - *(g) Y runs on an air-gapped machine*

- Enforced modularity = module boundaries provided by a mechanism external to the modules trusted by all modules

# Virtual Memory & Page Tables

VPN     offset

Virtual Address

| 0 | 1 | 0 | 1 | 0 | 1 |
|---|---|---|---|---|---|

Address Translation

Physical Address

| 1 | 1 | 1 | 0 | 1 | 0 | 1 |
|---|---|---|---|---|---|---|

PFN     offset

*Virtual namespace*

*Name translation controls visibility*

*"if you can't name it, you can't use it"*

*Physical namespace*

63    56 | 55    48 | 47    40 | 39    32 | 31    24 | 23    16 | 15    8 | 7    0

9          9          9          9          12

P4 index   P3 index   P2 index   P1 index   Offset

48 | 47    40 | 39    32 | 31    24 | 23    16 | 15    8 | 7    0

9          9          9          9          12

P4 index    P3 index    P2 index    P1 index    Offset

47    40 | 39    32 | 31    24 | 23    16 | 15    8 | 7    0

$2^9$=512 entries
8 bytes / entry

9          9          9          9          12

P4 table
P3 table
P2 table
P1 table

P4 entry
P3 entry
P2 entry
P1 entry

4K memory page
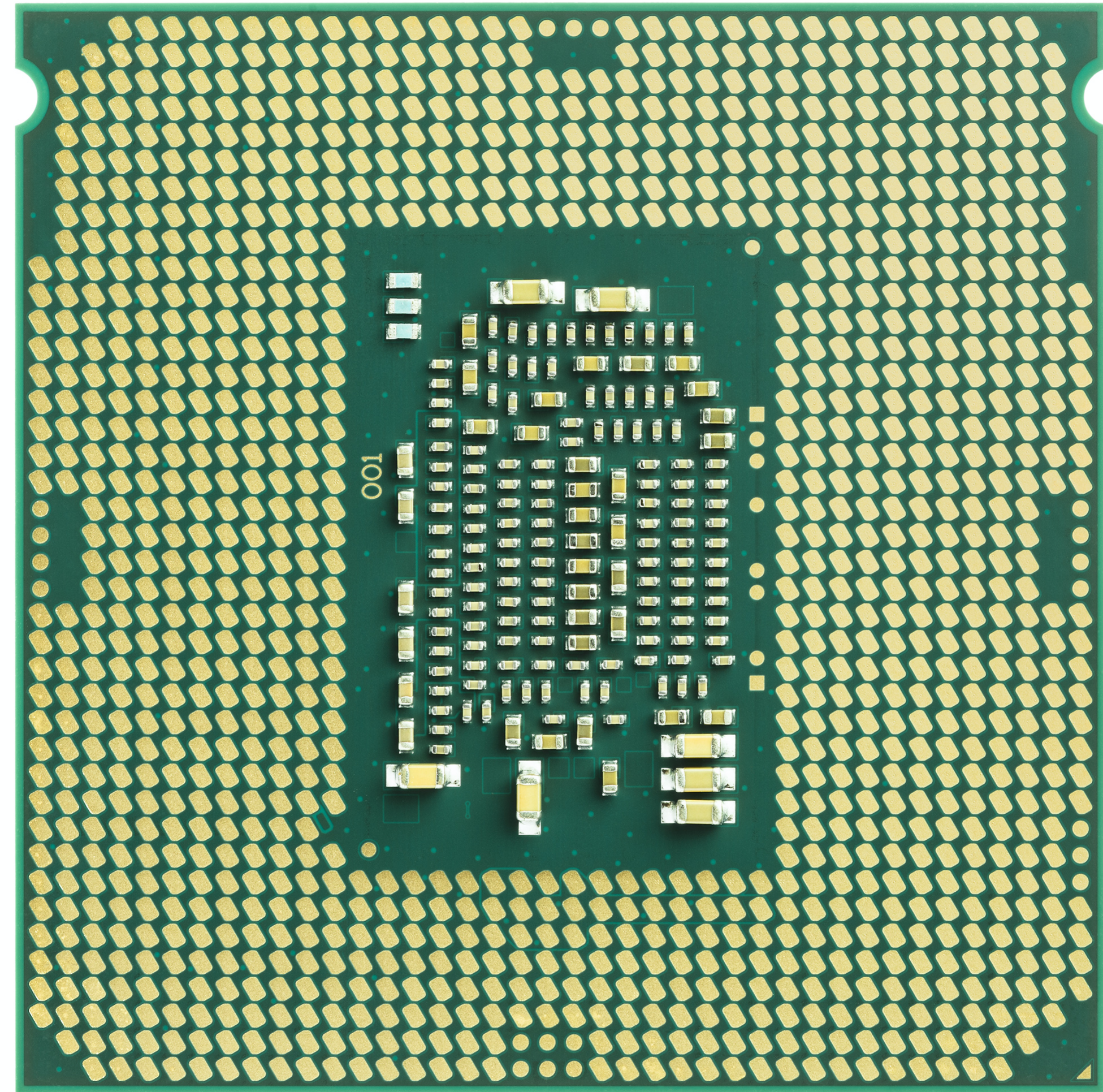
CR3 register

*Hierarchy of name directories:*
*• scalability for sparse namespaces*
*• locality of reference*

Diagrams courtesy of https://os.phil-opp.com/page-tables/

# Linear Address

| 47 | 39 | 38 | 30 | 29 | 21 | 20 | 12 | 11 | 0 |
|----|----|----|----|----|----|----|----|----|---|
| PML4 | | Directory Ptr | | Directory | | Table | | Offset | |

PML4 / 9

Directory Ptr

Directory / 9

Table / 9

Offset / 12

4-KiB Page

Physical Addr

PDE with PS=0 / 40

Page-Directory

Page Table

PTE / 40

Page-Directory-Pointer Table

PDPTE / 40

/ 9

PML4E / 40

**Hierarchy of name directories:**
- *scalability for sparse namespaces*
- *locality of reference*

*Root of hierarchy*

/ 40

CR3

Linear Address

| 47 | 39 38 | 30 29 | 21 20 | 12 11 | 0 |
|---|---|---|---|---|---|
| PML4 | Directory Ptr | Directory | Table | Offset | |

4-KiB Page

Physical Addr

PTE

Page Table

PDE with PS=0

Page-Directory

Page-Directory-Pointer Table

PDPTE

*Hierarchy of name directories:*
- *scalability for sparse namespaces*
- *locality of reference*

PML4E

*Root of hierarchy*

CR3

Linear Address

| 47 | 39 38 | 30 29 | 21 20 | 12 11 | 0 |
|---|---|---|---|---|---|
| PML4 | Directory Ptr | Directory | Table | Offset | |

4-KiB Page

Physical Addr

PTE

Page Table

PDE with PS=0

Page-Directory

Page-Directory-
Pointer Table

PDPTE

Hierarchy of name directories:
• scalability for sparse namespaces
• locality of reference

PML4E

Root of hierarchy

CR3

Linear Address

| 47 | 39 38 | 30 29 | 21 20 | 12 11 | 0 |
|---|---|---|---|---|---|
| PML4 | Directory Ptr | Directory | Table | Offset | |

PML4 /9

Directory Ptr /9

Directory /9

Table /9

Offset /12

4-KiB Page

Physical Addr

PTE

40

Page Table

PDE with PS=0

40

Page-Directory

Page-Directory-Pointer Table

PDPTE
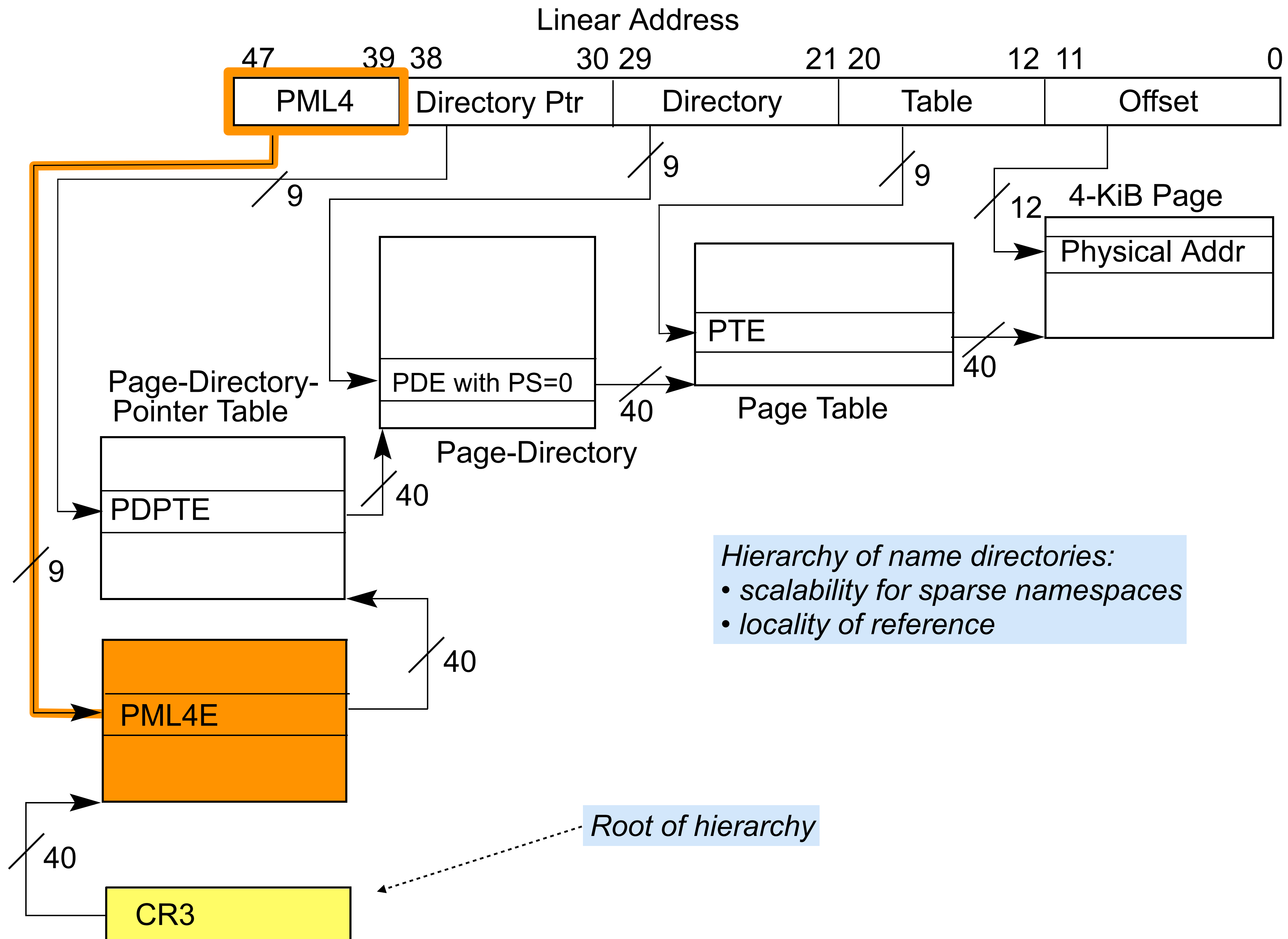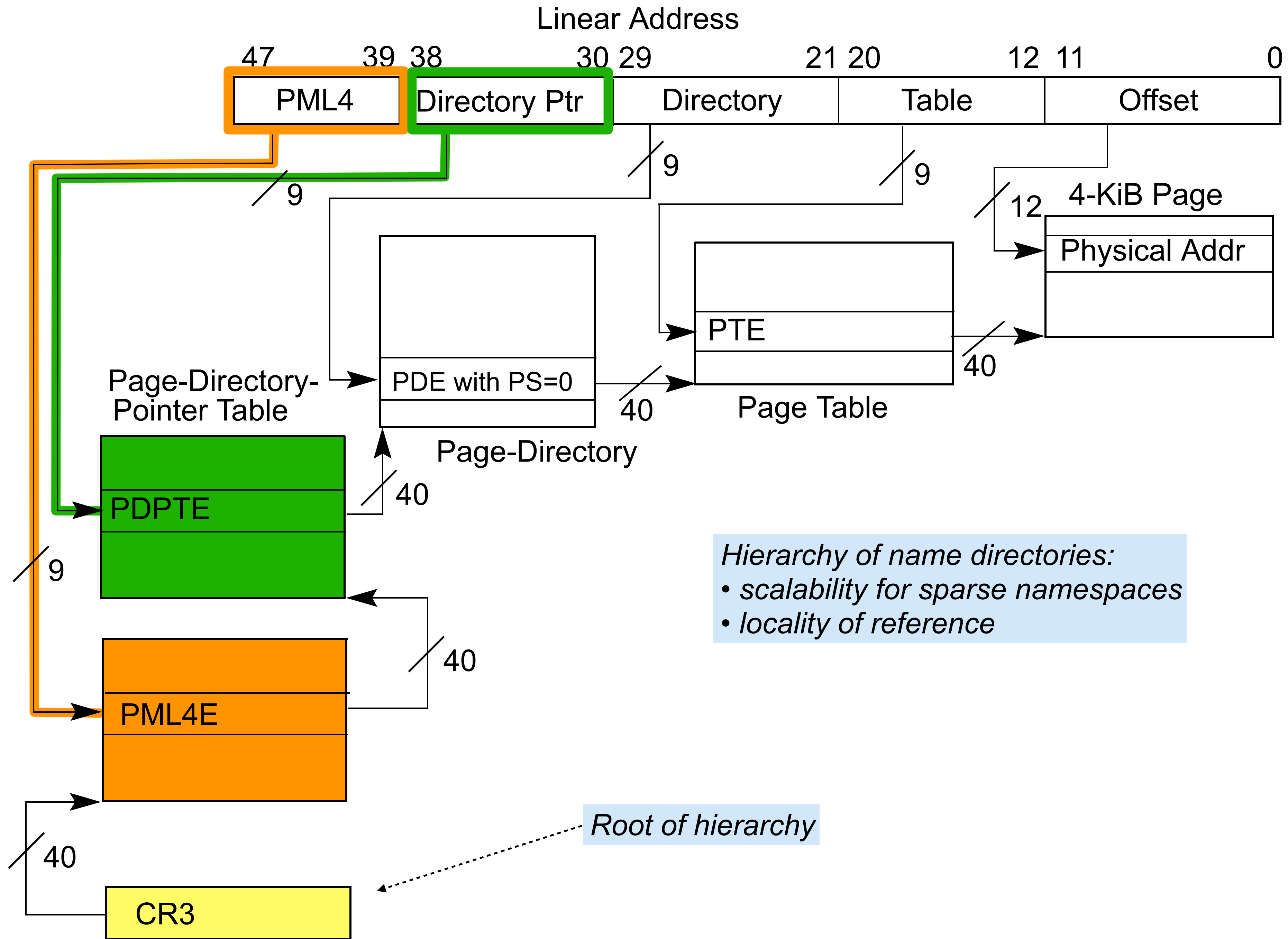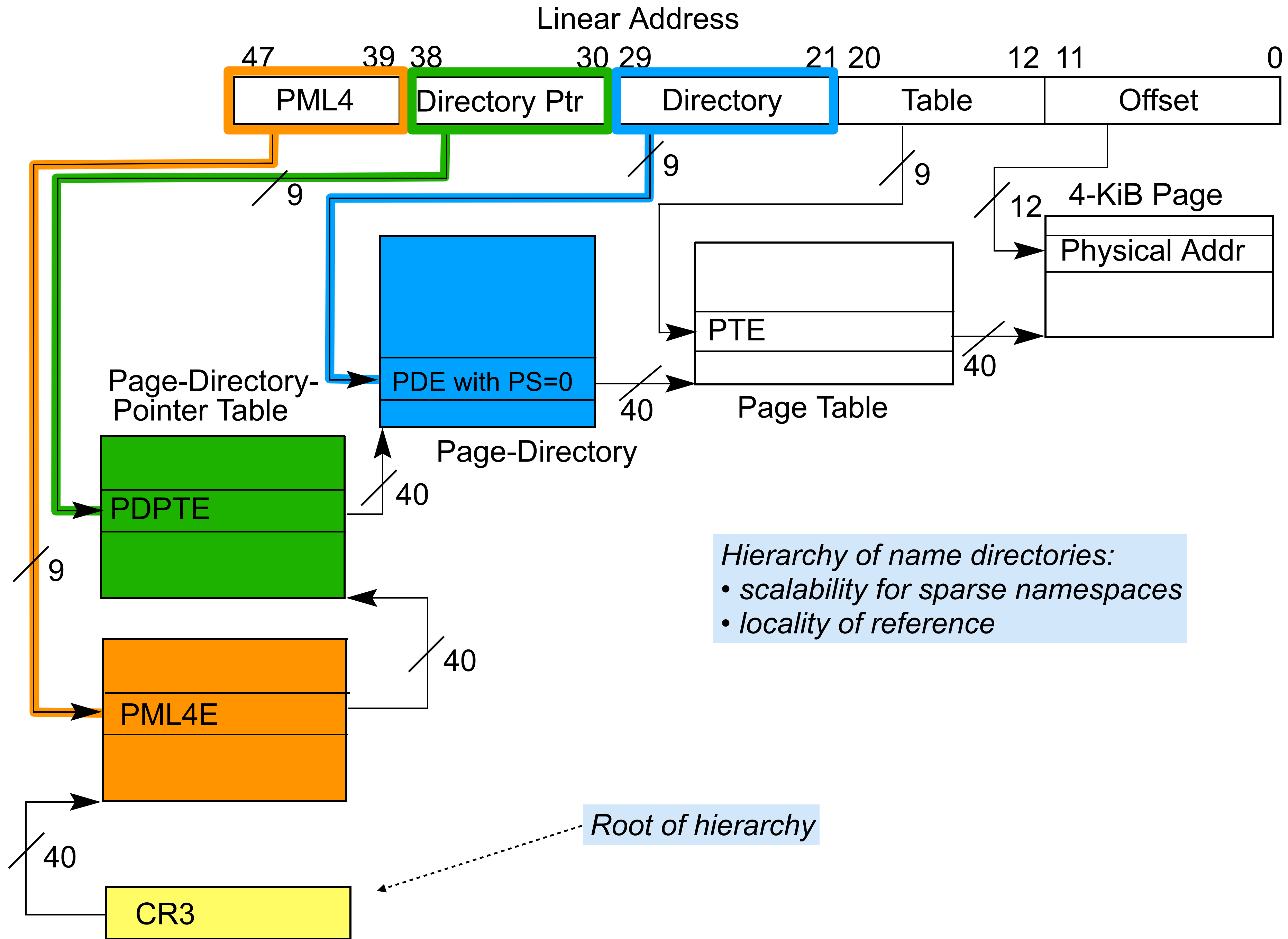
40

PML4E

40

CR3

40

*Hierarchy of name directories:*
- *scalability for sparse namespaces*
- *locality of reference*

*Root of hierarchy*

Linear Address

| 47 | 39 38 | 30 29 | 21 20 | 12 11 | 0 |
|---|---|---|---|---|---|
| PML4 | Directory Ptr | Directory | Table | Offset | |

9

9

9

12

4-KiB Page

Physical Addr

PTE

40

Page Table

PDE with PS=0

40

Page-Directory

Page-Directory-
Pointer Table

PDPTE

40

9

Hierarchy of name directories:
• scalability for sparse namespaces
• locality of reference

PML4E

40

Root of hierarchy

40

CR3

Linear Address

| 47 | 39 | 38 | 30 | 29 | 21 | 20 | 12 | 11 | 0 |
|---|---|---|---|---|---|---|---|---|---|
| PML4 | | Directory Ptr | | Directory | | Table | | Offset | |

9  9  9  9  12

4-KiB Page

Physical Addr

PTE

Page Table

40

PDE with PS=0

40

Page-Directory

Page-Directory-Pointer Table

PDPTE

40

PML4E

40

CR3

40

*Hierarchy of name directories:*
• *scalability for sparse namespaces*
• *locality of reference*

*Root of hierarchy*

Bit positions: 63 62 61 60 59 58 57 56 55 54 53 52 51 … M M−1 … 32 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0

| Fields | | | | | | | Flags | | | | | | | | | | Entry |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Reserved | | | Address of PML4 table | | | | | | | | PCD | PWT | Ign. | | | | CR3 |
| XD[3] | Ignored | Rsvd. | Address of page-directory-pointer table | | | Ign. | Rsvd | Ign | A | PCD | PWT | U/S | R/W | 1 | | | PML4E: present |
| Ignored | | | | | | | | | | | | | | 0 | | | PML4E: not present |
| XD | Prot. Key | Ignored | Rsvd. | Address of 1GB page frame | Reserved | PAT | Ign. | G | 1 | D | A | PCD | PWT | U/S | R/W | 1 | PDPTE: 1GB page |
| XD | Ignored | Rsvd. | Address of page directory | | | Ign. | 0 | Ign | A | PCD | PWT | U/S | R/W | 1 | | | PDPTE: page directory |
| Ignored | | | | | | | | | | | | | | 0 | | | PDTPE: not present |
| XD | Prot. Key | Ignored | Rsvd. | Address of 2MB page frame | Reserved | PAT | Ign. | G | 1 | D | A | PCD | PWT | U/S | R/W | 1 | PDE: 2MB page |
| XD | Ignored | Rsvd. | Address of page table | | | Ign. | 0 | Ign | A | PCD | PWT | U/S | R/W | 1 | | | PDE: page table |
| Ignored | | | | | | | | | | | | | | 0 | | | PDE: not present |
| XD | Prot. Key | Ignored | Rsvd. | Address of 4KB page frame | | Ign. | G | PAT | D | A | PCD | PWT | U/S | R/W | 1 | | PTE: 4KB page |
| Ignored | | | | | | | | | | | | | | 0 | | | PTE: not present |

Linear Address

| 47 | 39 38 | 30 29 | 21 20 | 12 11 | 0 |
|---|---|---|---|---|---|
| PML4 | Directory Ptr | Directory | Table | Offset | |

4-KByte Page

Physical Addr

PTE

Page Table

PDE with PS=0

Page-Directory

Page-Directory-Pointer Table

PDPTE

PML4E

CR3

## Linear Address

| 47 | 39 38 | 30 29 | 21 20 | 12 11 | 0 |
|---|---|---|---|---|---|
| PML4 | Directory Ptr | Directory | Table | Offset | |

Page-Directory-Pointer Table

PDE with PS=0

Page-Directory

PTE

Page Table

4-KByte Page

Physical Addr

PDPTE

PML4E

CR3

## Linear Address

| 47 | 39 38 | 30 29 | 0 |
|---|---|---|---|
| PML4 | Directory Ptr | Offset | |

Page-Directory-Pointer Table

PDPTE with PS=1

PML4E

CR3

1-GByte Page

Physical Addr

Page size trade-offs:

Bigger pages => fewer PTEs =>
less memory for PTs => fewer TLB misses

Smaller pages => less internal fragmentation

(UltraSPARC up to 16 GiB pages)

Global · Page attribute type · Dirty · Accessed · Caching disabled · Write-through caching · User / Supervisor · Write / Read-only · Present

| 6 6 6 6 5 5 5 5 5 5 5 5 5 | M¹ | M-1 | 3 3 3 2 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 2 1 0 9 8 7 6 5 4 3 2 1 | | | 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 | | | | | | | | | | | |
| Reserved² | | | Address of PML4 table | Ignored | | | PCD | PWT | Ign. | | | | **CR3** |
| XD³ Ignored | | Rsvd. | Address of page-directory-pointer table | Ign. | Rsvd | Ign | A | PCD | PWT | U/S | R/W | **1** | **PML4E: present** |
| Ignored | | | | | | | | | | | | **0** | **PML4E: not present** |
| XD Prot. Key⁴ Ignored | | Rsvd. | Address of 1GB page frame · Reserved | PAT | Ign. | G 1 D A | PCD | PWT | U/S | R/W | **1** | **PDPTE: 1GB page** |
| XD Ignored | | Rsvd. | Address of page directory | Ign. | **0** Ign A | PCD | PWT | U/S | R/W | **1** | | **PDPTE: page directory** |
| Ignored | | | | | | | | | | | | **0** | **PDTPE: not present** |
| XD Prot. Key⁴ Ignored | | Rsvd. | Address of 2MB page frame · Reserved | PAT | Ign. | G 1 D A | PCD | PWT | U/S | R/W | **1** | **PDE: 2MB page** |
| XD Ignored | | Rsvd. | Address of page table | Ign. | **0** Ign A | PCD | PWT | U/S | R/W | **1** | | **PDE: page table** |
| Ignored | | | | | | | | | | | | **0** | **PDE: not present** |
| XD Prot. Key⁴ Ignored | | Rsvd. | Address of 4KB page frame | Ign. | G PAT D A | PCD | PWT | U/S | R/W | **1** | | **PTE: 4KB page** |
| Ignored | | | | | | | | | | | | **0** | **PTE: not present** |

Execute disabled

Global
Page attribute type
Dirty
Accessed
Caching disabled
Write-through caching
User / Supervisor
Write / Read-only
Present

| 6 6 6 6 5 5 5 5 5 5 5 5 5 | M[1] | M-1 | 3 3 3 2 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 | | | | | | | | | | | | |
| 3 2 1 0 9 8 7 6 5 4 3 2 1 | | | 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 | | | | | | | | | | | | |

| Bits 63–52 / M | M–1 … | Address | PAT/Reserved | Ign. | G | D | A | PCD | PWT | U/S | R/W | P | Entry |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Reserved[2] | | Address of PML4 table | | Ignored | | | | PCD | PWT | | | Ign. | **CR3** |
| XD[3] | Ignored | Rsvd. — Address of page-directory-pointer table | | Ign. | Rsvd | Ign | A | PCD | PWT | U/S | R/W | **1** | **PML4E: present** |
| Ignored | | | | | | | | | | | | **0** | **PML4E: not present** |
| XD | Prot. Key[4] — Ignored | Rsvd. — Address of 1GB page frame — Reserved | PAT | Ign. | G | **1** D | A | PCD | PWT | U/S | R/W | **1** | **PDPTE: 1GB page** |
| XD | Ignored | Rsvd. — Address of page directory | | Ign. | **0** | Ign | A | PCD | PWT | U/S | R/W | **1** | **PDPTE: page directory** |
| Ignored | | | | | | | | | | | | **0** | **PDTPE: not present** |
| XD | Prot. Key[4] — Ignored | Rsvd. — Address of 2MB page frame — Reserved | PAT | Ign. | G | **1** D | A | PCD | PWT | U/S | R/W | **1** | **PDE: 2MB page** |
| XD | Ignored | Rsvd. — Address of page table | | Ign. | **0** | Ign | A | PCD | PWT | U/S | R/W | **1** | **PDE: page table** |
| Ignored | | | | | | | | | | | | **0** | **PDE: not present** |
| XD | Prot. Key[4] — Ignored | Rsvd. — Address of 4KB page frame | | Ign. | G | PAT D | A | PCD | PWT | U/S | R/W | **1** | **PTE: 4KB page** |
| Ignored | | | | | | | | | | | | **0** | **PTE: not present** |

Diagonal column labels (left to right): Execute disabled, Global, Page attribute type, Dirty, Accessed, Caching disabled, Write-through caching, User / Supervisor, Write / Read-only, Present

| Bits 63–M | M | M-1 ... 12 | 11–9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Reserved² | | Address of PML4 table | Ignored | | | | | PCD | PWT | Ign. | | | **CR3** |
| XD³ — Ignored — Rsvd. | | Address of page-directory-pointer table | Ign. | Rsvd | Ign | | A | PCD | PWT | U/S | R/W | **1** | **PML4E: present** |
| Ignored | | | | | | | | | | | | **0** | **PML4E: not present** |
| XD — Prot. Key⁴ — Ignored — Rsvd. | | Address of 1 GB page frame — Reserved | PAT — Ign. | G | **1** | D | A | PCD | PWT | U/S | R/W | **1** | **PDPTE: 1GB page** |
| XD — Ignored — Rsvd. | | Address of page directory | Ign. | **0** | Ign | | A | PCD | PWT | U/S | R/W | **1** | **PDPTE: page directory** |
| Ignored | | | | | | | | | | | | **0** | **PDTPE: not present** |
| XD — Prot. Key⁴ — Ignored — Rsvd. | | Address of 2 MB page frame — Reserved | PAT — Ign. | G | **1** | D | A | PCD | PWT | U/S | R/W | **1** | **PDE: 2MB page** |
| XD — Ignored — Rsvd. | | Address of page table | Ign. | **0** | Ign | | A | PCD | PWT | U/S | R/W | **1** | **PDE: page table** |
| Ignored | | | | | | | | | | | | **0** | **PDE: not present** |
| XD — Prot. Key⁴ — Ignored — Rsvd. | | Address of 4KB page frame | Ign. | G | PAT | D | A | PCD | PWT | U/S | R/W | **1** | **PTE: 4KB page** |
| Ignored | | | | | | | | | | | | **0** | **PTE: not present** |

Execute disabled

Global
Page attribute type
Dirty
Accessed
Caching disabled
Write-through caching
User / Supervisor
Write / Read-only
Present

| 6 6 6 6 5 5 5 5 5 5 5 5 5 5 | M¹ | M-1 | 3 3 3 2 2 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Reserved² | | Address of PML4 table | | Ignored | | P C D | P W T | Ign. | | | | **CR3** |
| XD³ Ignored | Rsvd. | Address of page-directory-pointer table | | Ign. | Rsvd | Ign | A | P C D | P W T | U/S | R/W | 1 → **PML4E: present** |
| Ignored | | | | | | | | | | | 0 | **PML4E: not present** |
| XD Prot. Key⁴ Ignored | Rsvd. | Address of 1GB page frame | Reserved | PAT | Ign. | G 1 | D A | P C D | P W T | U/S | R/W | 1 → **PDPTE: 1GB page** |
| XD Ignored | Rsvd. | Address of page directory | | Ign. | 0 Ign | A | P C D | P W T | U/S | R/W | 1 → **PDPTE: page directory** |
| Ignored | | | | | | | | | | | 0 | **PDTPE: not present** |
| XD Prot. Key⁴ Ignored | Rsvd. | Address of 2MB page frame | Reserved | PAT | Ign. | G 1 | D A | P C D | P W T | U/S | R/W | 1 → **PDE: 2MB page** |
| XD Ignored | Rsvd. | Address of page table | | Ign. | 0 Ign | A | P C D | P W T | U/S | R/W | 1 → **PDE: page table** |
| Ignored | | | | | | | | | | | 0 | **PDE: not present** |
| XD Prot. Key⁴ Ignored | Rsvd. | Address of 4KB page frame | | Ign. | G PAT | D A | P C D | P W T | U/S | R/W | 1 → **PTE: 4KB page** |
| Ignored | | | | | | | | | | | 0 | **PTE: not present** |

4 x 100 ns table lookup +
1 x 100 ns access =
500 ns =
~1,500 instructions @ 3 GHz clock

I could be doing a lot of useful work
while waiting to read memory !

What to do ?

# Caching

# Key Issues in Caching

- Locality

  - *temporal (reuse within short amount of time)*

  - *spatial (use shortly physically nearby data)*

- Replacement algorithm

  - *LRU, Time-aware LRU in CDNs, Least-frequent LRU in CDNs, MRU in scanning big data*

- Write-through vs. write-back

- Working set

- Direct-mapped vs. partially associative vs. fully associative

- Size vs. speed trade-off

# Intel Skylake Microarchitecture



Gen9

Core    Core

L3$ Slice    L3$ Slice

L3$ Slice    L3$ Slice

Ring

Core    Core

**System Agent**

Display Controller

PCIe

eDRAM Controller (optional)

Memory Controller

**Front End**

Instruction Cache Tag

μOP Cache Tag

L1 Instruction Cache 32KiB 8-Way

Instruction TLB

**16 Bytes/cycle**

Branch Predictor (BPU)

Instruction Fetch & PreDecode (16 B window)

MOP MOP MOP MOP MOP MOP

Instruction Queue (50, 2x25 entries)    Macro-Fusion

MOP MOP MOP MOP MOP

MicroCode Sequencer ROM (MS ROM)

5-Way Decode

Complex Decoder | Simple Decoder | Simple Decoder | Simple Decoder | Simple Decoder

1-4 μOPs    μOP    μOP    μOP    μOP

Stack Engine (SE)

Adder Adder Adder

**4 μOPs**    5 μOPs

Decoded Stream Buffer (DSB) (μOP Cache) (1.5k μOPs; 8-Way) (64 B window)

**6 μOPs**

MUX

Loop Stream Detector (LSD)    Allocation Queue (IDQ) (128, 2x64 μOPs)    Micro-Fusion

**64B/cycle**

μOP μOP μOP μOP μOP μOP

Register Alias Table (RAT)

Branch Order Buffer (BOB) (48-entry)

4 μOP

Load    FP    Int Vect    Int    Store

Common Data Buses (CDBs)

Move Elimination    Rename / Allocate / Retirement ReOrder Buffer (224 entries)    Ones Idioms    Zeroing Idioms

μOP μOP μOP μOP μOP μOP μOP μOP

Integer Physical Register File (180 Registers)

Scheduler Unified Reservation Station (RS) (97 entries)

Vector Physical Register File (168 Registers)

Port 0 | Port 1 | Port 5 | Port 6 | Port 2 | Port 3 | Port 4 | Port 7

μOP μOP μOP μOP μOP μOP μOP μOP

**EUs**

INT ALU / INT DIV / INT Vect ALU / INT Vect MUL / FP FMA / AES / Vect String / FP DIV / Branch

INT ALU / INT MUL / INT Vect ALU / INT Vect MUL / FP FMA / Bit Scan

INT ALU / Vect Shuffle / INT Vect ALU / LEA

INT ALU / Branch

AGU / Load Data

AGU / Load Data

Store Data

AGU

256bit/cycle

**Execution Engine**

Unified STLB

L2 Cache 256KiB 4-Way

**32B/cycle** To L3

**Memory Subsystem**

Store Buffer & Forwarding (56 entries)

32B/cycle

Load Buffer (72 entries)

32B/cycle    32B/cycle

L1 Data Cache 32KiB 8-Way

Data TLB

Line Fill Buffers (LFB) (10 entries)

**64B/cycle**

**Front End**

Instruction Cache Tag

μOP Cache Tag

L1 Instruction Cache
32KiB 8-Way

Instruction TLB

**16 Bytes/cycle**

Branch Predictor (BPU)

Instruction Fetch & PreDecode
(16 B window)

MOP MOP MOP MOP MOP MOP

Instruction Queue
(50, 2x25 entries)

Macro-Fusion

MOP MOP MOP MOP MOP

MicroCode Sequencer ROM (MS ROM)

5-Way Decode

Complex Decoder | Simple Decoder | Simple Decoder | Simple Decoder | Simple Decoder

1-4 μOPs   μOP   μOP   μOP   μOP

Stack Engine (SE)

Adder Adder Adder

**4 μOPs**

Decoded Stream Buffer (DSB)
(μOP Cache)
(1.5k μOPs; 8-Way)
(64 B window)

**6 μOPs**

**5 μOPs**

MUX

Loop Stream Detector (LSD)

Allocation Queue (IDQ) (128, 2x64 μOPs)

Micro-Fusion

μOP μOP μOP μOP μOP

Branch Order Buffer (BOB) (48-entry)

ate / Retirement
r (224 entries)

Ones Idioms | Zeroing Idioms

P    μOP    μOP    μOP    μOP

Scheduler
Reservation Station (RS)
(97 entries)

Vector Physical Register File
(168 Registers)

Port 6 | Port 2 | Port 3 | Port 4 | Port 7

μOP    μOP    μOP    μOP    μOP

NT ALU Branch | AGU Load Data | AGU Load Data | Store Data | AGU

FP FMA | FP FMA
AES | Bit Scan
Vect String
FP DIV
Branch

**EUs**

**Execution Engine**

Unified STLB

L2 Cache
256KiB 4-Way

**32B/cycle**

**256bit/cycle**

Store Buffer & Forwarding
(56 entries)

**32B/cycle**

**64B/cycle**

Load Buffer
(72 entries)

**32B/cycle**

L1 Data Cache
32KiB 8-Way

Data TLB

Line Fill Buffers (LFB)
(10 entries)

**Memory Subsystem**

L3 cache

RAM

*George Candea*

*Principles of Computer Systems*

**Front End**

Instruction Cache Tag
μOP Cache Tag

L1 Instruction Cache
32KiB 8-Way

Instruction TLB

**16 Bytes/cycle**

Branch Predictor (BPU)

Instruction Fetch & PreDecode
(16 B window)

MOP MOP MOP MOP MOP MOP

Instruction Queue
(50, 2x25 entries)    Macro-Fusion

MOP MOP MOP MOP MOP

MicroCode Sequencer ROM (MS ROM)

5-Way Decode

Complex Decoder | Simple Decoder | Simple Decoder | Simple Decoder | Simple Decoder

1-4 μOPs   μOP   μOP   μOP   μOP

Stack Engine (SE)

Adder Adder Adder

4 μOPs

5 μOPs

Decoded Stream Buffer (DSB)
(μOP Cache)
(1.5k μOPs; 8-Way)
(64 B window)

6 μOPs

MUX

Loop Stream Detector (LSD)    Allocation Queue (IDQ) (128, 2x64 μOPs)    Micro-Fusion

μOP μOP μOP μOP μOP    Branch Order Buffer (BOB) (48-entry)

...ate / Retirement ...r (224 entries)    Ones Idioms    Zeroing Idioms

μOP   μOP   μOP   μOP   μOP

Scheduler Reservation Station (RS) (97 entries)    Vector Physical Register File (168 Registers)

Port 6   Port 2   Port 3   Port 4   Port 7

μOP   μOP   μOP   μOP   μOP

...NT ALU Branch    AGU Load Data    AGU Load Data    Store Data    AGU

FP FMA | FP FMA
AES | Bit Scan
Vect String
FP DIV
Branch

**EUs**

**Execution Engine**

Unified STLB

L2 Cache 256KiB 4-Way

**32B/cycle**

256bit/cycle

L3 cache

RAM

Store Buffer & Forwarding (56 entries)

32B/cycle

64B/cycle

Load Buffer (72 entries)

32B/cycle

32B/cycle

L1 Data Cache 32KiB 8-Way    Data TLB
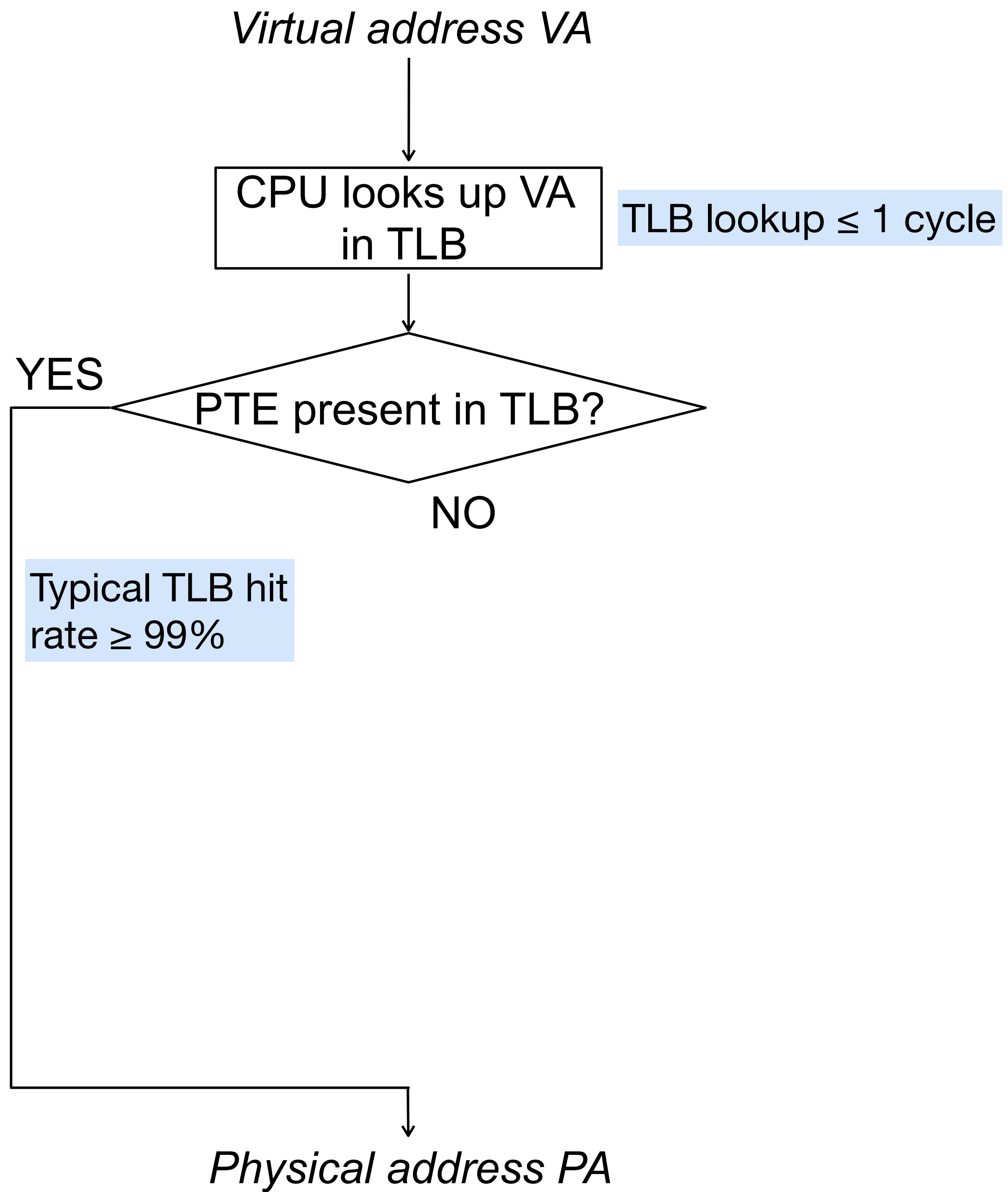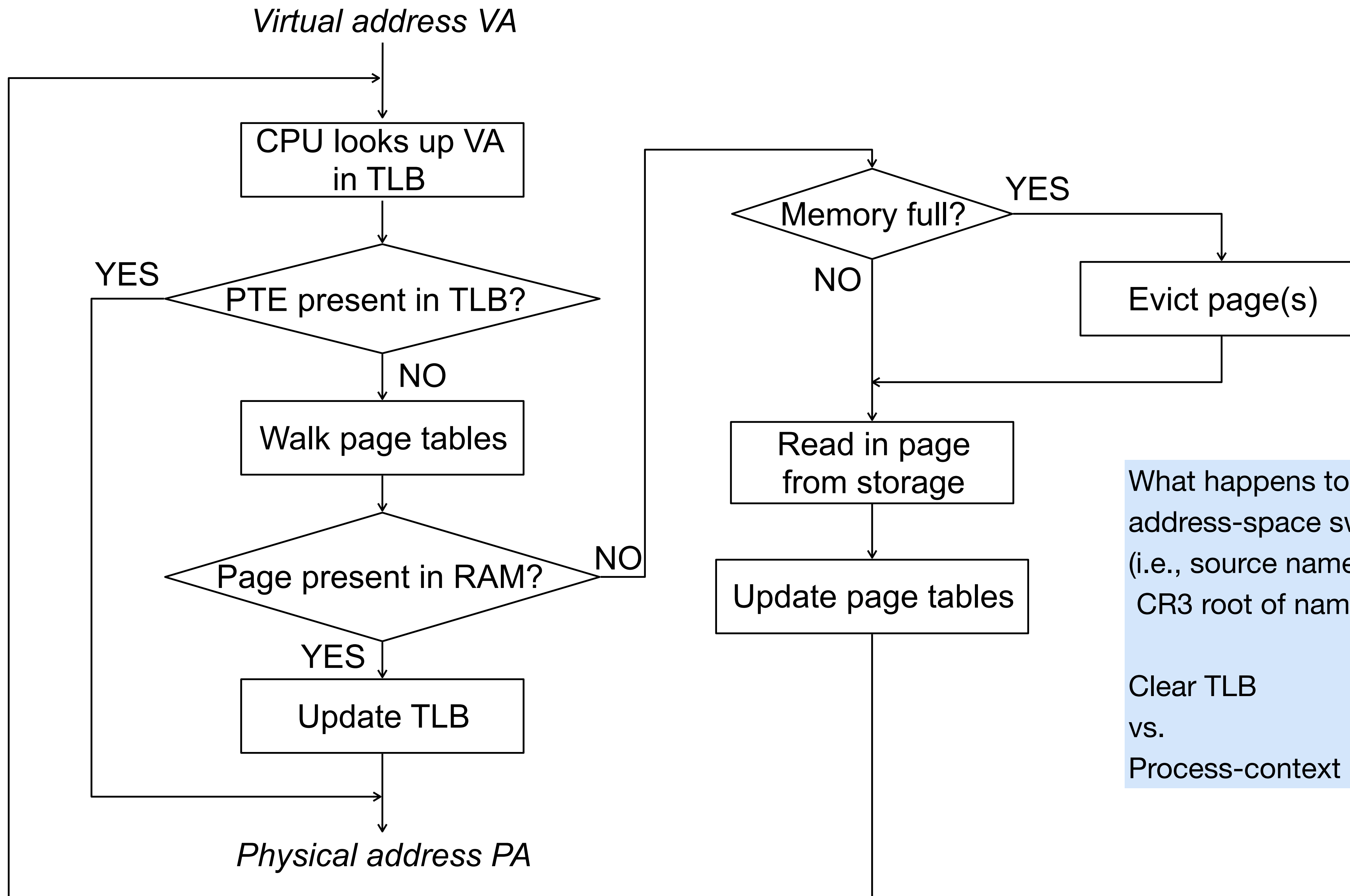
Line Fill Buffers (LFB) (10 entries)

**Memory Subsystem**

Cache the translations of names using TLB

Split TLB to exploit access locality and on-chip location
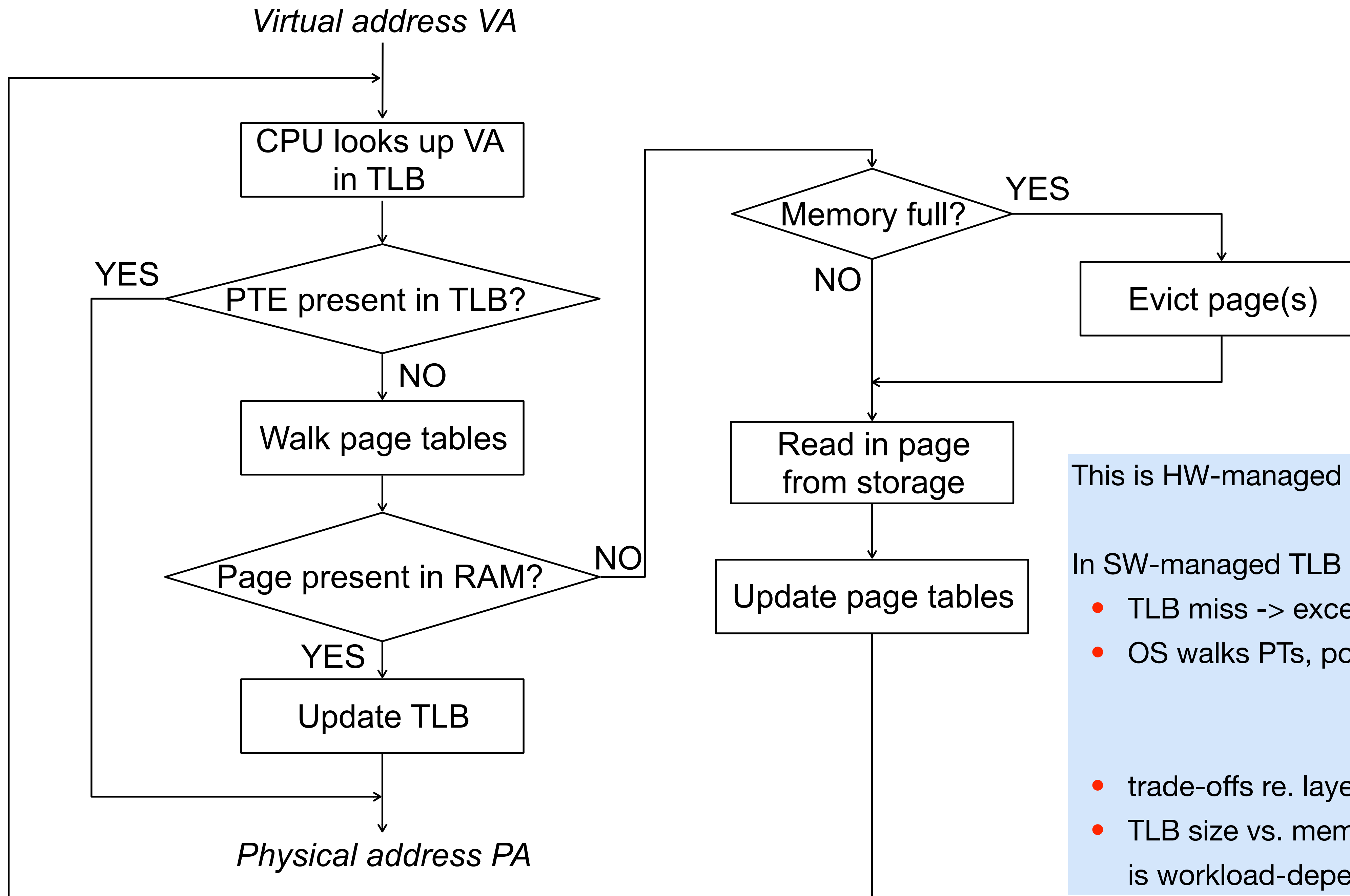
Technology trade-offs: speed vs. cost

*Virtual address VA*

CPU looks up VA in TLB

TLB lookup ≤ 1 cycle

PTE present in TLB?

YES

NO

Typical TLB hit rate ≥ 99%

Read from L1 = 3-4 cycles
Read from L2 = 10-12 cycles
Read from L3 = 30-80 cycles
Read from RAM = 100-150 cycles

*Physical address PA*

*Virtual address VA*

CPU looks up VA in TLB

PTE present in TLB?
- YES
- NO

Walk page tables

Page present in RAM?
- YES
- NO

Update TLB

*Physical address PA*

Memory full?
- YES → Evict page(s)
- NO

Read in page from storage

Update page tables

What happens to TLB on an address-space switch?
(i.e., source namespace differs =>
 CR3 root of namespace changes)

Clear TLB
vs.
Process-context identifiers

*Virtual address VA*

CPU looks up VA in TLB

PTE present in TLB?

YES

NO

Walk page tables

Page present in RAM?

NO

YES

Update TLB

*Physical address PA*

Memory full?

YES

NO

Evict page(s)

Read in page from storage

Update page tables

This is HW-managed TLB !

In SW-managed TLB (MIPS, SPARC, …)
- TLB miss -> exception to OS
- OS walks PTs, populates TLB

- trade-offs re. layer to delegate to
- TLB size vs. memory size trade-off is workload-dependent => hard!

# Indexing and Tagging in Caches

- Index into cache → check the tag

- Types of caches

  - *PIPT (Physically indexed / Physically tagged)*

  - *VIVT (Virtually indexed / Virtually tagged)*

  - *VIPT (Virtually indexed / Physically tagged)*

  - ~~*PIVT (Physically indexed / Virtually tagged)*~~

- Caches on x86

  - *L1 is VIPT*

  - *L2 and L3 are PIPT*

# Constants in System Design

# Systems "Constants" (CPU cycles)

- Register-register ADD/OR/etc. < 1 CPU cycle

- Memory write ~1 cycle

- Correctly / incorrectly predicted "if" branch = 1 cycle / 10-20 cycles

- L1 / L2 / L3 / main RAM read = 3-4 cycles / 10-12 cycles / 30-80 cycles / 100-150 cycles

- TLB hit / miss = 0.5 - 1 cycle / 7-21 cycles

- CAS = 15-30 cycles

- C function direct / indirect call = 15-30 cycles / 20-50 cycles

- Kernel syscall = 1,000 - 1,500 cycles

- Thread context switch (direct costs) = 2,000 cycles

- On NUMA, different-socket mem hierarchy access is 3 - 10x that of non-NUMA

# Systems "Constants" (absolute time)

- L1 / L2 / main RAM reference = 1 / 4 / 100 ns

- Mutex lock or unlock = 17 ns

- Branch misprediction = 3 ns

- Send 2KB (2,000 bytes) over commodity network = 88 ns

- Compress 1 KB with Snappy = 2,000 ns = 2 microsec

- SSD random read = 16,000 ns = 16 microsec

- Read 1 MB sequentially from RAM = 5,000 ns = 5 microsec

- Packet roundtrip in same datacenter < 50 microsec

- Read 1 MB sequentially from SSD = 78 microsec

- Seek on magnetic disk = 3 millisec

- Read 1 MB sequentially from magnetic disk = 1 millisec

- Packet roundtrip CA -> Amsterdam -> CA = 150 millisec