



# Modularity through Memory Virtualization

---

Prof. George Canea

*School of Computer & Communication Sciences*

# Objectives

---

- Understand memory virtualization
  - *the coolest form of modularization we have in operating systems*
- Back-of-the-envelope calculations
  - *sanity-check a design before writing any code*

# Outline

---

---

- Names and Page Tables
  - $PT = \text{directory}$  for mapping memory names (VA) to memory locations (PA)
  - Specific example: Intel Skylake microarchitecture
- Caching
  - Generalities
  - Caching in name translation
- Key reference values for system design

# Examples of Names

---

# Examples of Names

---

- Memory address

# Examples of Names

---

- Memory address
  - *names a memory location*

# Examples of Names

---

- Memory address
  - *names a memory location*
- Pointer

# Examples of Names

---

- Memory address
  - *names a memory location*
- Pointer
  - `void*` names the location of a memory word

# Examples of Names

---

- Memory address
  - *names a memory location*
- Pointer
  - `void*` names the location of a memory word
  - `int*` names the location of an integer

# Examples of Names

---

- Memory address
  - *names a memory location*
- Pointer
  - `void*` names the location of a memory word
  - `int*` names the location of an integer
- Virtual memory address

# Examples of Names

---

- Memory address
  - *names a memory location*
- Pointer
  - `void*` names the location of a memory word
  - `int*` names the location of an integer
- Virtual memory address
  - *names (from userspace) a physical memory address*

# **Names & Page Tables**

# Names

---

- Scope
  - *Private*: unique within a context (e.g., a private IP address)
  - *Global*: unique across contexts (e.g., a global IP address)

# Names

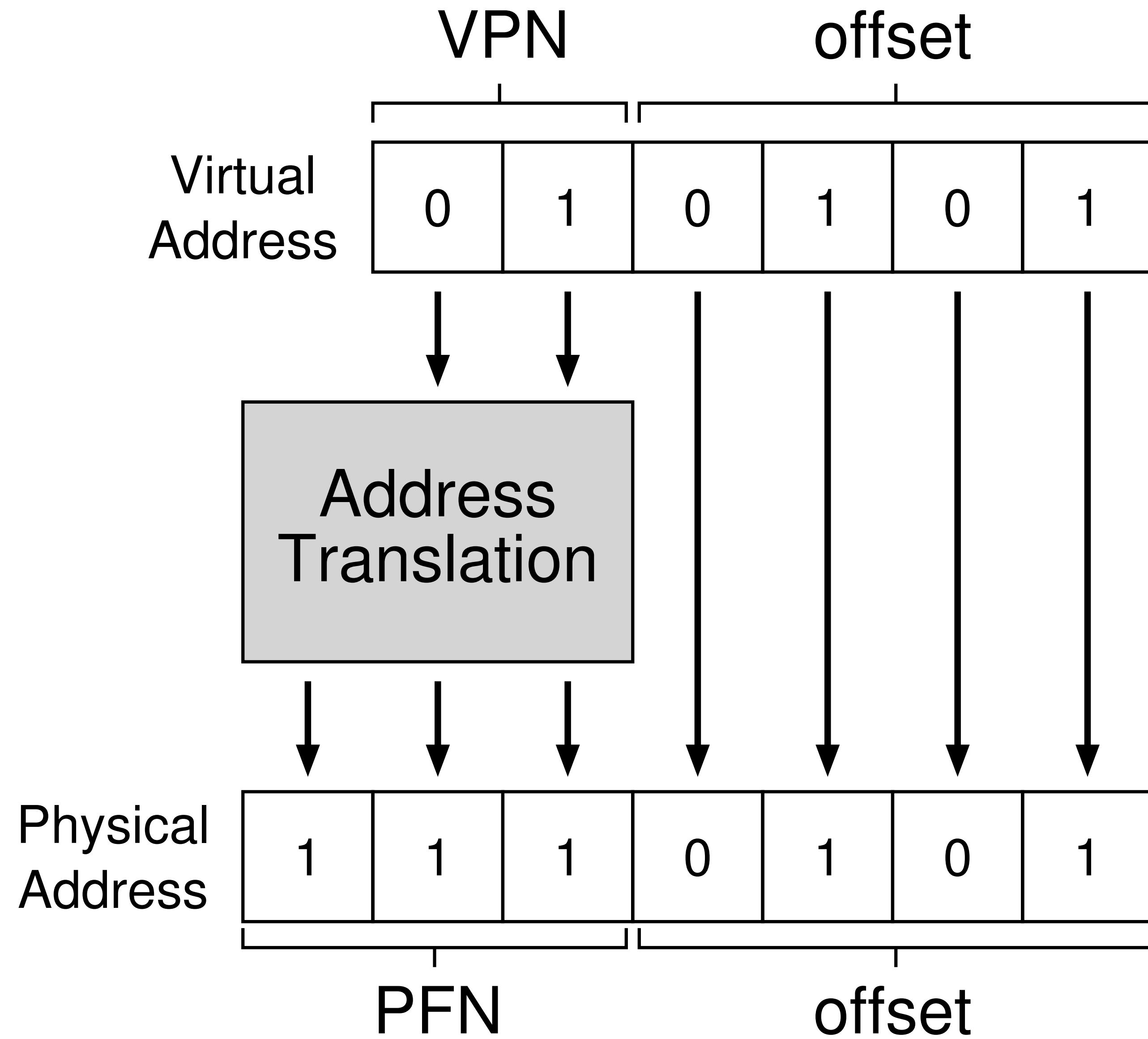
---

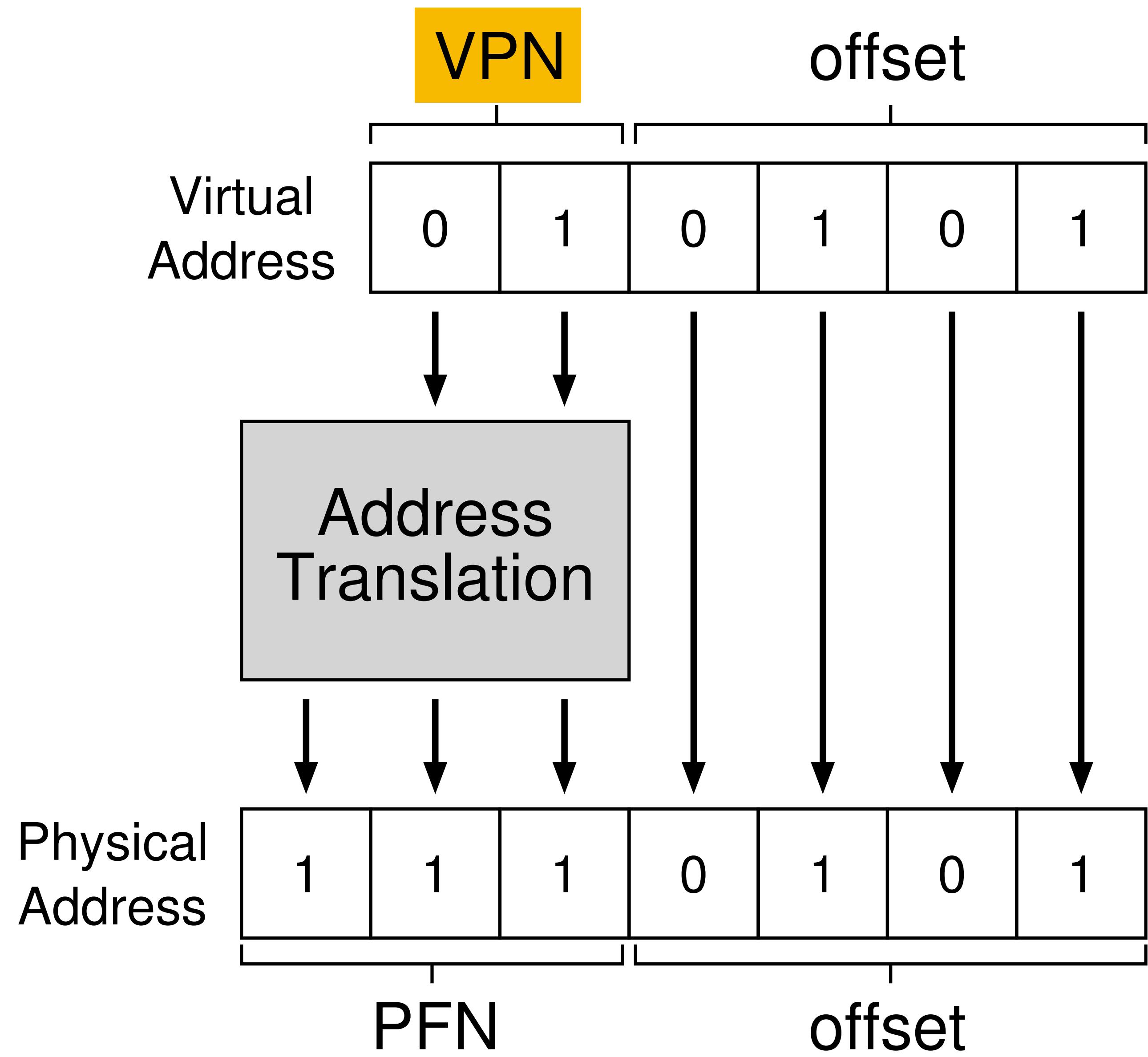
- Scope
  - *Private*: unique within a context (e.g., a private IP address)
  - *Global*: unique across contexts (e.g., a global IP address)
- Structure
  - *Hierarchical*: name relationship implies object relationship (e.g., two IP addresses sharing the same prefix)
  - *Flat*: name relationship implies nothing (e.g., content IDs in Peer-to-Peer networks)

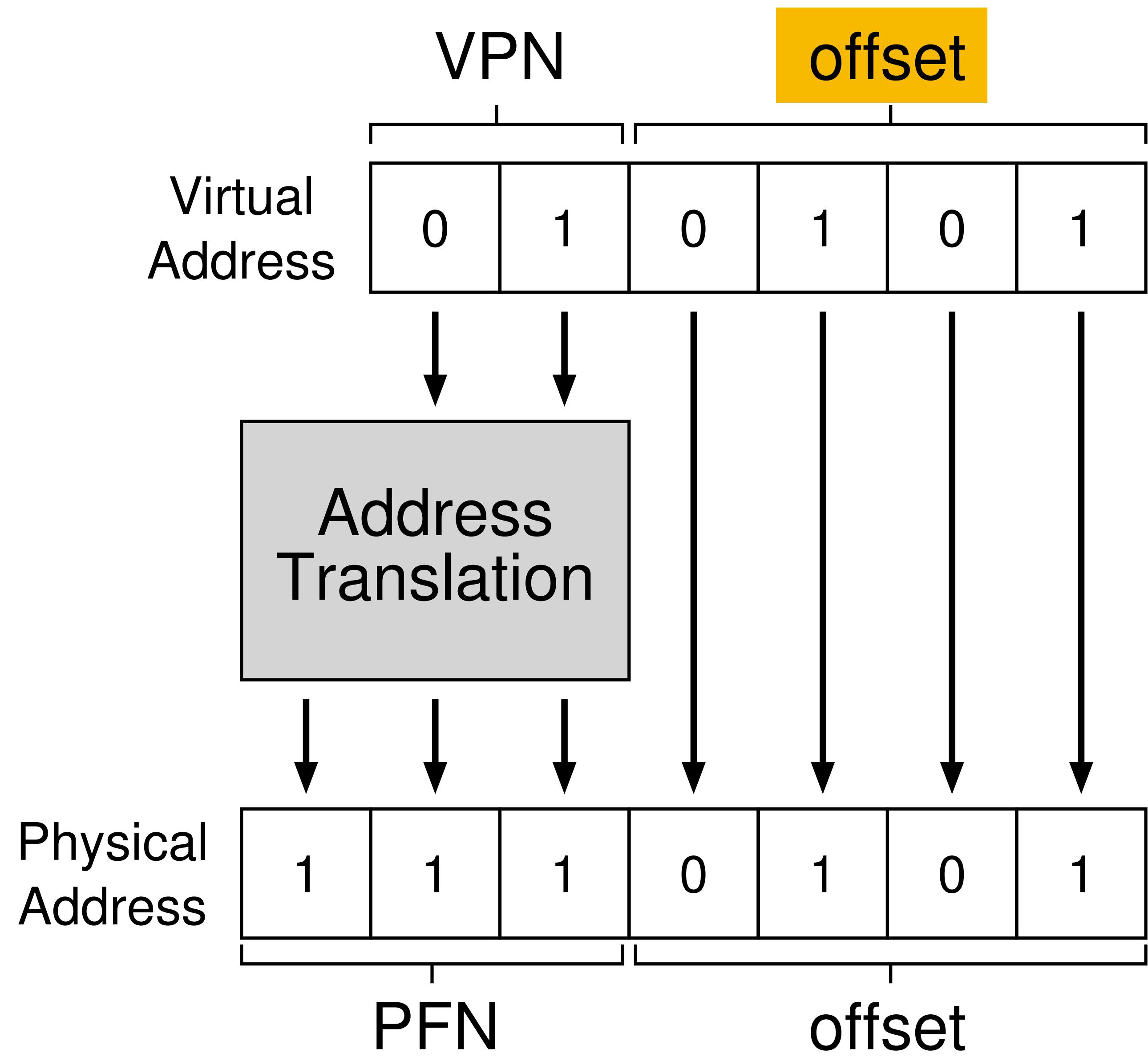
# Names

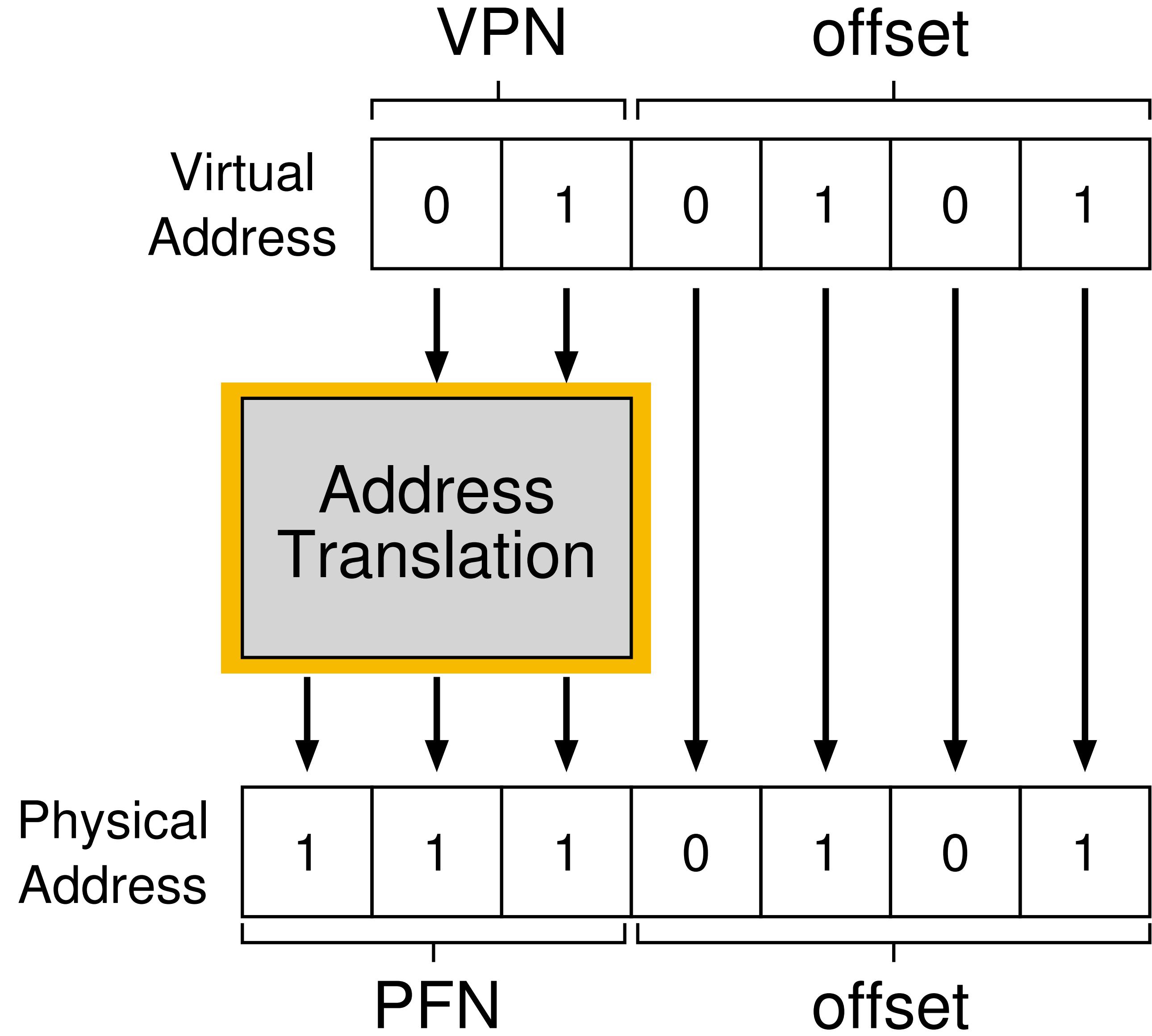
---

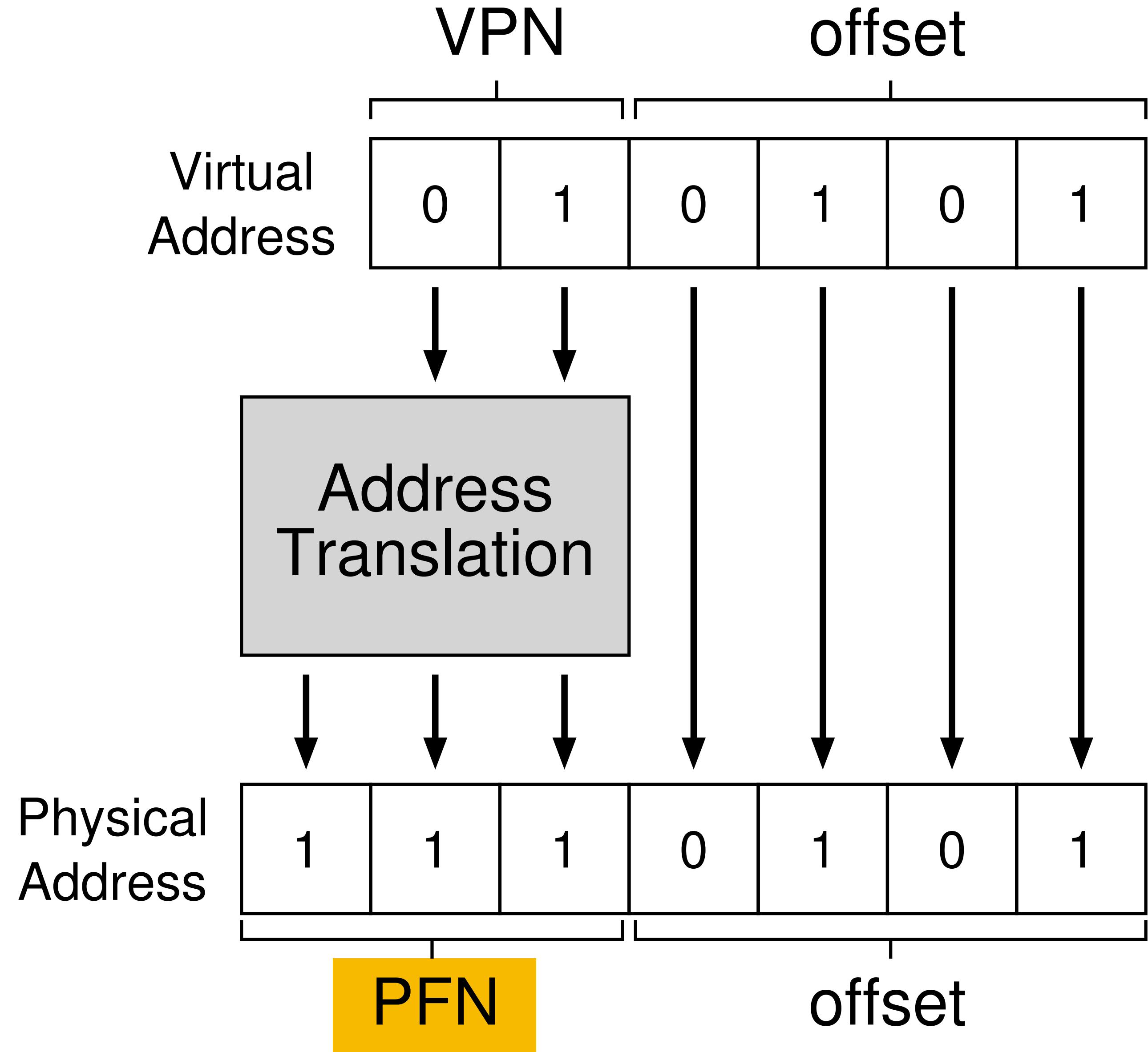
- Scope
  - *Private*: unique within a context (e.g., a private IP address)
  - *Global*: unique across contexts (e.g., a global IP address)
- Structure
  - *Hierarchical*: name relationship implies object relationship (e.g., two IP addresses sharing the same prefix)
  - *Flat*: name relationship implies nothing (e.g., content IDs in Peer-to-Peer networks)
- Naming system
  - *Directories* of name→value mappings, support name lookups and updates

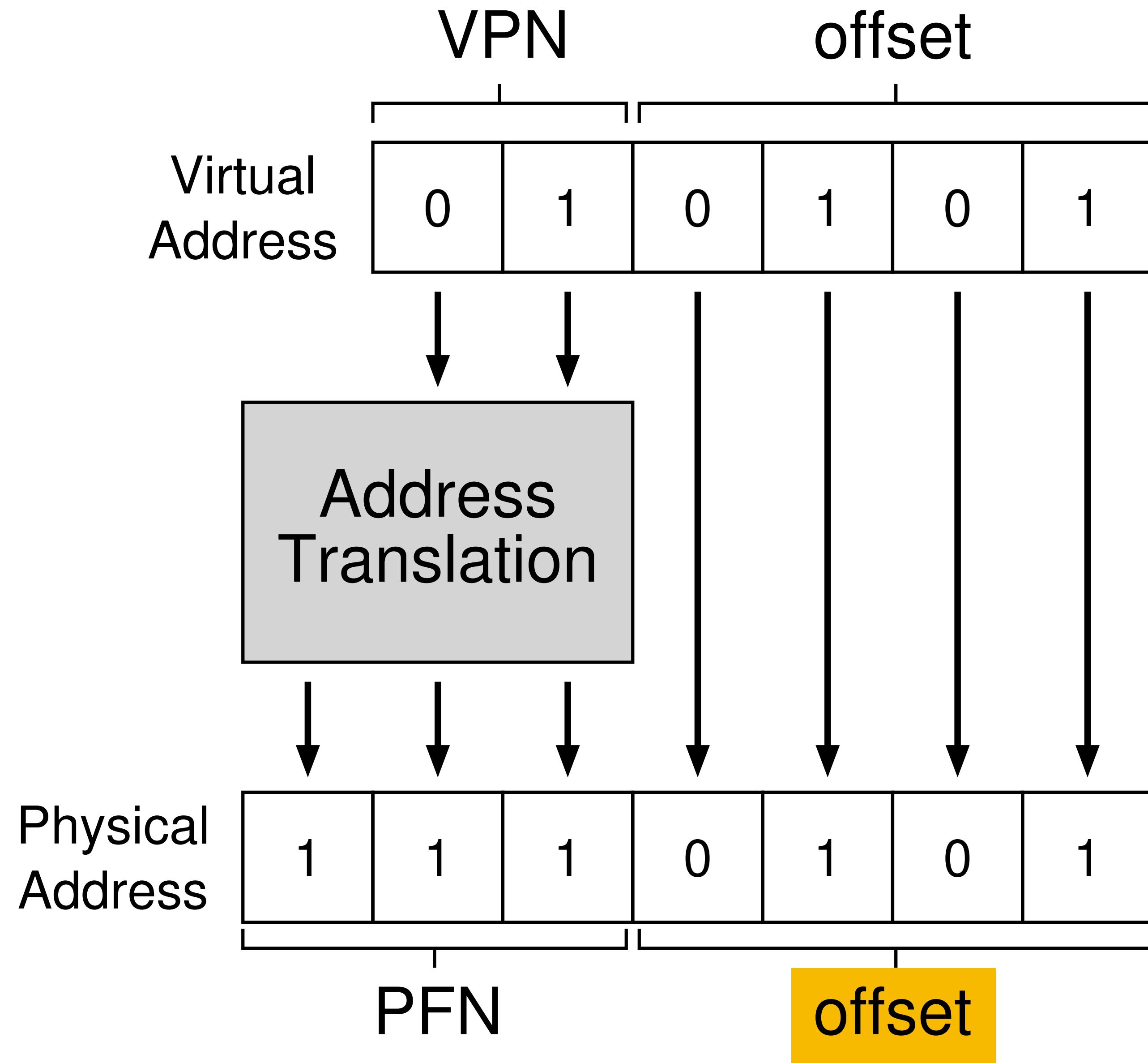


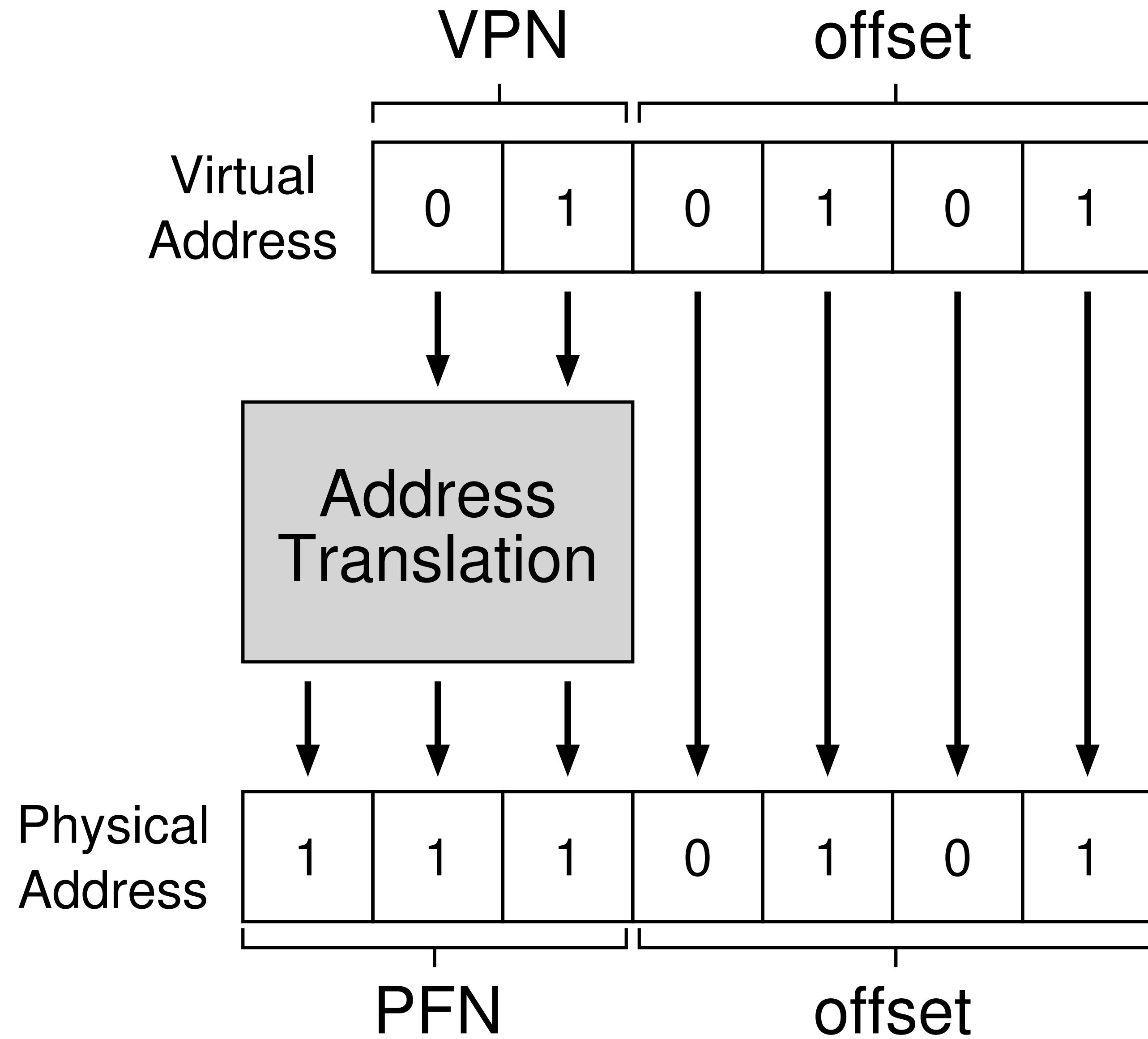




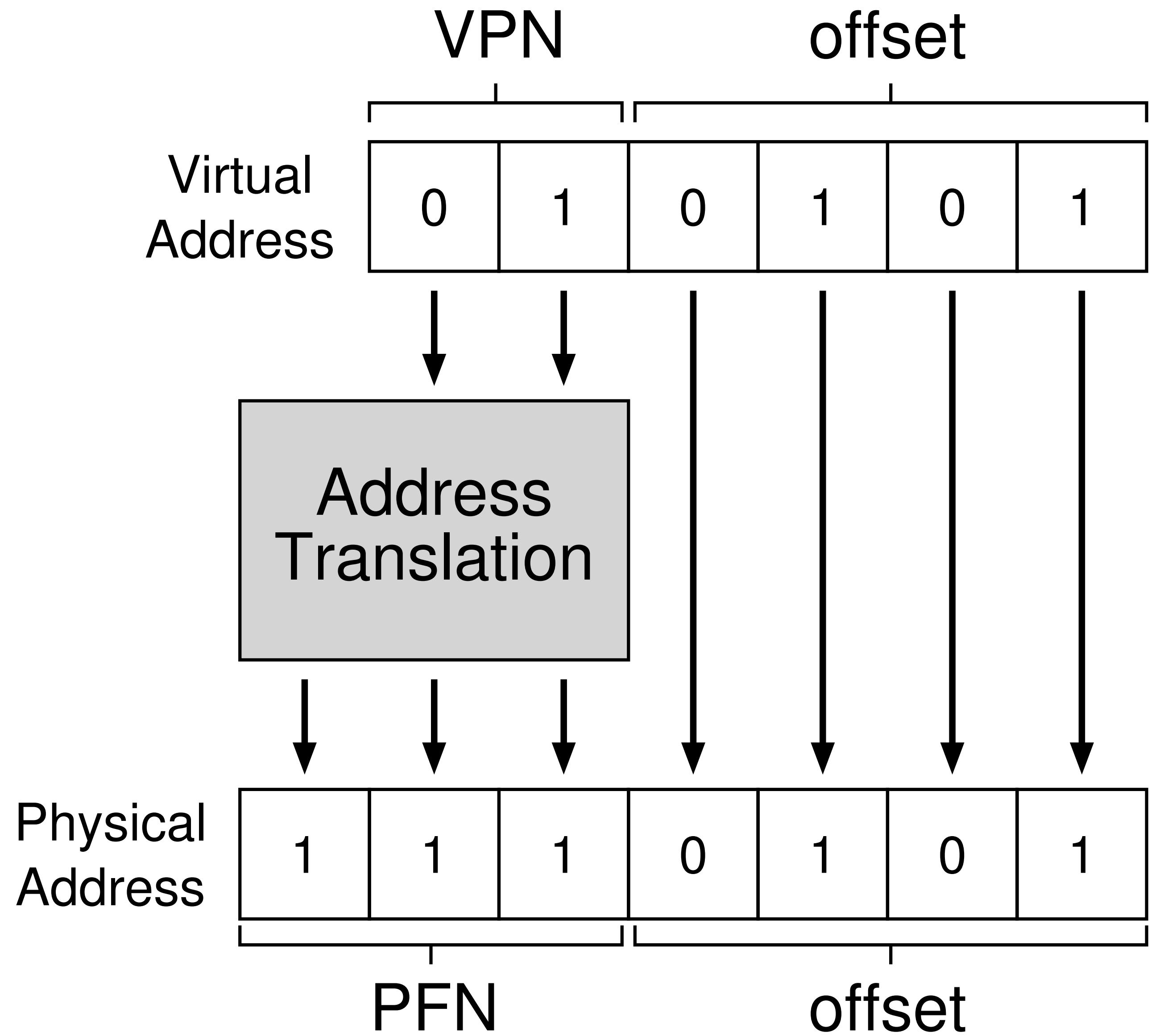








# Indirection

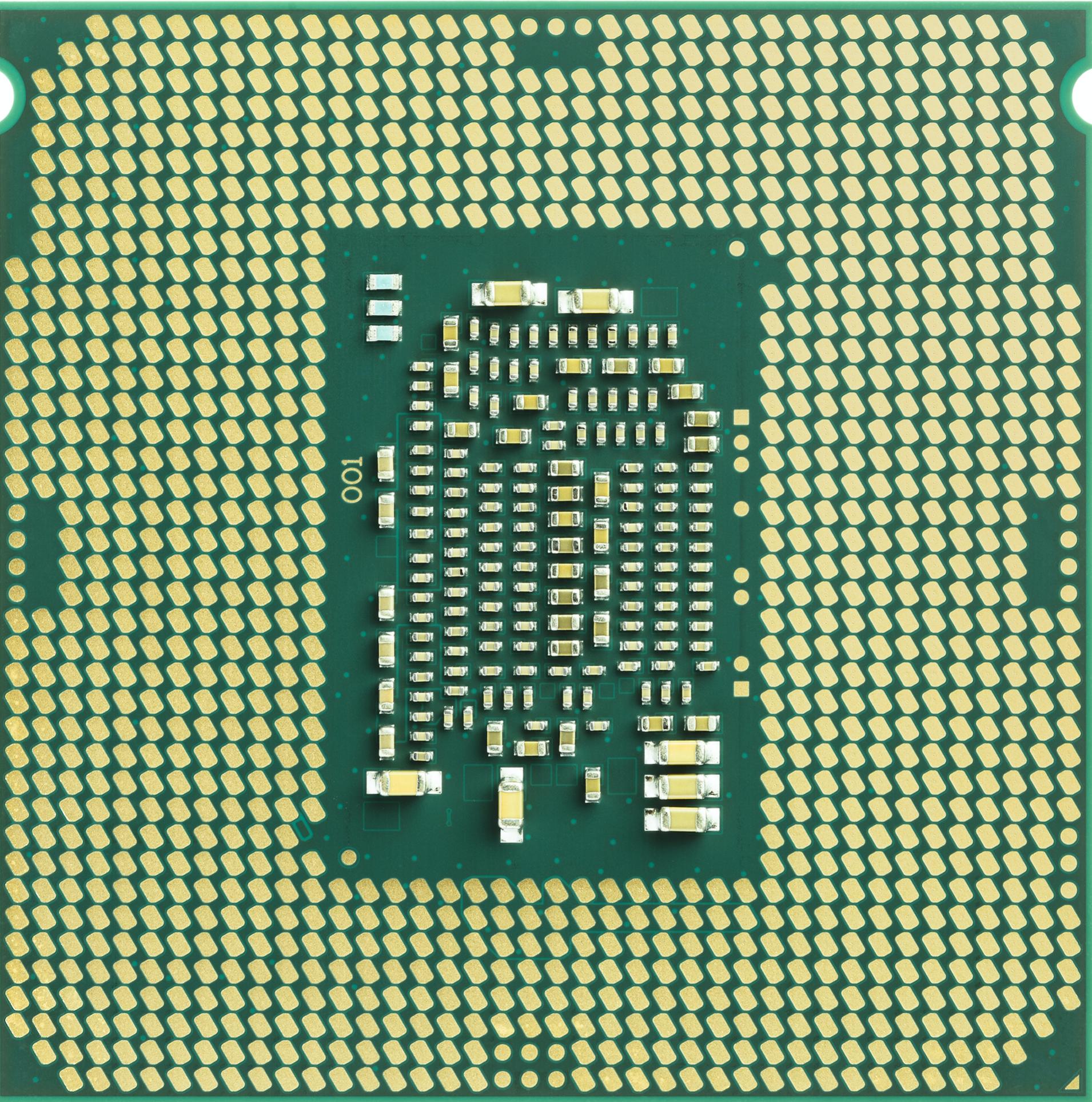


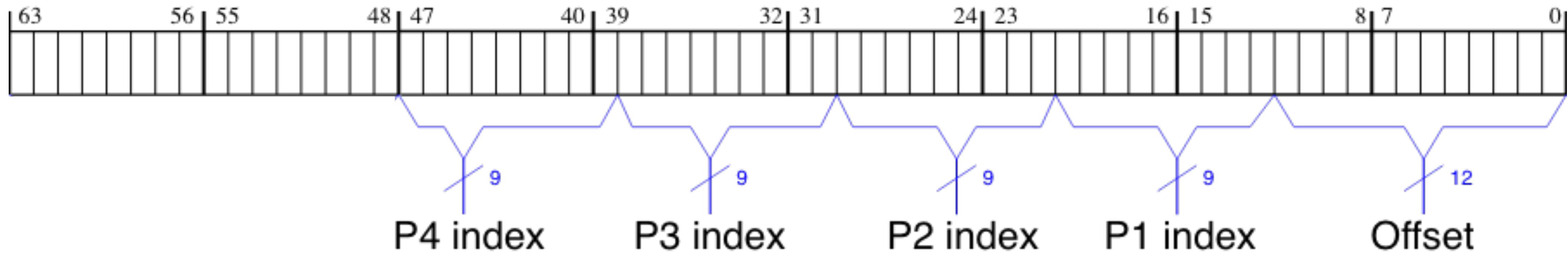
*Virtual  
namespace*

*Name  
translation  
controls  
visibility*

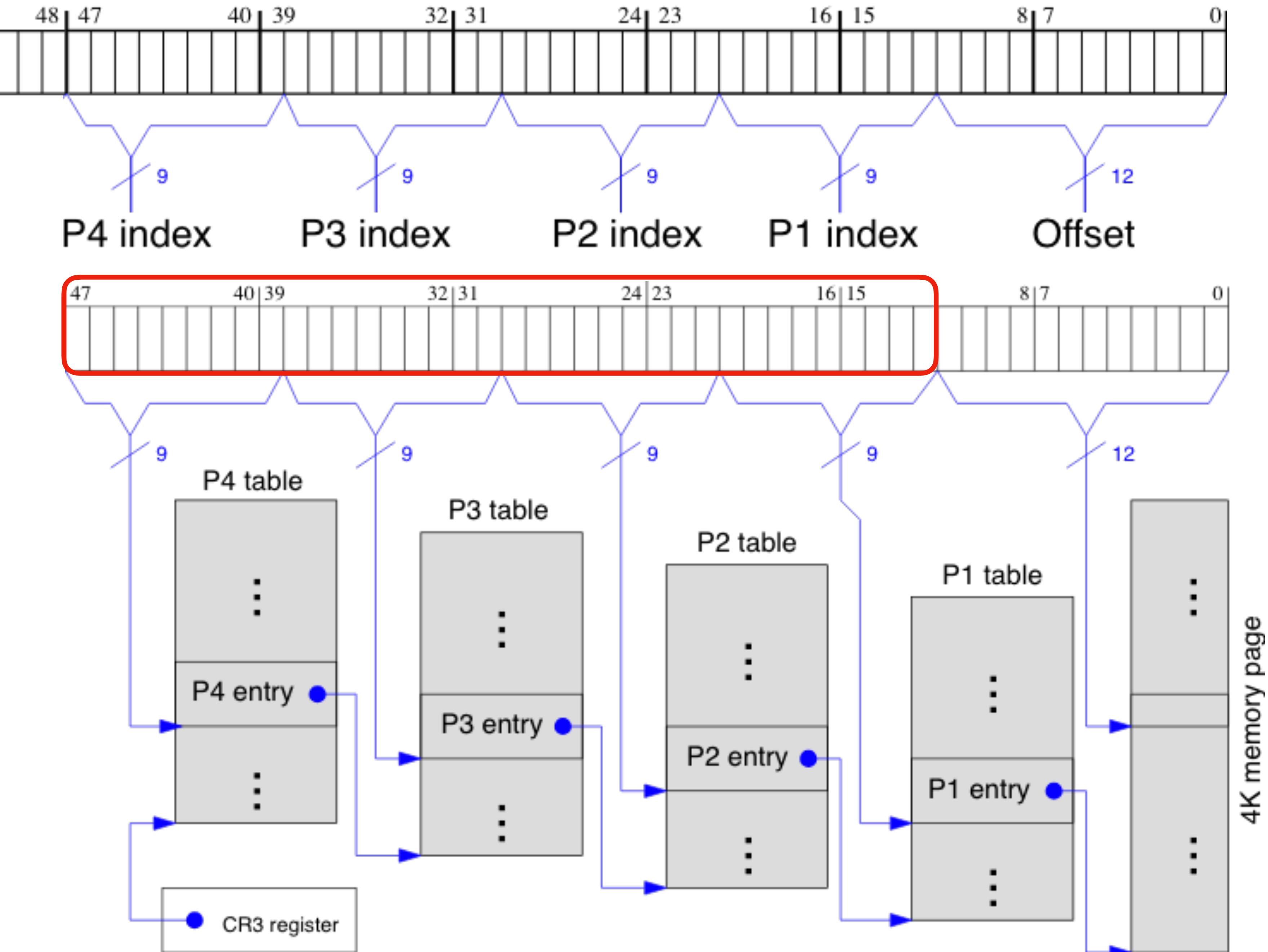
*"if you can't  
name it, you  
can't use it"*

*Physical  
namespace*

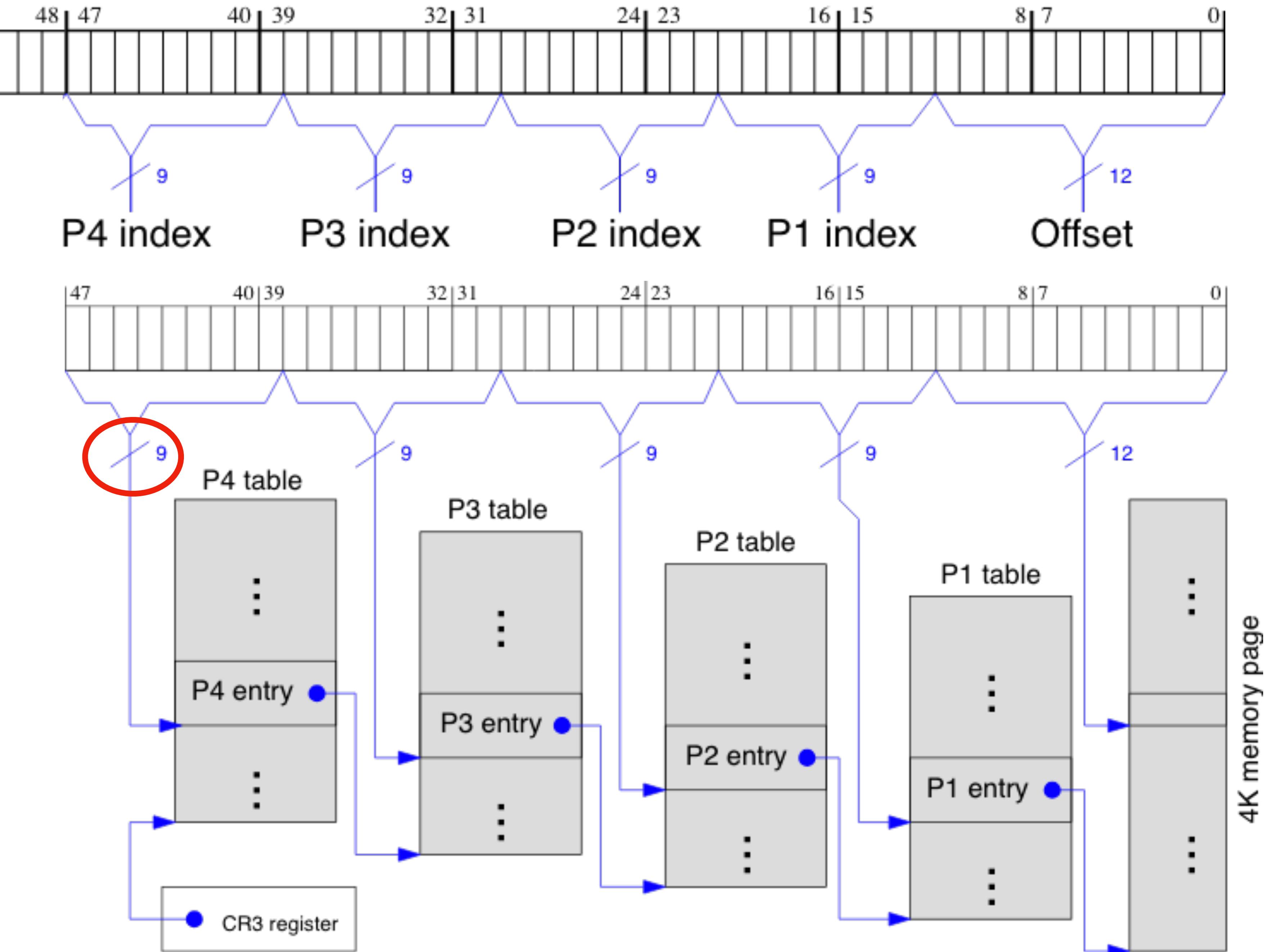




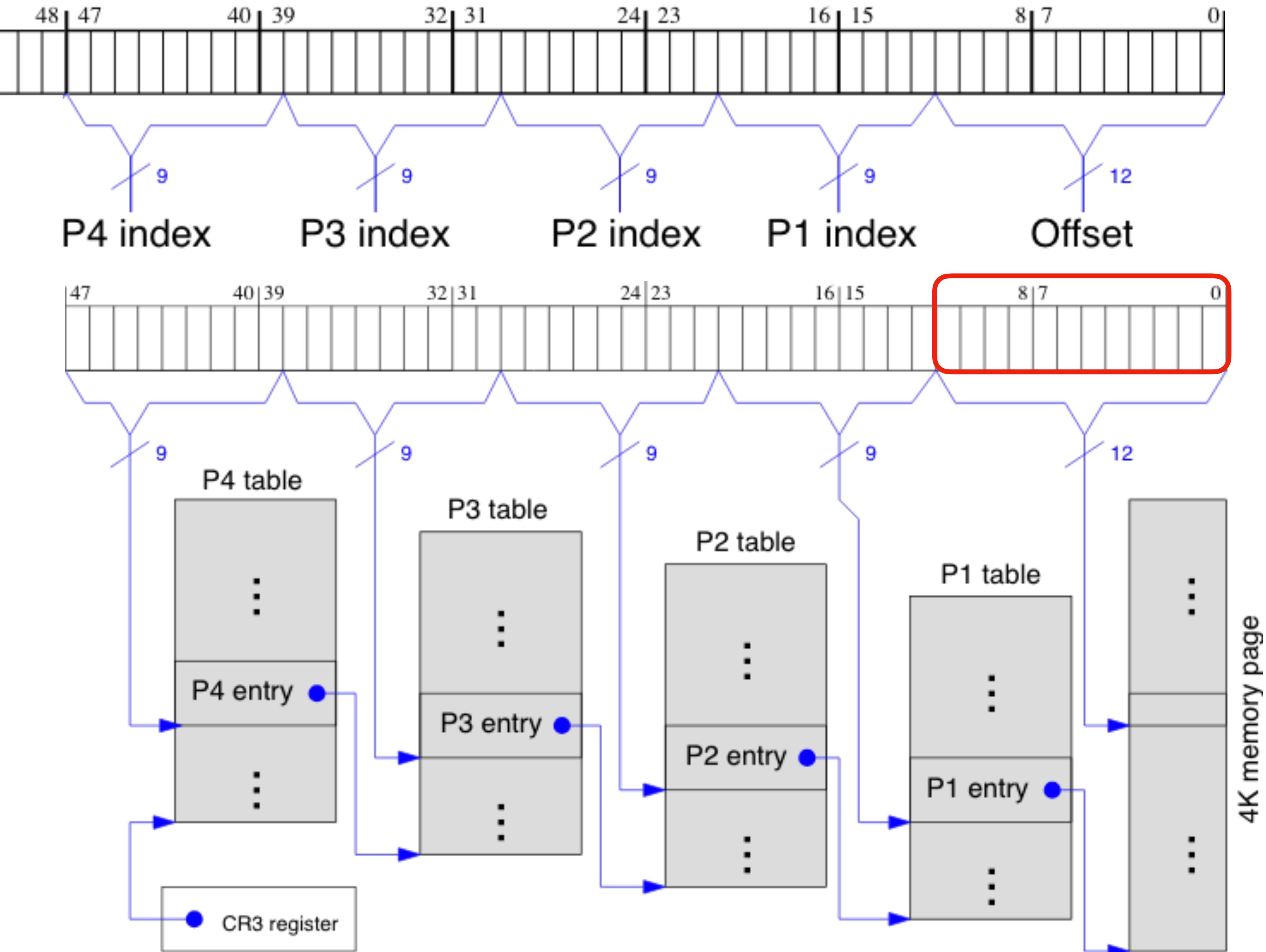
HW walks page tables  
OS updates them

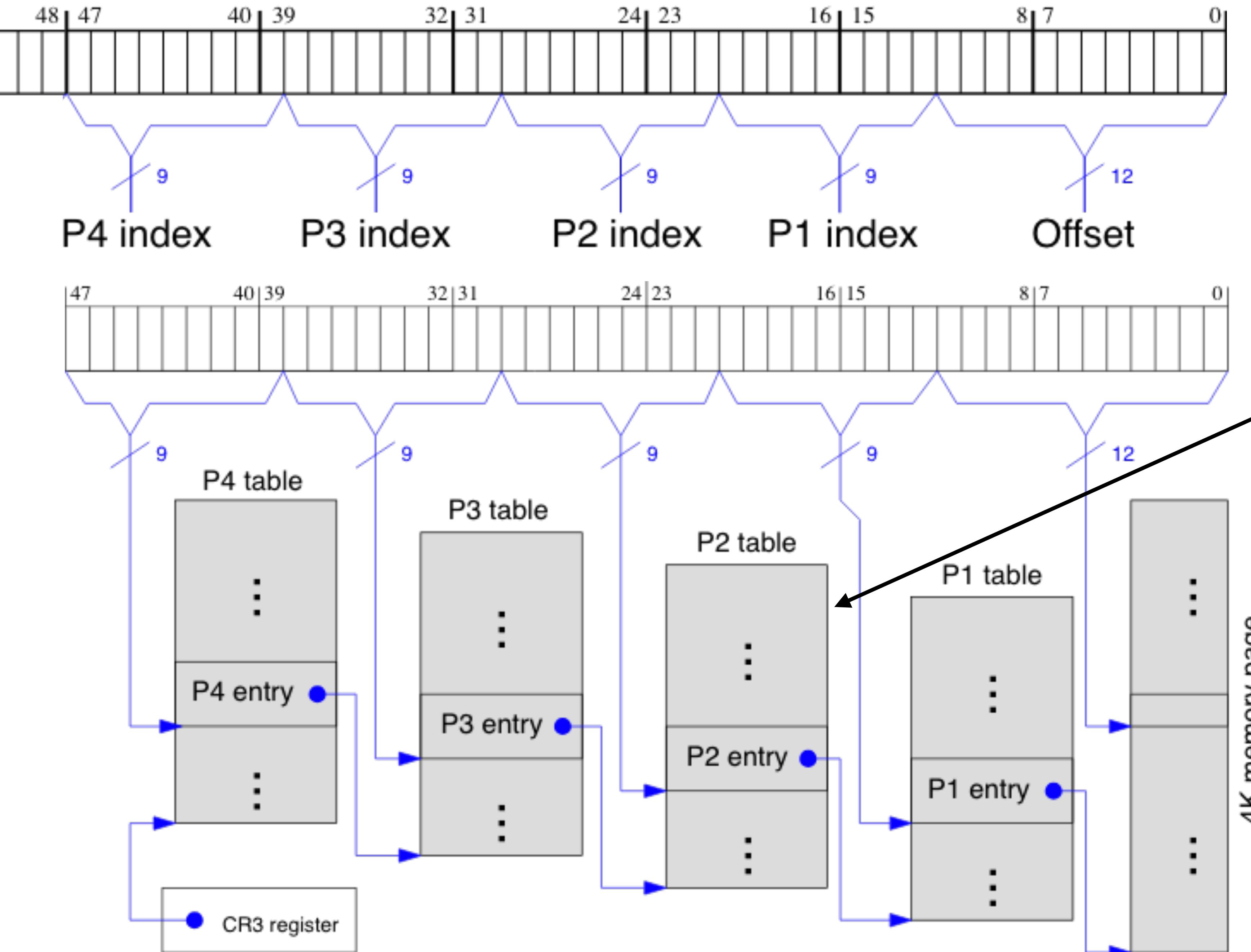


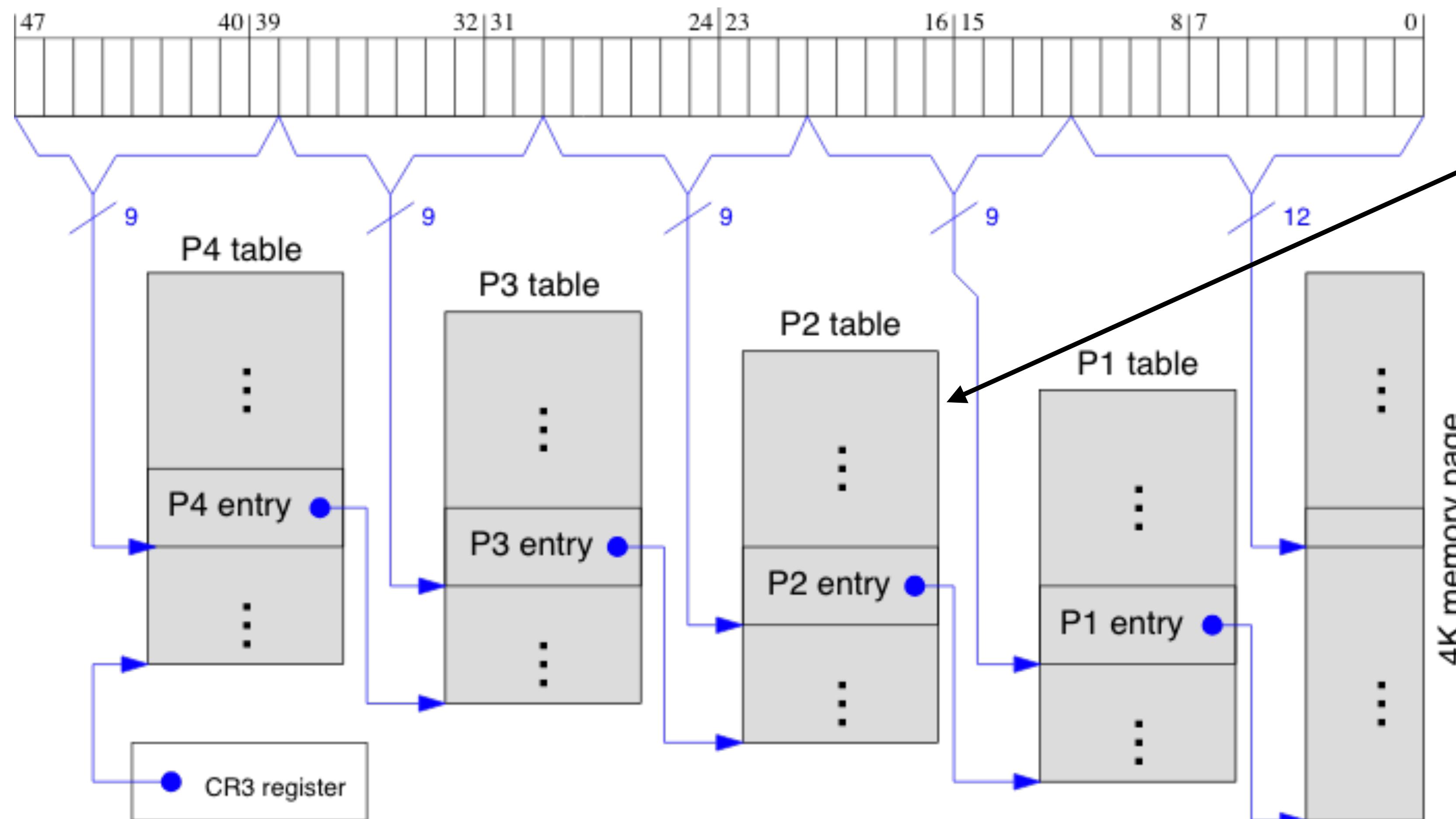
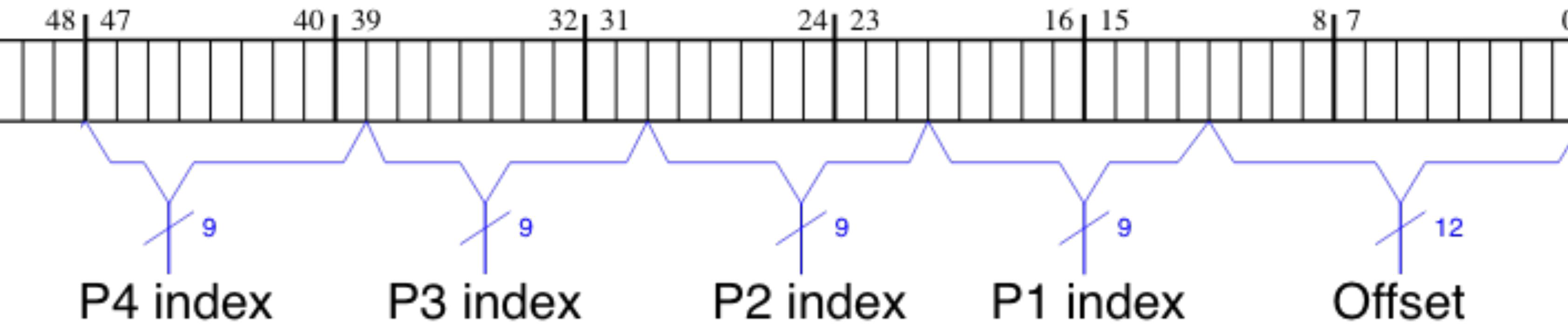
HW walks page tables  
OS updates them



HW walks page tables  
OS updates them







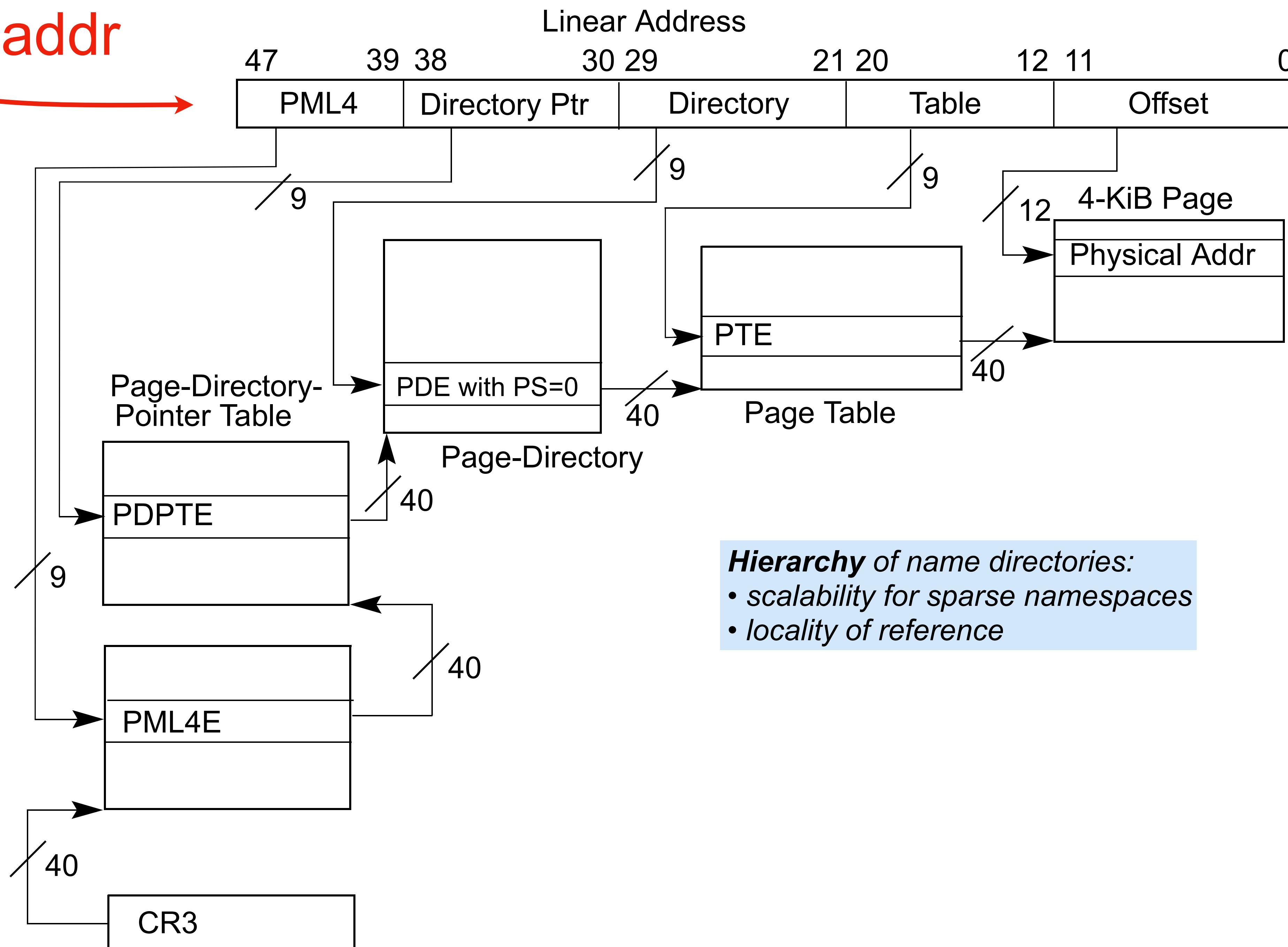
$2^9 = 512$  entries  
8 bytes / entry

**Hierarchy** of name directories:  

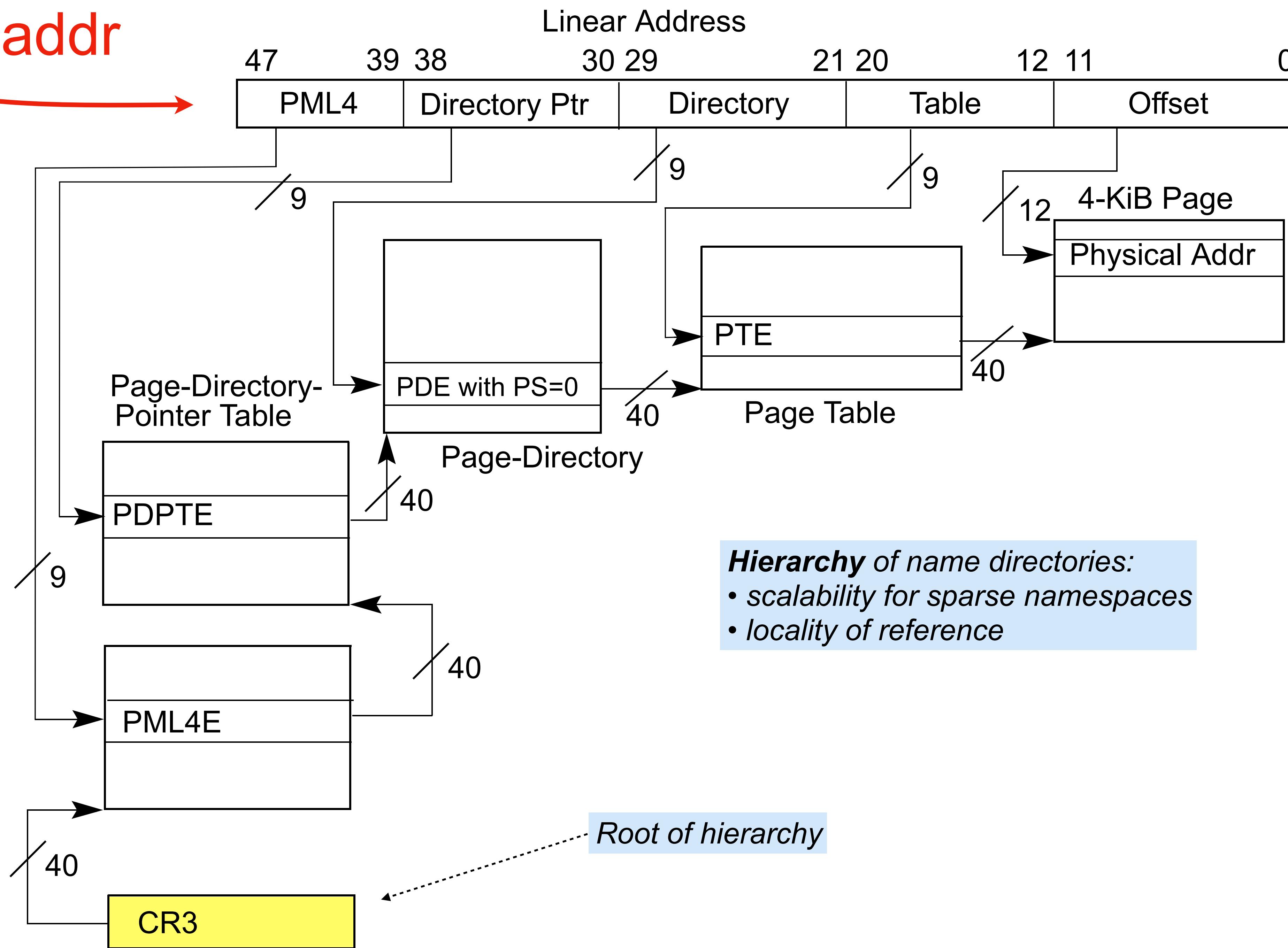
- scalability for sparse namespaces
- locality of reference

HW walks page tables  
OS updates them

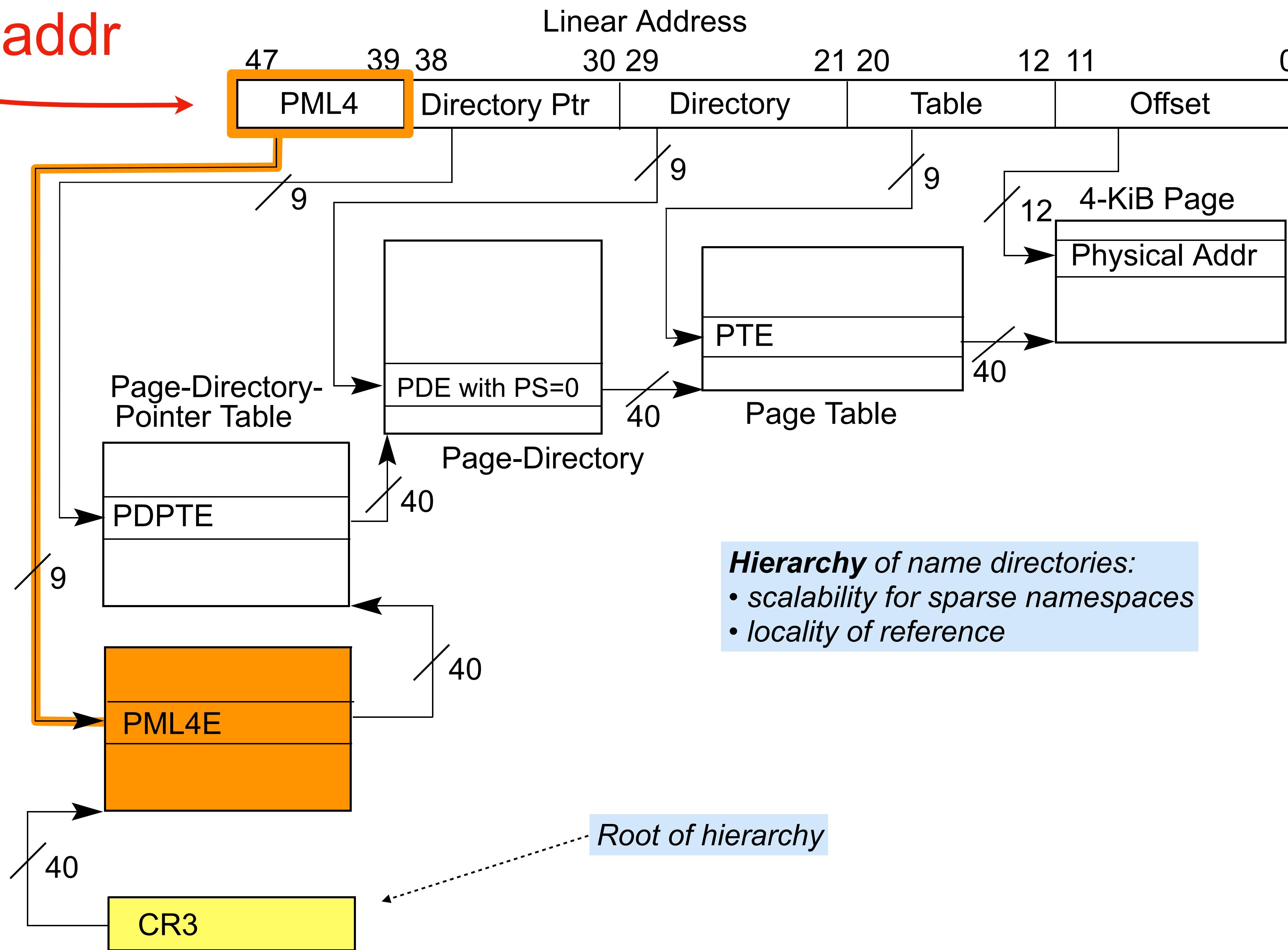
# Virtual addr



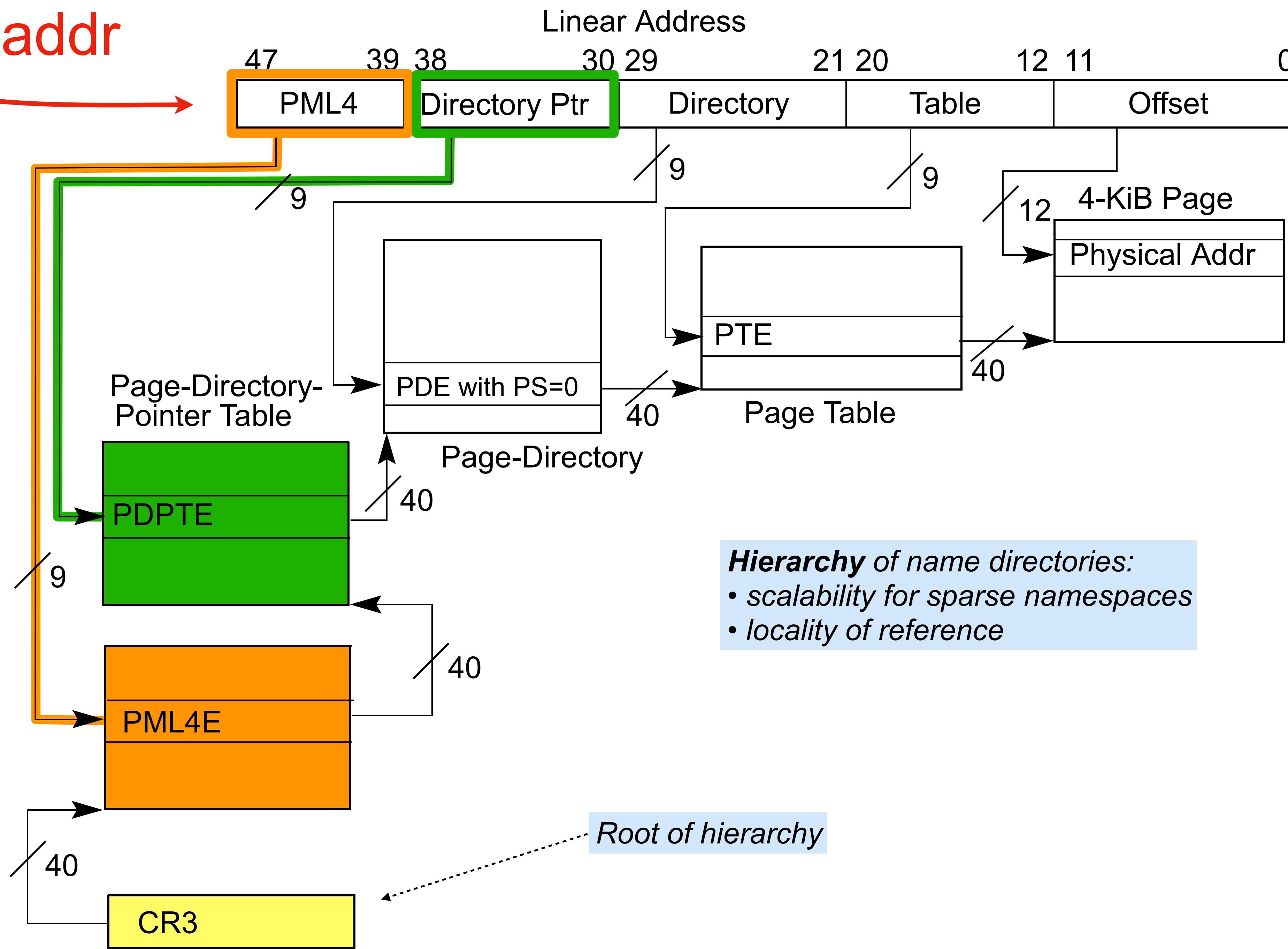
# Virtual addr



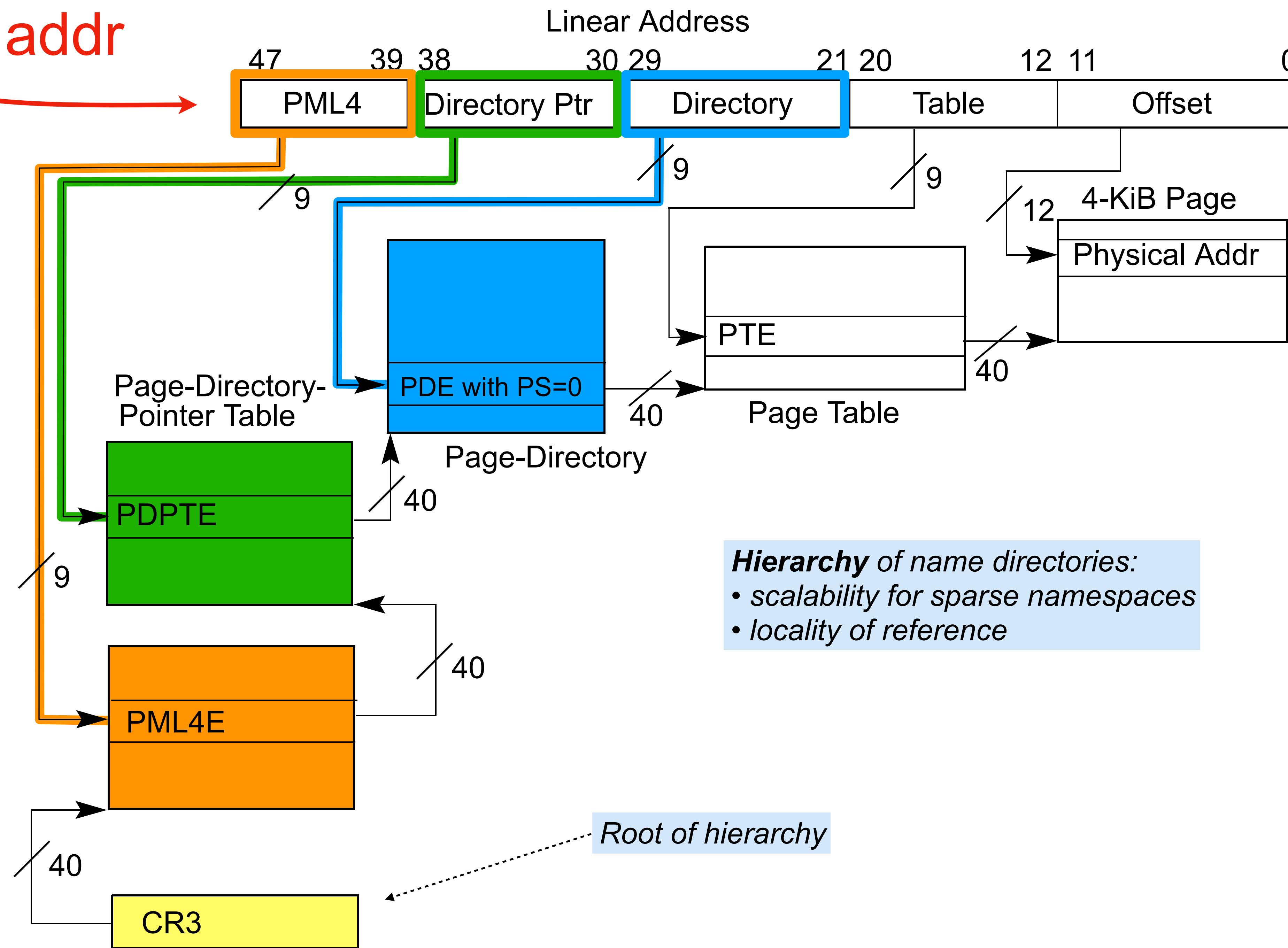
# Virtual addr



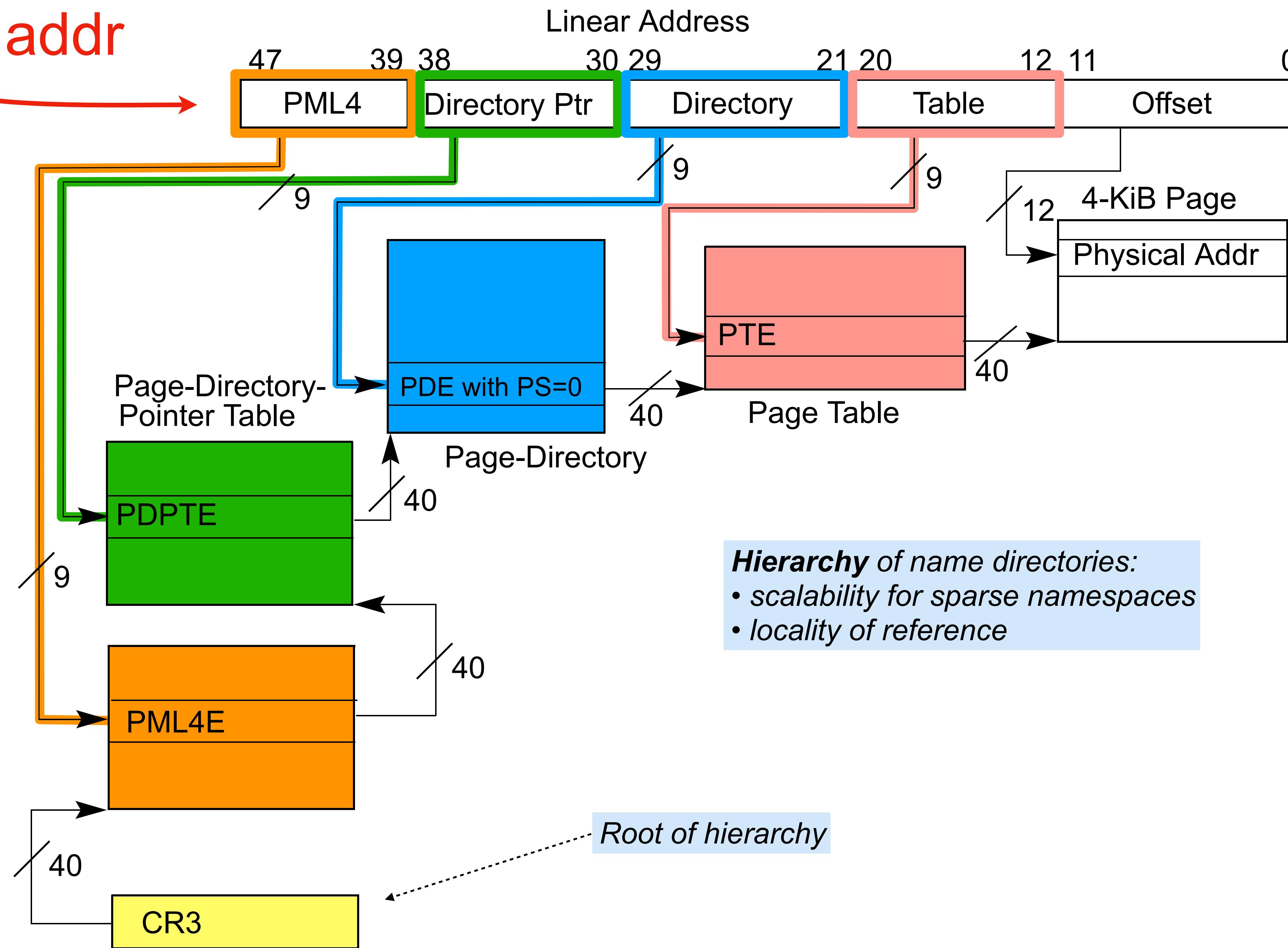
# Virtual addr



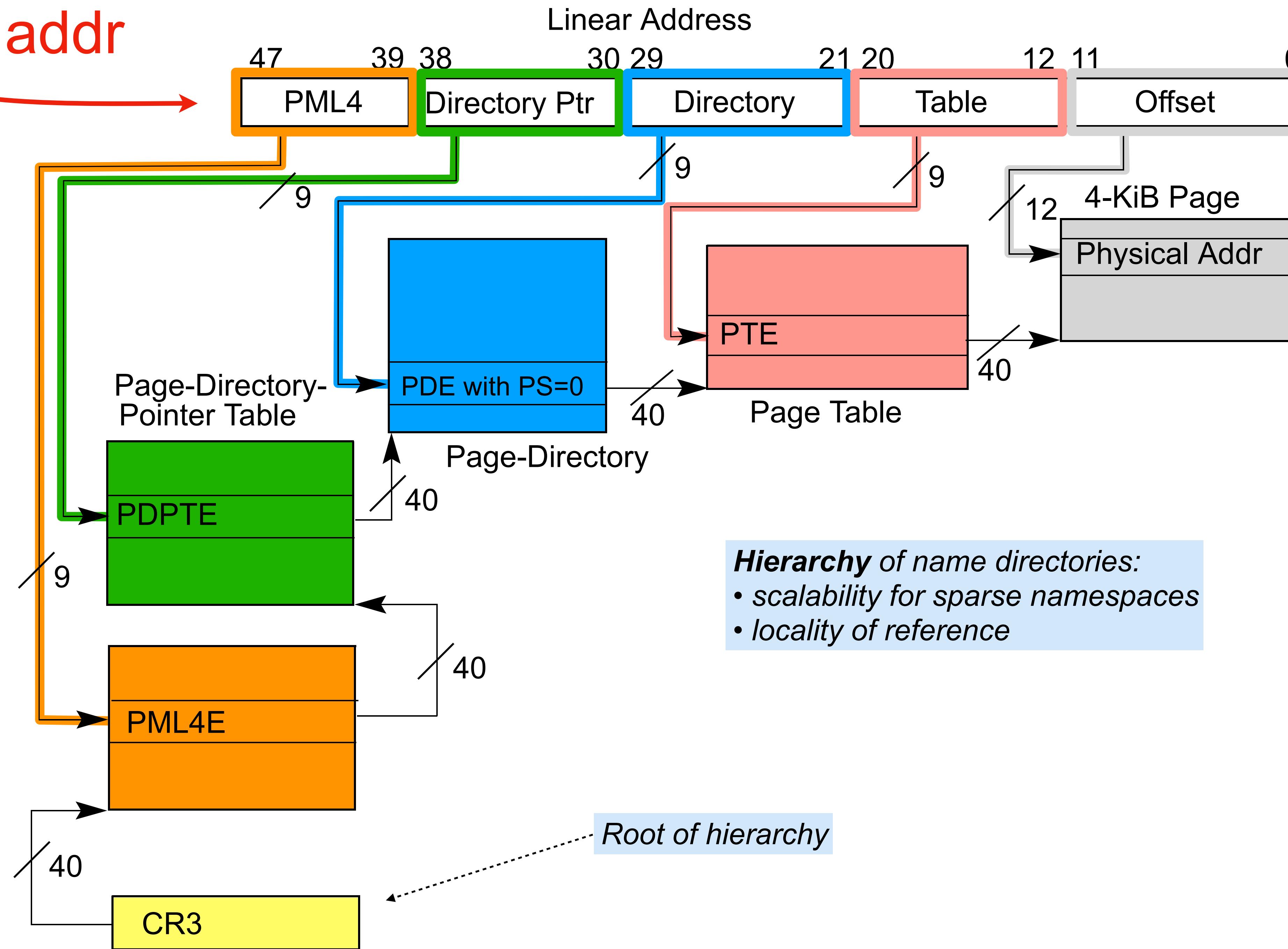
# Virtual addr



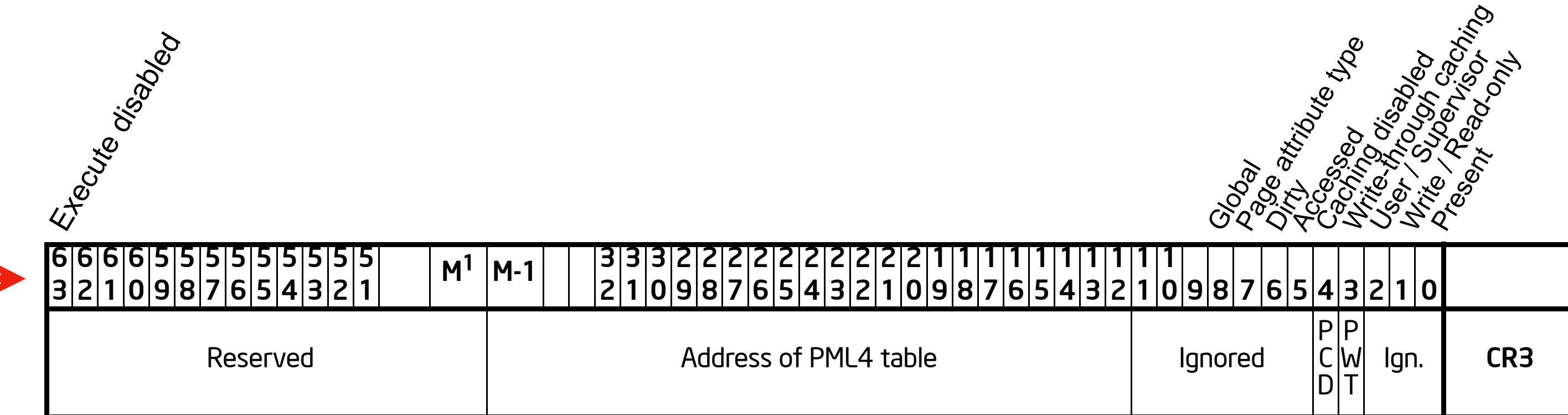
# Virtual addr



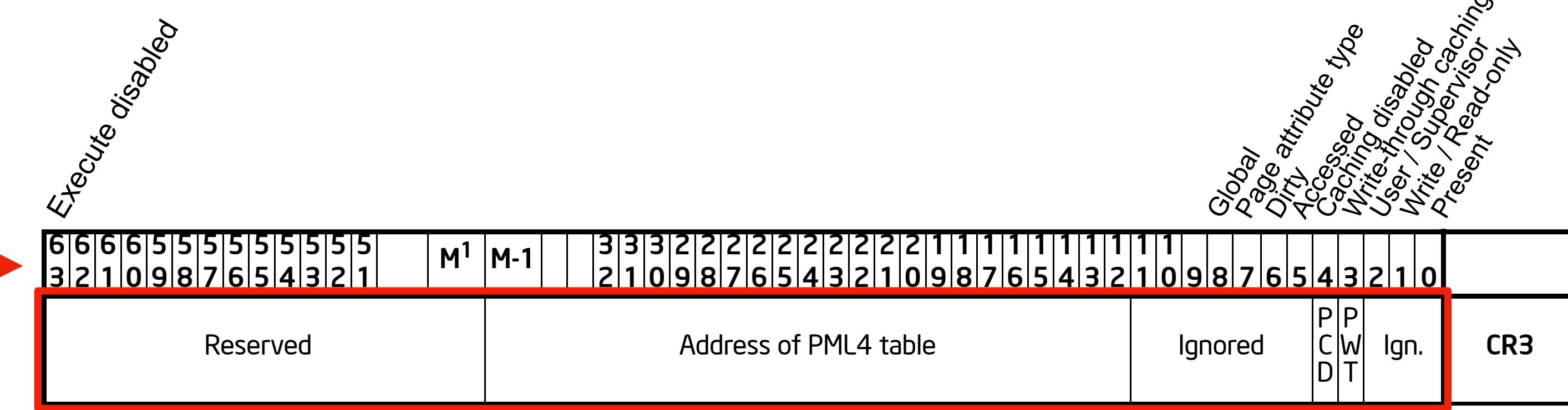
# Virtual addr



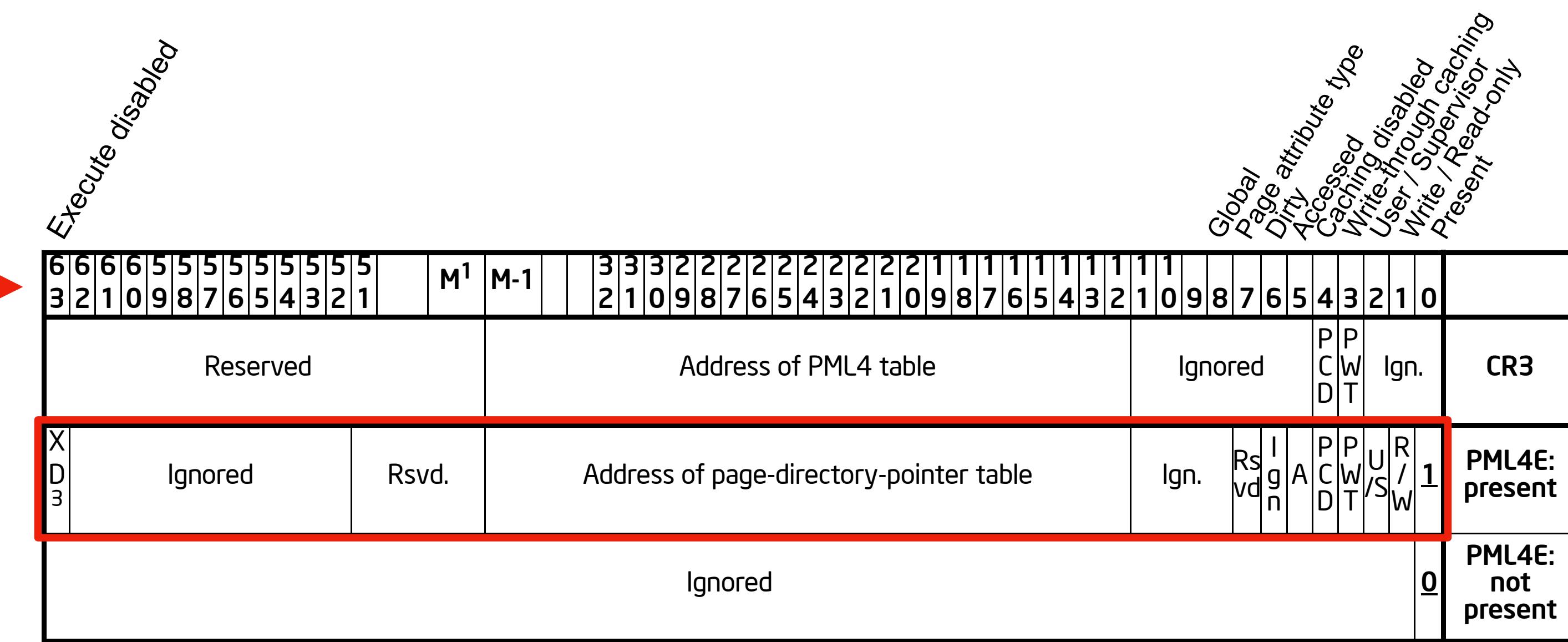
# VA bit index



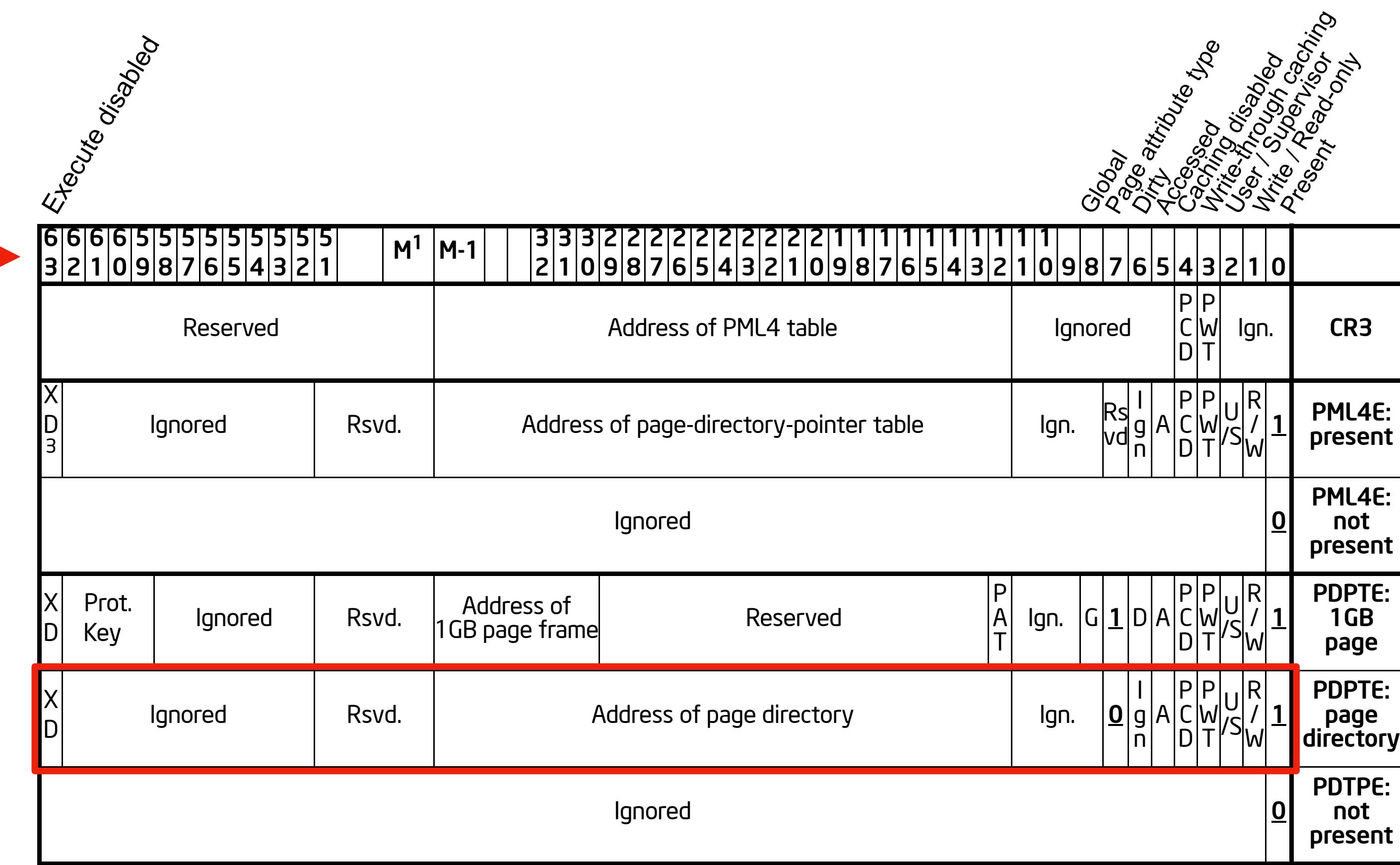
# VA bit index



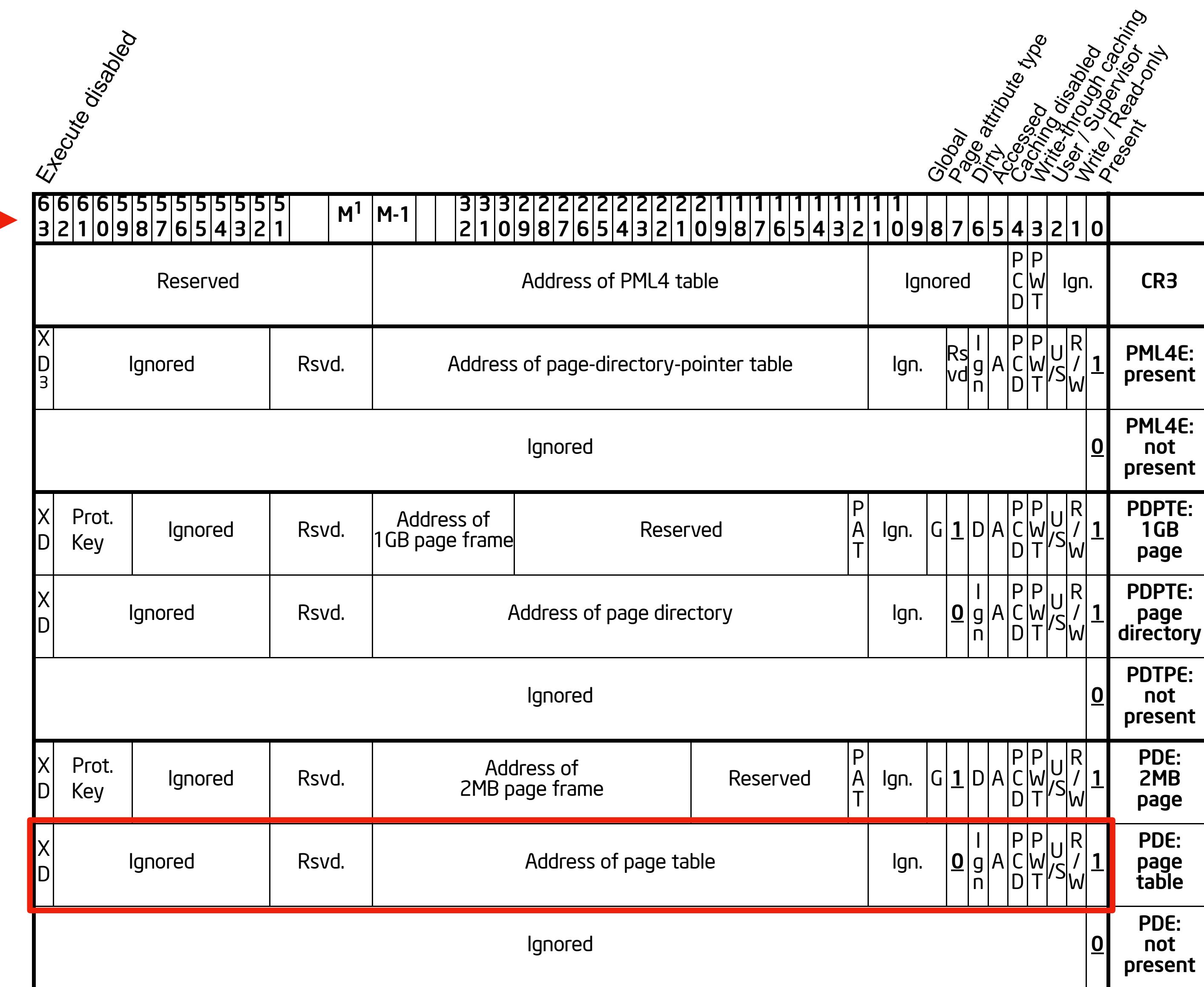
# VA bit index



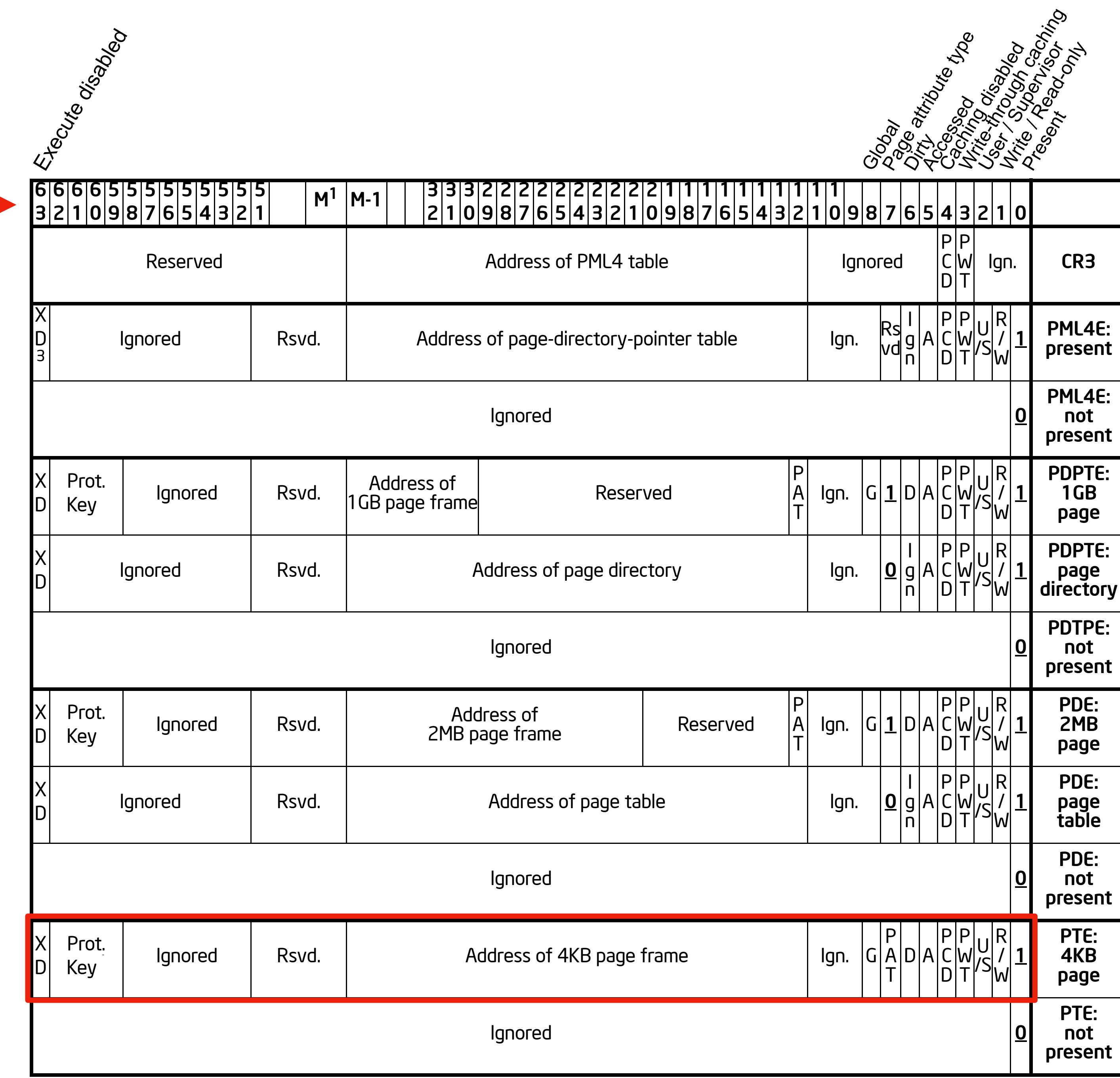
# VA bit index



# VA bit index



# VA bit index



CR3 (root)

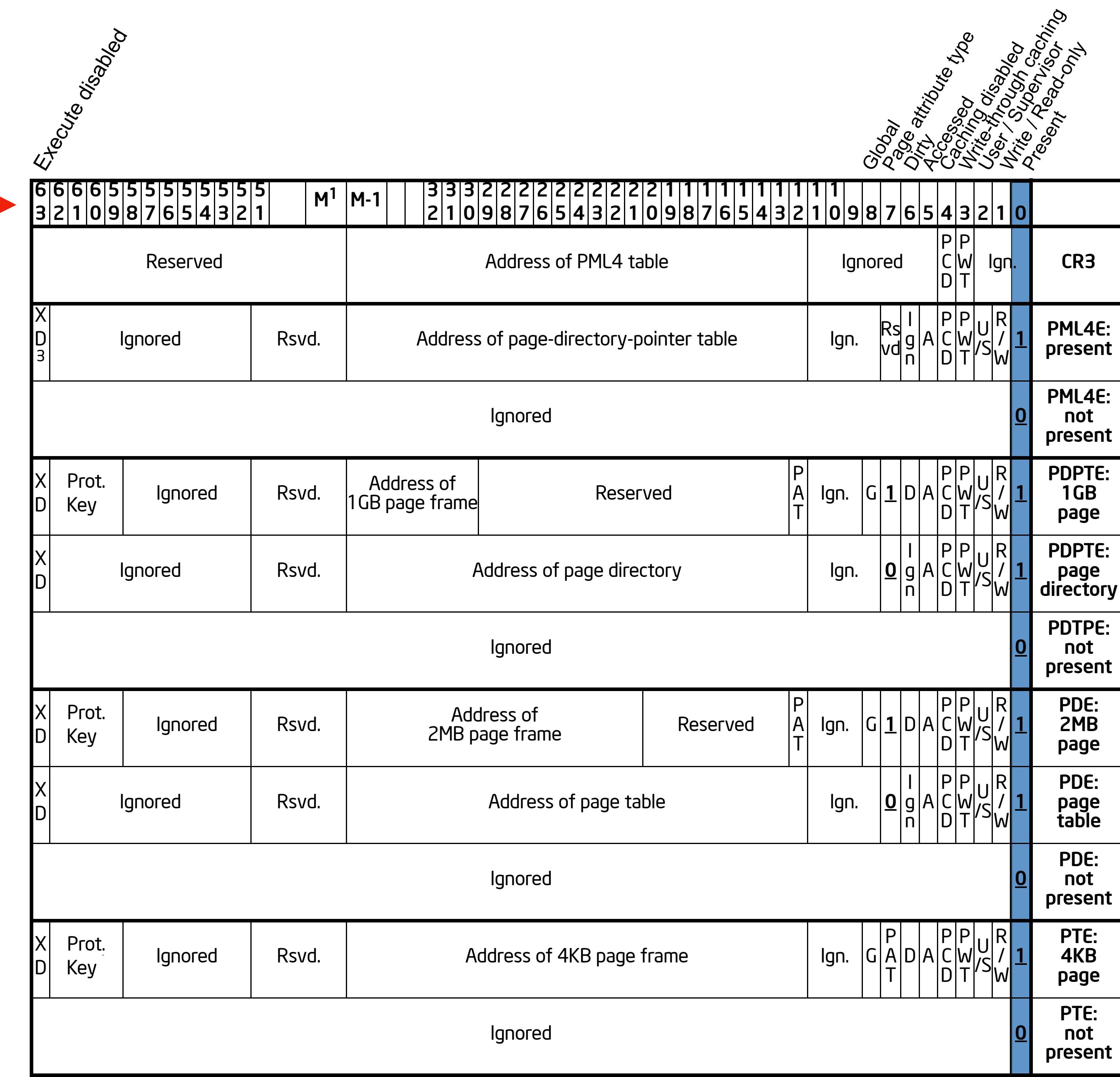
Top-level page table

2nd-level page table

3rd-level page table

Last-level page table

# VA bit index



CR3 (root)

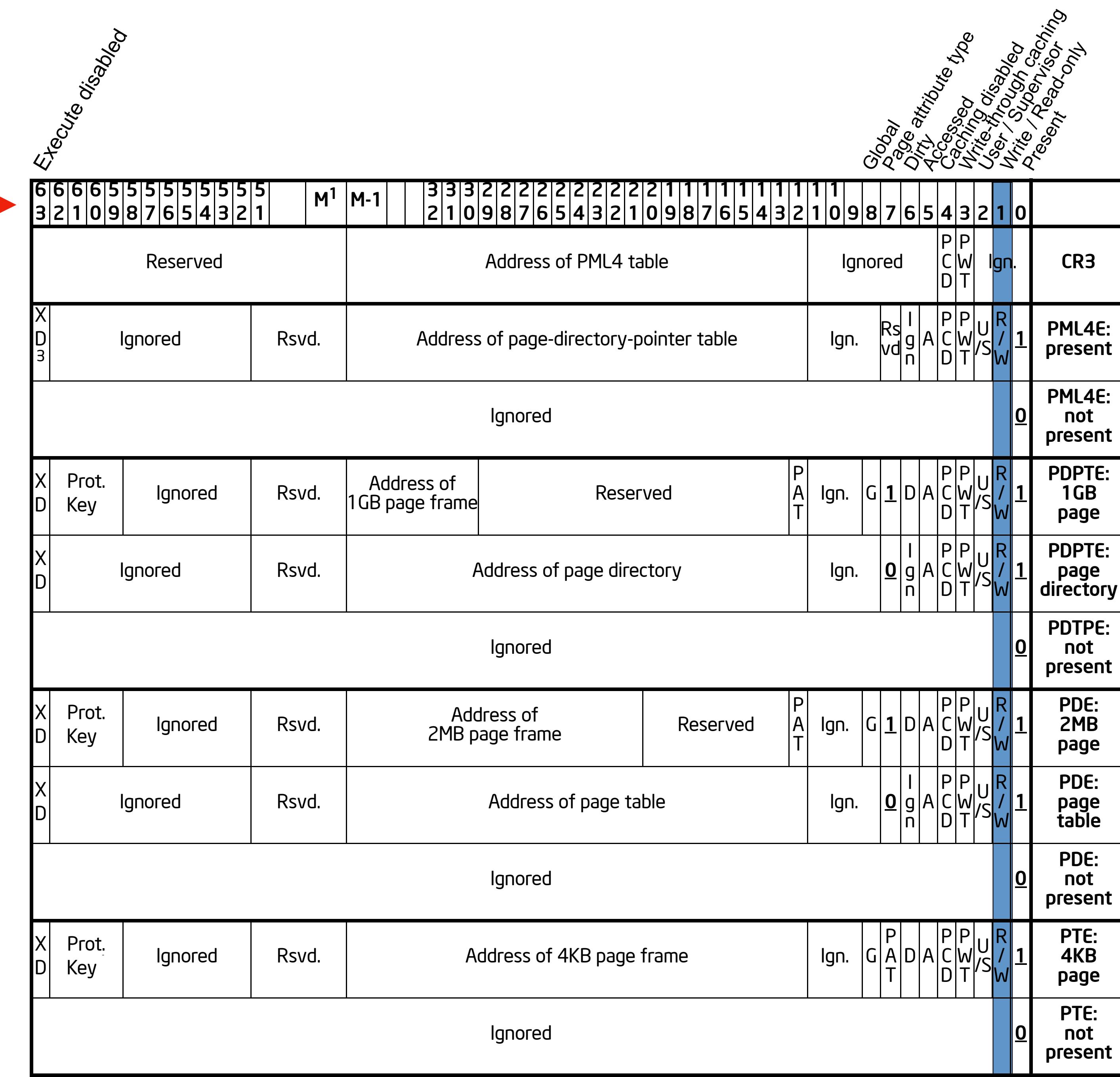
Top-level page table

2nd-level page table

3rd-level page table

Last-level page table

# VA bit index



CR3 (root)

Top-level page table

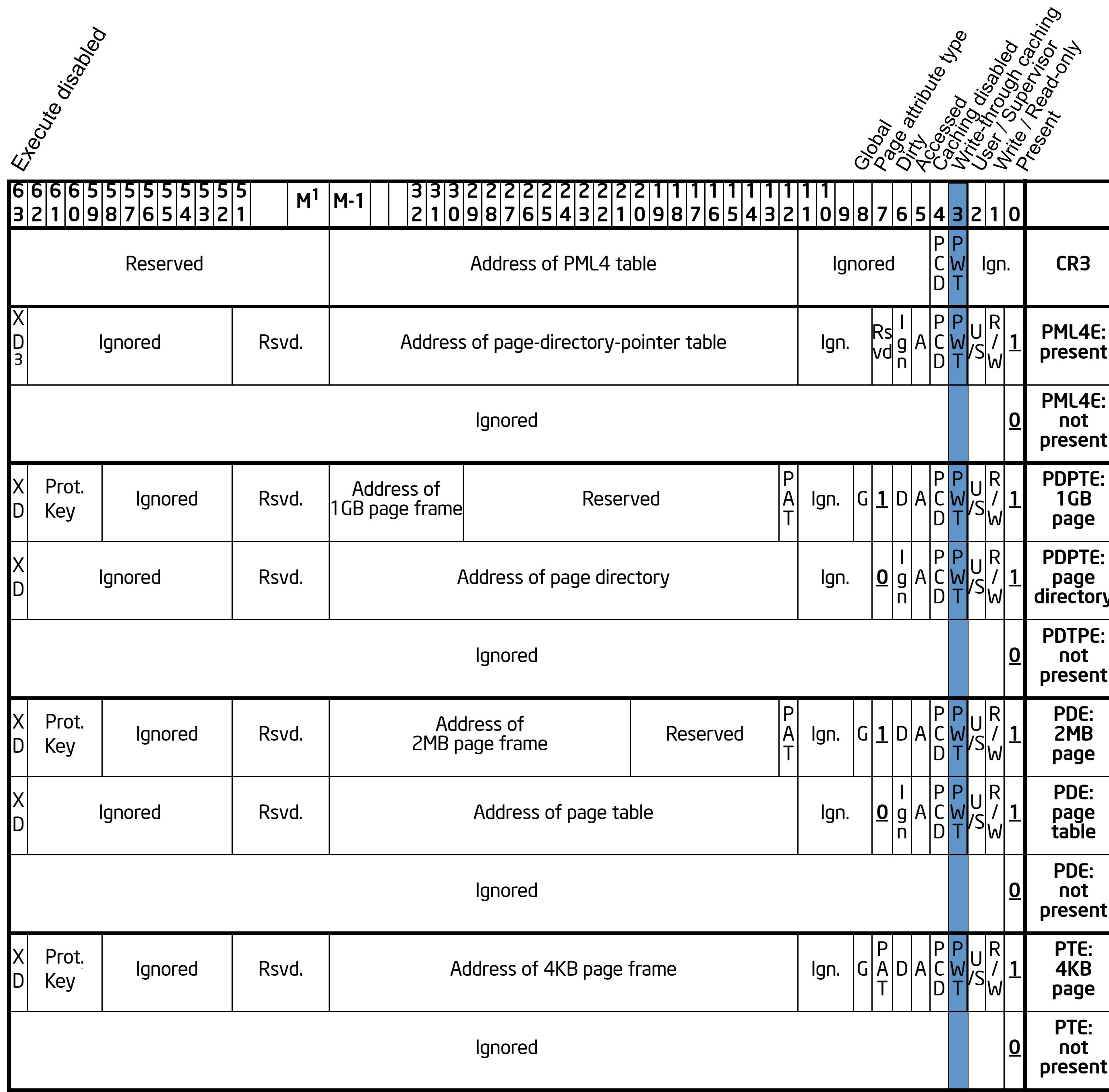
2nd-level page table

3rd-level page table

Last-level page table

# VA bit index

**Execute disabled**



# CR3 (root)

# Top-level page table



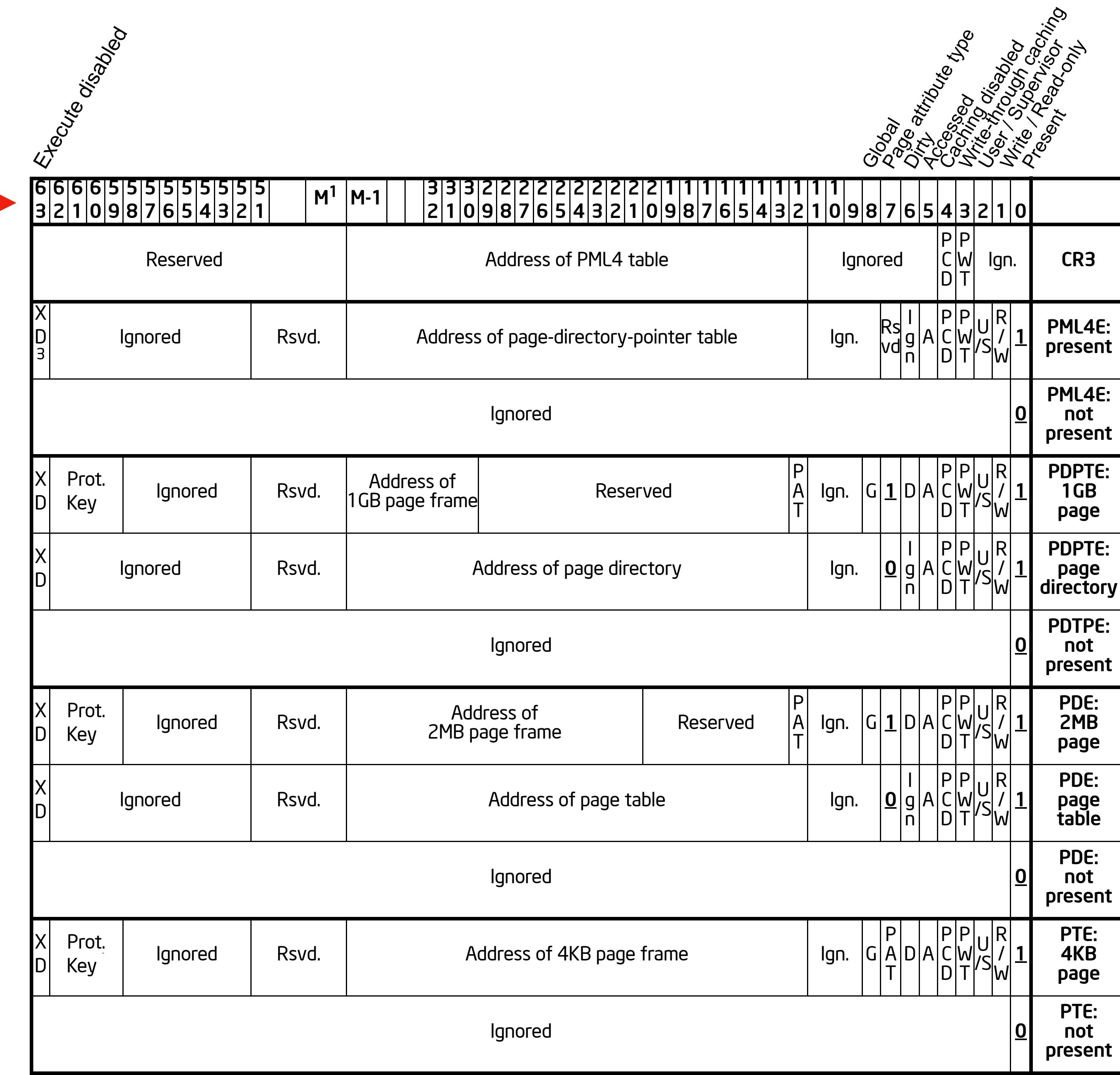
# 2nd-level page table



# 3rd-level page table

# Last-level page table

# VA bit index



CR3 (root)

Top-level page table

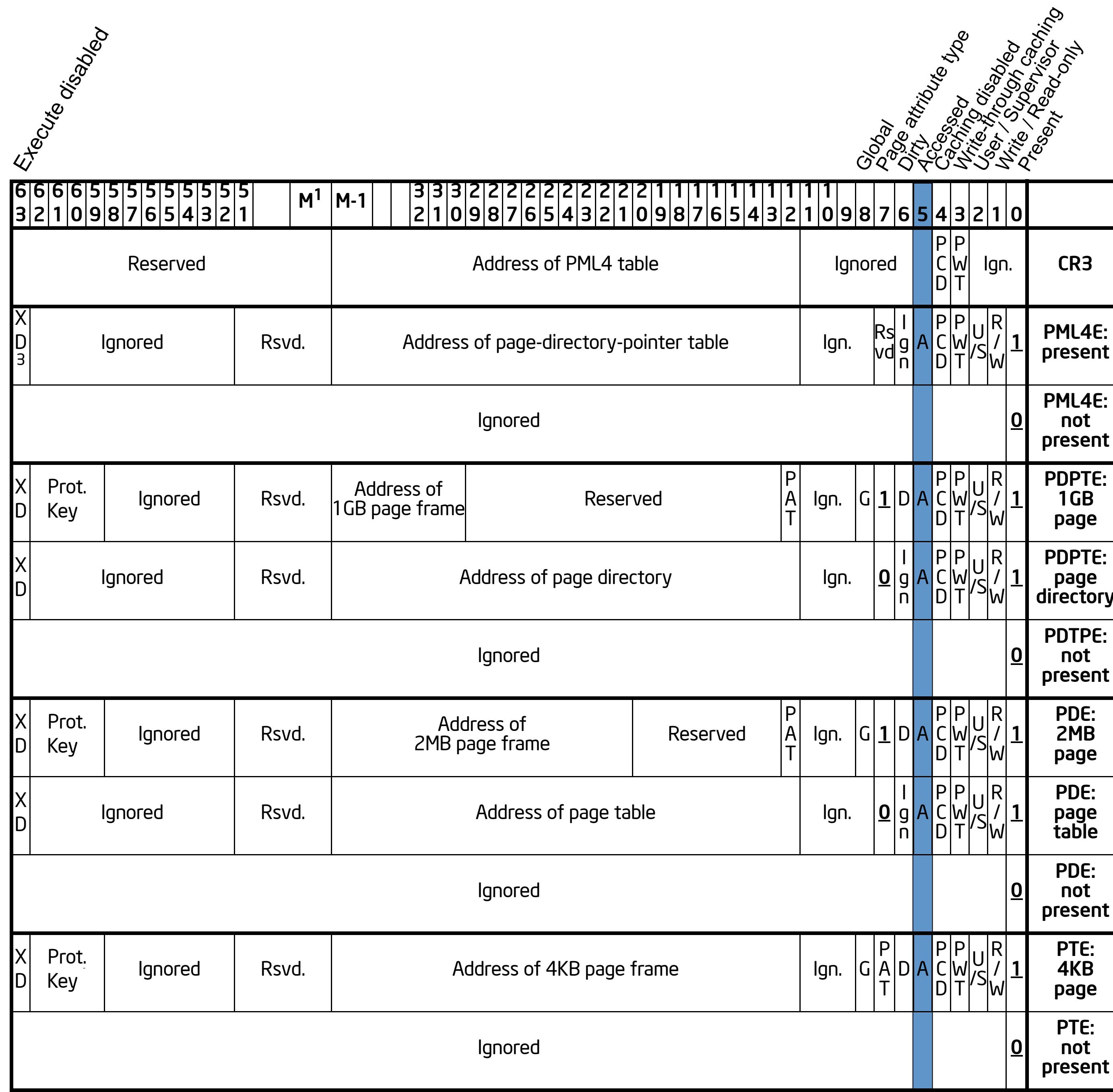
2nd-level page table

3rd-level page table

Last-level page table

# VA bit index

**Execute disabled**



# CR3 (root)

# Top-level page table

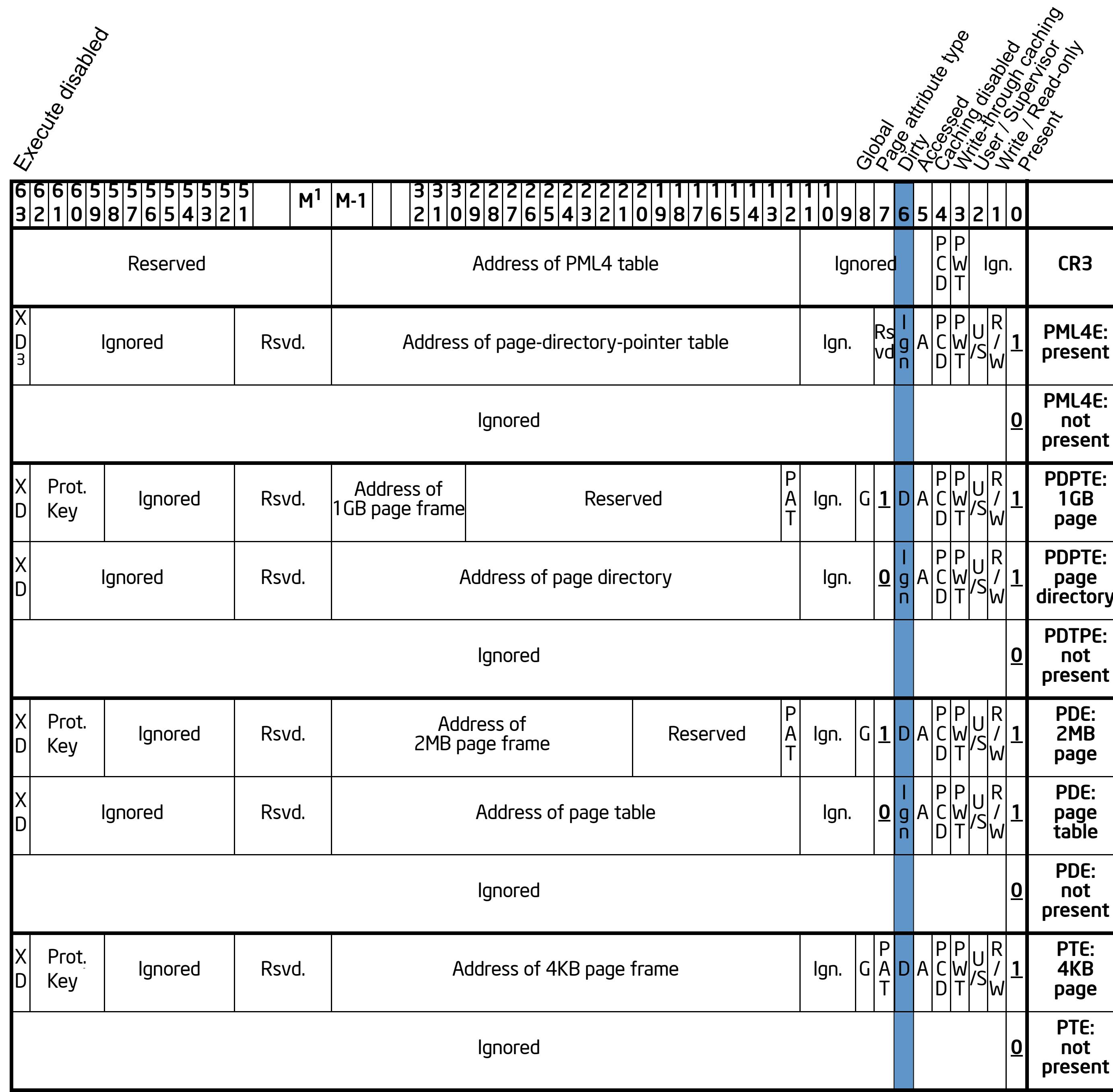
# 2nd-level page table

# 3rd-level page table

# Last-level page table

# VA bit index

**Execute disabled**



# CR3 (root)

# Top-level page table



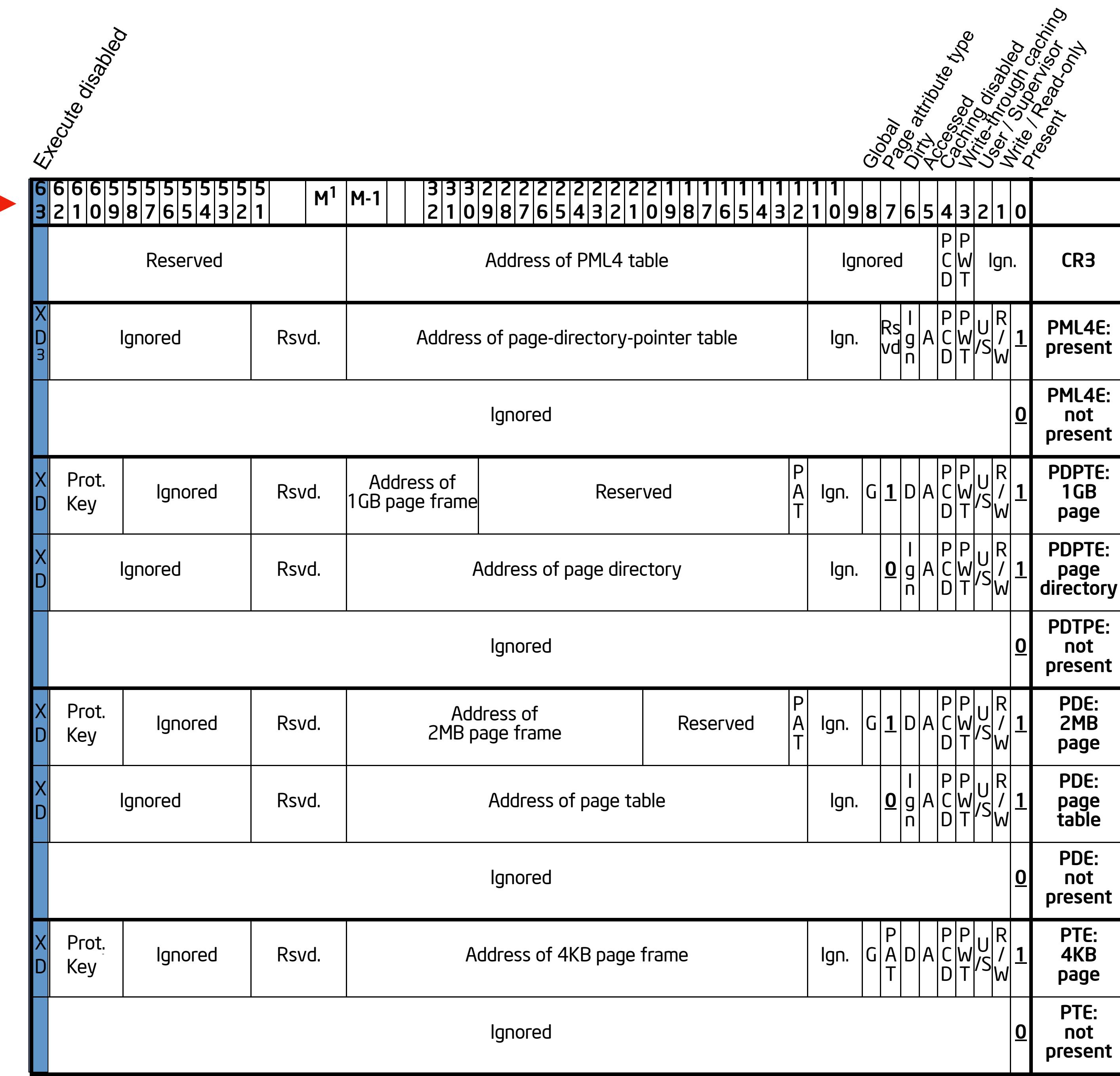
# 2nd-level page table



# 3rd-level page table

# Last-level page table

# VA bit index



# CR3 (root)

# Top-level page table



# 2nd-level page table



# 3rd-level page table

# Last-level page table

Global Page attribute type											
Accessed											
Caching disabled											
Write-through caching											
User / Supervisor Present											
<b>Execute disabled</b>											
6 6 6 6 5 5 5 5 5 5   M <sup>1</sup>   M-1   3 3 3 2 2 2 2 2 2   1 1 1 1 1 1 1 1   1 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											
Reserved <sup>2</sup>			Address of PML4 table				Ignored	P C W D T	Ign.		
X D 3	Ignored	Rsvd.	Address of page-directory-pointer table				Ign.	R s v d I g n A P C W D T U S R W 1	PML4E: present		
Ignored											
Ignored	Rsvd.	Address of 1GB page frame				Reserved	P A T I g n D A P C W D T U S R W 1	PDPTE: 1GB page			
X D	Ignored	Rsvd.	Address of page directory				Ign.	O I g n A P C W D T U S R W 1	PDPTE: page directory		
Ignored											
X D	Prot. Key <sup>4</sup>	Ignored	Rsvd.	Address of 2MB page frame		Reserved	P A T I g n D A P C W D T U S R W 1	PDE: 2MB page			
X D	Ignored	Rsvd.	Address of page table				Ign.	O I g n A P C W D T U S R W 1	PDE: page table		
Ignored											
Ignored	Rsvd.	Address of 4KB page frame				Ign.	G P A D A P C W D T U S R W 1	PTE: 4KB page			
Ignored											

Global Page attribute type											
Accessed Caching disabled Write-through caching User / Supervisor Present											
Dirty											
Write / Read-only											
User / Supervisor Present											
<b>CR3</b>											
Ignored											
P P C W D T											
Ign.											
Rsvd.											
Address of PML4 table											
X D 3 Ignored Rsvd.											
PML4E: present											
Ignored											
0											
PDPTE: 1GB page											
X D Ignored Rsvd.											
Address of page directory											
Ignored											
0											
PDTPE: not present											
Ignored											
0											
PDE: 2MB page											
X D Ignored Rsvd.											
Address of 2MB page frame											
Ignored											
0											
PDE: page table											
X D Ignored Rsvd.											
Address of page table											
Ignored											
0											
PDE: not present											
X D Ignored Rsvd.											
Address of 4KB page frame											
Ignored											
0											
PTE: 4KB page											
X D Ignored Rsvd.											
Address of 4KB page frame											
Ignored											
0											
PTE: not present											

Execute disabled

Ignored

Rsvd.

Ignored

Rsvd.

Ignored

Rsvd.

Ignored

Rsvd.

Ignored

Rsvd.

Address of PML4 table

Address of page-directory-pointer table

Address of page directory

Address of page table

Address of 4KB page frame

Address of 2MB page frame

Address of 1GB page frame

Address of M

Address of 32-bit memory location

Address of 64-bit memory location

Address of 4KB page frame

Address of 2MB page frame

Address of 1GB page frame

Address of M

Address of 32-bit memory location

Address of 64-bit memory location

Address of 4KB page frame

Address of 2MB page frame

Address of 1GB page frame

Address of M

Address of 32-bit memory location

Address of 64-bit memory location

Address of 4KB page frame

Address of 2MB page frame

Address of 1GB page frame

Address of M

Address of 32-bit memory location

Address of 64-bit memory location

Address of 4KB page frame

Address of 2MB page frame

Address of 1GB page frame

Address of M

Address of 32-bit memory location

Address of 64-bit memory location

Address of 4KB page frame

Address of 2MB page frame

Address of 1GB page frame

Address of M

Address of 32-bit memory location

Address of 64-bit memory location

Address of 4KB page frame

Address of 2MB page frame

Address of 1GB page frame

Address of M

Address of 32-bit memory location

Address of 64-bit memory location

Address of 4KB page frame

Address of 2MB page frame

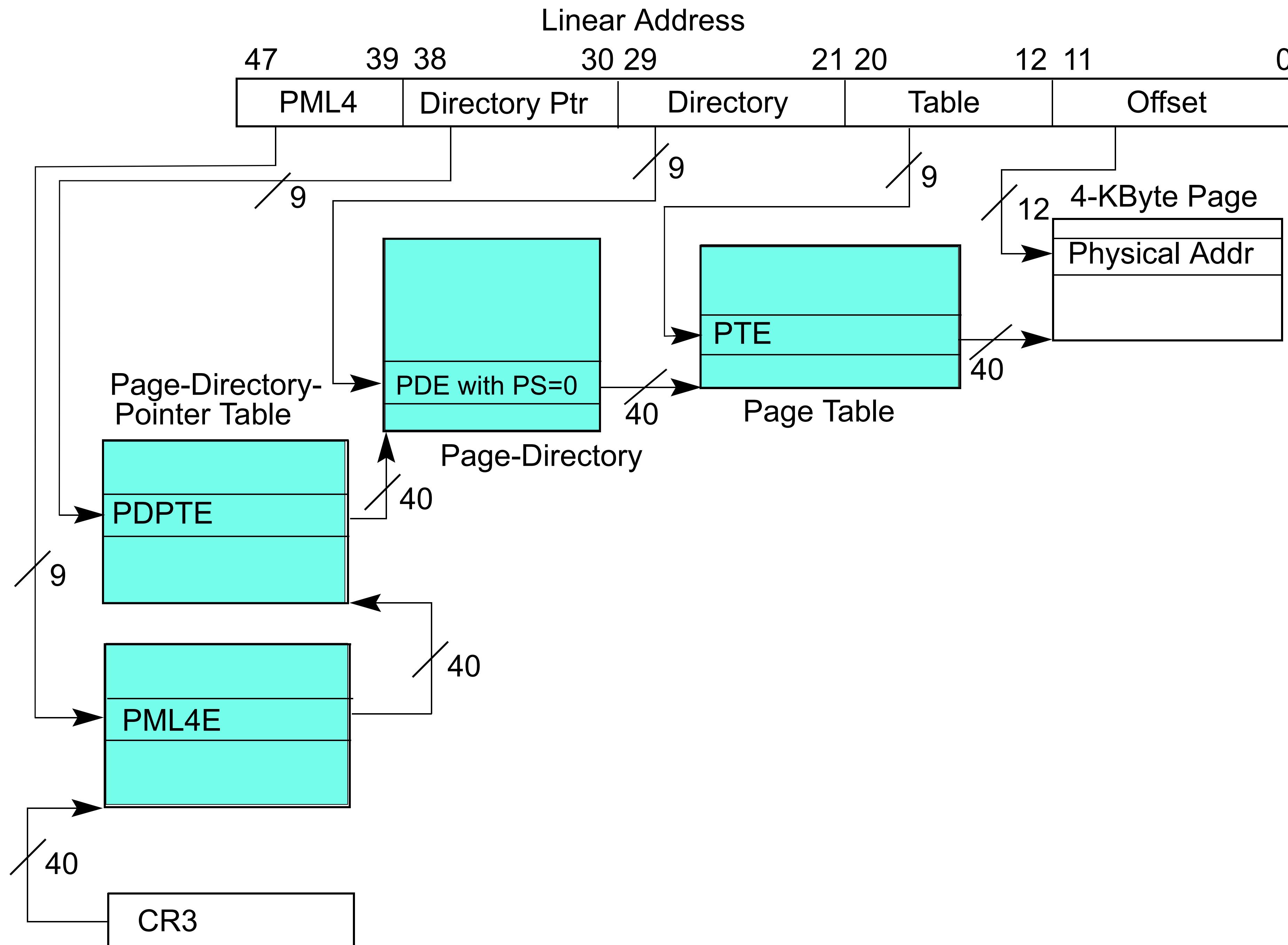
Address of 1GB page frame

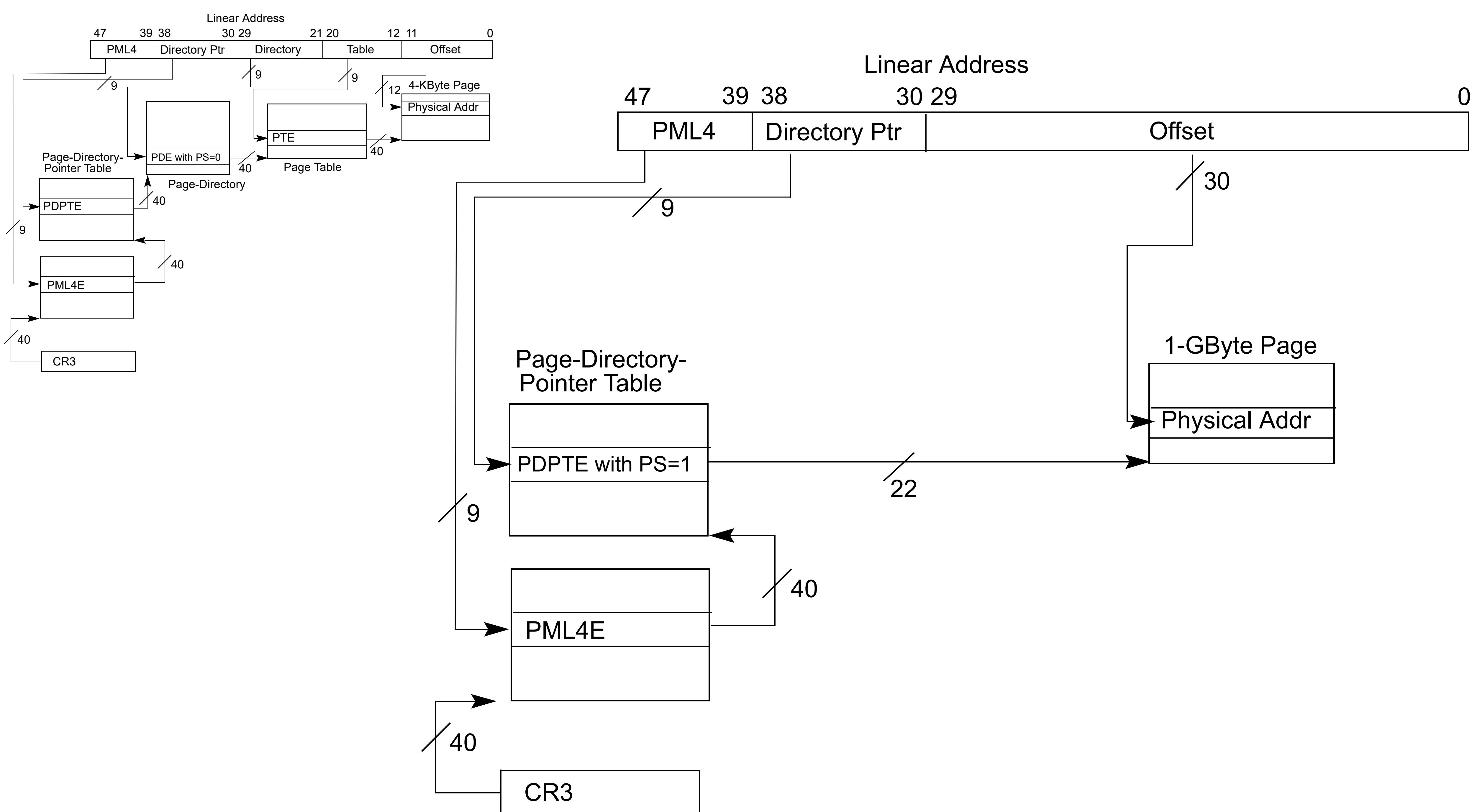
Address of M

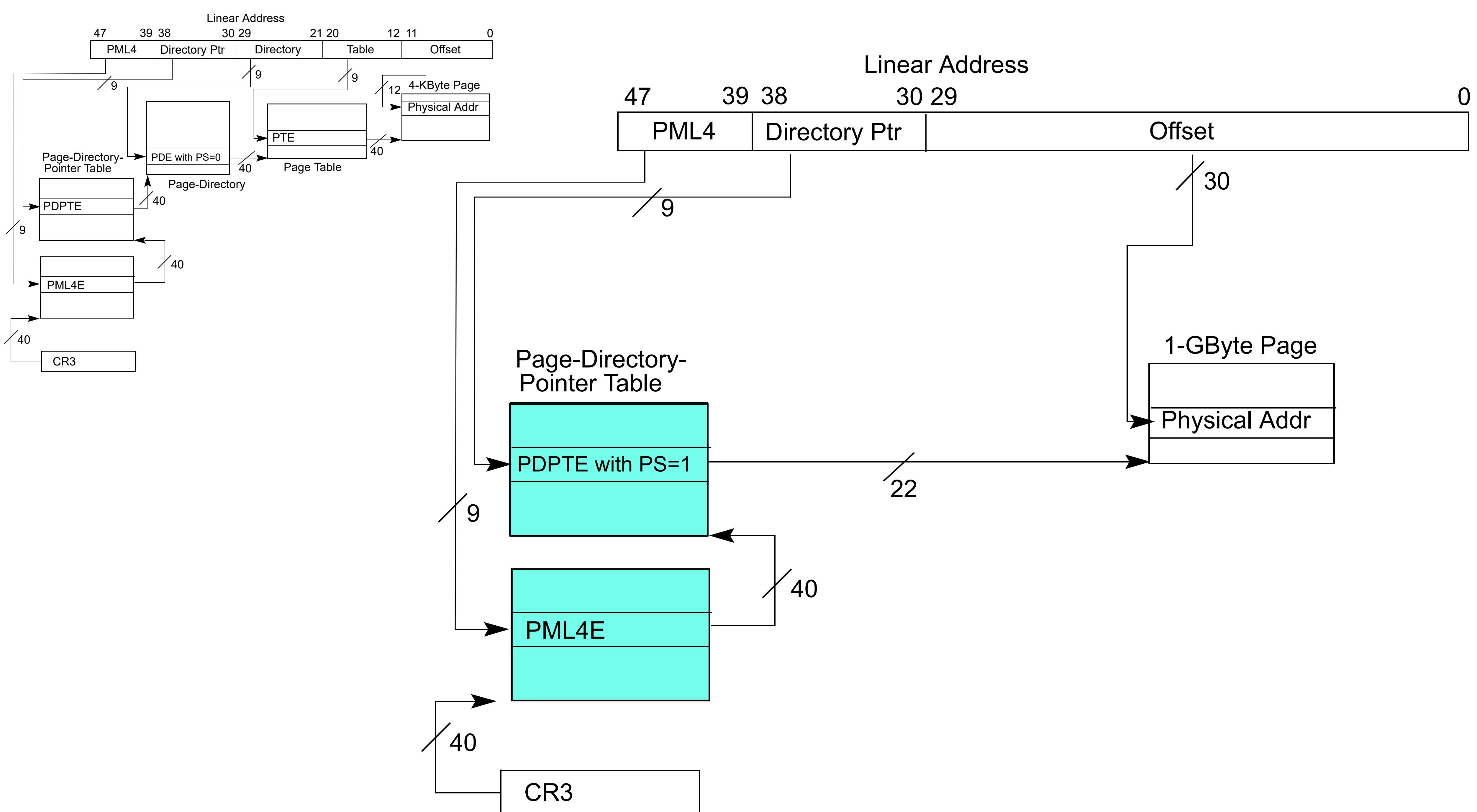
Address of 32-bit memory location

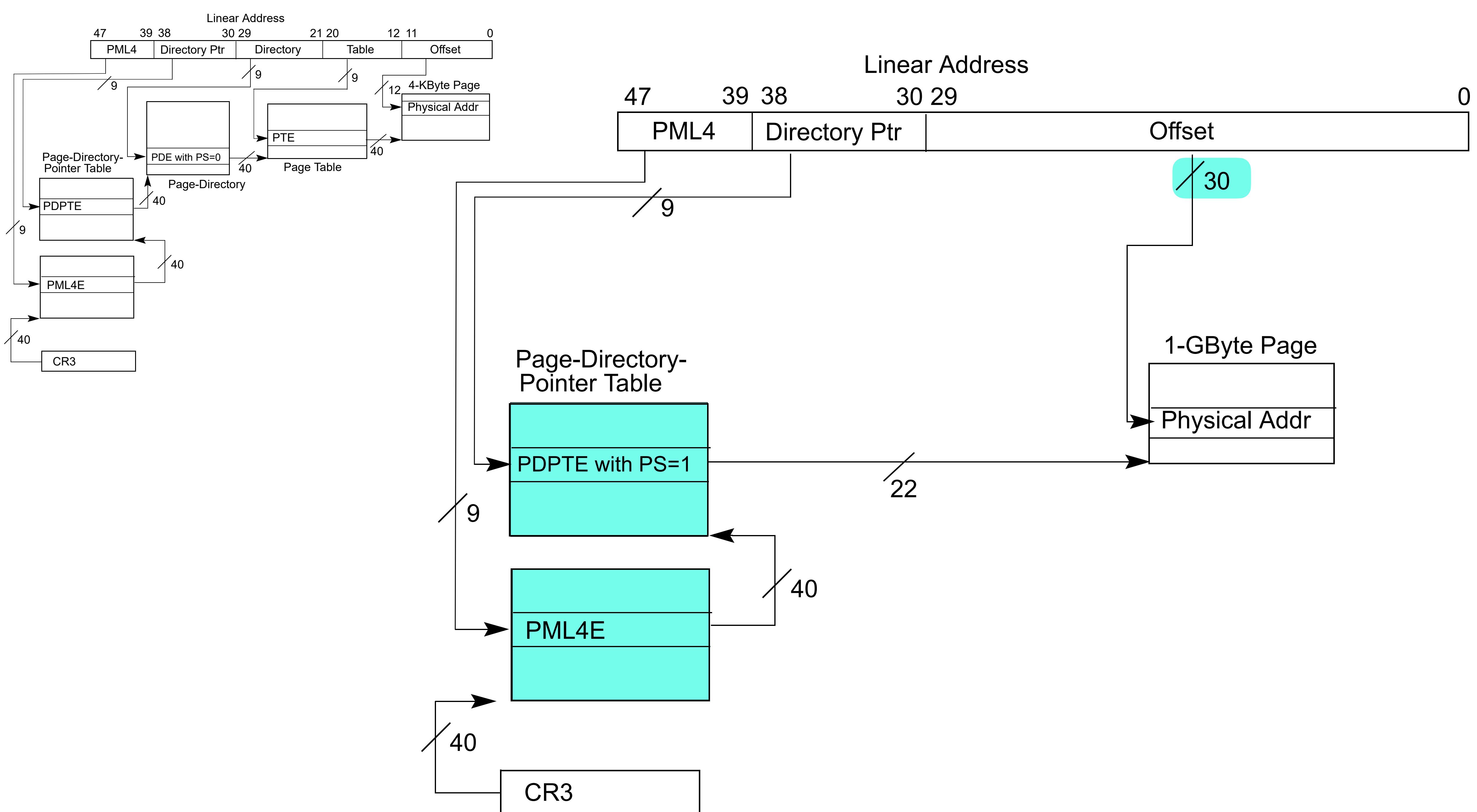
Global Page attribute type											
Accessed											
Caching disabled											
Write-through caching											
User / Supervisor Present											
<b>Execute disabled</b>											
6 6 6 6 5 5 5 5 5 5   M <sup>1</sup>   M-1   3 3 3 2 2 2 2 2 2   1 1 1 1 1 1 1 1   1 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0											
Reserved <sup>2</sup>			Address of PML4 table				Ignored	P C W D T	Ign.		
X D 3	Ignored	Rsvd.	Address of page-directory-pointer table				Ign.	R s v d I g n A P C W D T U S R W 1	PML4E: present		
Ignored									0		
X D	Prot. Key <sup>4</sup>	Ignored	Rsvd.	Address of 1GB page frame	Reserved		P A T	Ign. G 1 D A P C W D T U S R W 1	PDPT <sup>E</sup> : 1GB page		
X D	Ignored	Rsvd.	Address of page directory				Ign.	O I g n A P C W D T U S R W 1	PDPT <sup>E</sup> : page directory		
Ignored									0		
X D	Prot. Key <sup>4</sup>	Ignored	Rsvd.	Address of 2MB page frame		Reserved	P A T	Ign. G 1 D A P C W D T U S R W 1	PDE: 2MB page		
X D	Ignored	Rsvd.	Address of page table				Ign.	O I g n A P C W D T U S R W 1	PDE: page table		
Ignored									0		
X D	Prot. Key <sup>4</sup>	Ignored	Rsvd.	Address of 4KB page frame			Ign. G P A D A P C W D T U S R W 1	PTE: 4KB page			
Ignored									0		

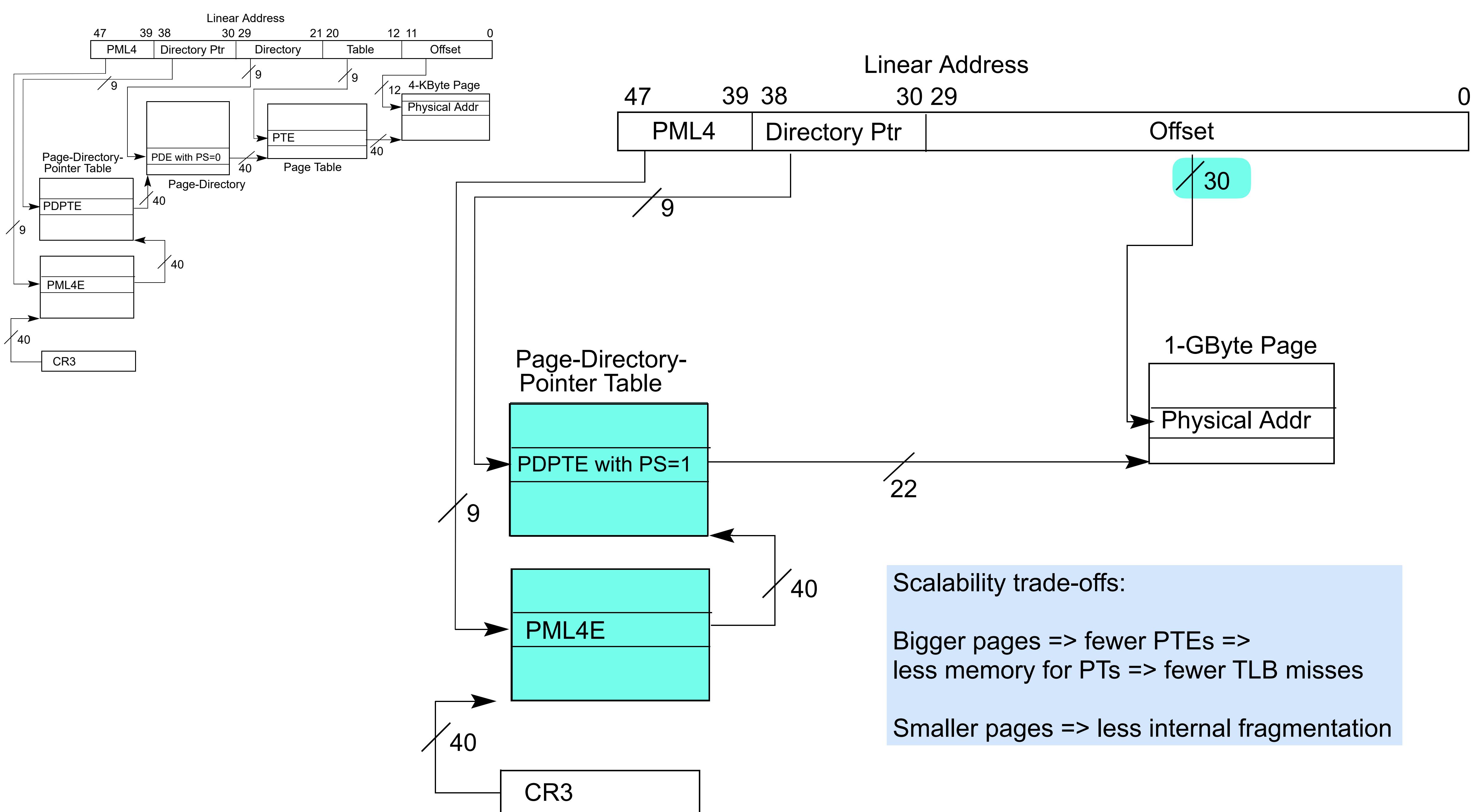
Execute disabled										Global Page attribute type																
										Dirty Accessed Caching disabled Write-through caching User / Supervisor Present Write / Read-only																
6 6 6 6 5 5 5 5 5 5 5 5 5   M <sup>1</sup> M-1   3 3 3 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1   1 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0																										
Reserved <sup>2</sup>										Address of PML4 table																
X D <sub>3</sub>	Ignored	Rsvd.	Address of page-directory-pointer table										Ign.	Rsvd	Ign.	A	P C W D T	U / S W	1	PML4E: present						
Ignored																				PML4E: not present						
X D	Prot. Key <sup>4</sup>	Ignored	Rsvd.	Address of 1GB page frame	Reserved					P A T	Ign.	G	1	D	A	P C W D T	U / S W	1	PDPTPE: 1GB page							
X D	Ignored	Rsvd.	Address of page directory										Ign.	0	I g n	A	P C W D T	U / S W	1	PDPTPE: page directory						
Ignored																				PDTPE: not present						
X D	Prot. Key <sup>4</sup>	Ignored	Rsvd.	Address of 2MB page frame	Reserved					P A T	Ign.	G	1	D	A	P C W D T	U / S W	1	PDE: 2MB page							
X D	Ignored	Rsvd.	Address of page table										Ign.	0	I g n	A	P C W D T	U / S W	1	PDE: page table						
Ignored																				PDE: not present						
X D	Prot. Key <sup>4</sup>	Ignored	Rsvd.	Address of 4KB page frame										Ign.	G	P A T	D A	P C W D T	U / S W	1	PTE: 4KB page					
Ignored																				PTE: not present						













Execute disabled												Global Page attribute type																	
												Accessed Caching disabled Write-through caching User / Supervisor Present																	
												Write / Read-only																	
6 6 6 6 5 5 5 5 5 5 5 5 3 2 1 0 9 8 7 6 5 4 3 2 1												M <sup>1</sup>	M-1	3 3 3 2 2 2 2 2 2 2 2 2 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	1 1	1 1	1 1	P P C W D T	P P C W D T	I g n.	Ign.	CR3							
X D <sub>3</sub>	Reserved <sup>2</sup>		Address of PML4 table												Ignored		P C D P C D U / S	P C D P C D U / S	R / W R / W	1	PML4E: present								
X D	Ignored		Rsvd.		Address of page-directory-pointer table												Ign.		R s v d I g n A	P C D P C D U / S	R / W R / W	0	PML4E: not present						
X D	Prot. Key <sup>4</sup>	Ignored		Rsvd.		Address of 1GB page frame		Reserved		P A T	Ign.		G 1	D A	P C D P C D U / S	P C D P C D U / S	R / W R / W	1	PDPTPE: 1GB page										
X D	Ignored		Rsvd.		Address of page directory												Ign.		O g n I g n A	P C D P C D U / S	R / W R / W	1	PDPTPE: page directory						
Ignored																								0	PDPTPE: not present				
X D	Prot. Key <sup>4</sup>	Ignored		Rsvd.		Address of 2MB page frame		Reserved		P A T	Ign.		G 1	D A	P C D P C D U / S	P C D P C D U / S	R / W R / W	1	PDE: 2MB page										
X D	Ignored		Rsvd.		Address of page table												Ign.		O g n I g n A	P C D P C D U / S	R / W R / W	1	PDE: page table						
Ignored																								0	PDE: not present				
X D	Prot. Key <sup>4</sup>	Ignored		Rsvd.		Address of 4KB page frame		Ignored		Ign.	G P A T	D A	P C D P C D U / S	P C D P C D U / S	R / W R / W	1	PTE: 4KB page												
Ignored																								0	PTE: not present				







Global attribute type										Dirty	Accessed	Caching disabled	Write-through caching	User / Supervisor	Write / Read-only	Present	
Execute disabled																	
6 3	6 2	6 1	6 0	5 9	5 8	5 7	5 6	5 5	5 4	5 3	5 2	5 1	1 0	1 0	1 0		
M <sup>1</sup>	M-1			3 2	3 1	3 0	3 9	3 8	2 7	2 6	2 5	2 4	2 3	2 2	2 1		
Reserved <sup>2</sup>		Address of PML4 table				Ignored		P C W D T	I g n o r e	I g n. v d n	A P C D T	P C W D T	U S R W /	R W /	CR3		
X D 3	Ignored	Rsvd.	Address of page-directory-pointer table				Ign.								PML4E: present		
Ignored										0	PML4E: not present						
X D	Prot. Key <sup>4</sup>	Ignored	Rsvd.	Address of 1GB page frame	Reserved		P A T	Ign.	G 1	D A	P C D T	P C W D T	U S R W /	R W /	PDPTE: 1GB page		
X D	Ignored	Rsvd.	Address of page directory				Ign.	0	I g n o r e	A P C D T	P C W D T	U S R W /	R W /	1	PDPTE: page directory		
Ignored										0	PDTPE: not present						
X D	Prot. Key <sup>4</sup>	Ignored	Rsvd.	Address of 2MB page frame		Reserved	P A T	Ign.	G 1	D A	P C D T	P C W D T	U S R W /	R W /	1	PDE: 2MB page	
X D	Ignored	Rsvd.	Address of page table				Ign.	0	I g n o r e	A P C D T	P C W D T	U S R W /	R W /	1	PDE: page table		
Ignored										0	PDE: not present						
X D	Prot. Key <sup>4</sup>	Ignored	Rsvd.	Address of 4KB page frame		Ign.	G 1	I D A 0	P C D T	P C W D T	U S R W /	R W /	1	PTE: 4KB page			
Ignored										0	PTE: not present						

Execute disabled										Global attribute type														
										Dirty Accessed Caching disabled User - Supervisor Write / Read-only Present														
66665555555555 3210987654321										M <sup>1</sup> M-1 3332222222211111111111 2109876543210987654321098765432109876543210														
Reserved <sup>2</sup>				Address of PML4 table						Ignored				P C W D	Ign.		CR3							
X D 3	Ignored		Rsvd.	Address of page-directory-pointer table						Ign.	R s v d n	I g n	A	P C D T	P C W T U S	R / W	1	PML4E: present						
Ignored										0														
X D	Prot. Key <sup>4</sup>	Ignored		Rsvd.	Address of 1GB page frame		Reserved				P A T	Ign.	G	1	D A	P C D T	P C W T U S	R / W	1	PDPTE: 1GB page				
X D	Ignored		Rsvd.	Address of page directory						Ign.	0	I g n	A	P C D T	P C W T U S	R / W	1	PDPTE: page directory						
Ignored										0														
X D	Prot. Key <sup>4</sup>	Ignored		Rsvd.	Address of 2MB page frame			Reserved		P A T	Ign.	G	1	D A	P C D T	P C W T U S	R / W	1	PDE: 2MB page					
X D	Ignored		Rsvd.	Address of page table						Ign.	0	I g n	A	P C D T	P C W T U S	R / W	1	PDE: page table						
Ignored										0														
X D	Prot. Key <sup>4</sup>	Ignored		Rsvd.	Address of 4KB page frame						Ign.	G	1	D A	P C D T	P C W T U S	R / W	1	PTE: 4KB page					
Ignored										0														

4 x 50 ns table lookup +  
1 x 50 ns access =  
250 ns =  
~750 instructions @ 3 GHz clock

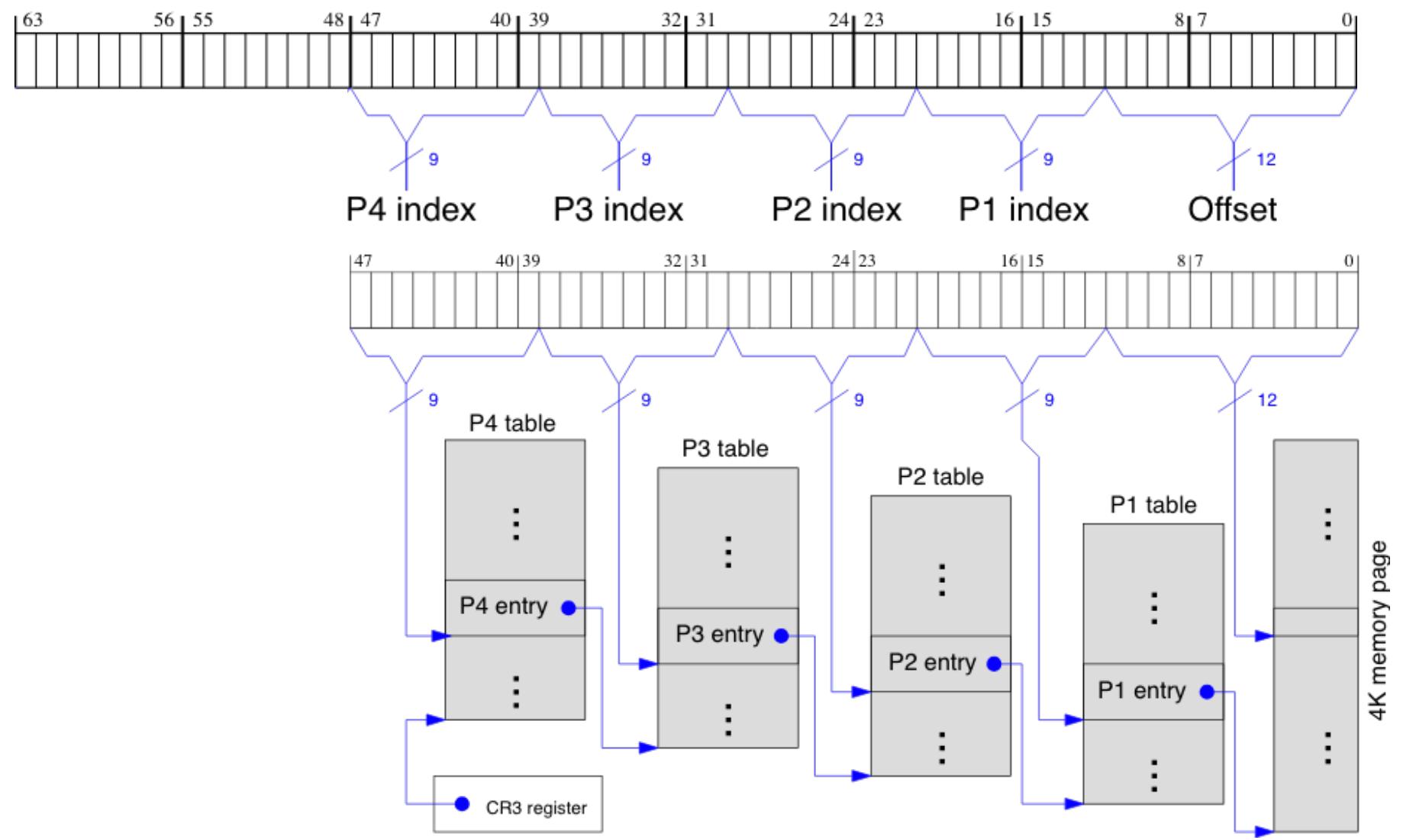
I could be doing a lot of useful work  
while waiting to read memory !

What to do ?

# **Caching Generalities**

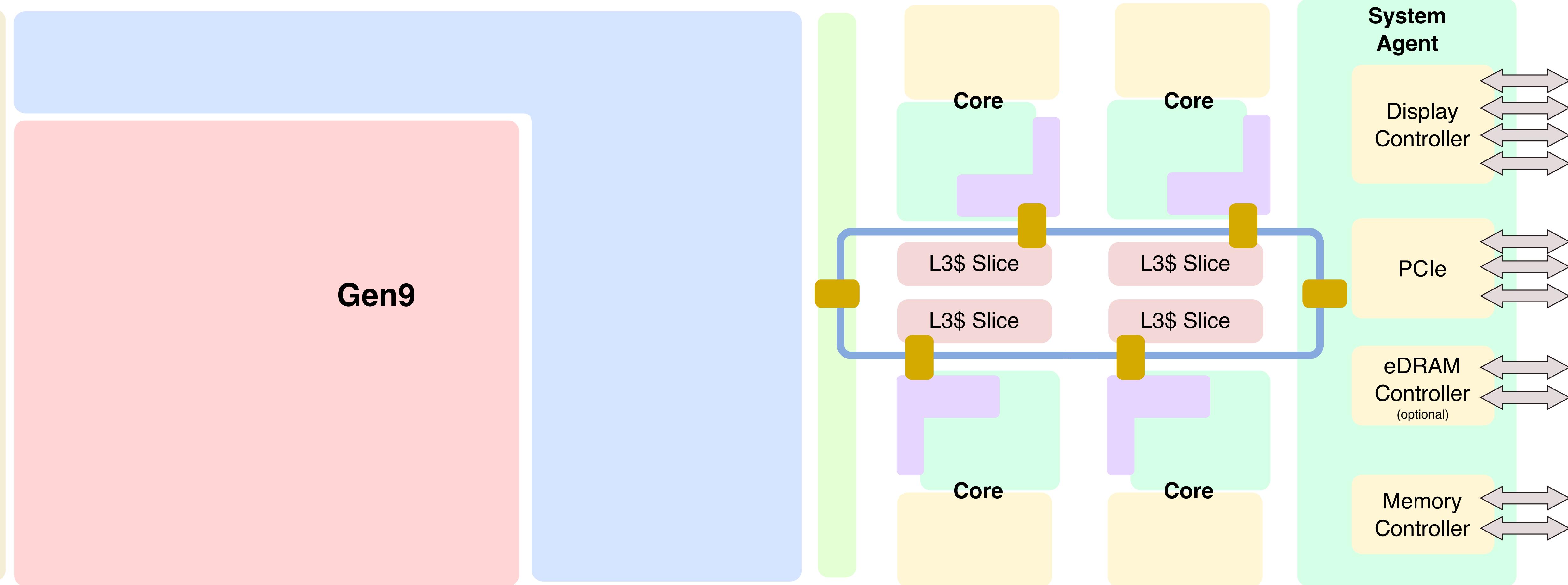
# Recap of Key Concepts in Caching

- 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
- Speed vs. size (and thus cost) trade-off

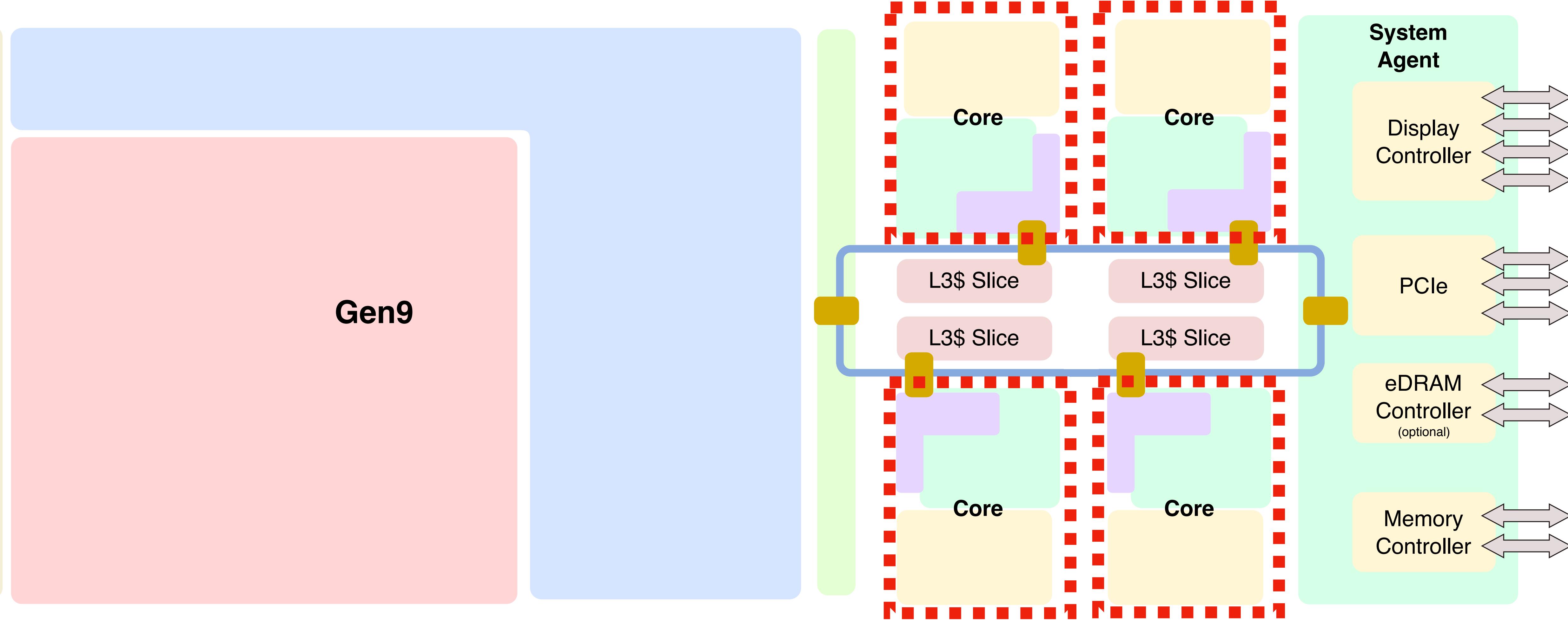


# Caching & TLBs

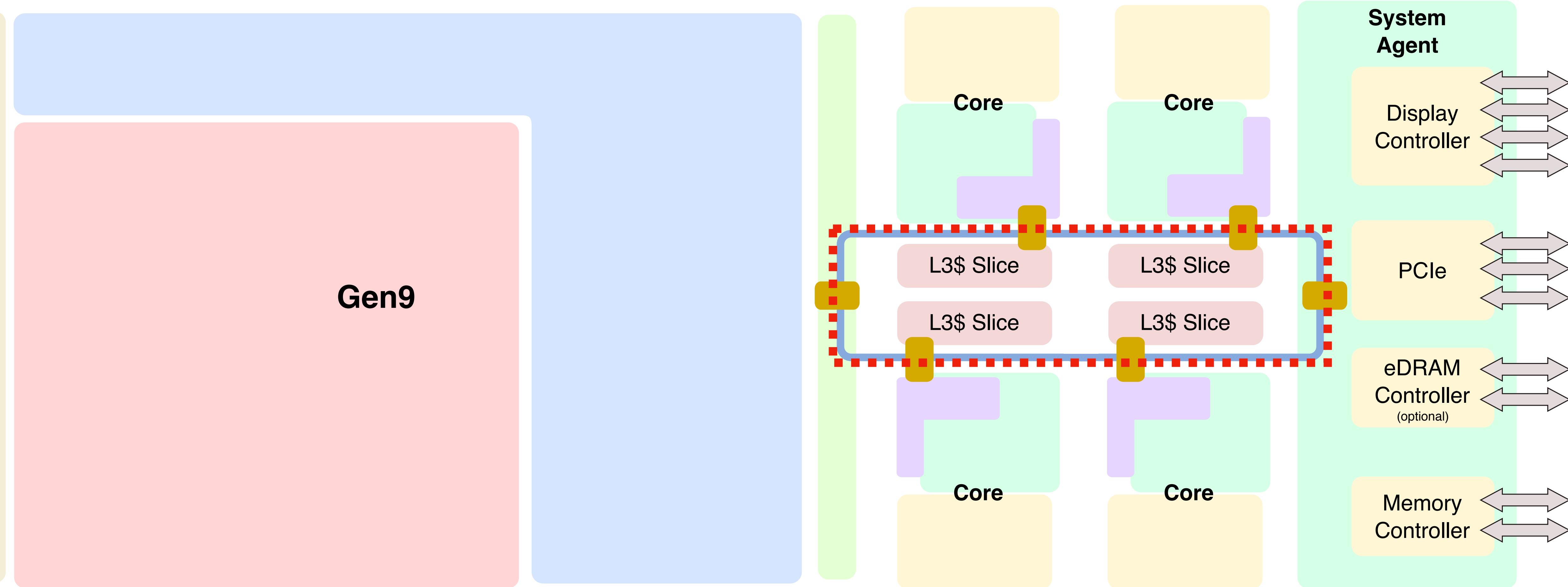
# Intel Skylake Microarchitecture



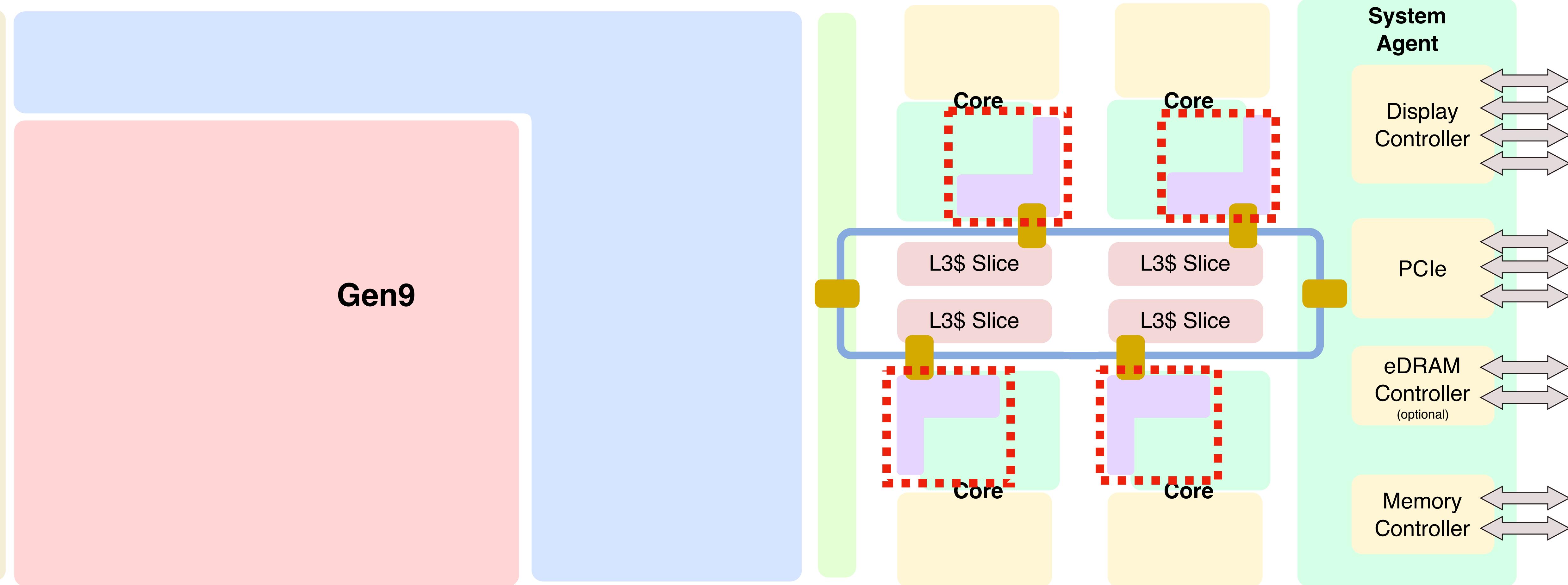
# Intel Skylake Microarchitecture



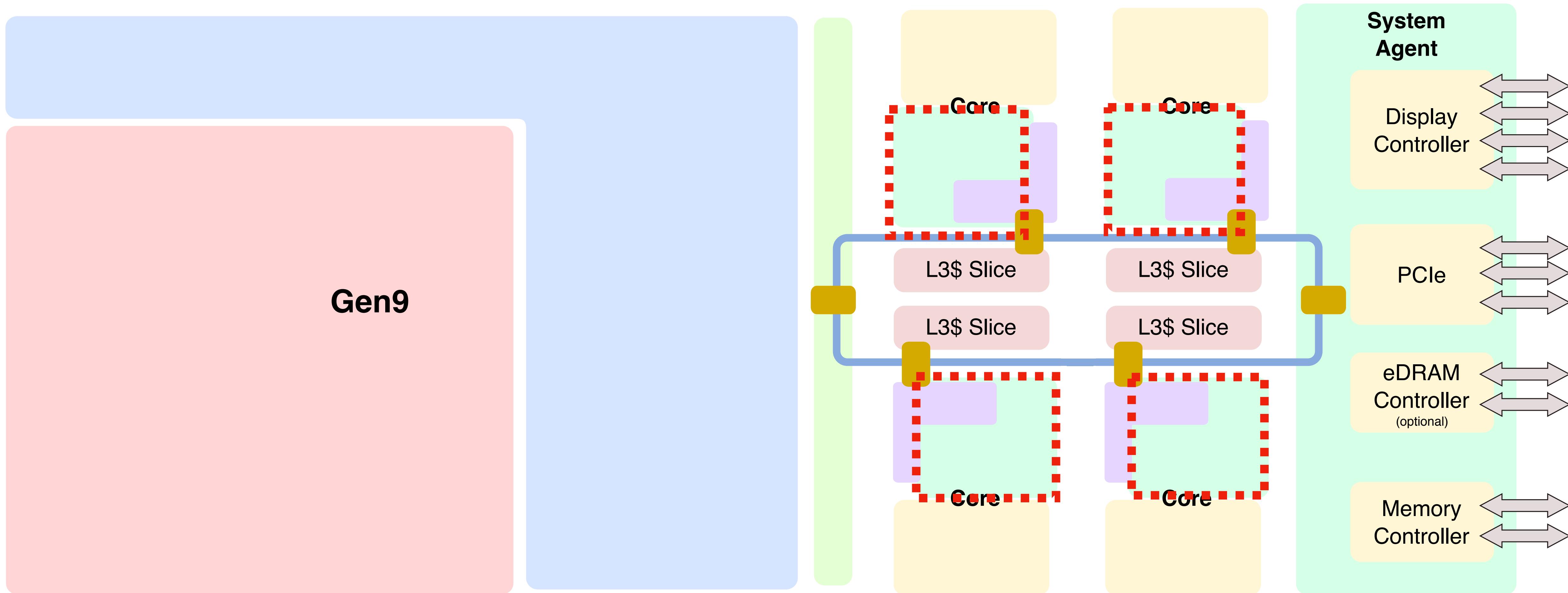
# Intel Skylake Microarchitecture



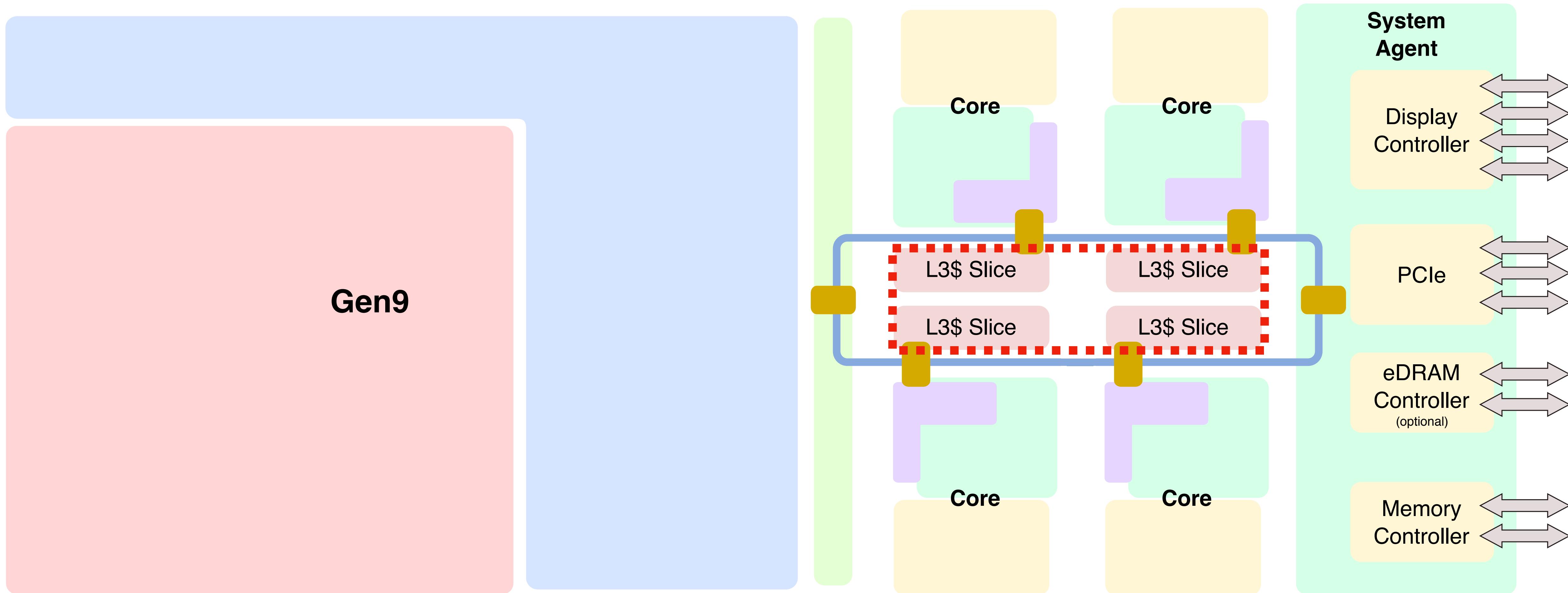
# Intel Skylake Microarchitecture



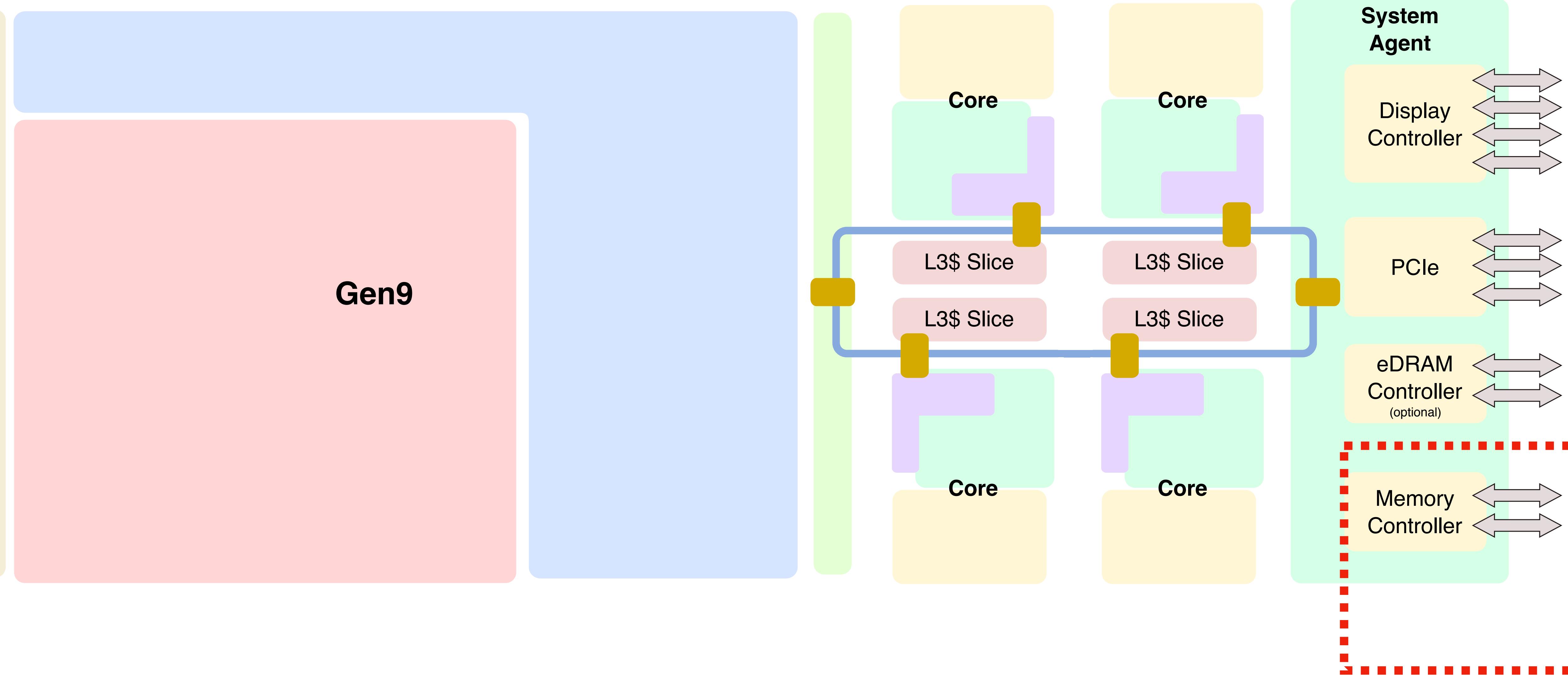
# Intel Skylake Microarchitecture



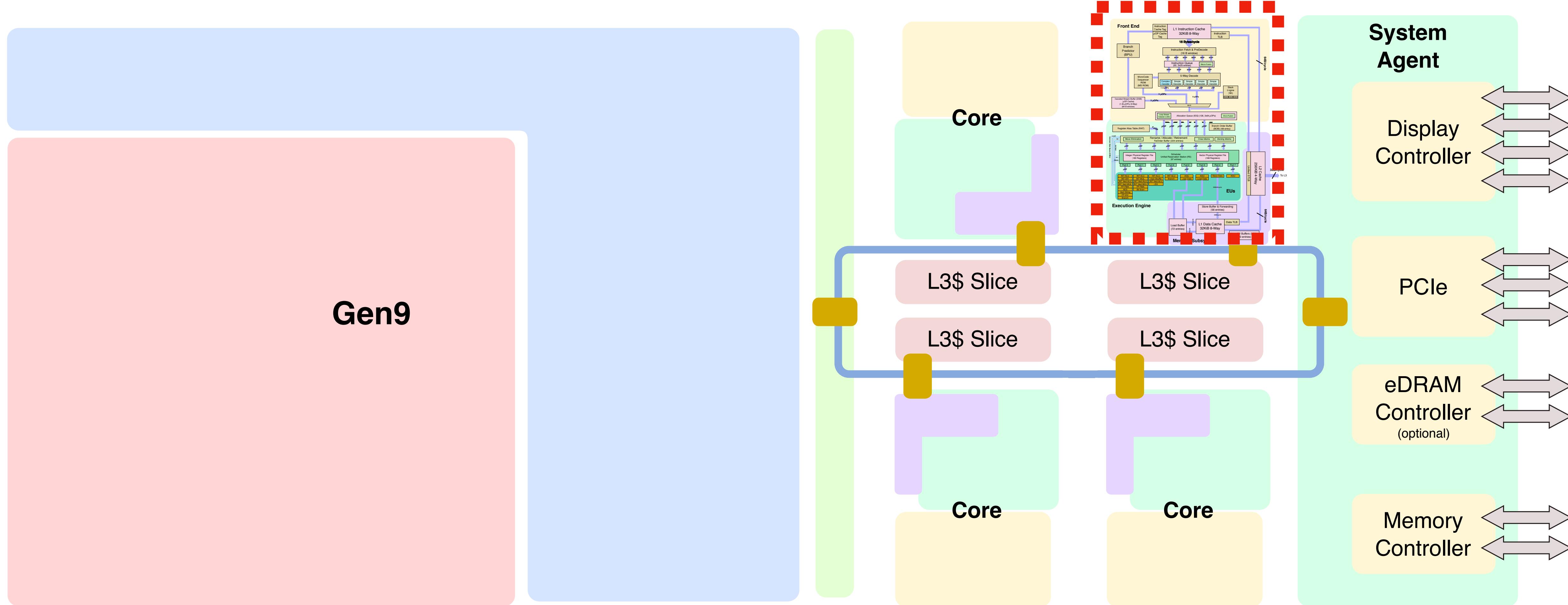
# Intel Skylake Microarchitecture

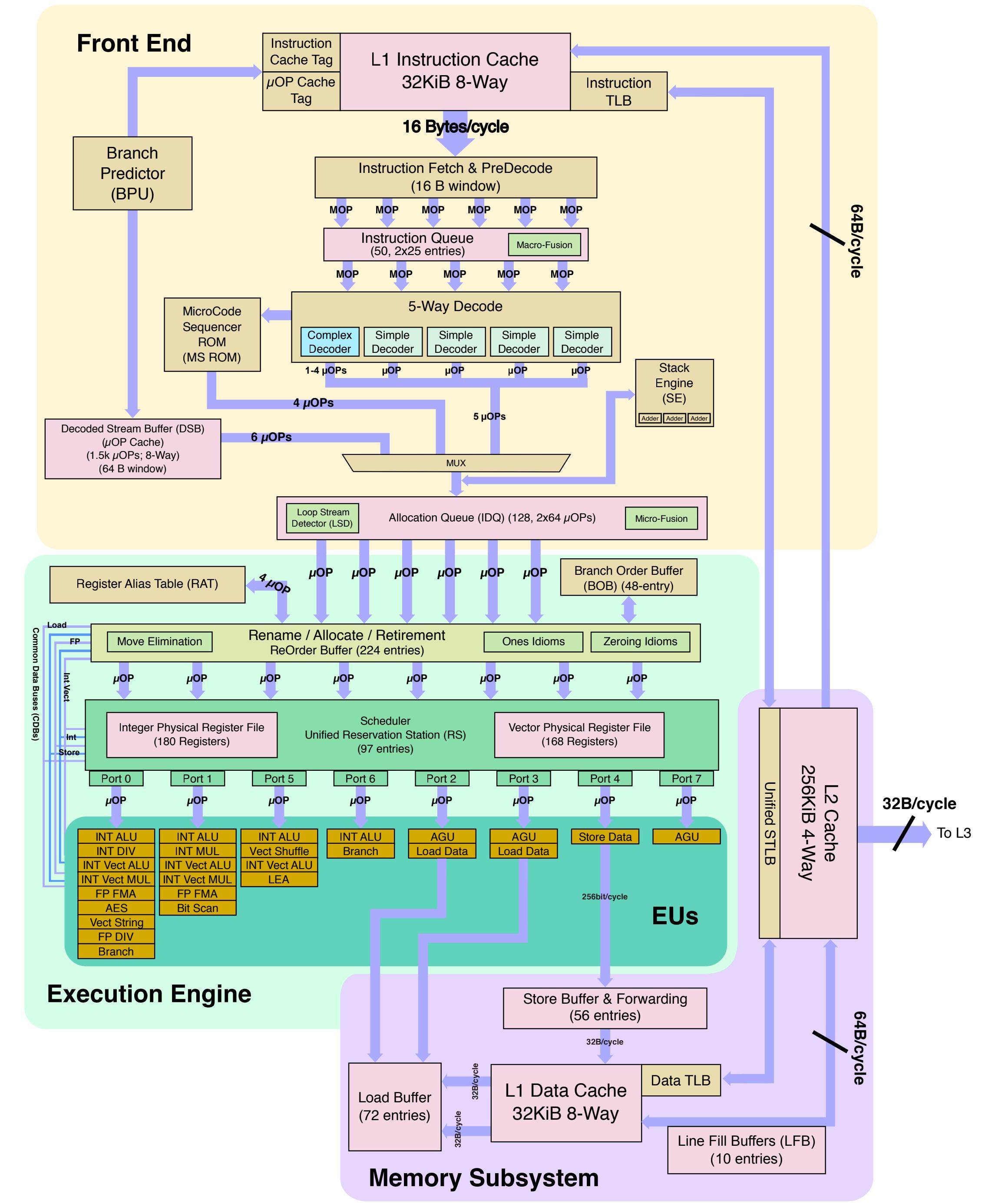


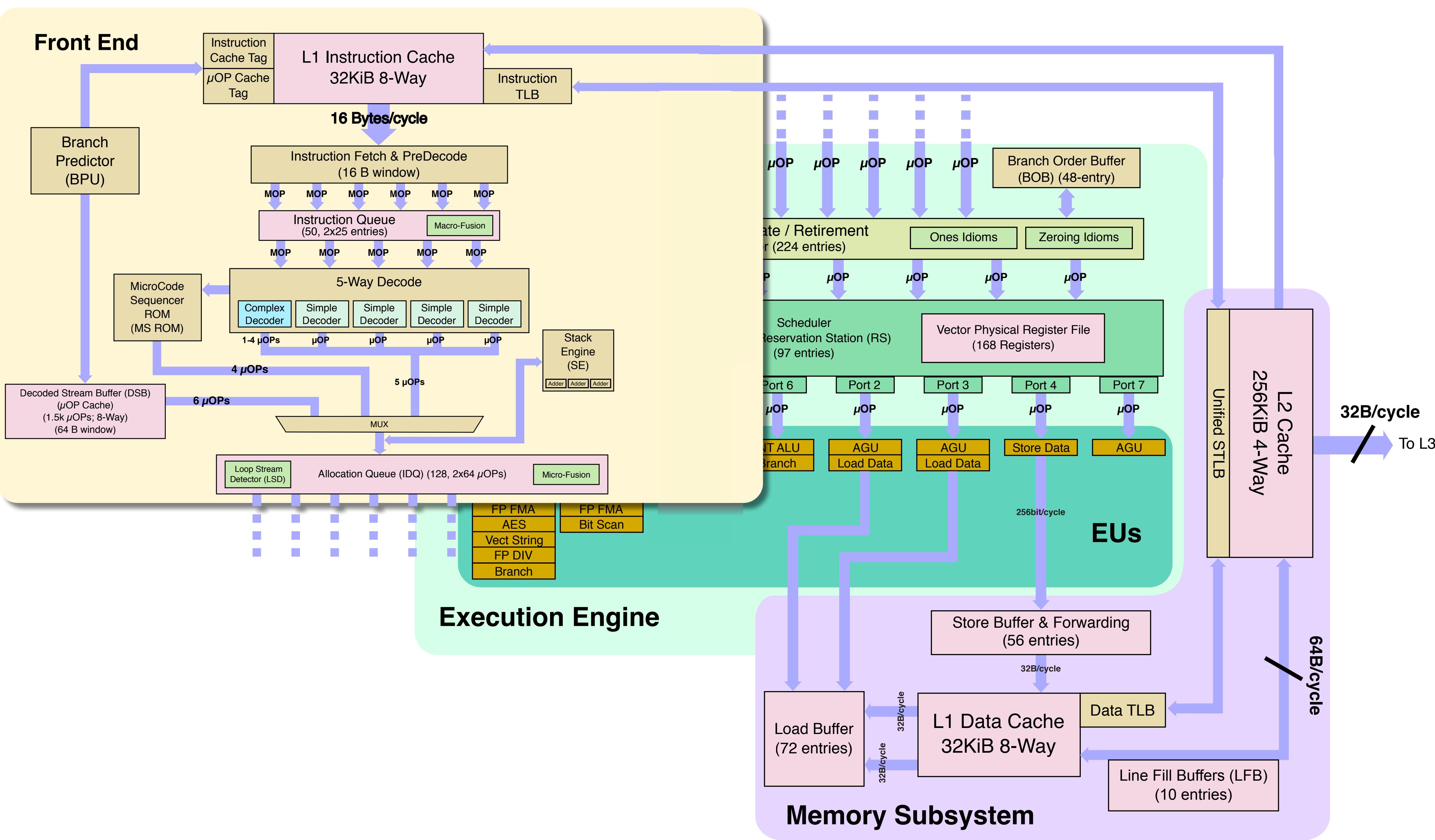
# Intel Skylake Microarchitecture

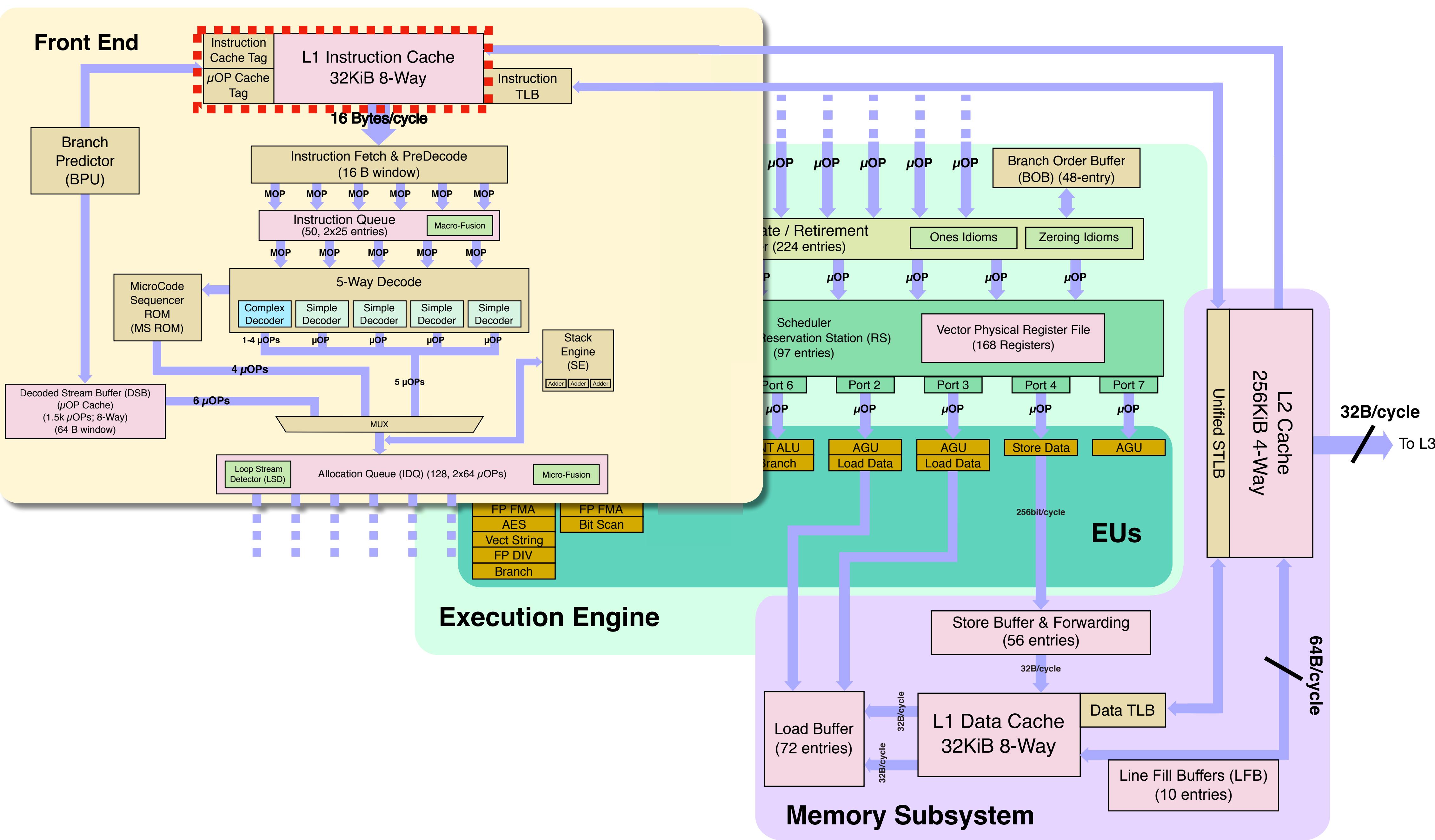


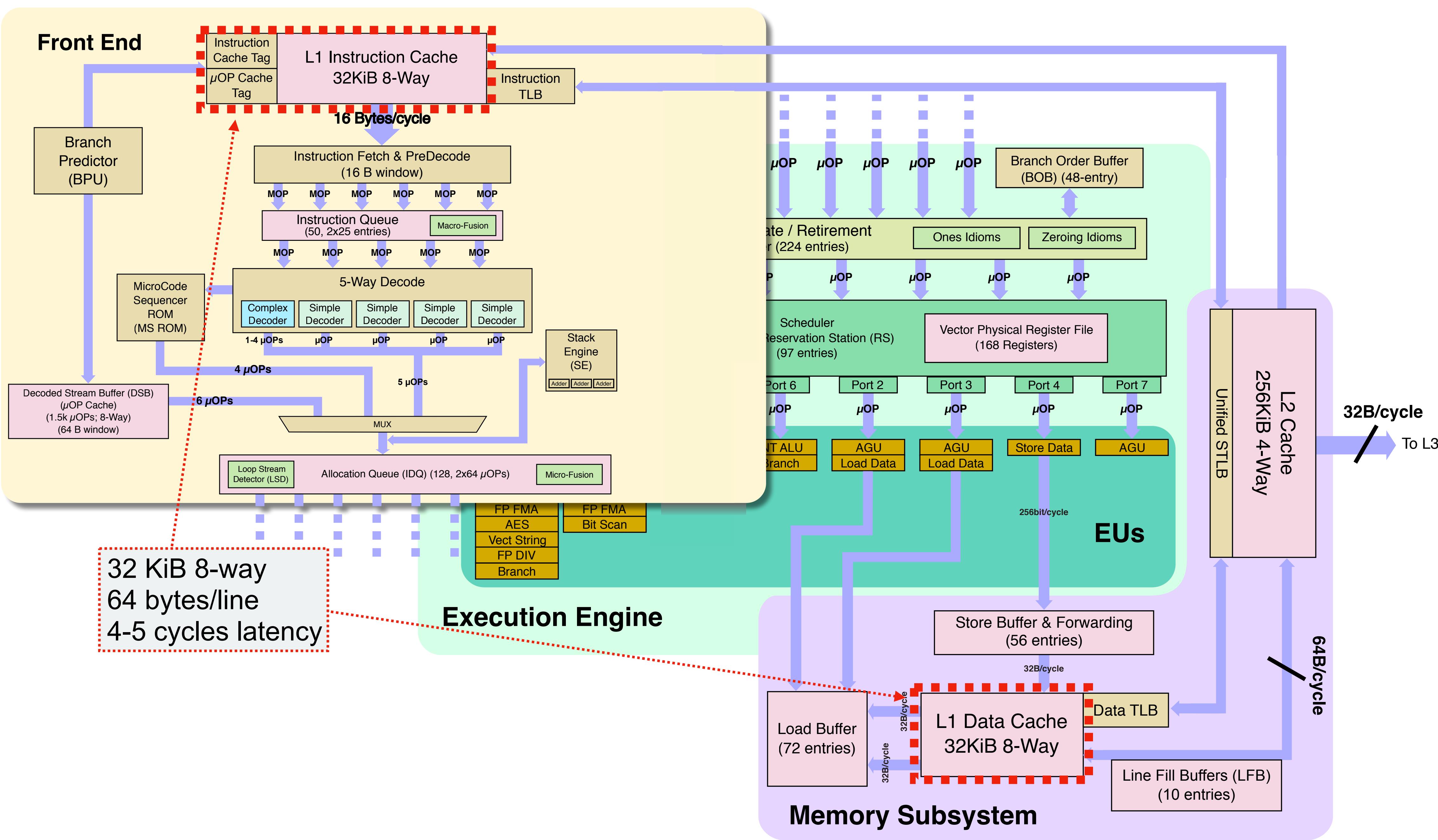
# Intel Skylake Microarchitecture

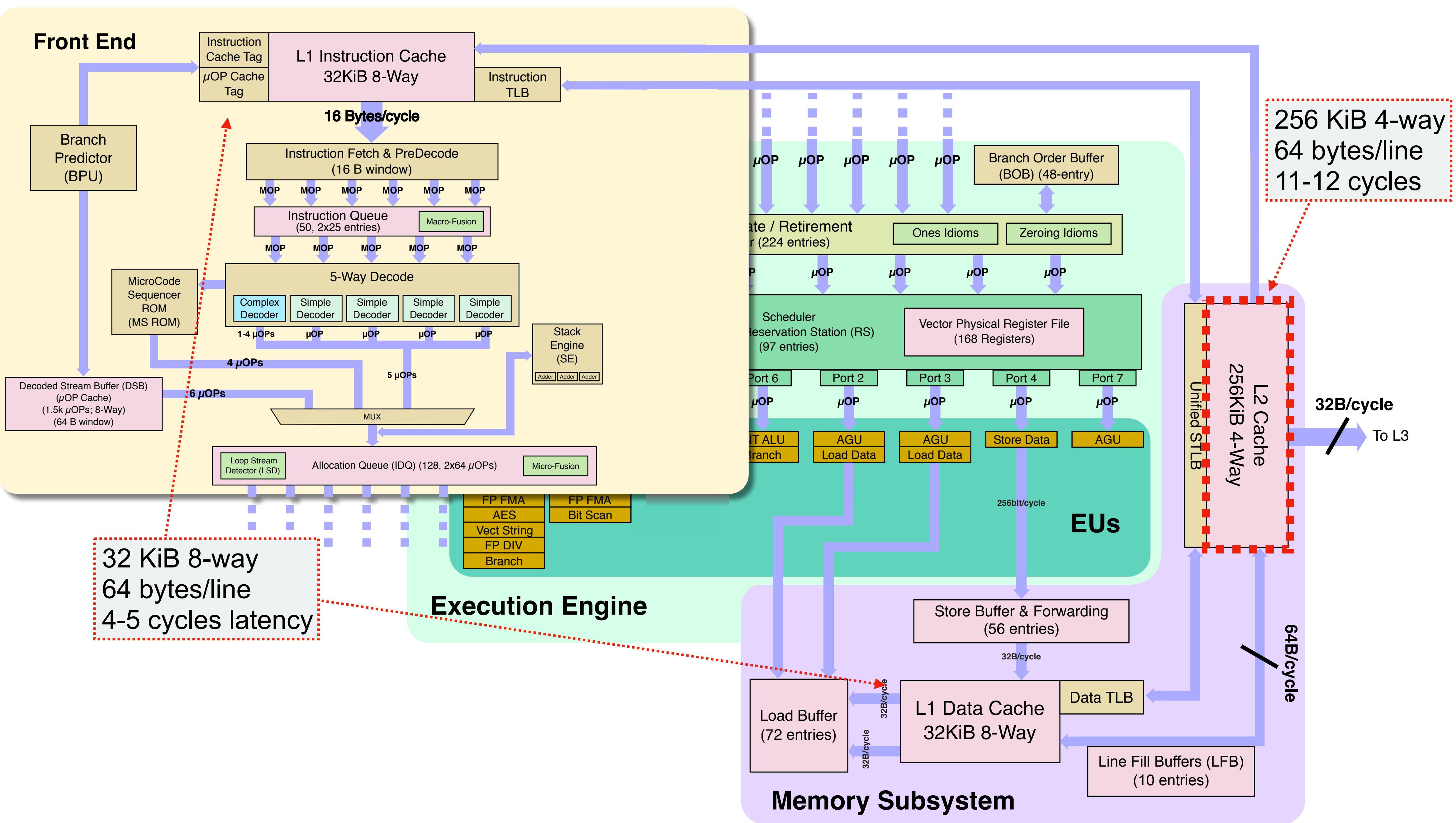


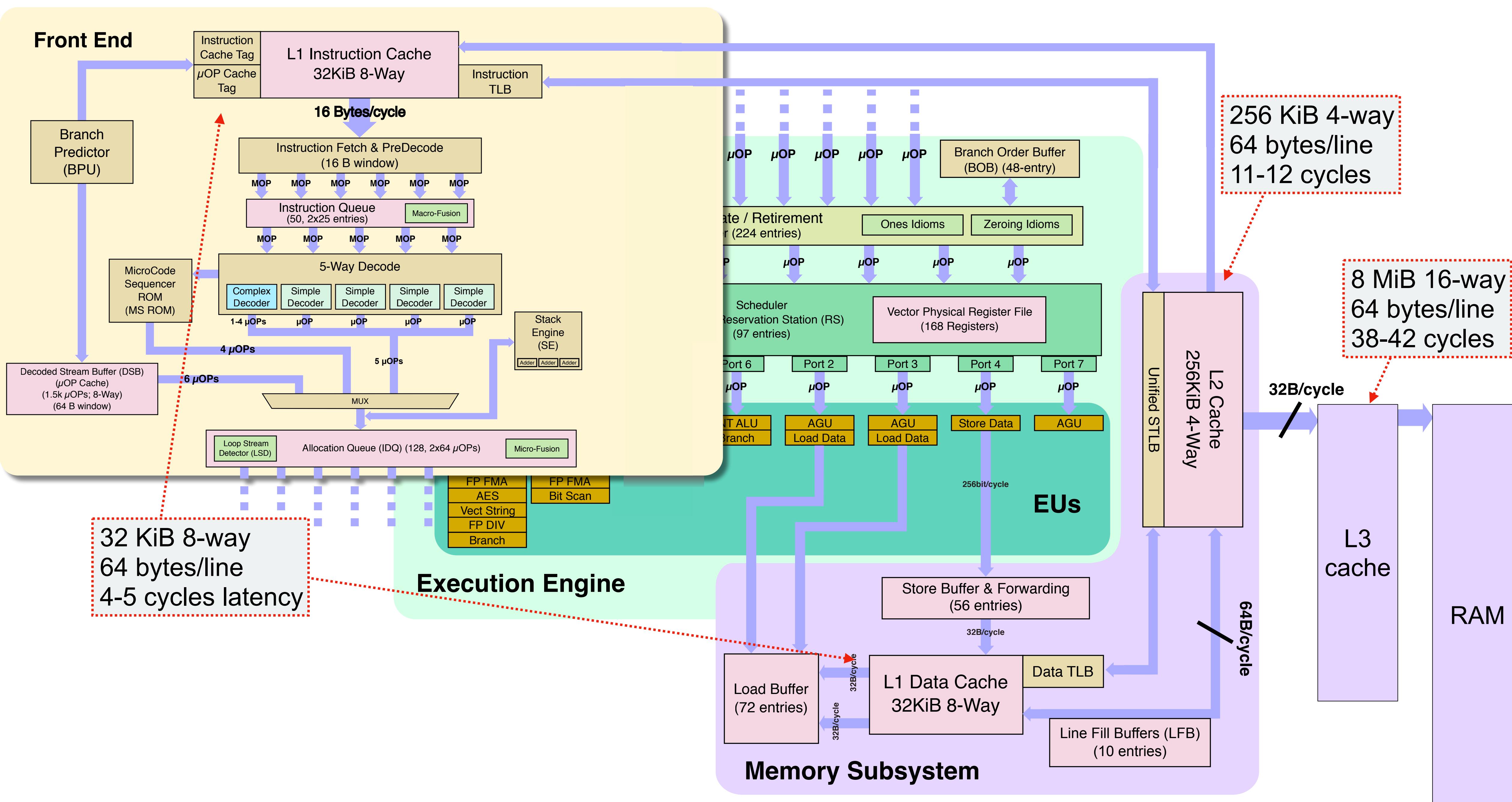


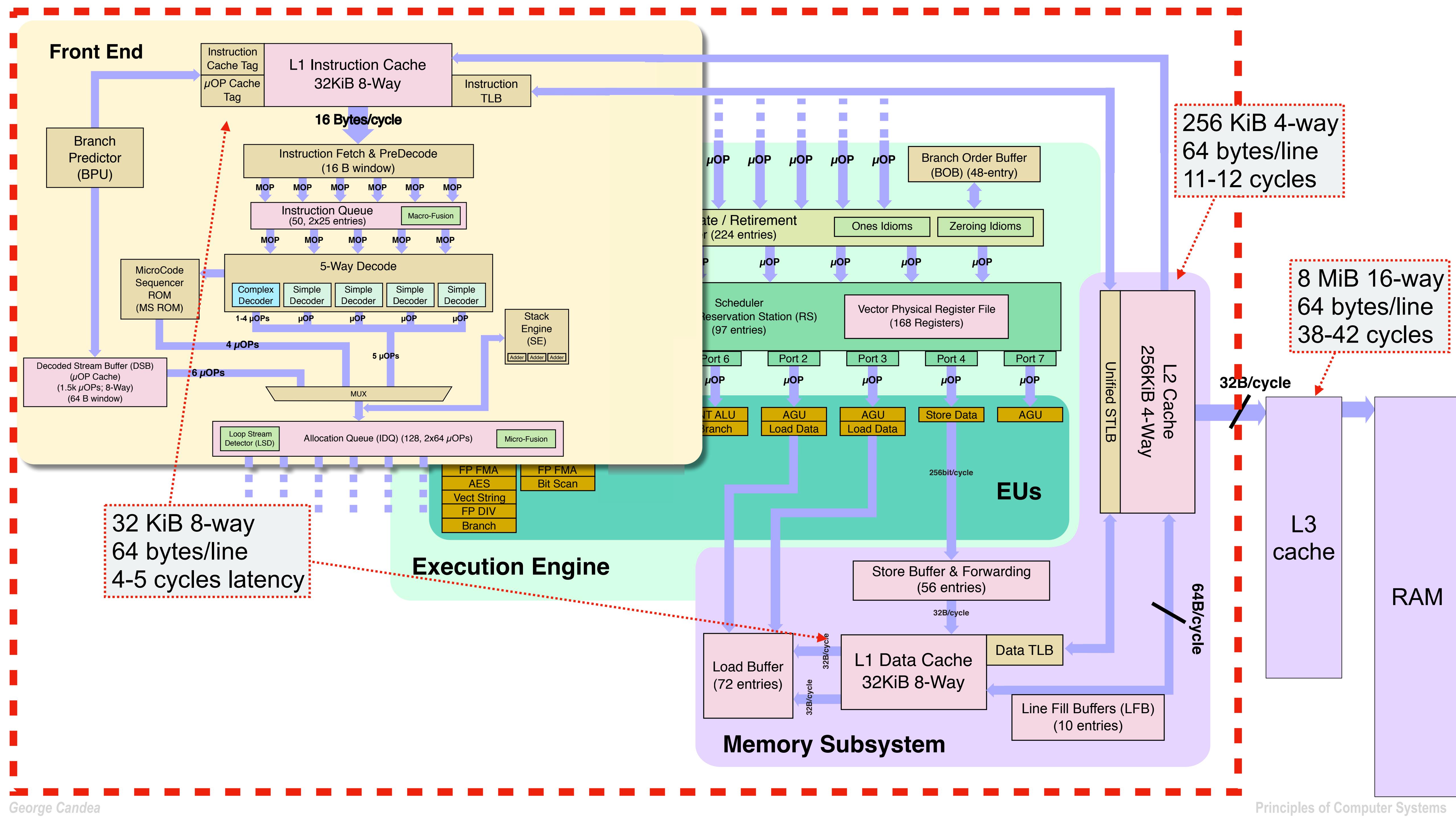


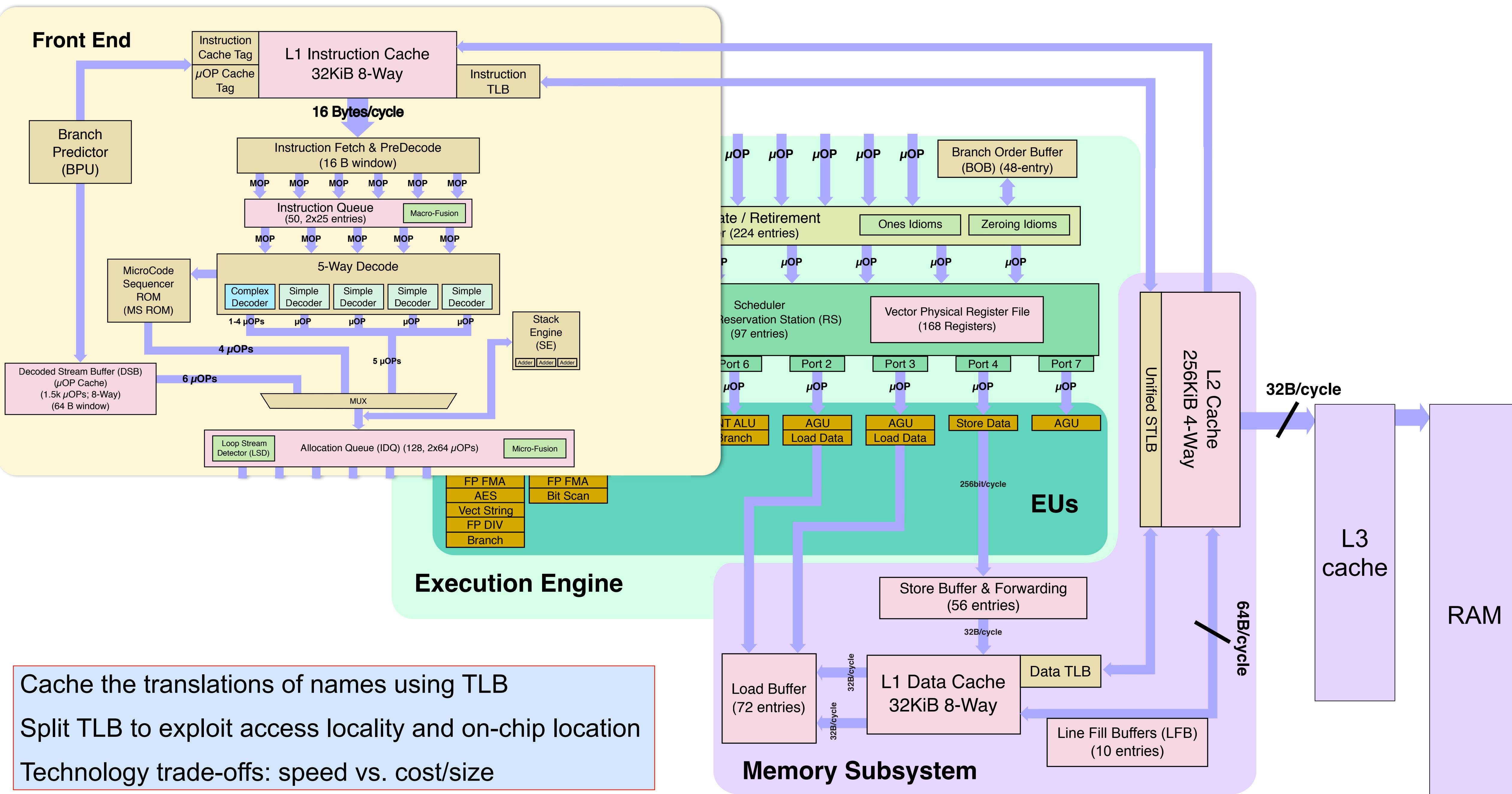


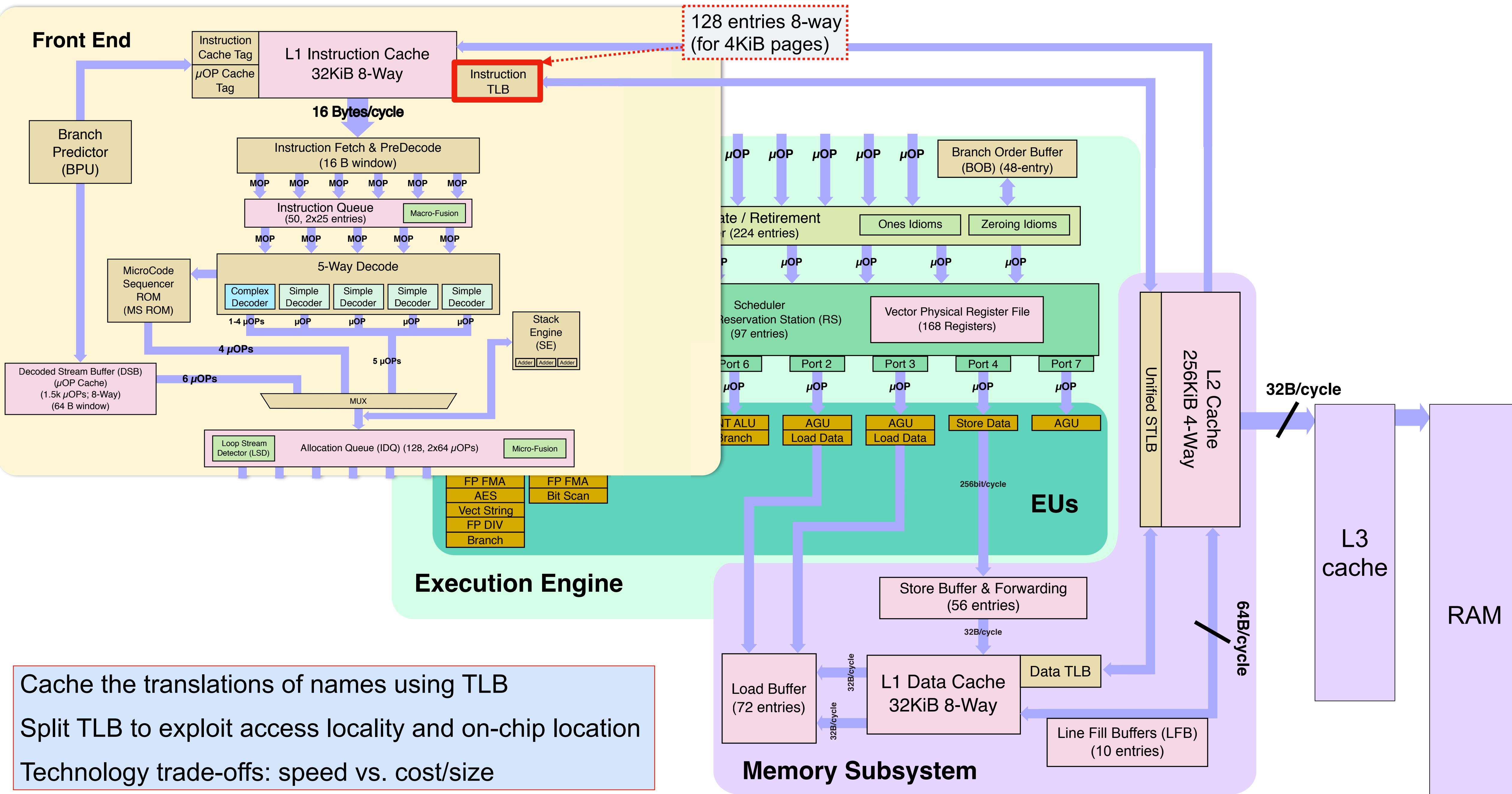


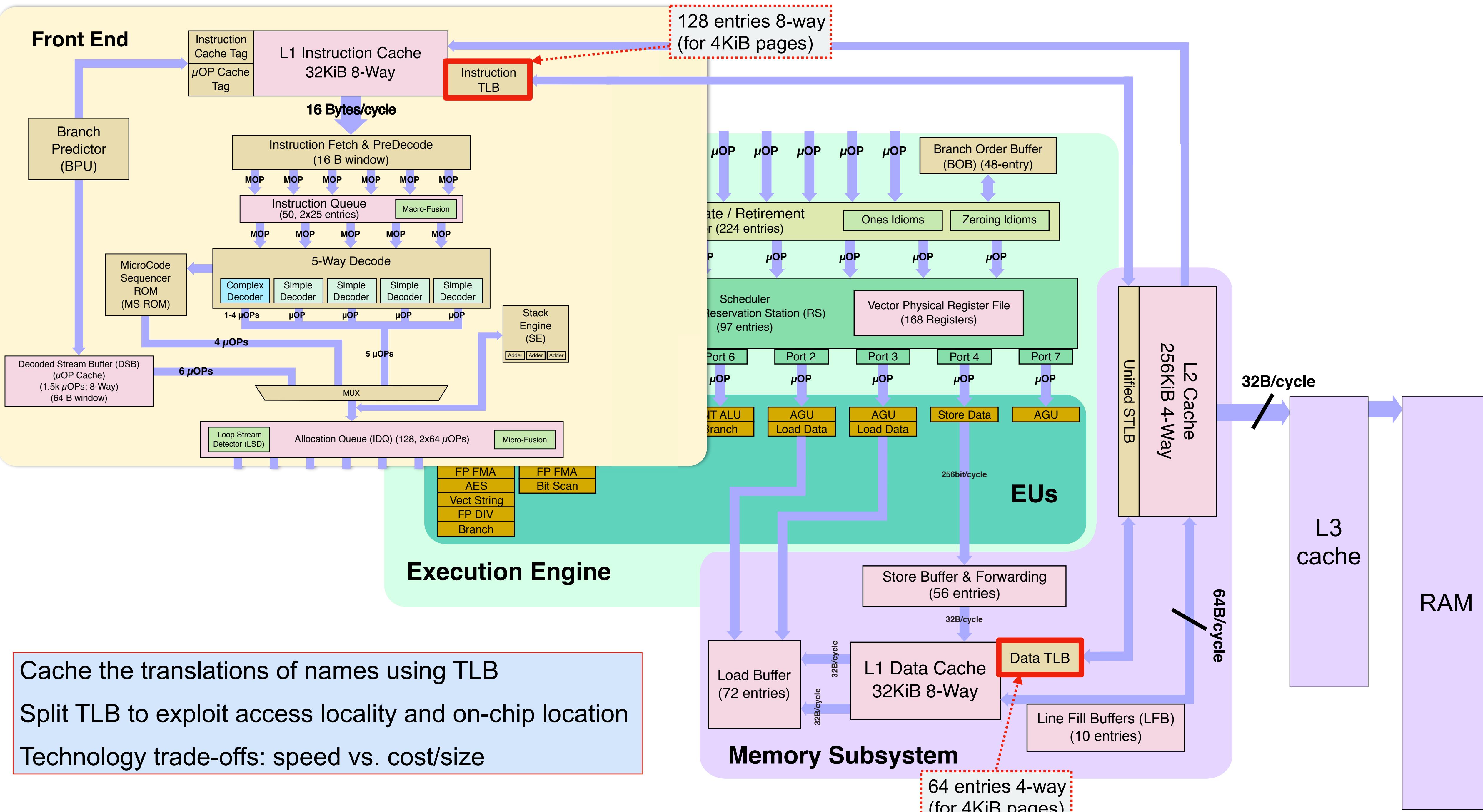


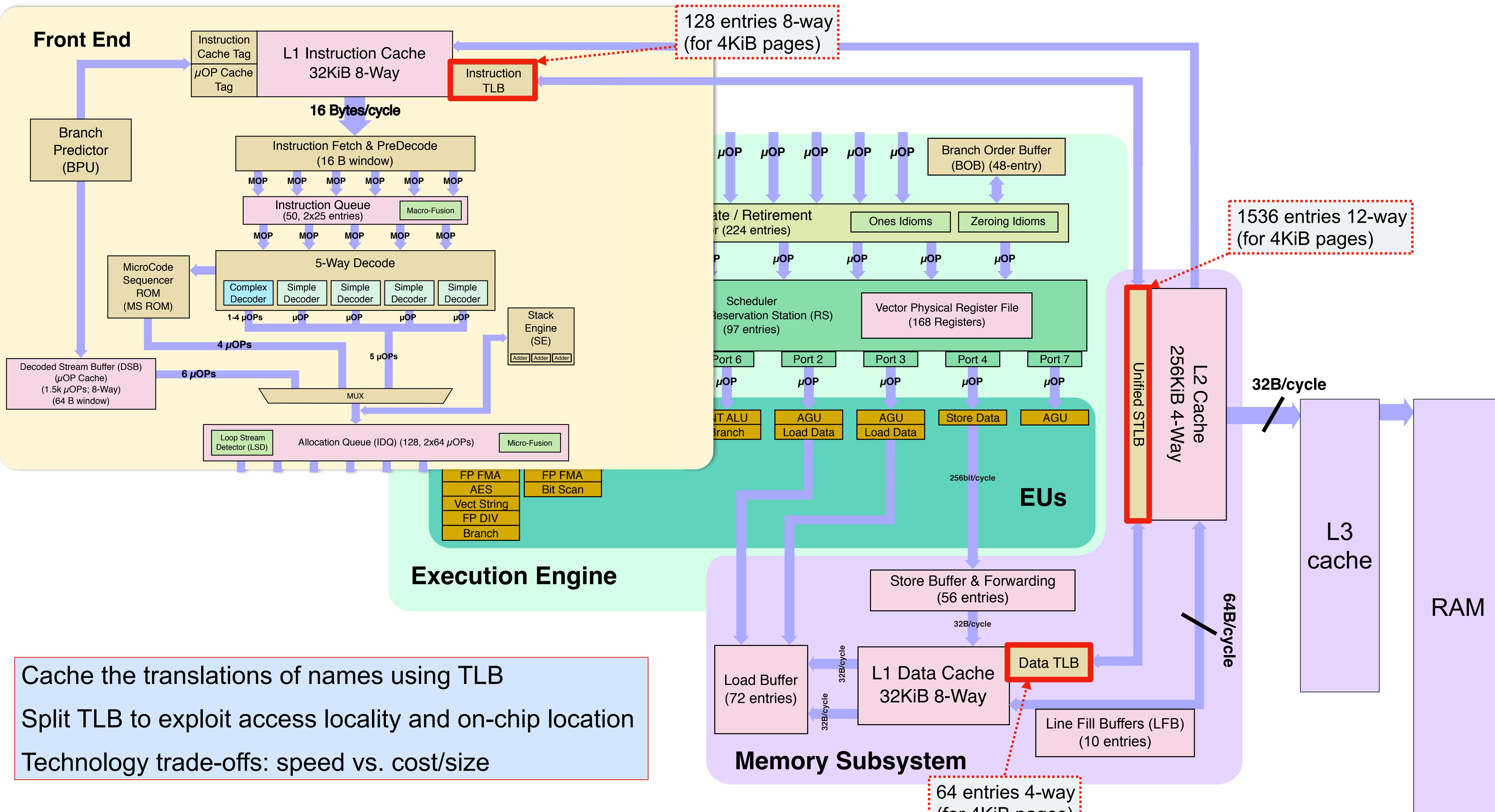






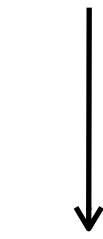




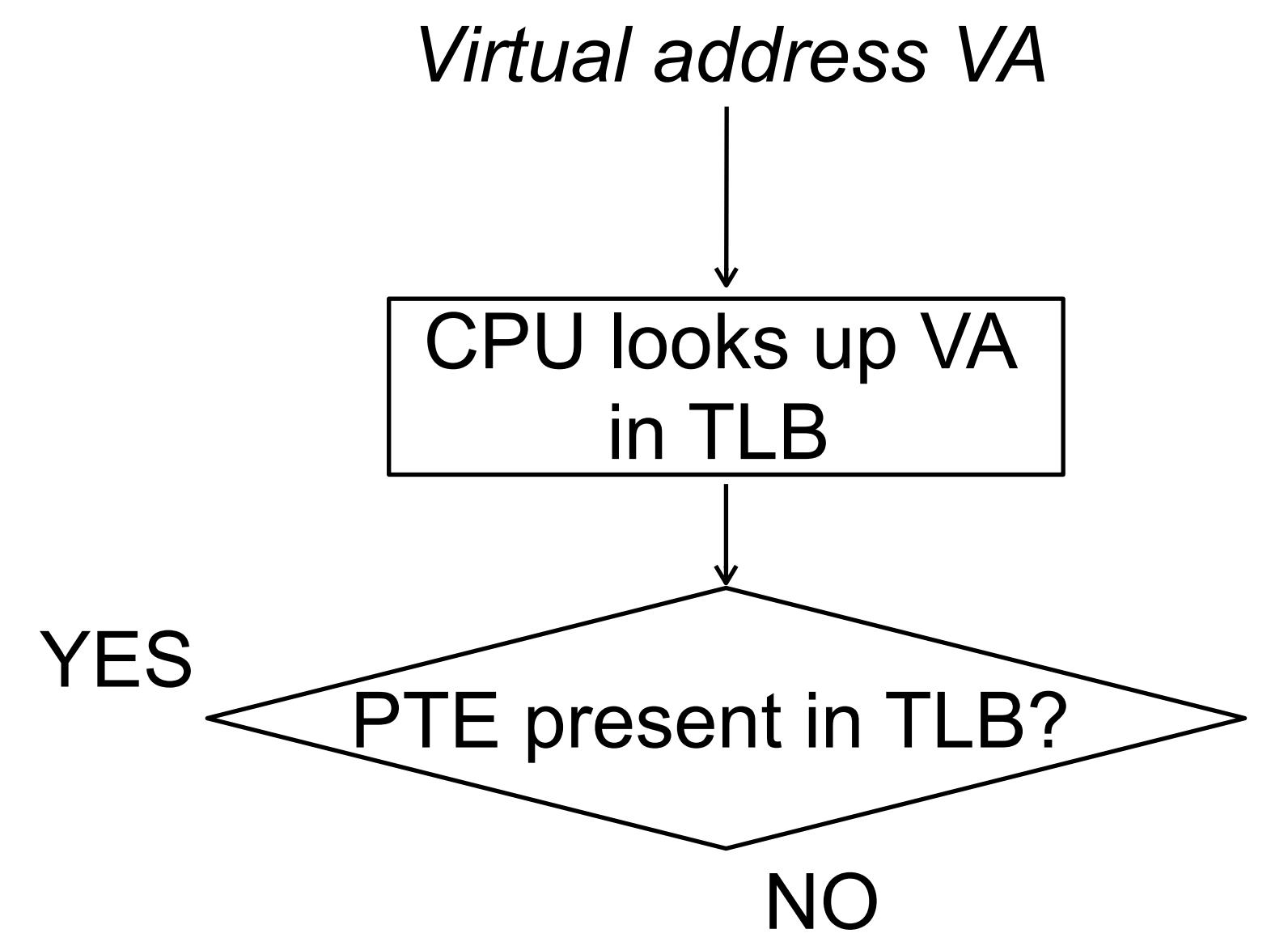


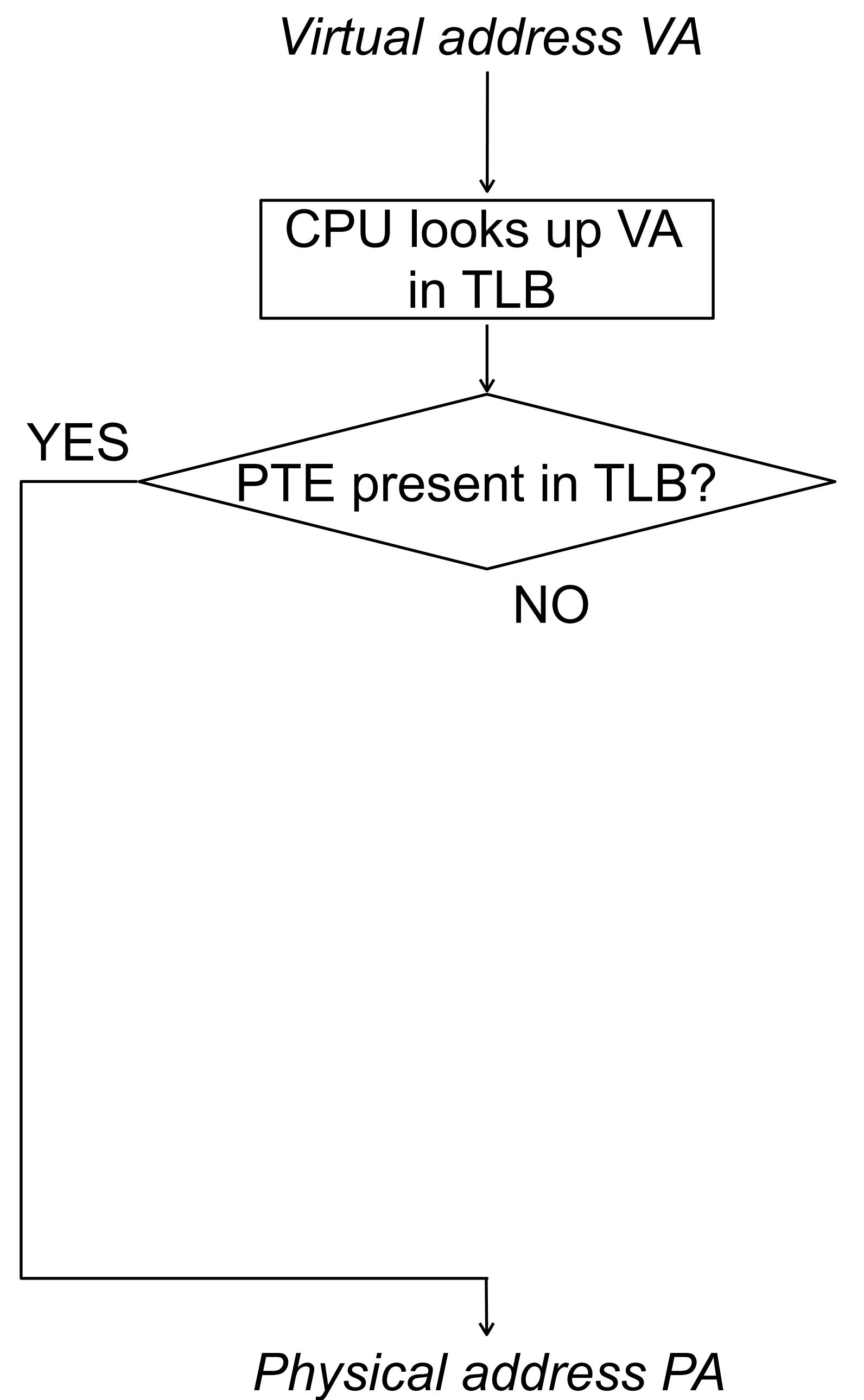
## *Virtual address VA*

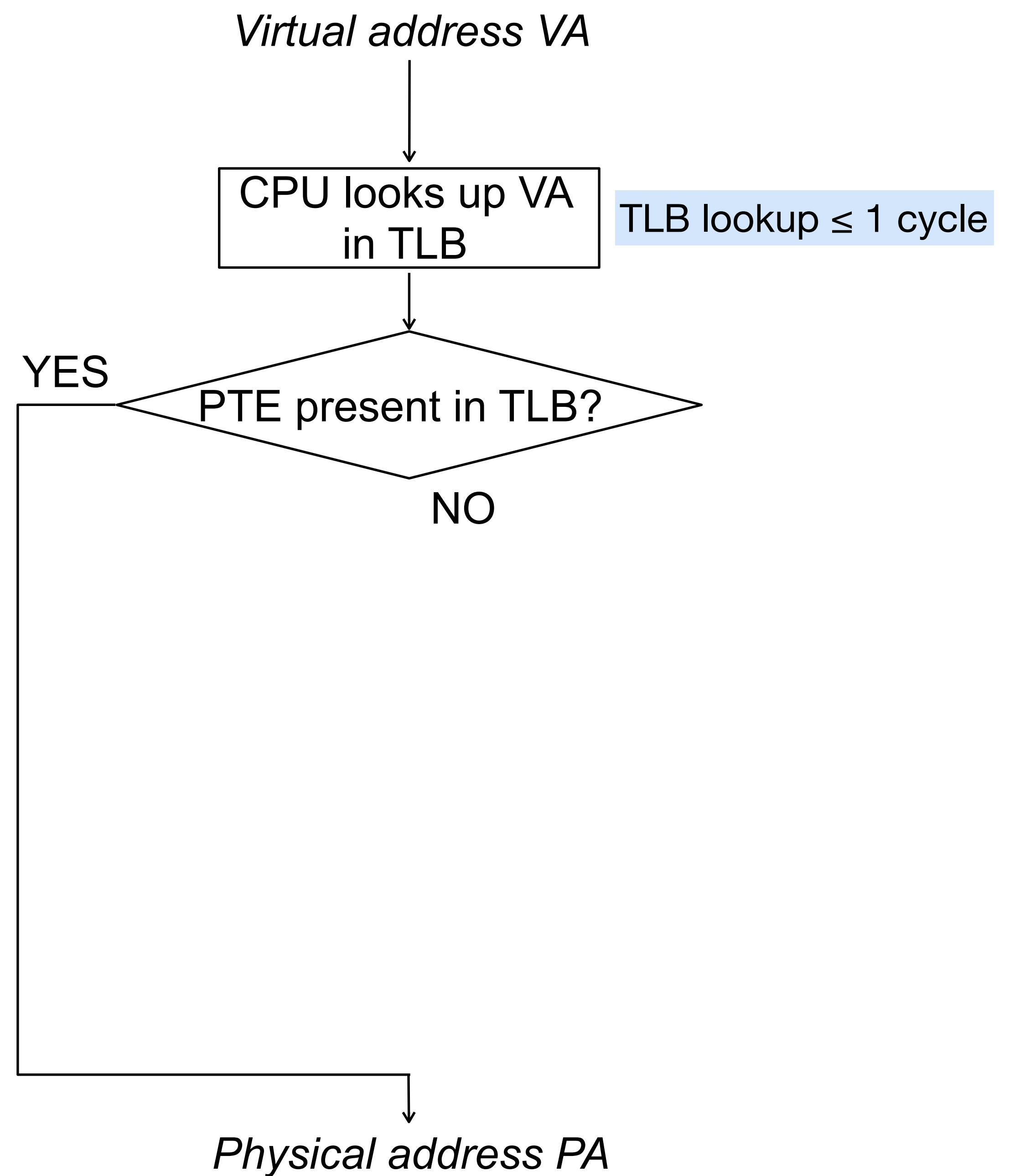
*Virtual address VA*

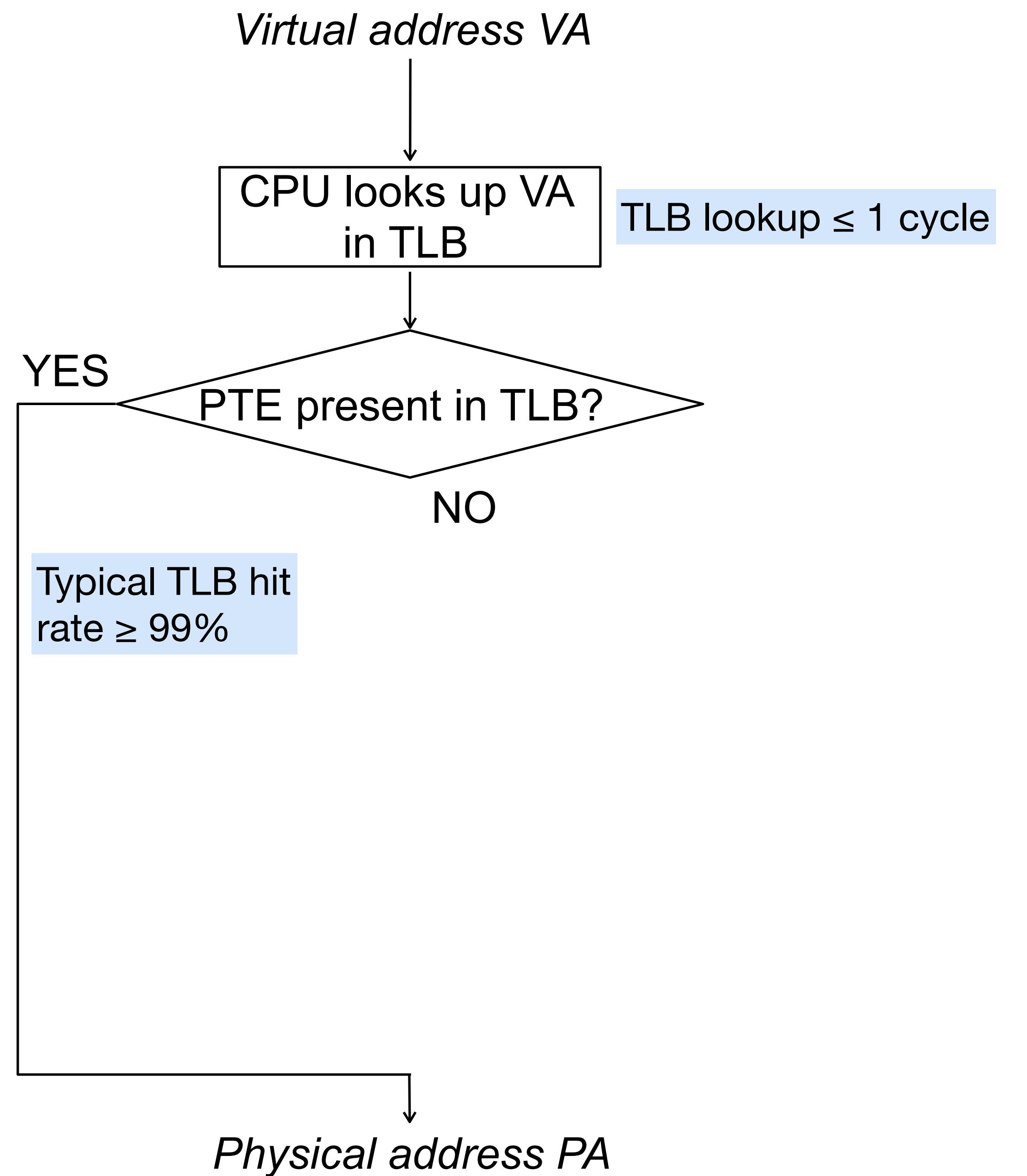


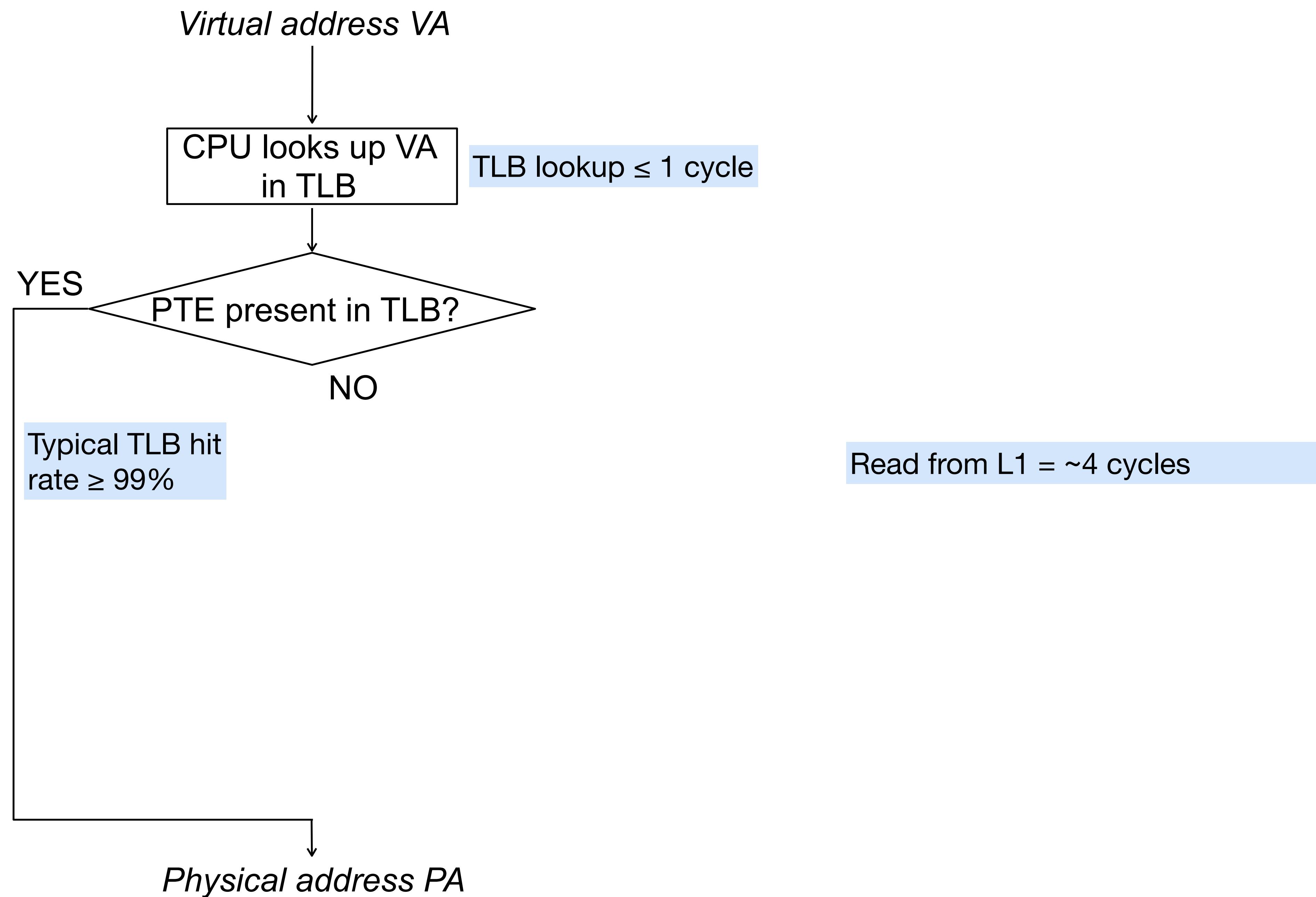
CPU looks up VA  
in TLB

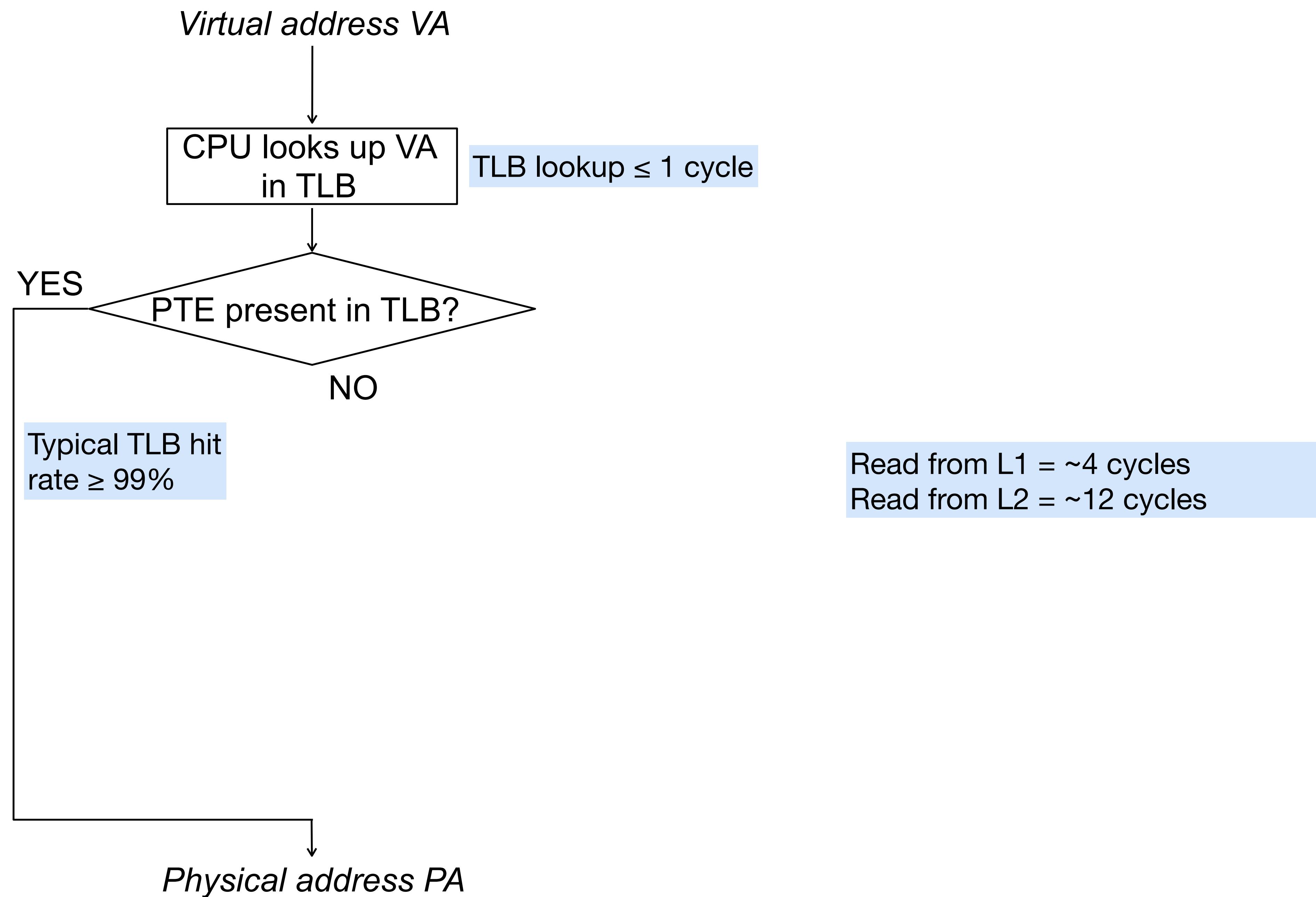


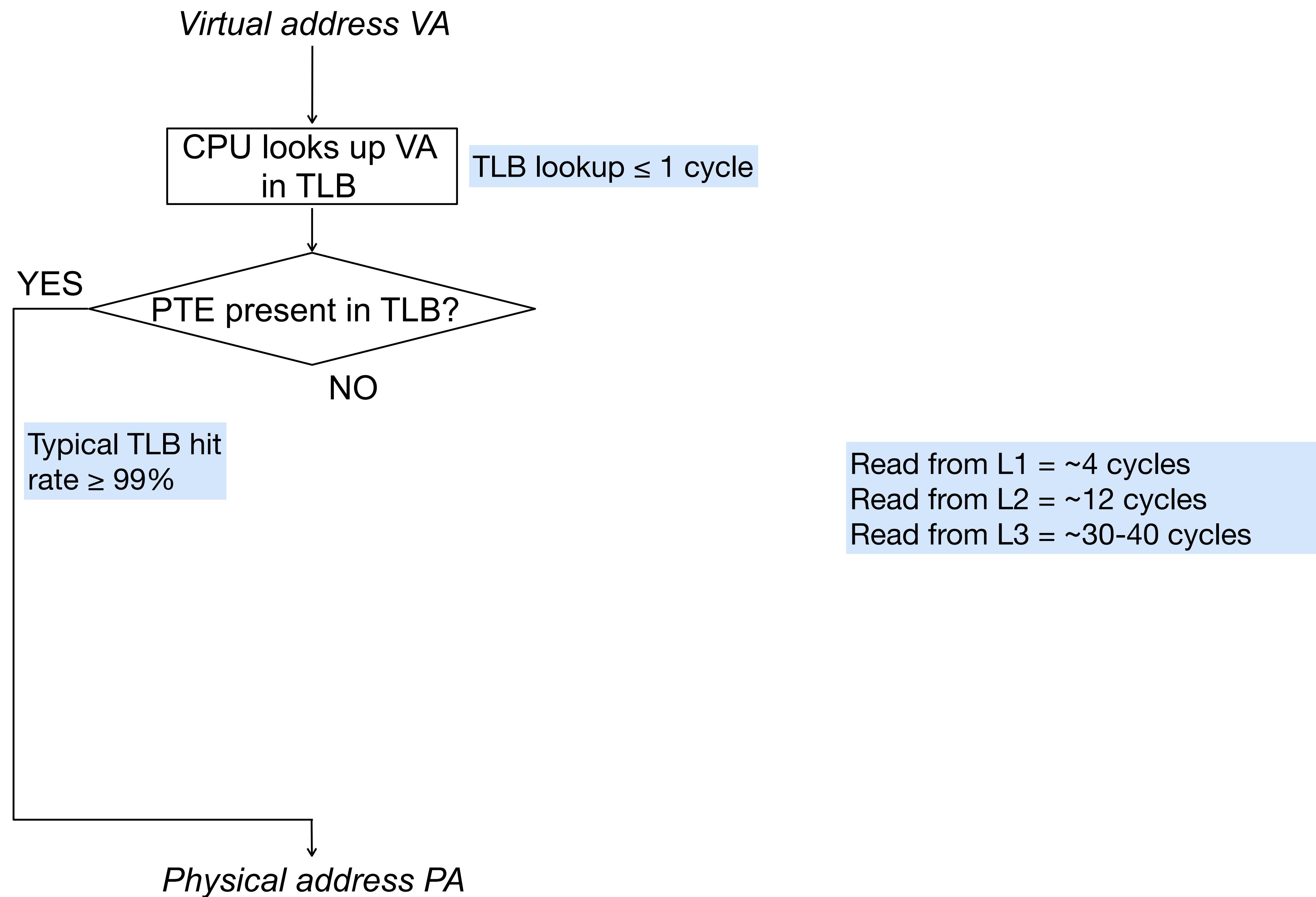


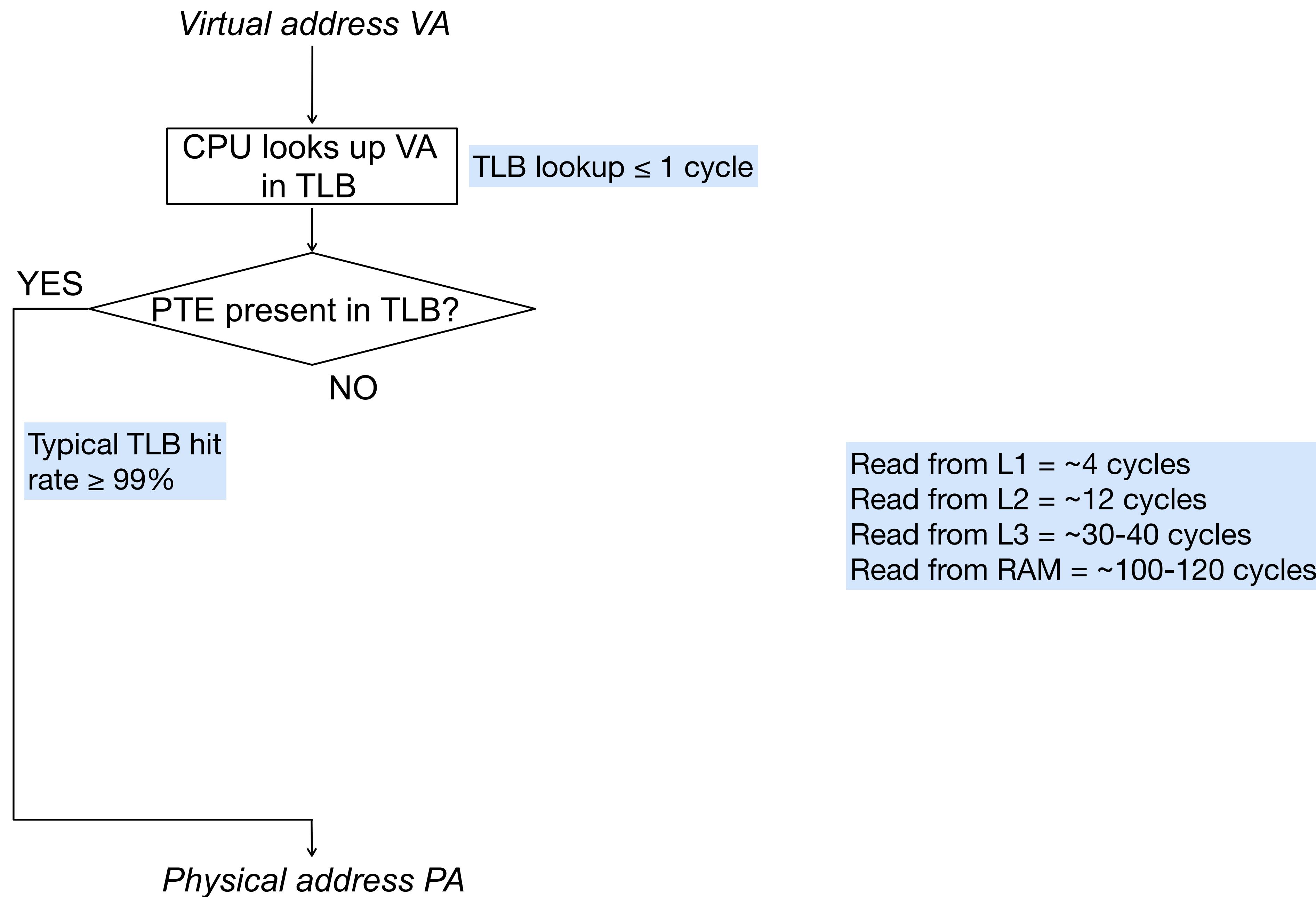


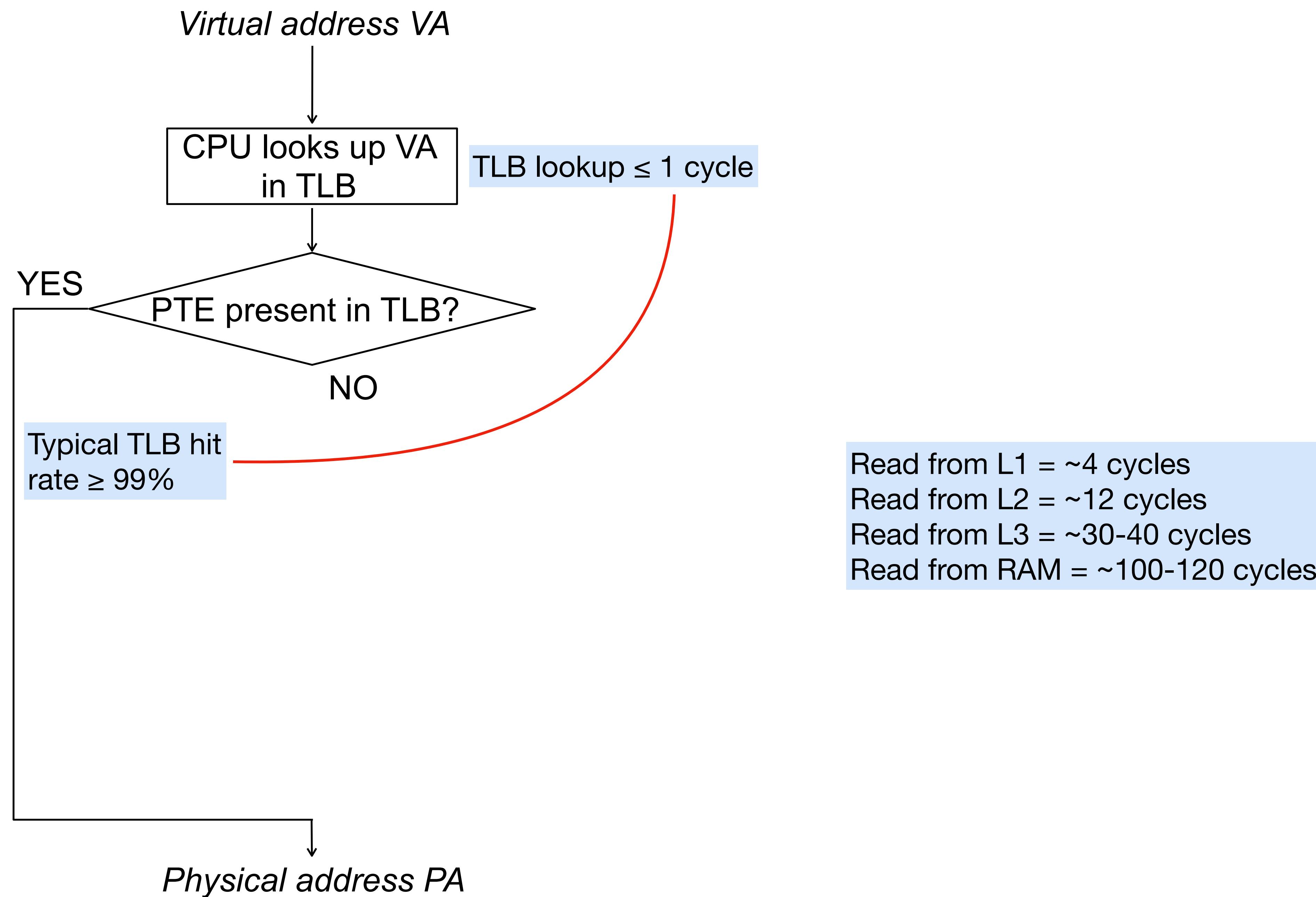


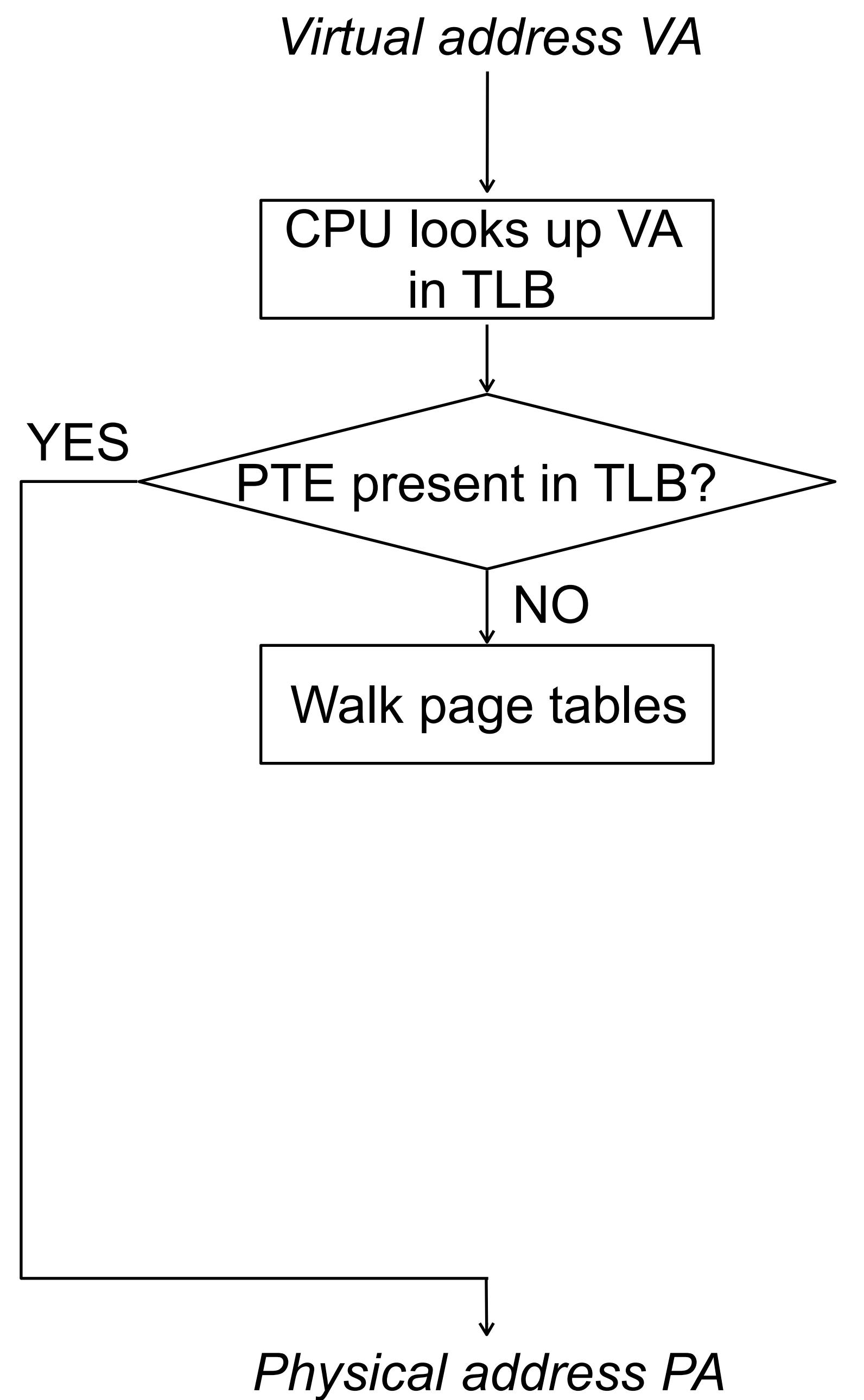


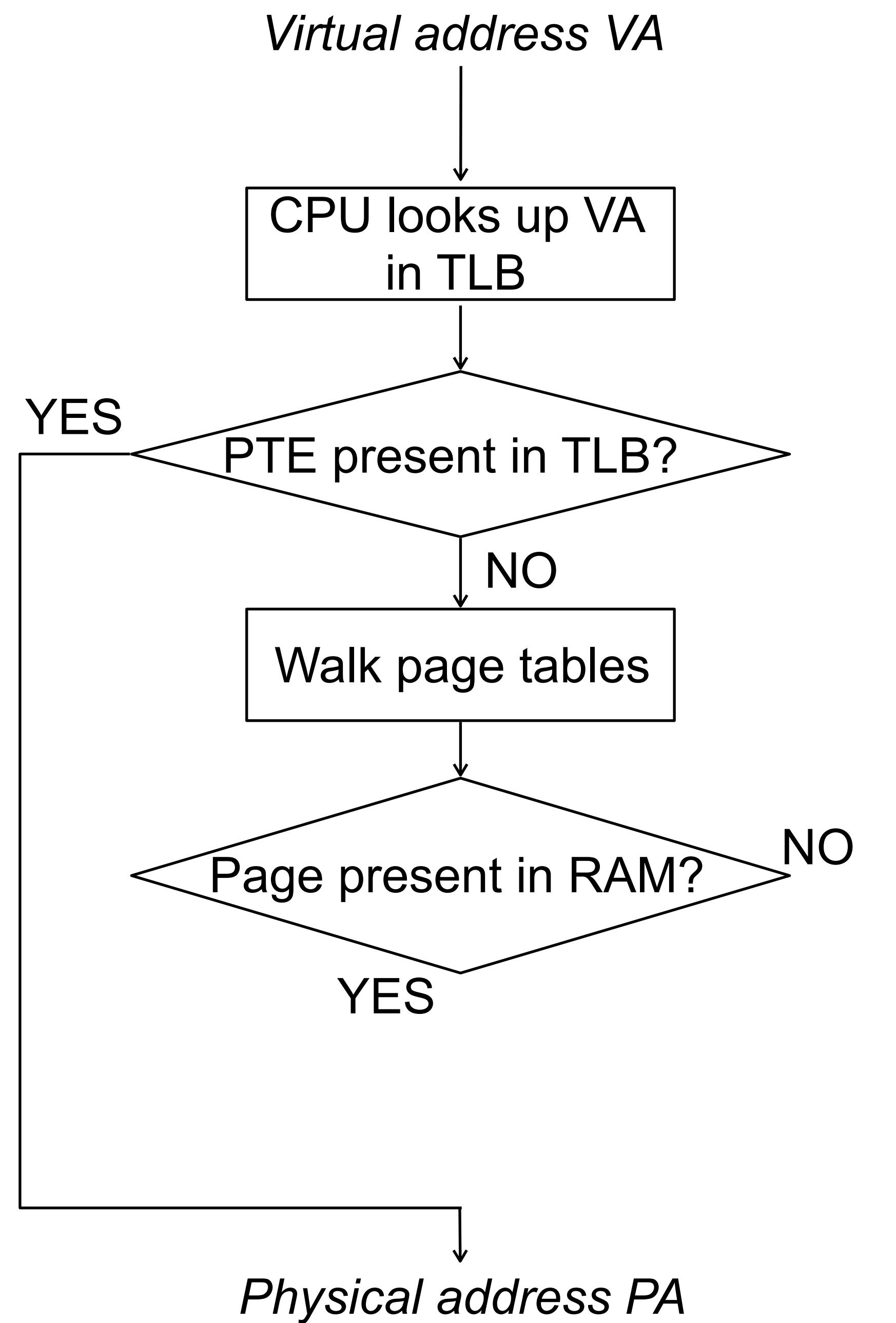


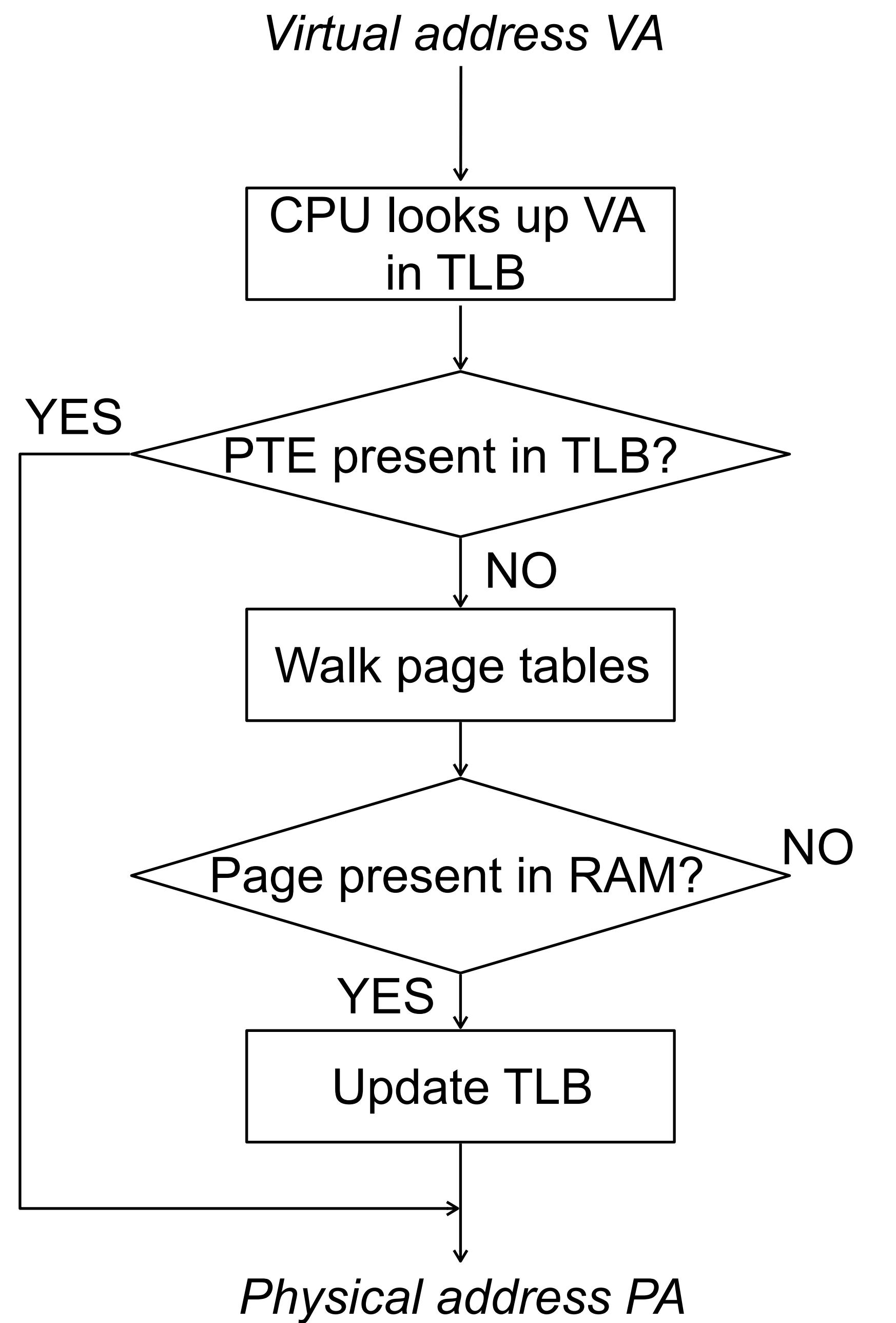


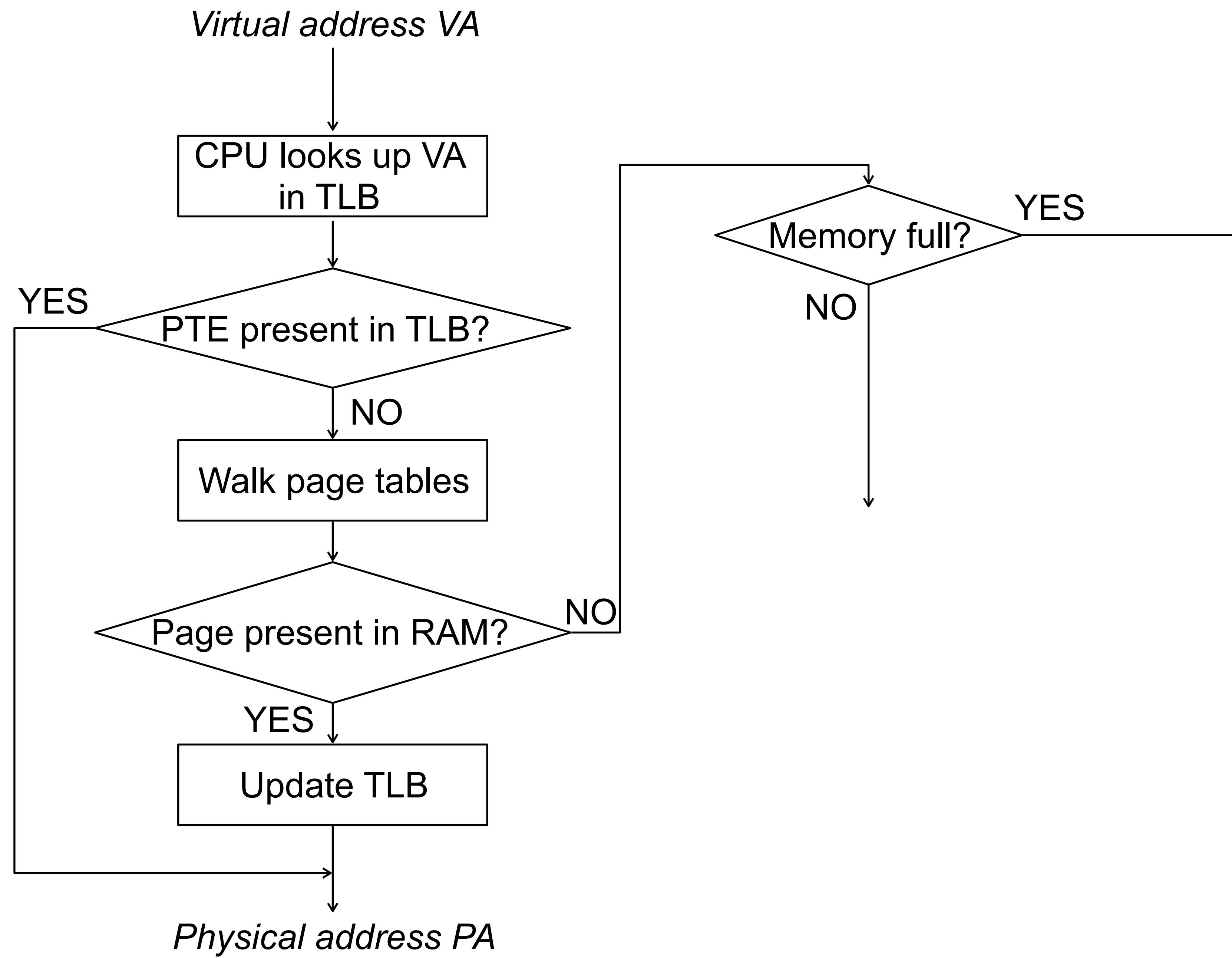


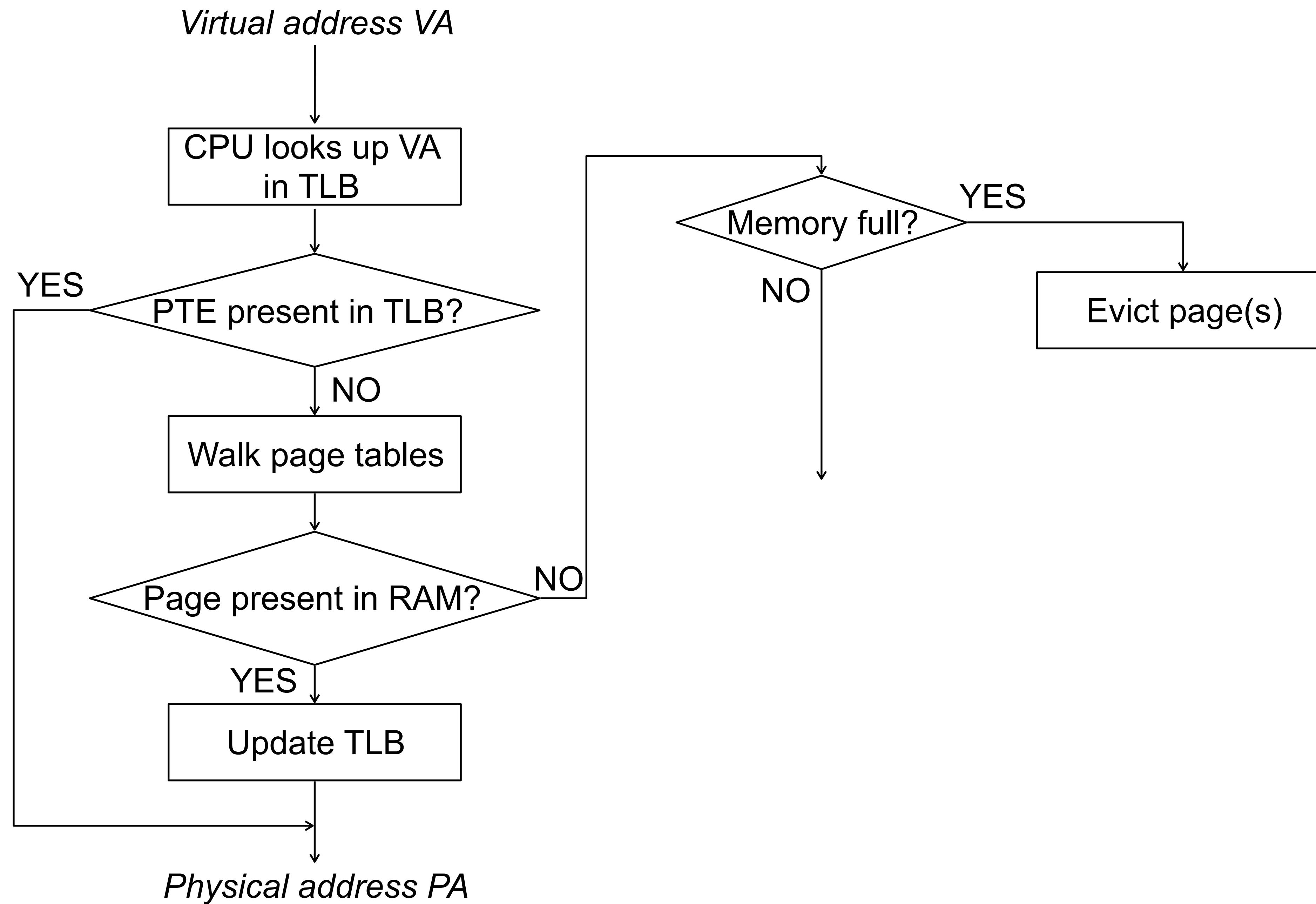


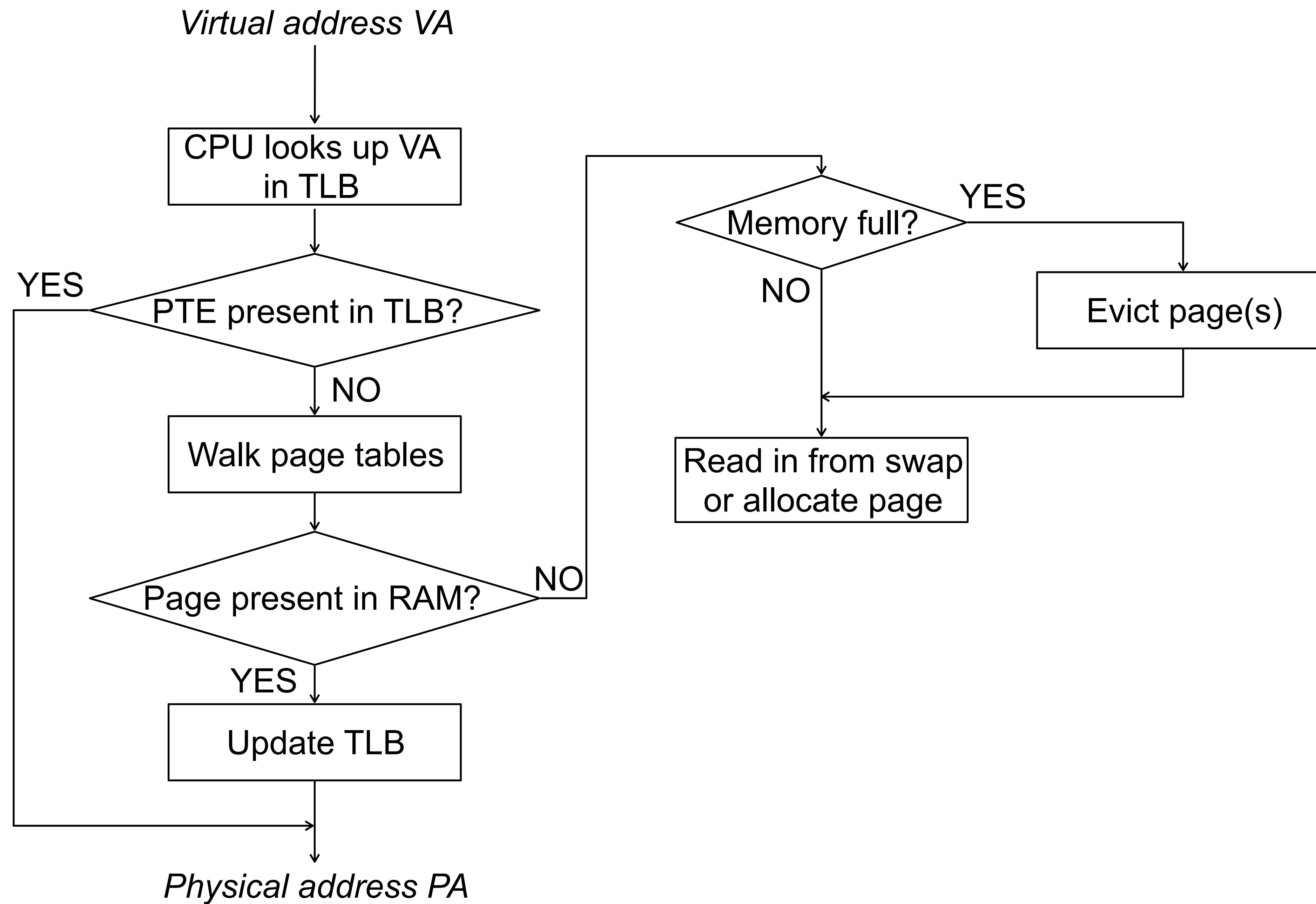


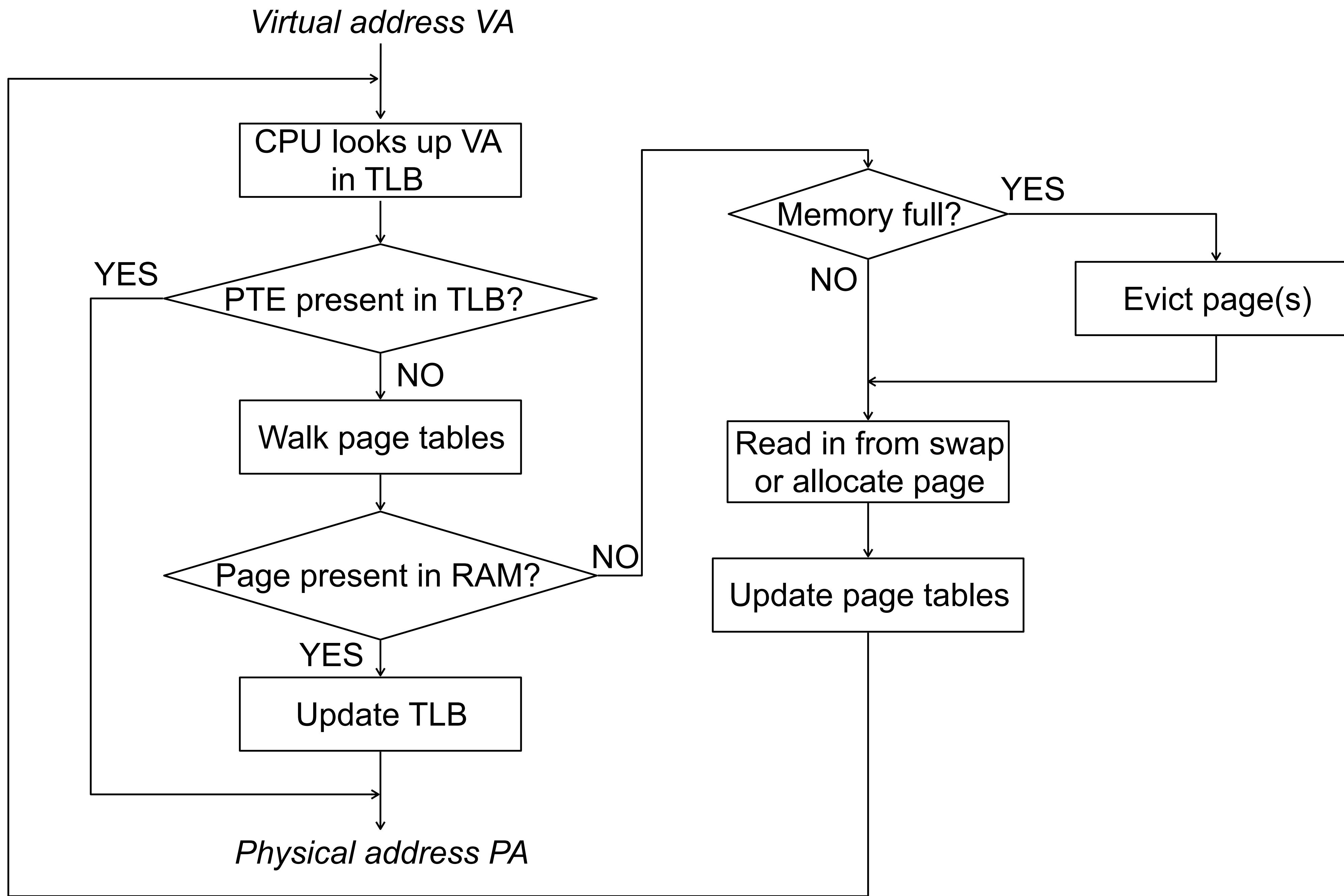


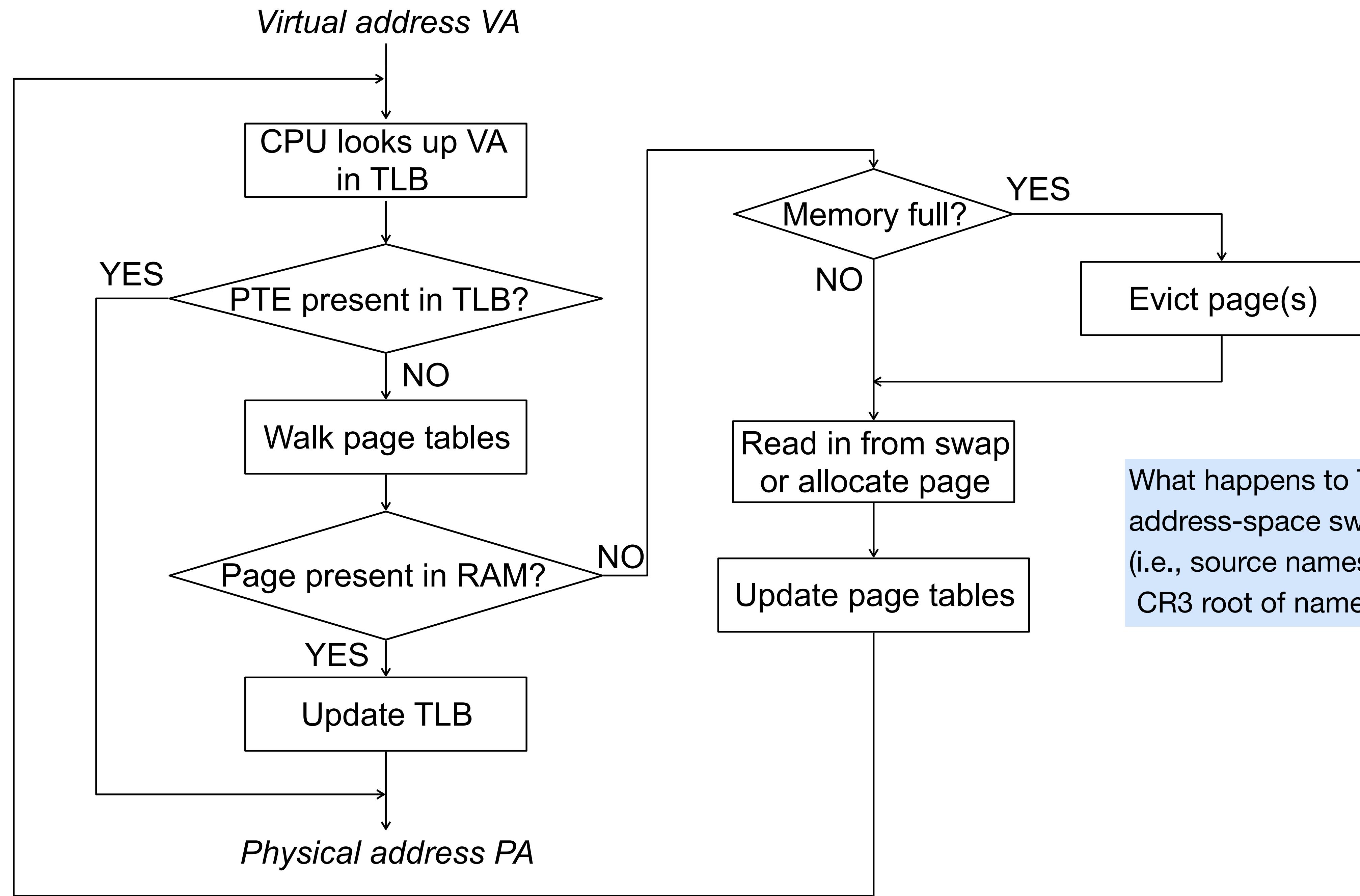




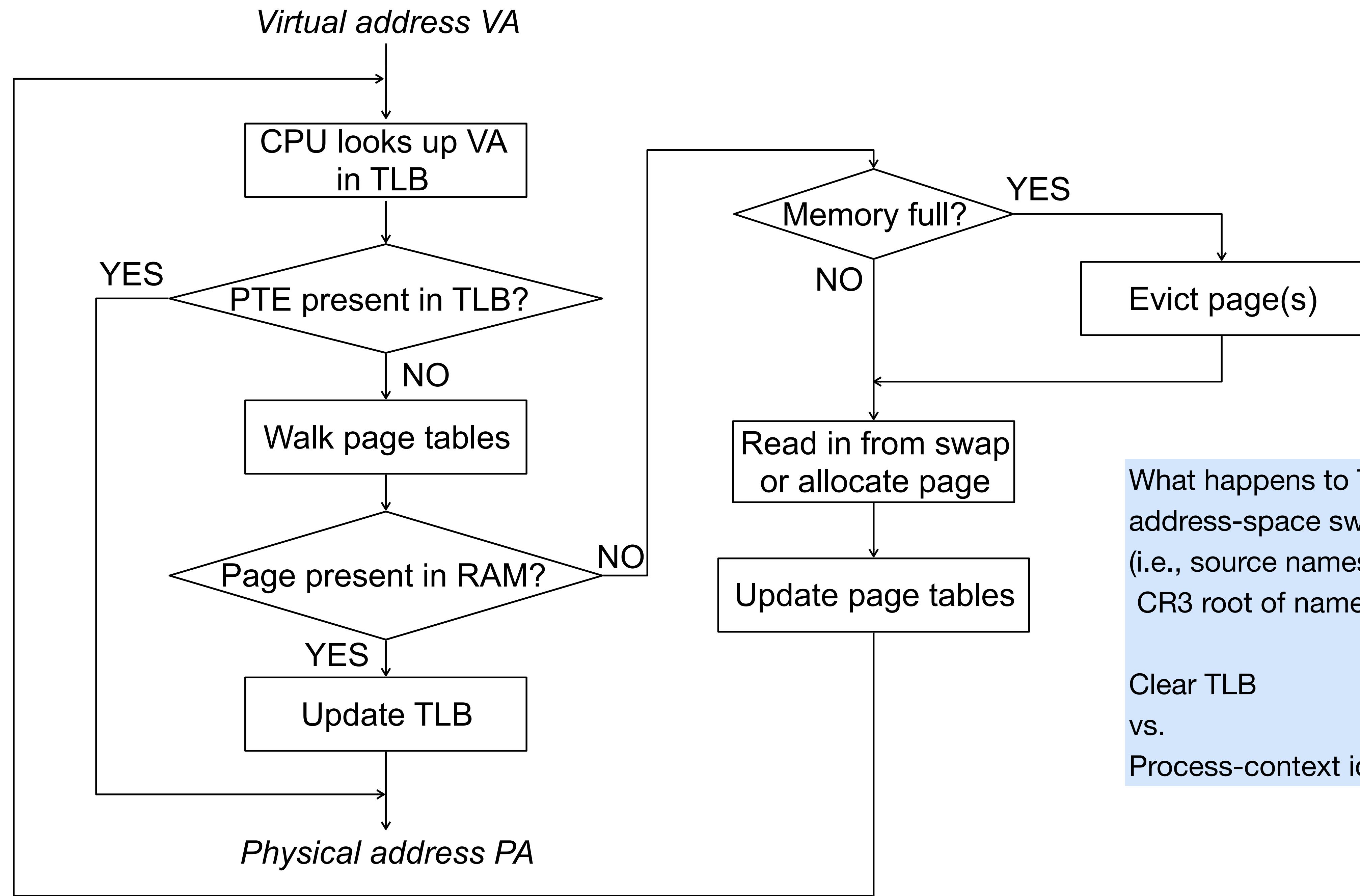






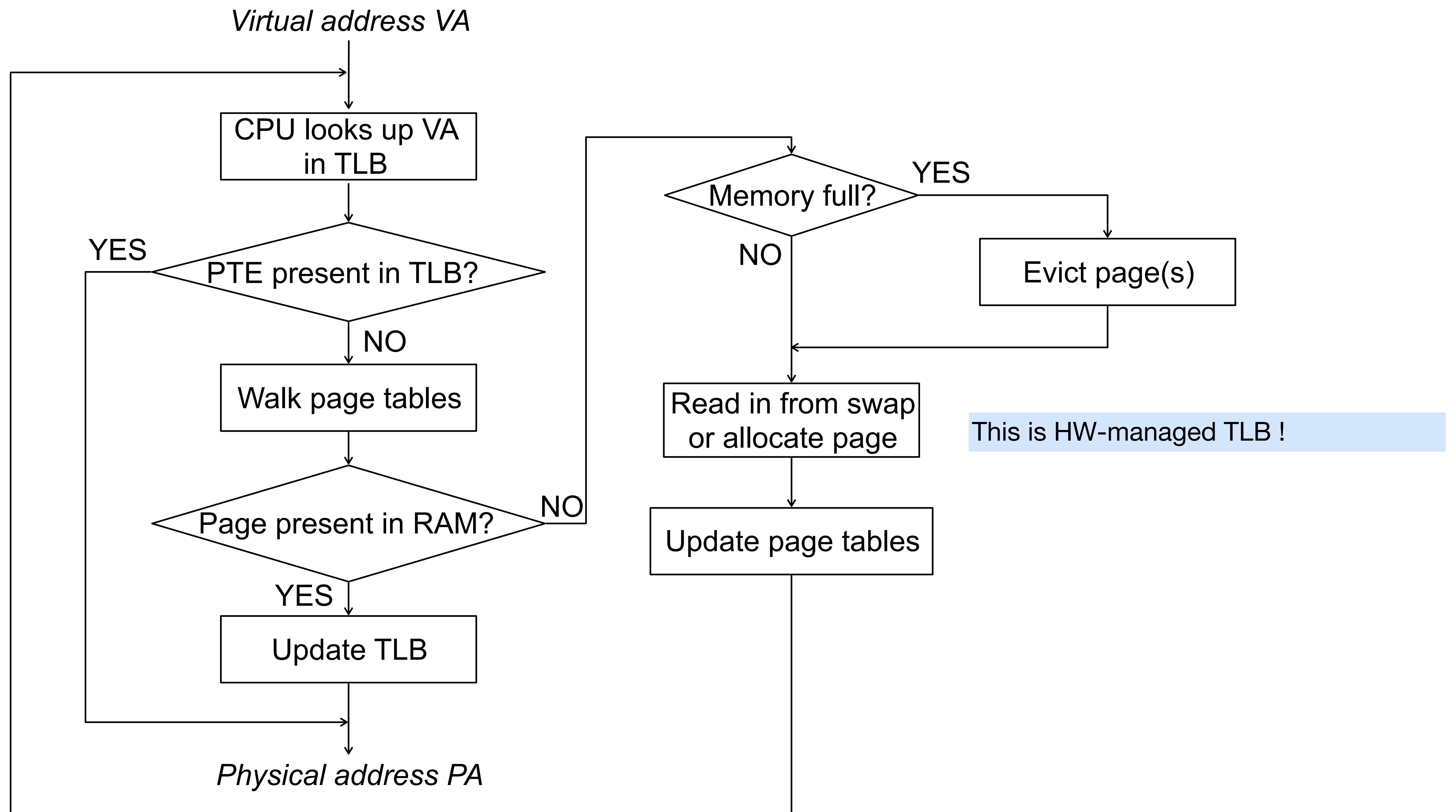


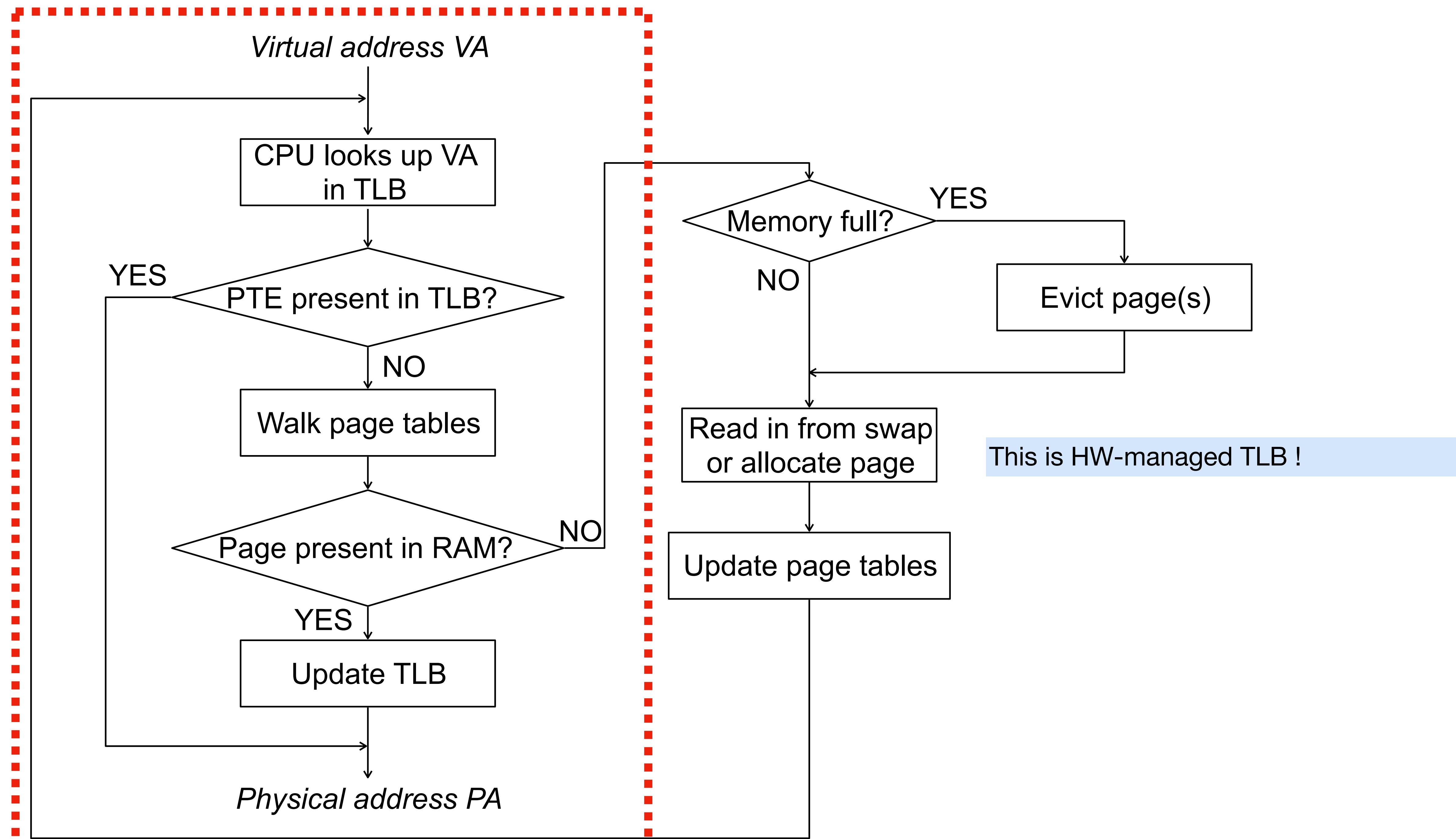
What happens to TLB on an address-space switch?  
 (i.e., source namespace differs =>  
 CR3 root of namespace changes)

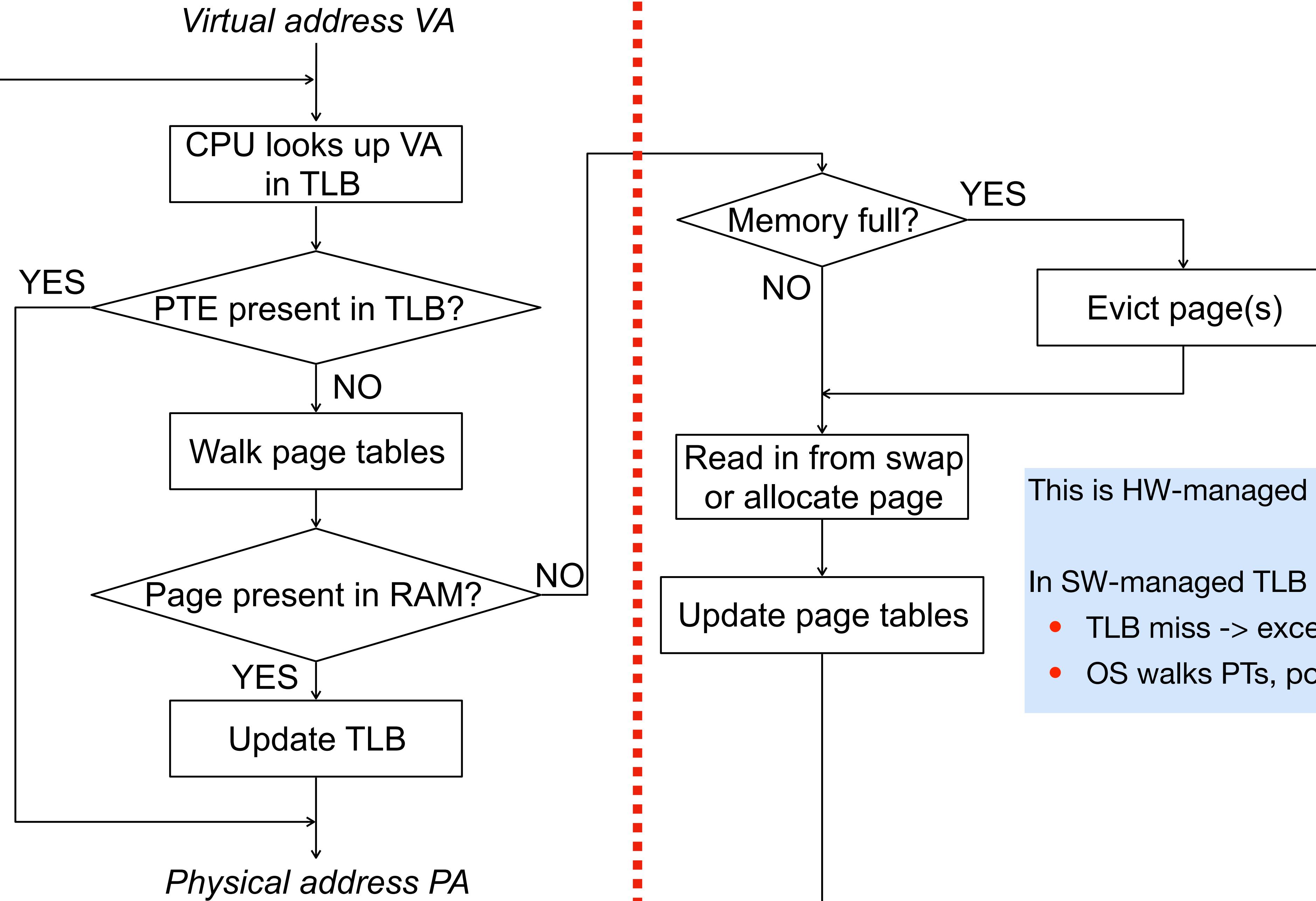


What happens to TLB on an address-space switch?  
 (i.e., source namespace differs =>  
 CR3 root of namespace changes)

Clear TLB  
 vs.  
 Process-context identifier (PCID)



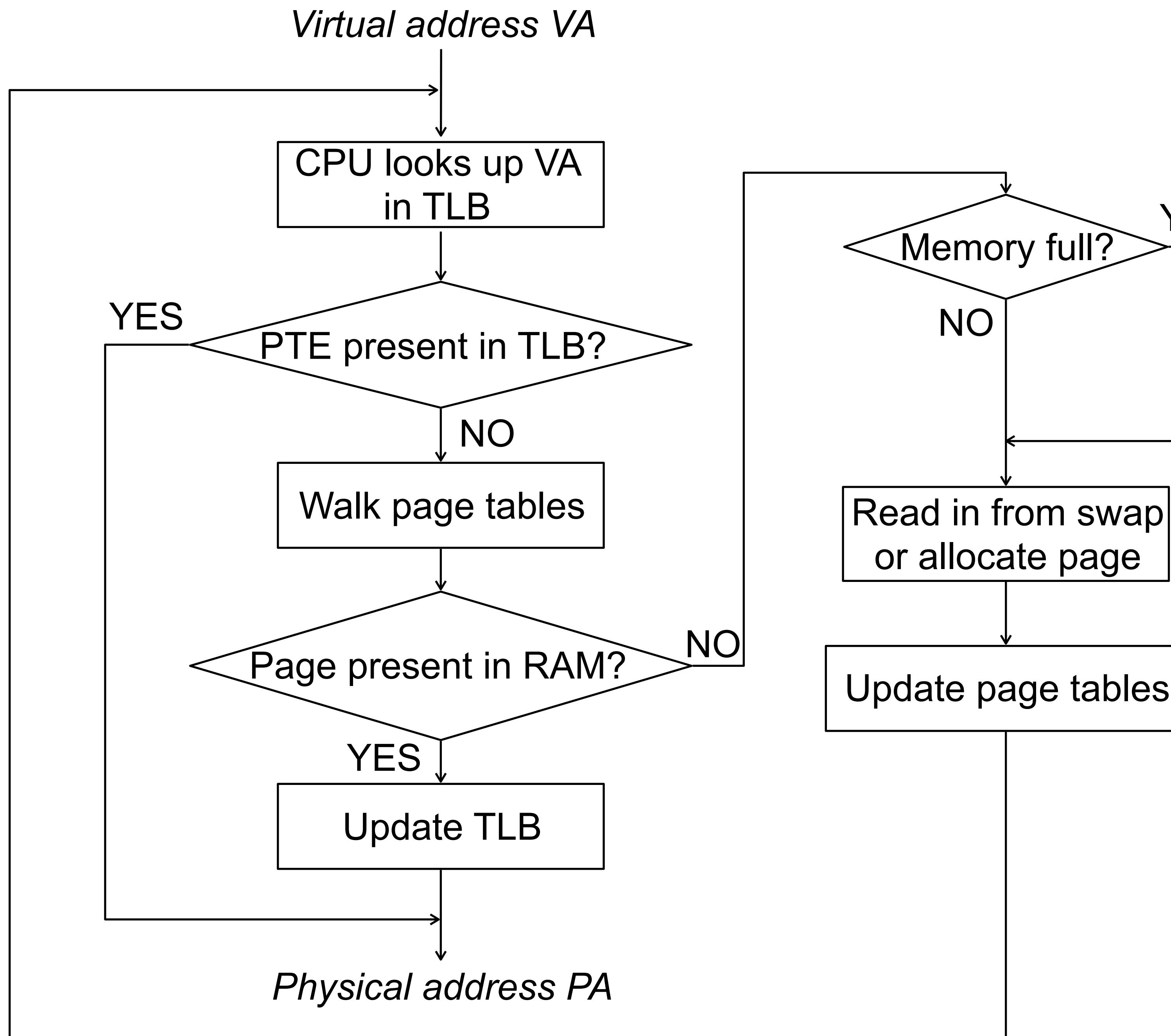




This is HW-managed TLB !

In SW-managed TLB (MIPS, SPARC, ...)

- TLB miss -> exception to OS
- OS walks PTs, populates TLB



This is HW-managed TLB !

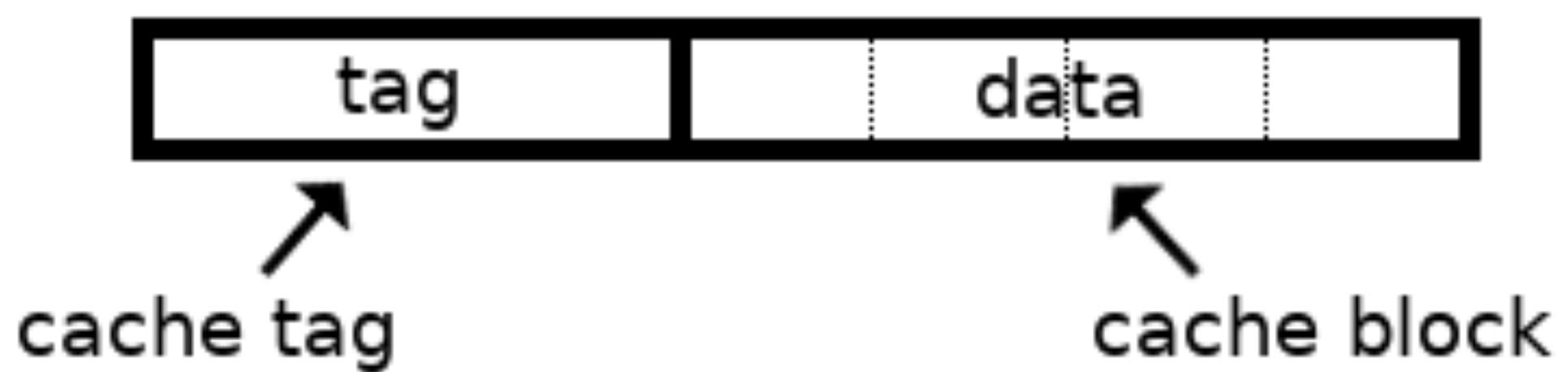
In SW-managed TLB (MIPS, SPARC, ...)

- 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!

# Cache structure

## Cache entry (line)

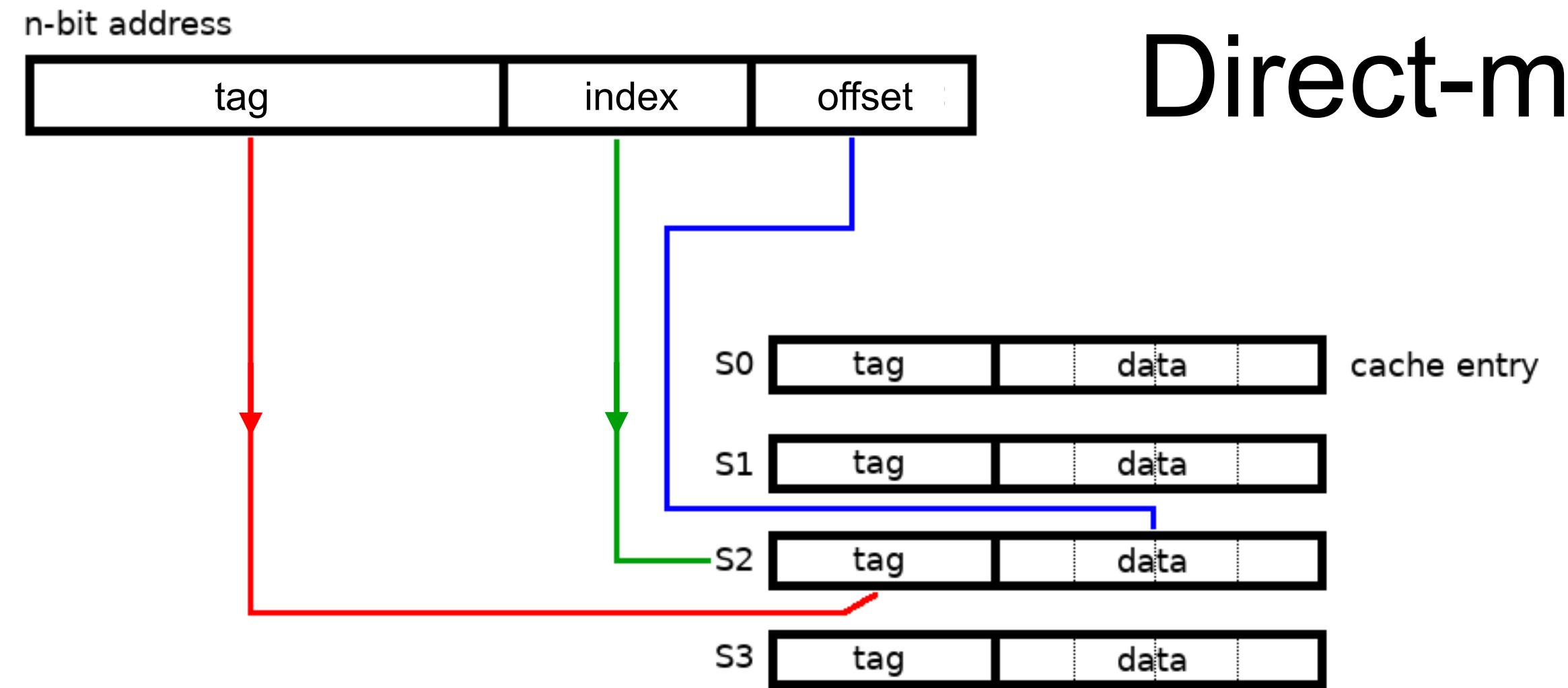


## Address (cache lookup key)

n-bit address

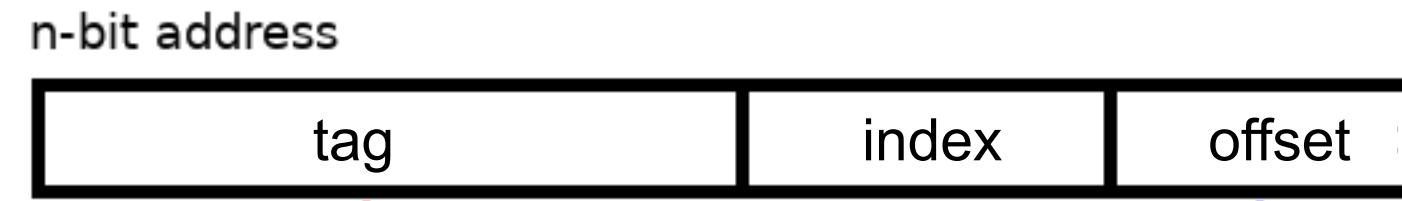


# Cache structure

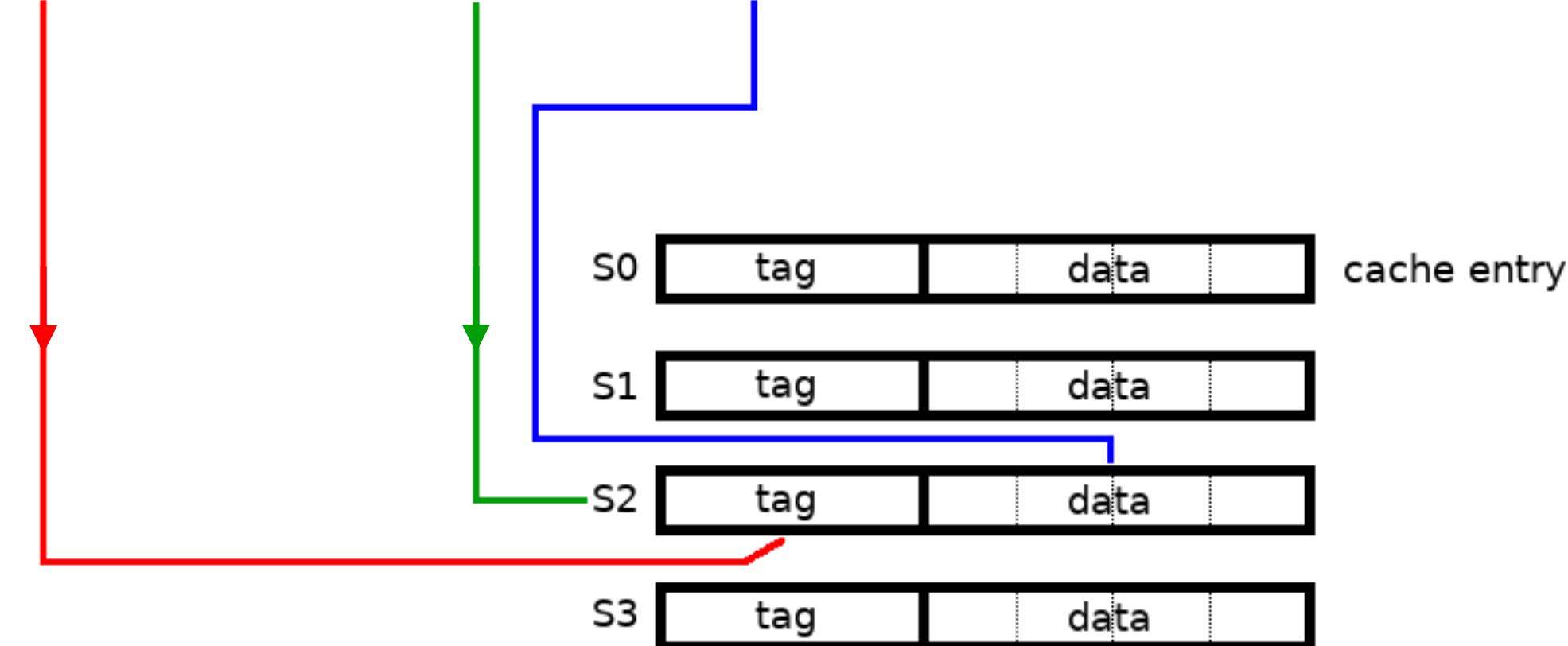


## Direct-mapped cache

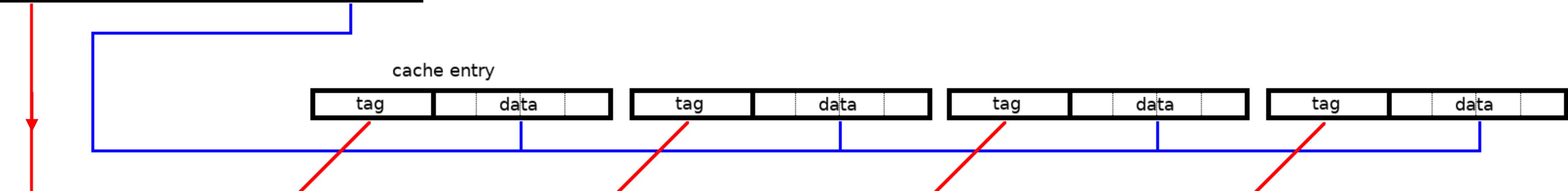
# Cache structure



## Direct-mapped cache



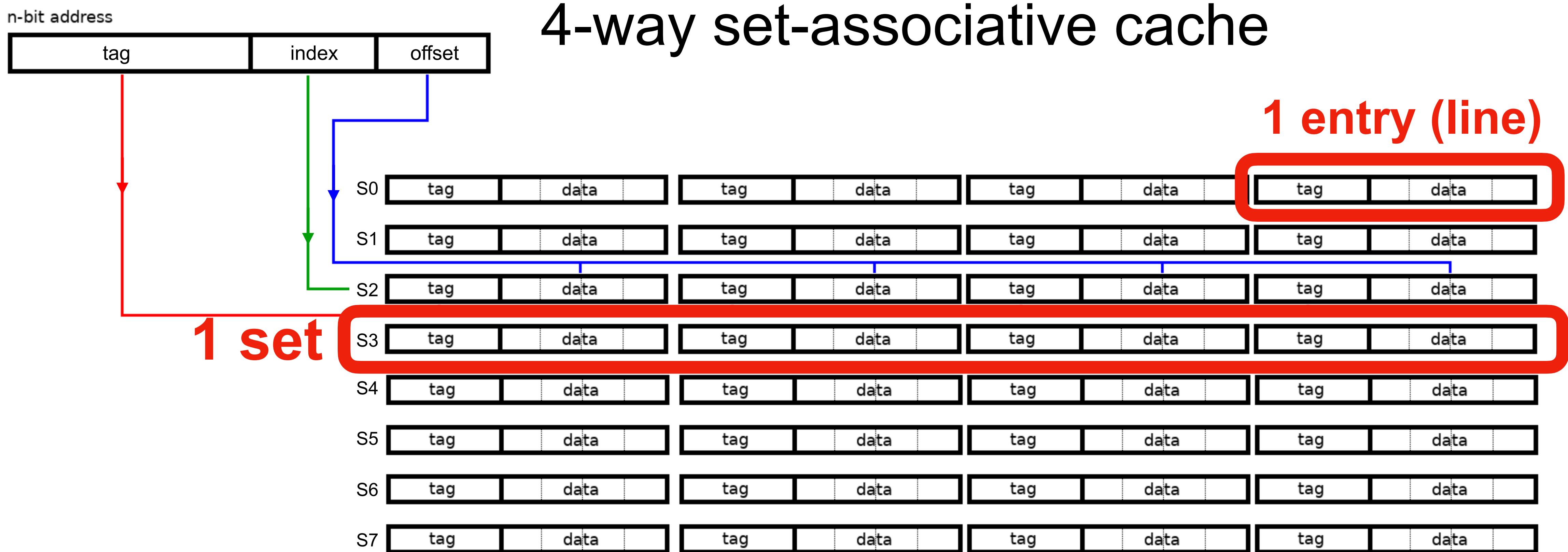
## Fully associative cache



# Cache structure

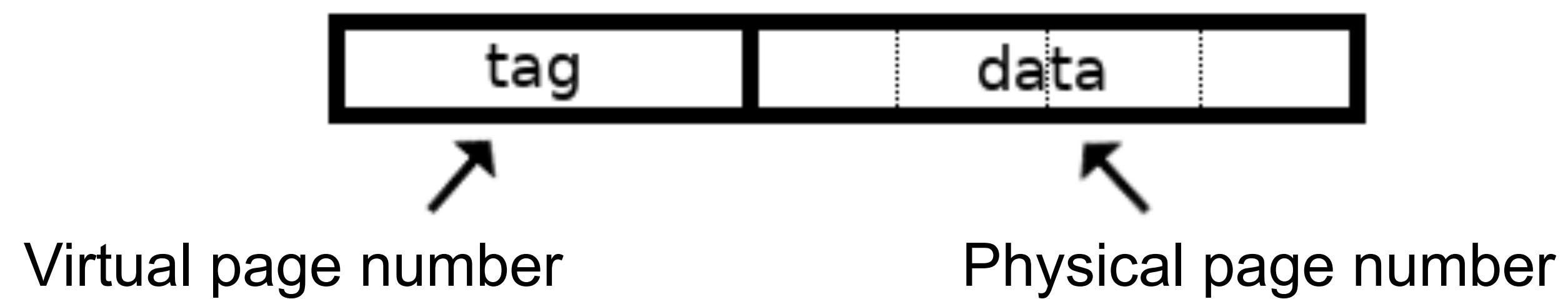


# Cache structure



# Overhead of memory-location name resolution

TLB cache line



# Overhead of memory-location name resolution

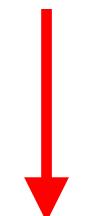
## TLB cache line



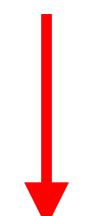
- Skylake (4 KiB pages)
  - $L1 i\text{-}TLB: 16 \text{ sets} * 8\text{-way set-associative} = 128 \text{ entries}$
  - $L1 d\text{-}TLB: 16 \text{ sets} * 4\text{-way set associative} = 64 \text{ entries}$
  - $L2 \text{ shared } TLB: 256 \text{ sets} * 12\text{-way set associative} = 1536 \text{ entries}$

# PIPT L1 d/i-cache ?

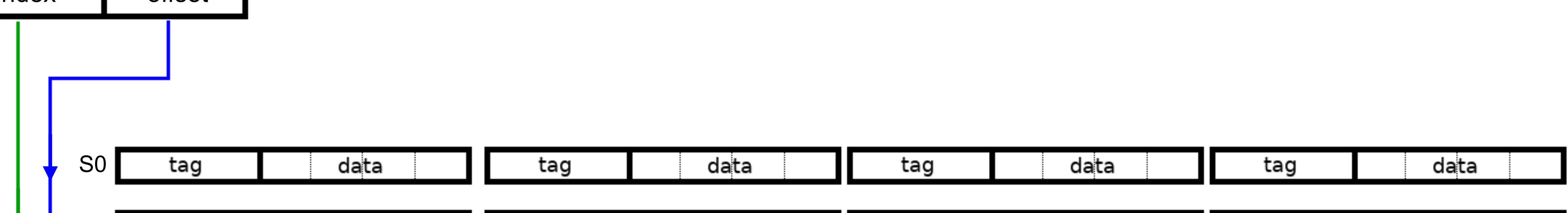
Virtual address



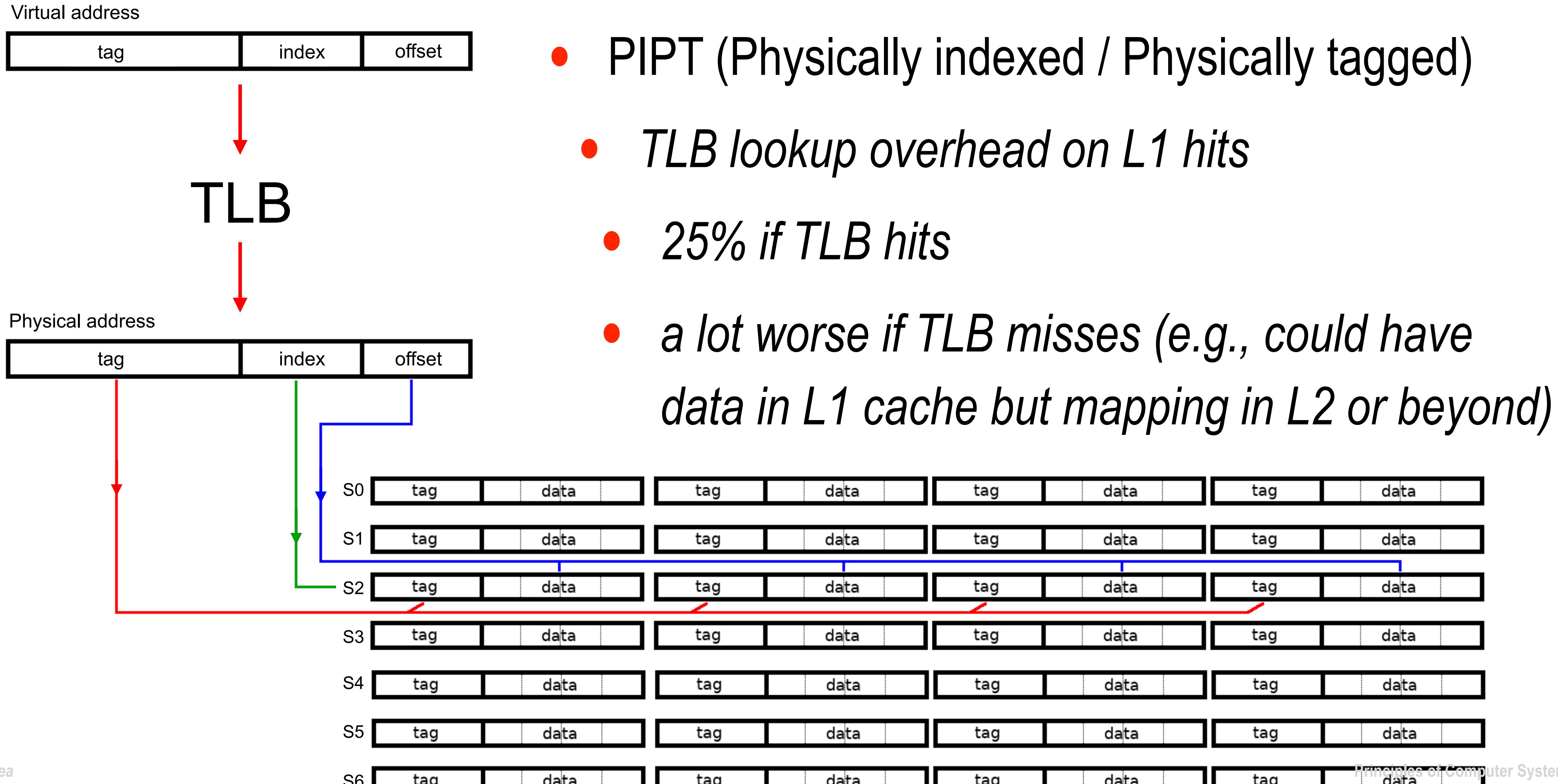
TLB



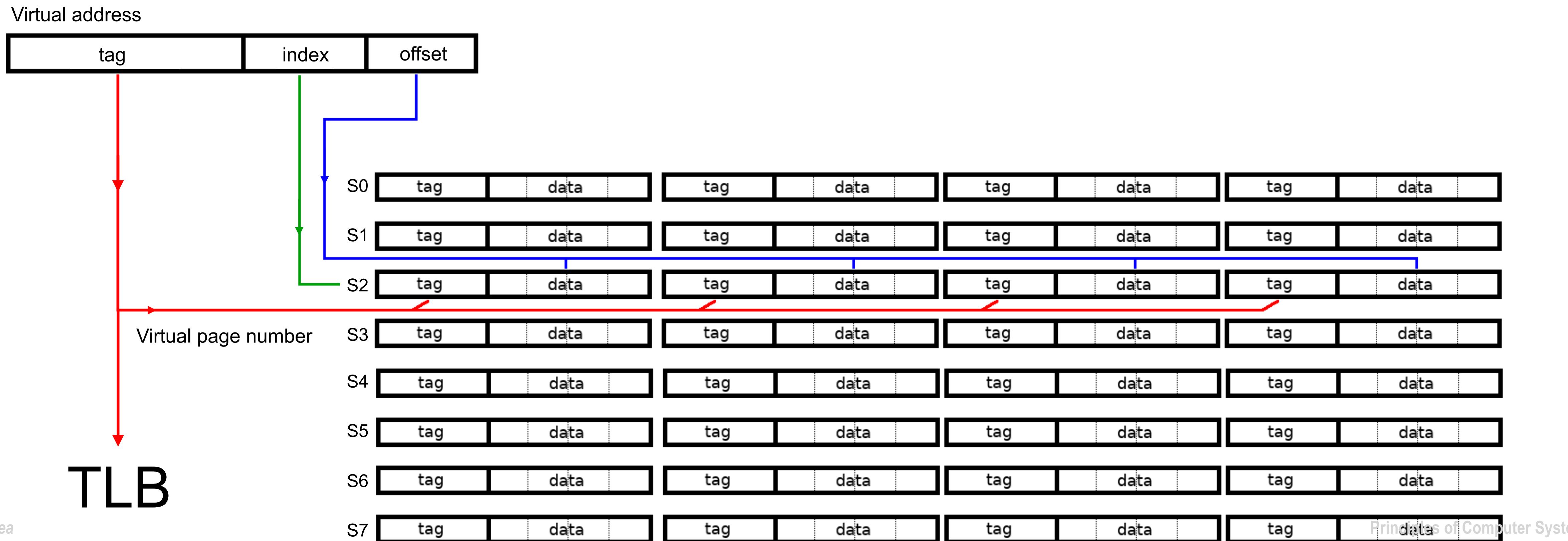
Physical address



# PIPT L1 d/i-cache ?

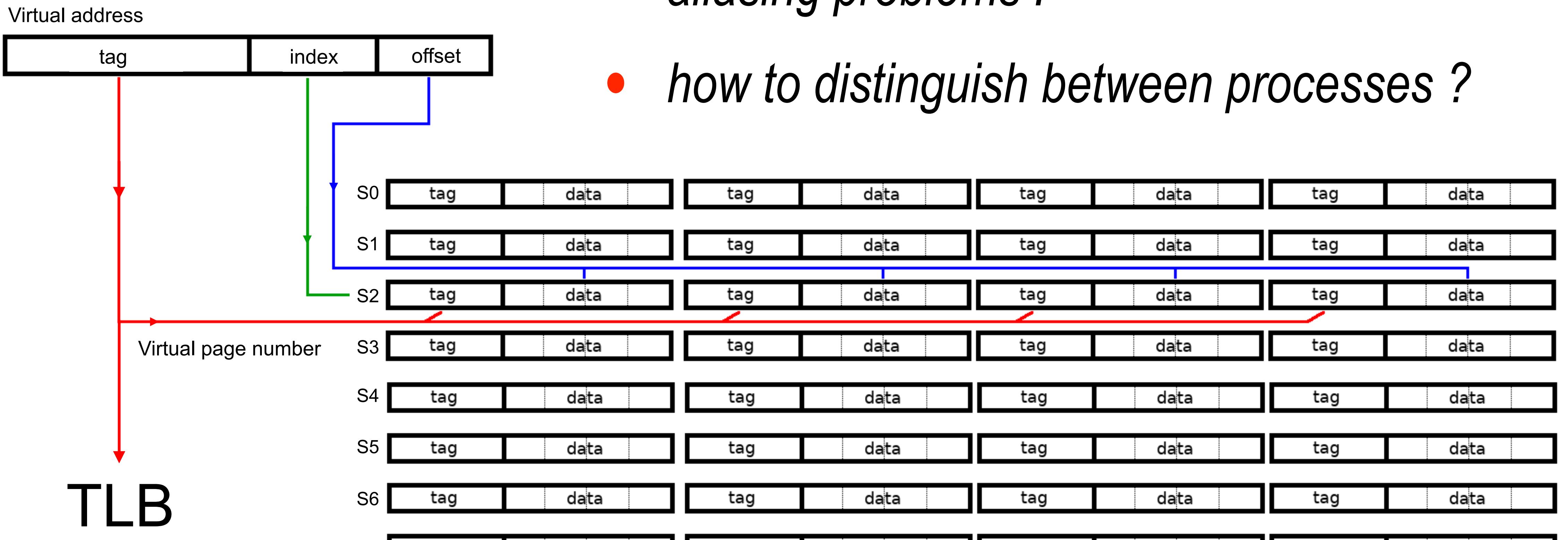


# VIRT L1 d/i-cache ?

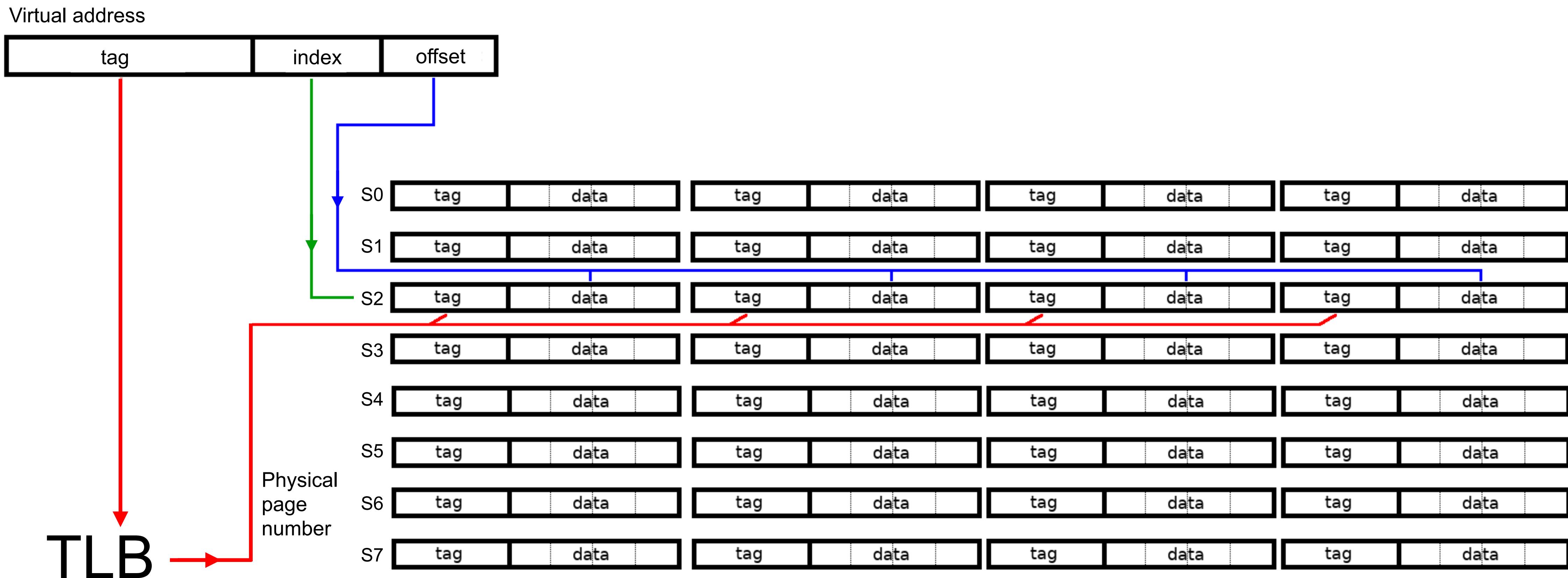


# VIVT L1 d/i-cache ?

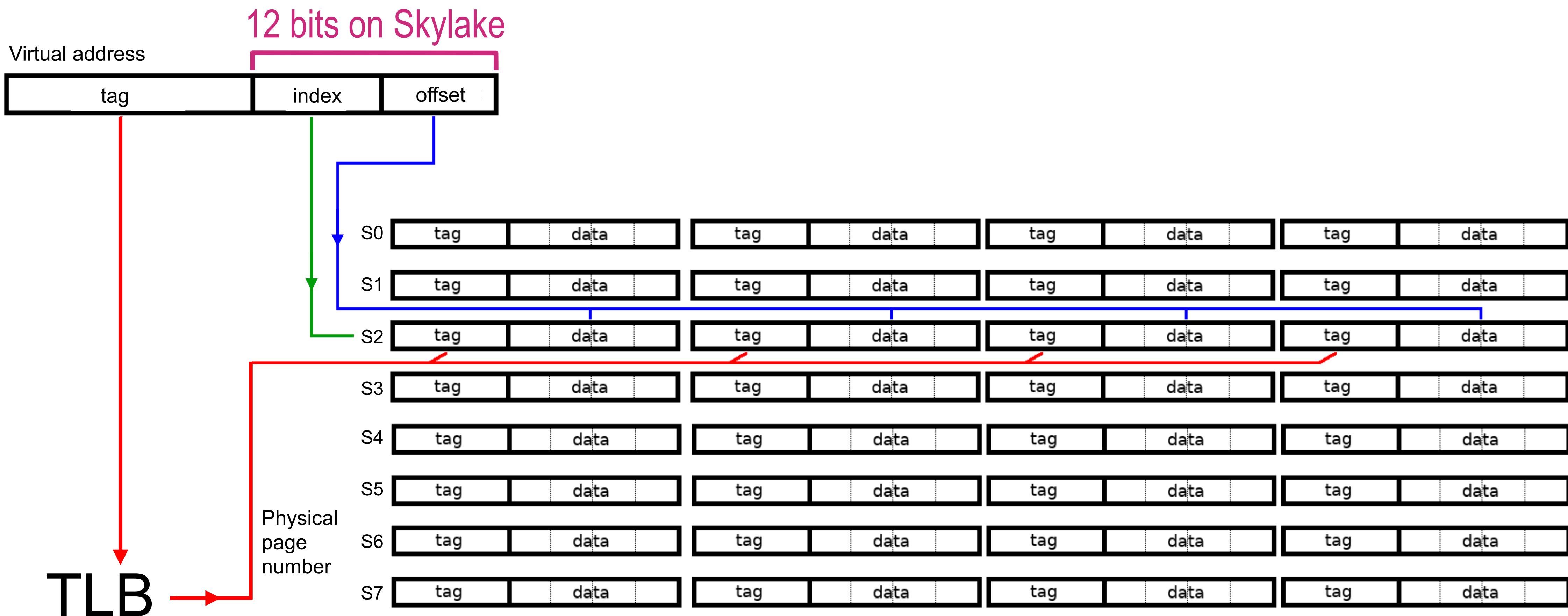
- VIVT (Virtually indexed / Virtually tagged)
  - *TLB lookups become free*
  - *aliasing problems !*
  - *how to distinguish between processes ?*



# VIPT L1 d/i-cache ?

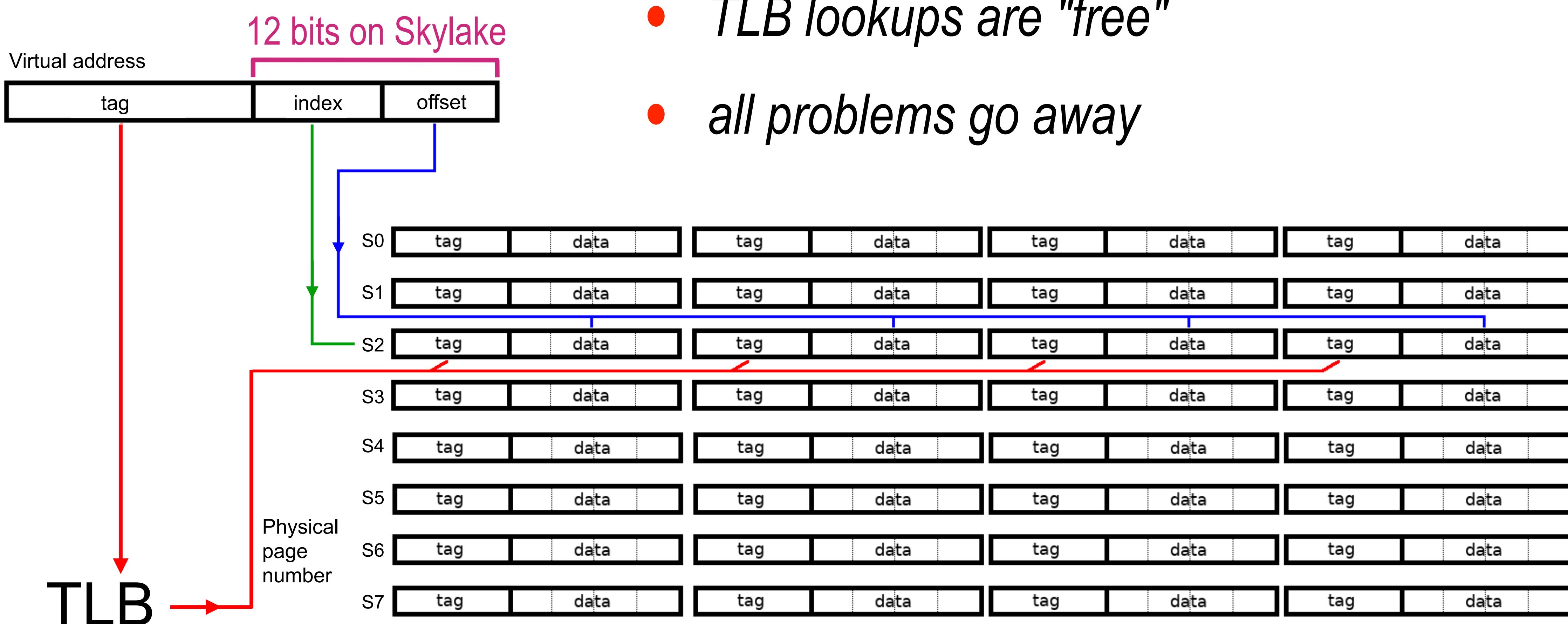


# VIPT L1 d/i-cache ?



# VIPT L1 d/i-cache ?

- VIPT (Virtually indexed / Physically tagged)



# zero-cost VA-to-PA resolution



- Types of main memory caches
  - *PIPT (Physically indexed / Physically tagged)*
    - TLB lookup introduces 25% overhead on L1 hits
  - *VIVT (Virtually indexed / Virtually tagged)*
    - Aliasing problems + how to distinguish between processes?
  - *VIPT (Virtually indexed / Physically tagged)*
    - Parallel lookup in TLB and L1 => mask entirely the TLB's (<1 cycle) as part of L1's 3-4 cycles
  - ~~*PIVT (Physically-indexed / Virtually-tagged)*~~
- Caches on x86
  - *L1 is VIPT*
  - *L2 and L3 are PIPT*

This slide is not about the TLB but about the memory cache (in particular L1, since for L2/L3 the TLB overhead is negligible)

# Summary

---

- MMU maps memory VAs (names) to PAs (values)
- Space scalability
  - *use pages and offsets as a way to avoid looking up every single location*
  - *name space is sparse => use hierarchical, multi-level page tables (directories)*
  - *reduce number name lookups by aggregating locations (huge pages)*
- Performance scalability
  - *cache mappings in TLBs*
  - *mask TLB lookup with VIPT memory caches (in the common case)*
- Make common case fast, make rare case merely correct

# **Reference Values for System Design**

# Reference Values (absolute time)

- L1 / L2 / main RAM reference ~ 1 / 4 / 40-50 ns
- DRAM bandwidth ~200 GB/sec
- Flash SSD random access 50-100 microsec
- Flash SSD bandwidth 5-10 GB/sec
- Seek on enterprise-grade HDD ~4 millisec
- Bandwidth of enterprise-grade HDD ~150 MB/sec
- Mutex lock or unlock ~ 15-100 ns (*depending on whether share cache line, are contended, are cache-aligned ...*)
- Packet RTT within a datacenter < 10 microsec (ignoring software stacks)
- Bandwidth between two hosts in a data center 10-400 Gbps
- Compress 1 KB of data (with Snappy/LZ77) ~ 2 microsec
- Packet roundtrip CA -> Amsterdam -> CA = 150 millisec (*note: speed of light in fiber ~5 msec/1,000 km*)

# Reference Values (CPU cycles)

- Register-register ADD/OR/etc. < 1 CPU cycle (*in an OO, pipelined processor*)
- Register write ~1 cycle
- Correctly / incorrectly predicted "if" branch = 1 cycle / 10-20 cycles
- L1 / L2 / L3 / main RAM read = ~4 cycles / ~20 cycles / ~30-50 cycles / ~100-150 cycles
- TLB hit = 0.5 - 1 cycle
- CAS = 15-30 cycles (increases with the number of cores !)
- C function direct / indirect call = 10-20 cycles / 20-50 cycles
- Kernel crossing round-trip = 1,000 cycles
- Thread context switch (direct costs) = 2,000 cycles
- On NUMA, different-socket mem hierarchy access is 3 - 10x that of non-NUMA