



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

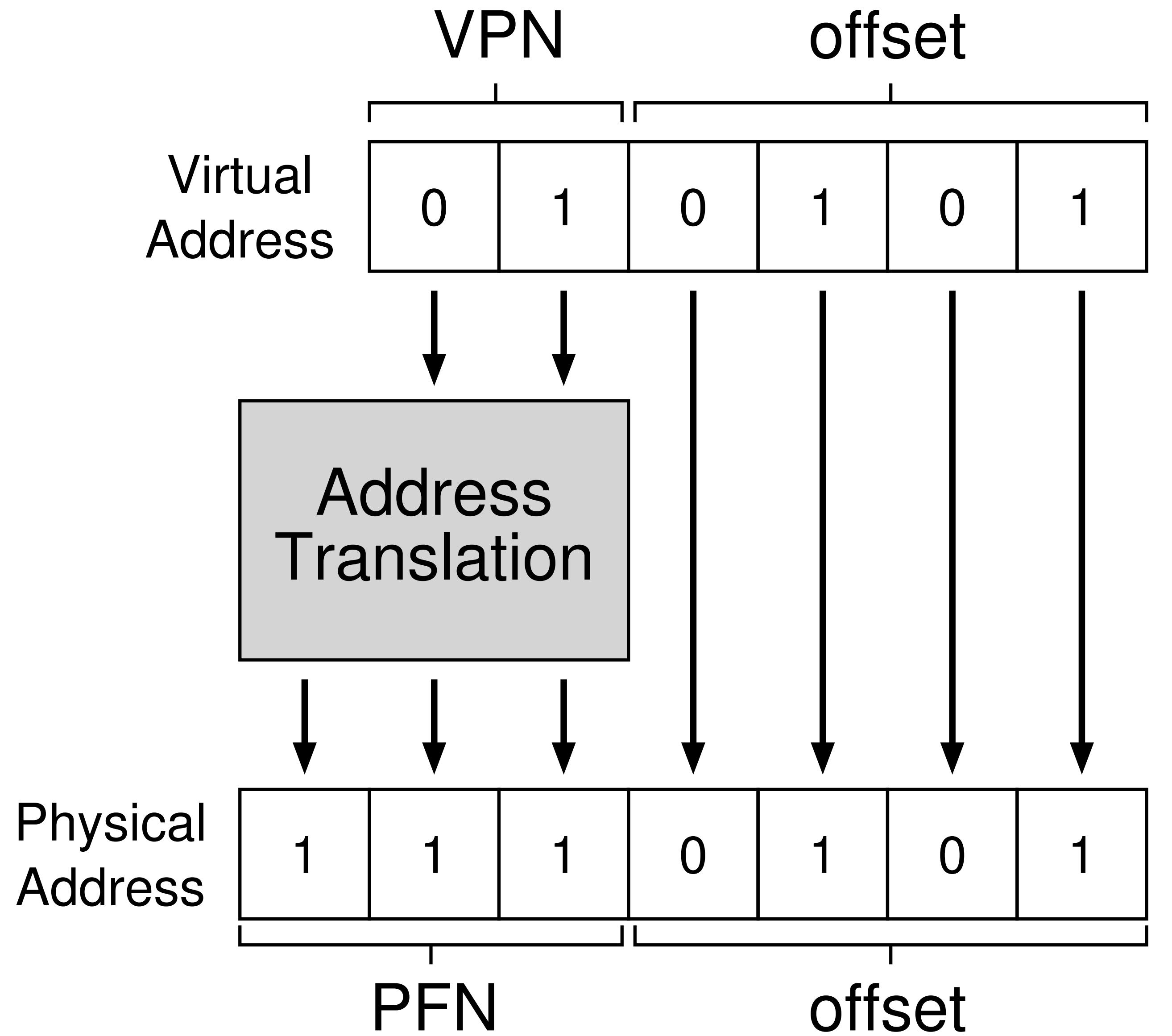
- Memory address
 - *names a memory location*
- Pointer
 - `void*` names the location of a memory word
 - `int*` names the location of an integer
- Virtual memory address
 - *names (from userspace) a physical memory address*

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)
- 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

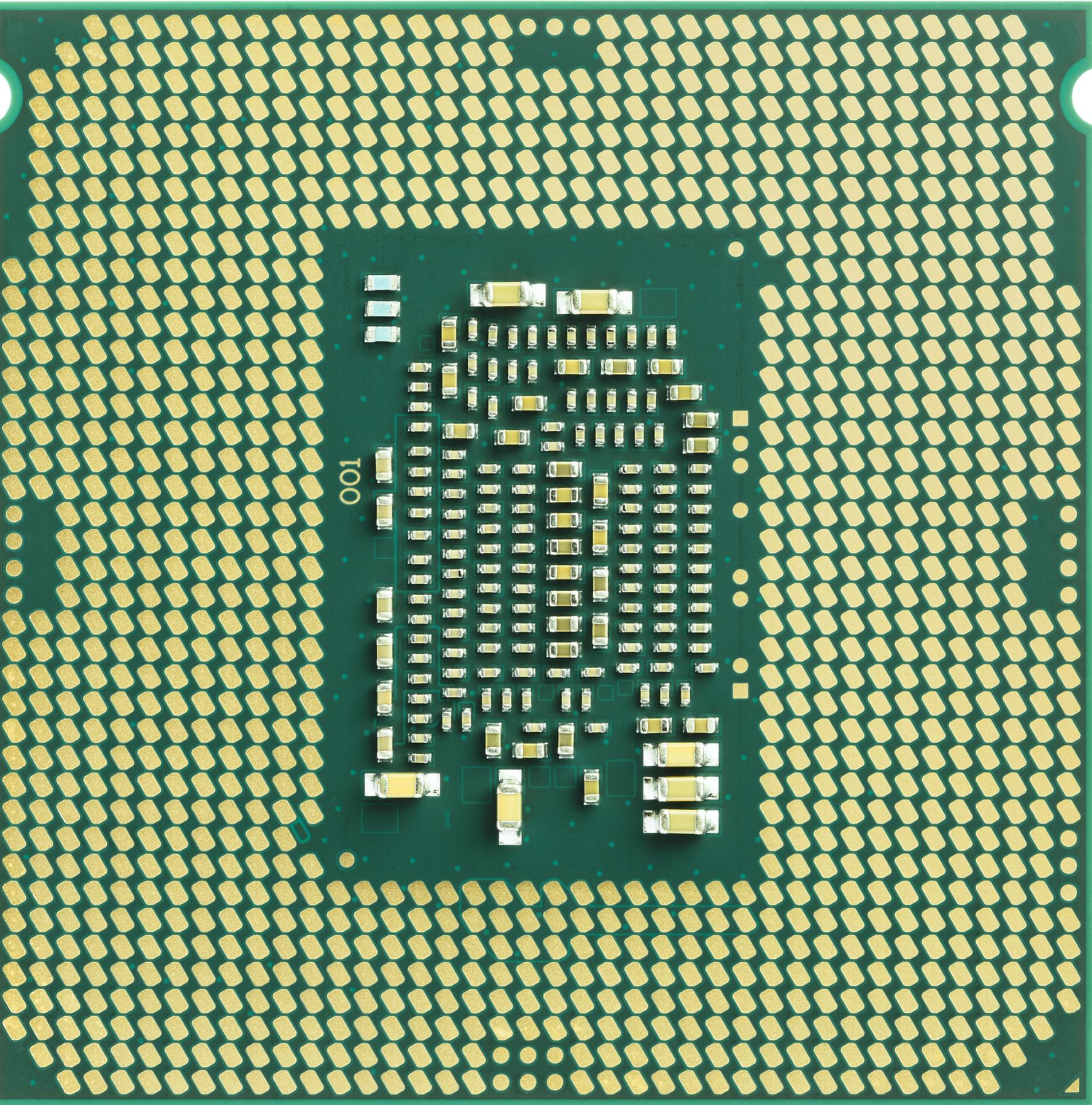


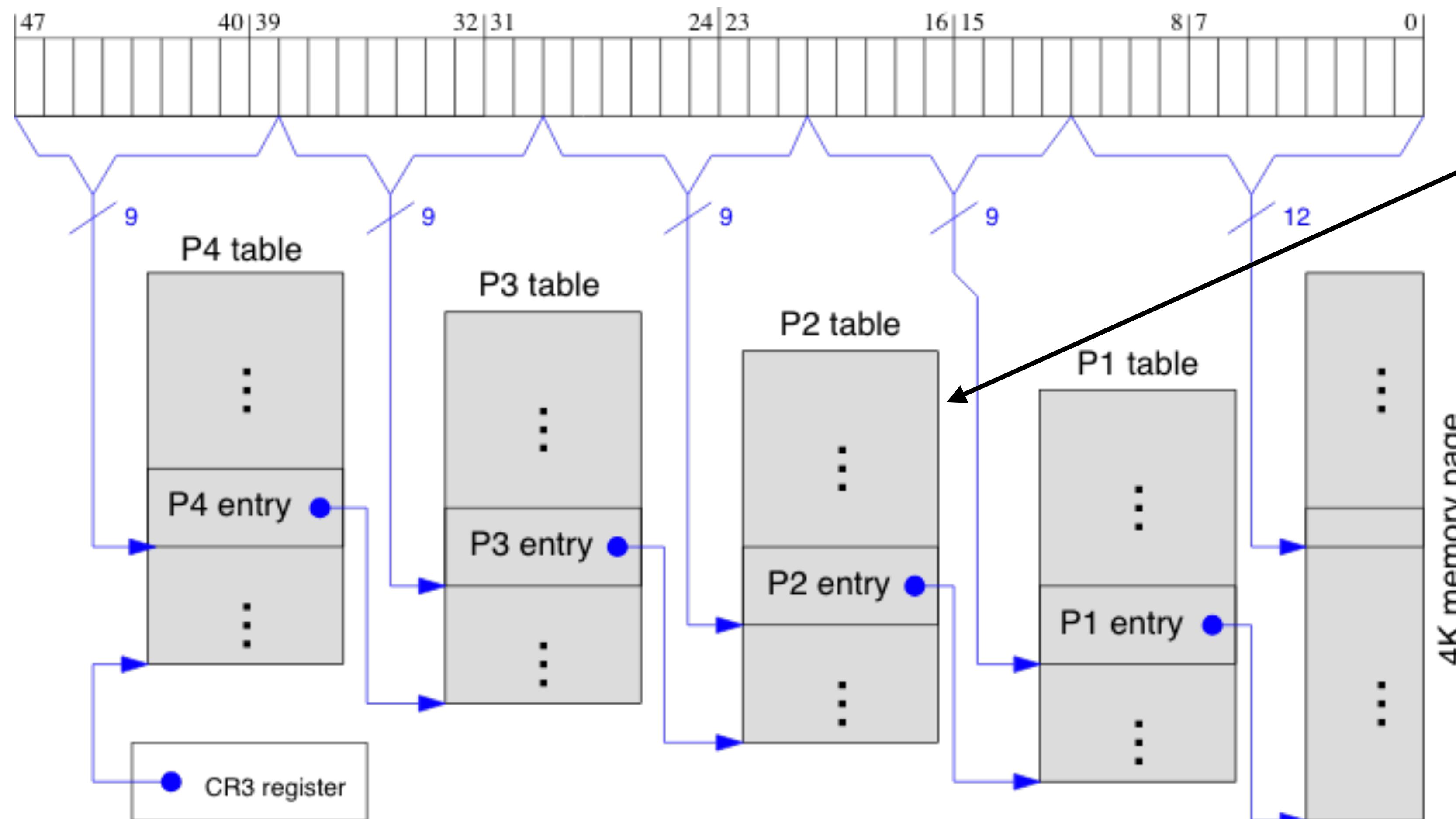
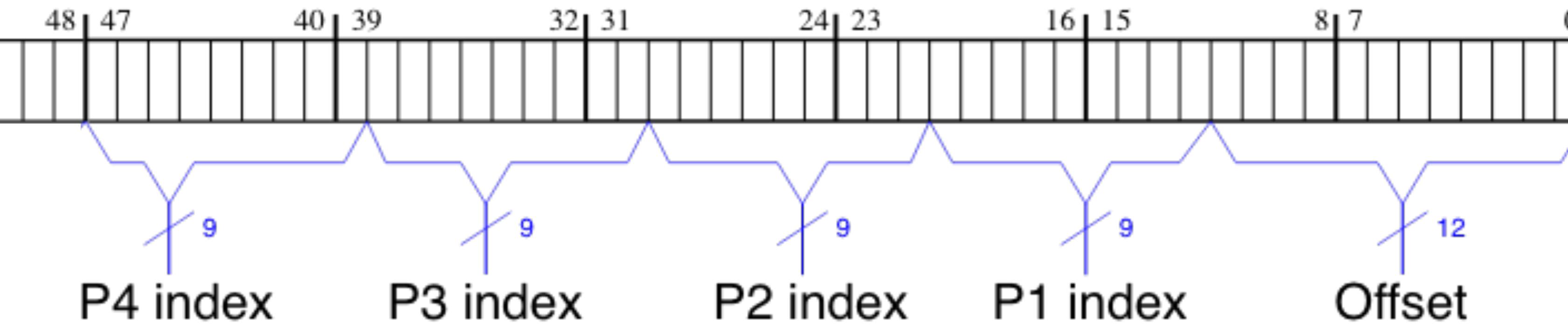
Virtual namespace

Name translation controls visibility

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

Physical namespace





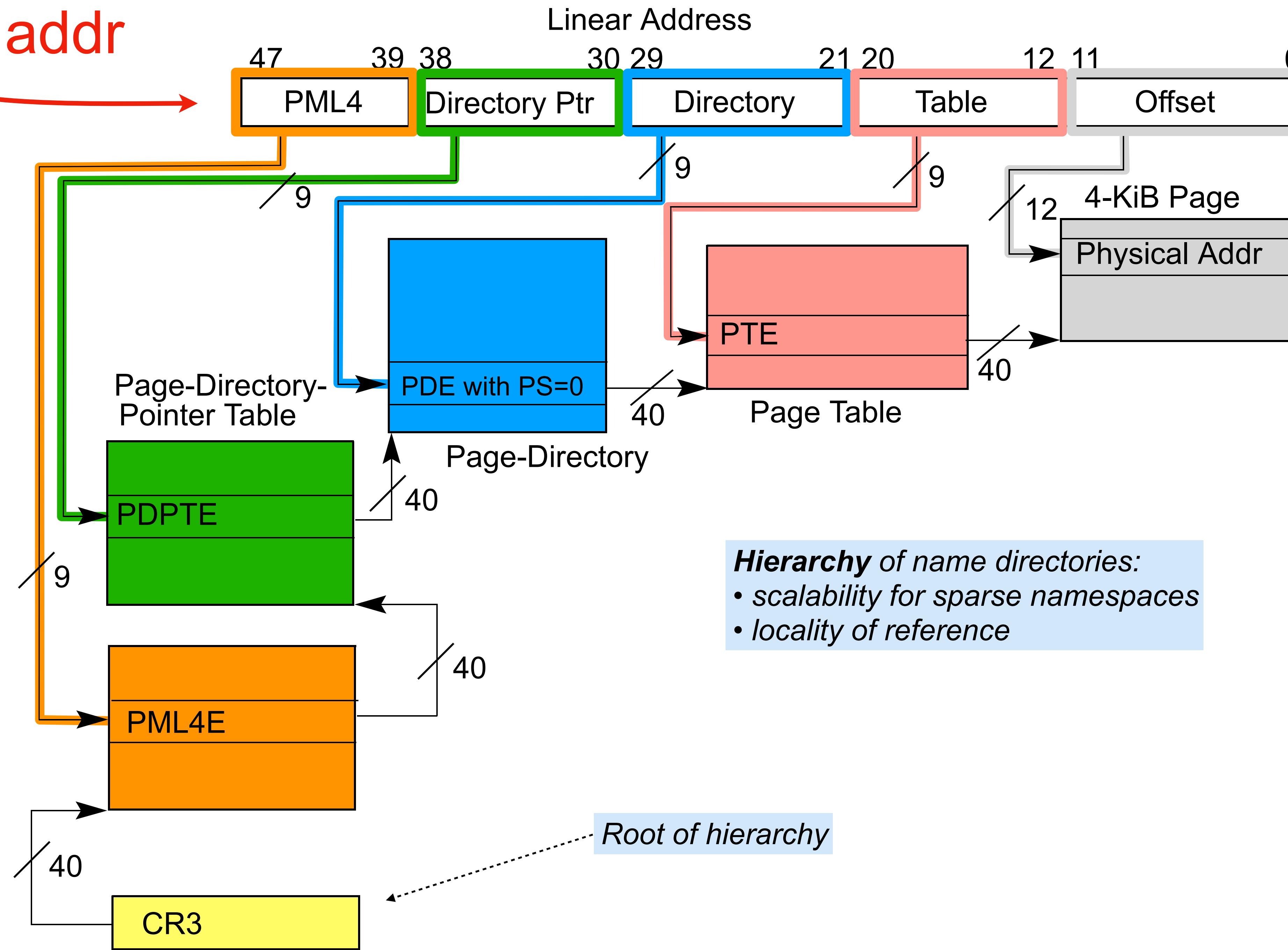
$2^9 = 512$ entries
8 bytes / entry

Hierarchy of name directories:

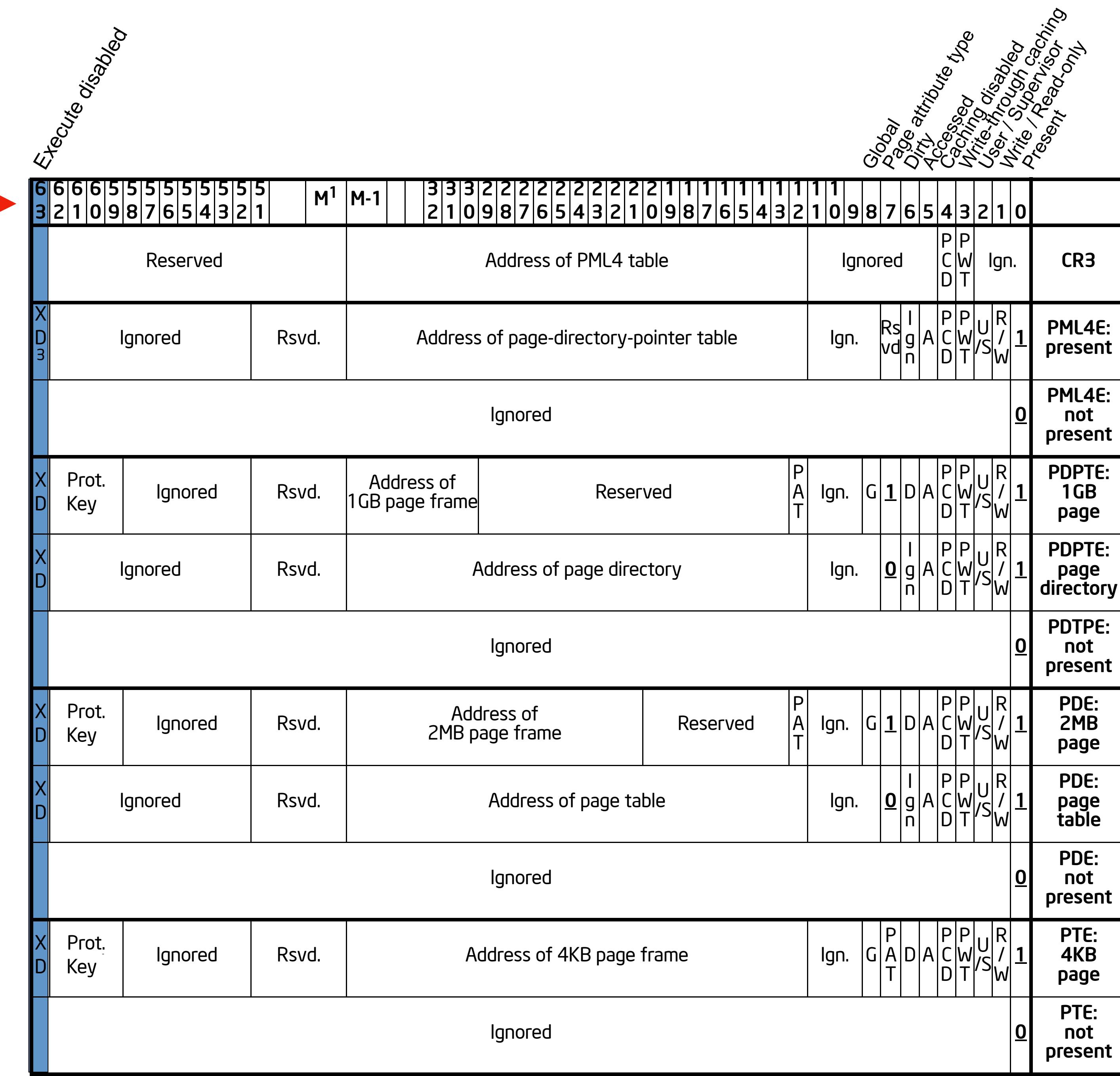
- scalability for sparse namespaces
- locality of reference

HW walks page tables
OS updates them

Virtual addr



VA bit index



CR3 (root)

Top-level page table



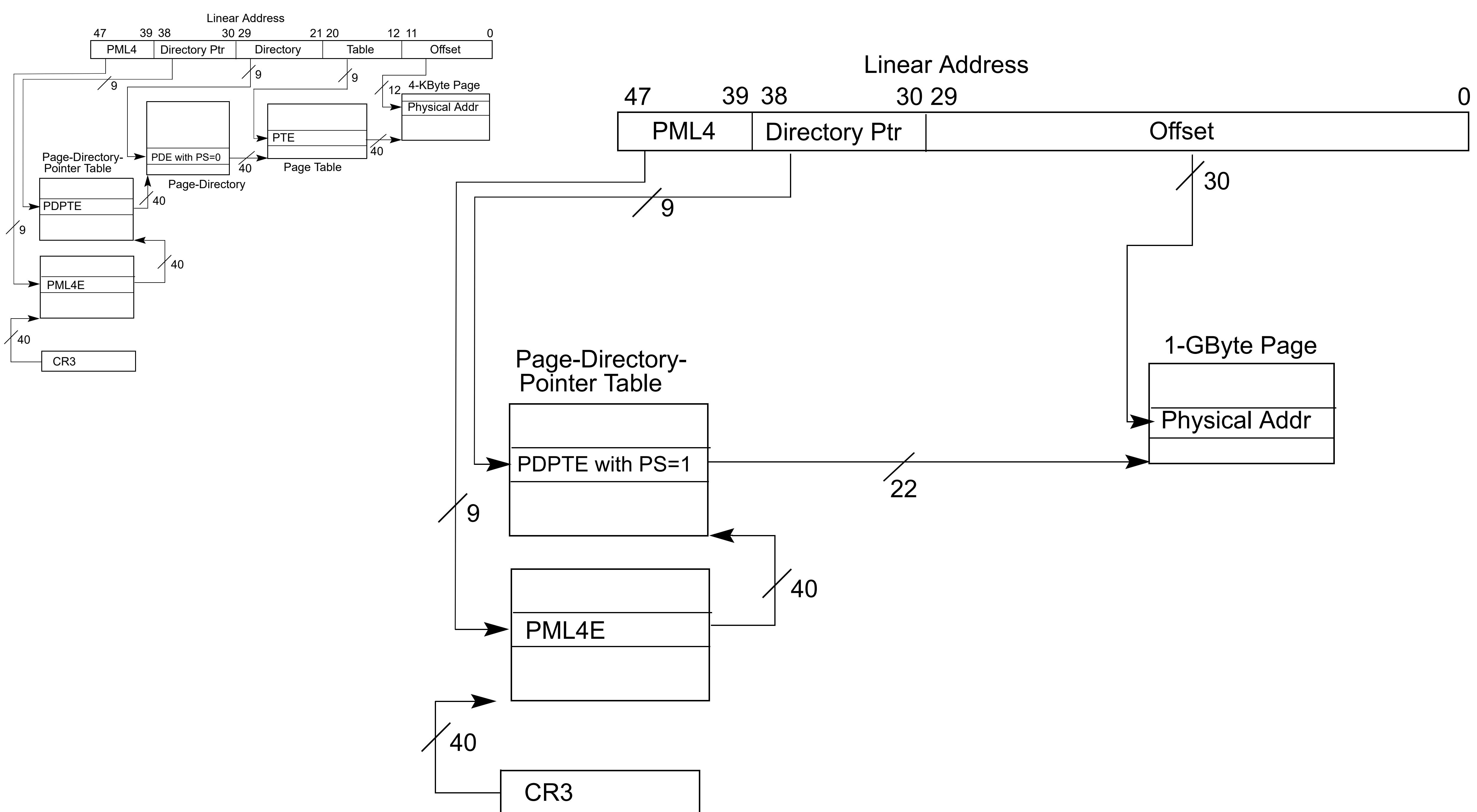
2nd-level page table

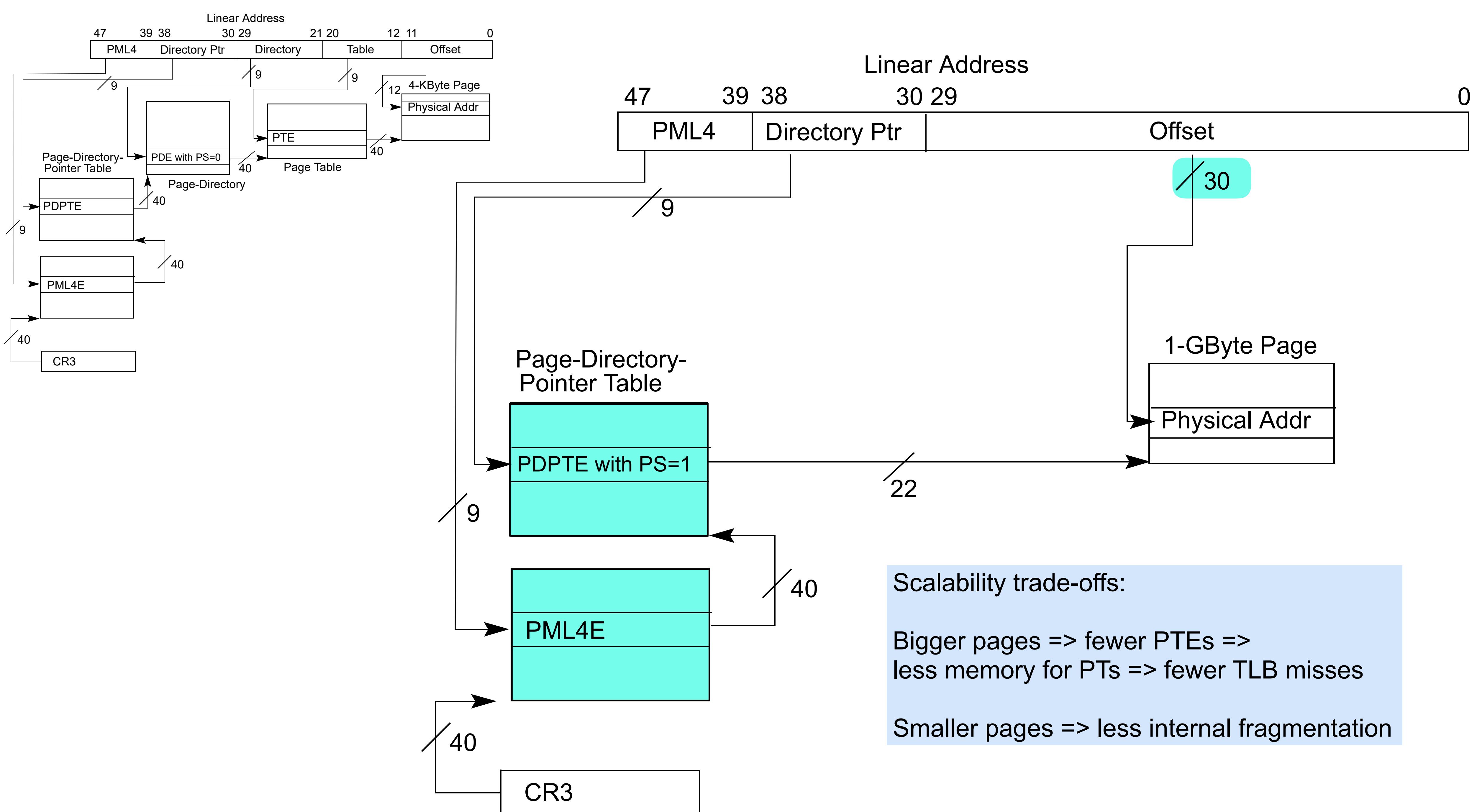


3rd-level page table

Last-level page table

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





Execute disabled										Global attribute type														
										Dirty Accessed Caching disabled Write-through caching User / Supervisor Present														
66665555555555 3210987654321										M ¹ M-1 3332222222211111111111 2109876543210987654321098765432109876543210														
Reserved ²				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 W D	P C W T	U S /S	R /W 1	PML4E: present						
Ignored										0														
X D	Prot. Key ⁴	Ignored		Rsvd.	Address of 1GB page frame		Reserved				P A T	Ign.	G 1	D A	P C W D	P C W T	U S /S	R /W 1	PDPTE: 1GB page					
X D	Ignored		Rsvd.	Address of page directory						Ign.	0	I g n	A	P C W D	P C W T	U S /S	R /W 1	PDPTE: page directory						
Ignored										0														
X D	Prot. Key ⁴	Ignored		Rsvd.	Address of 2MB page frame			Reserved		P A T	Ign.	G 1	D A	P C W D	P C W T	U S /S	R /W 1	PDE: 2MB page						
X D	Ignored		Rsvd.	Address of page table						Ign.	0	I g n	A	P C W D	P C W T	U S /S	R /W 1	PDE: page table						
Ignored										0														
X D	Prot. Key ⁴	Ignored		Rsvd.	Address of 4KB page frame						Ign.	G 1	D A	P C W D	P C W T	U S /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

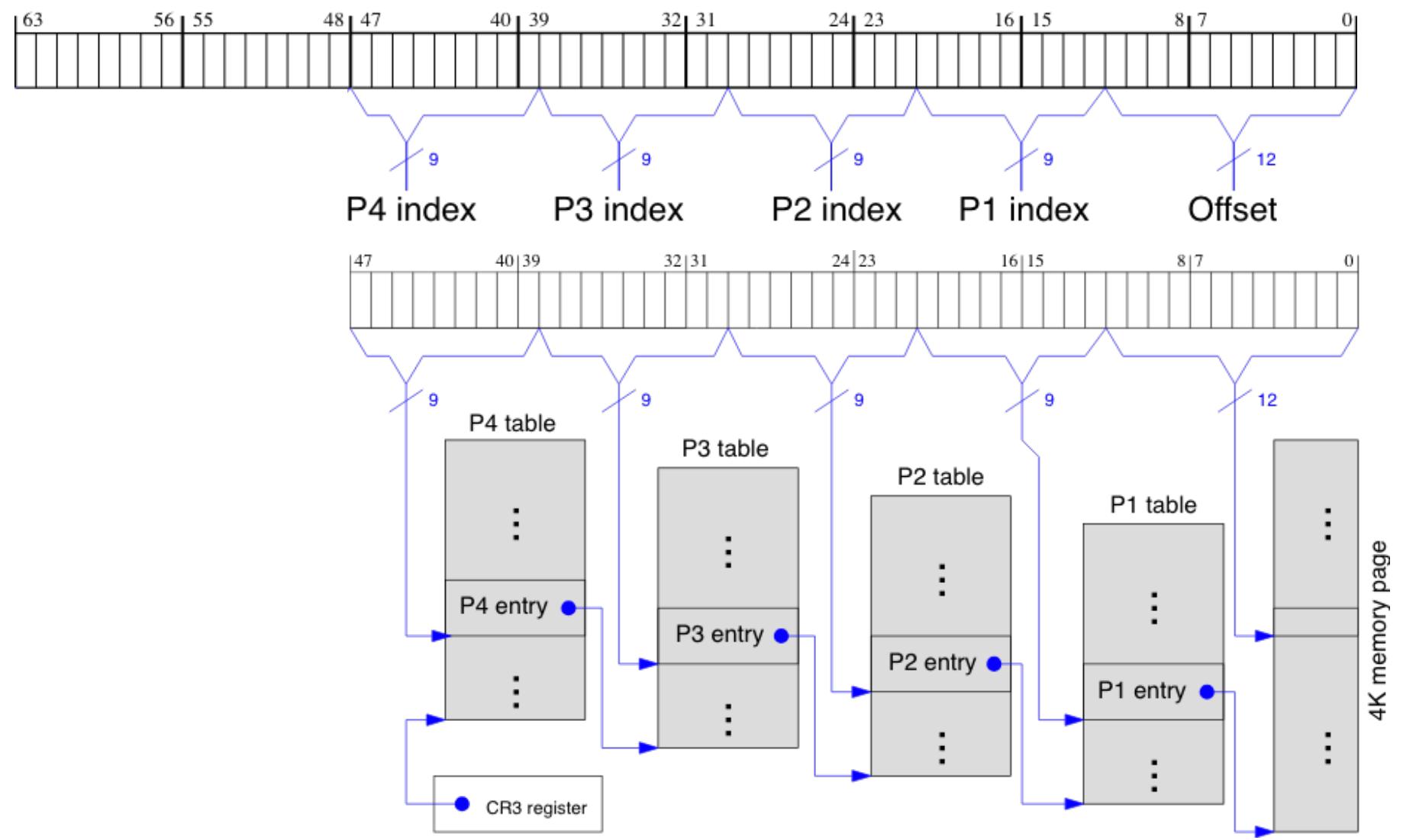
I could be doing a lot of useful work
while waiting to read memory !

What to do ?

Caching Generalities

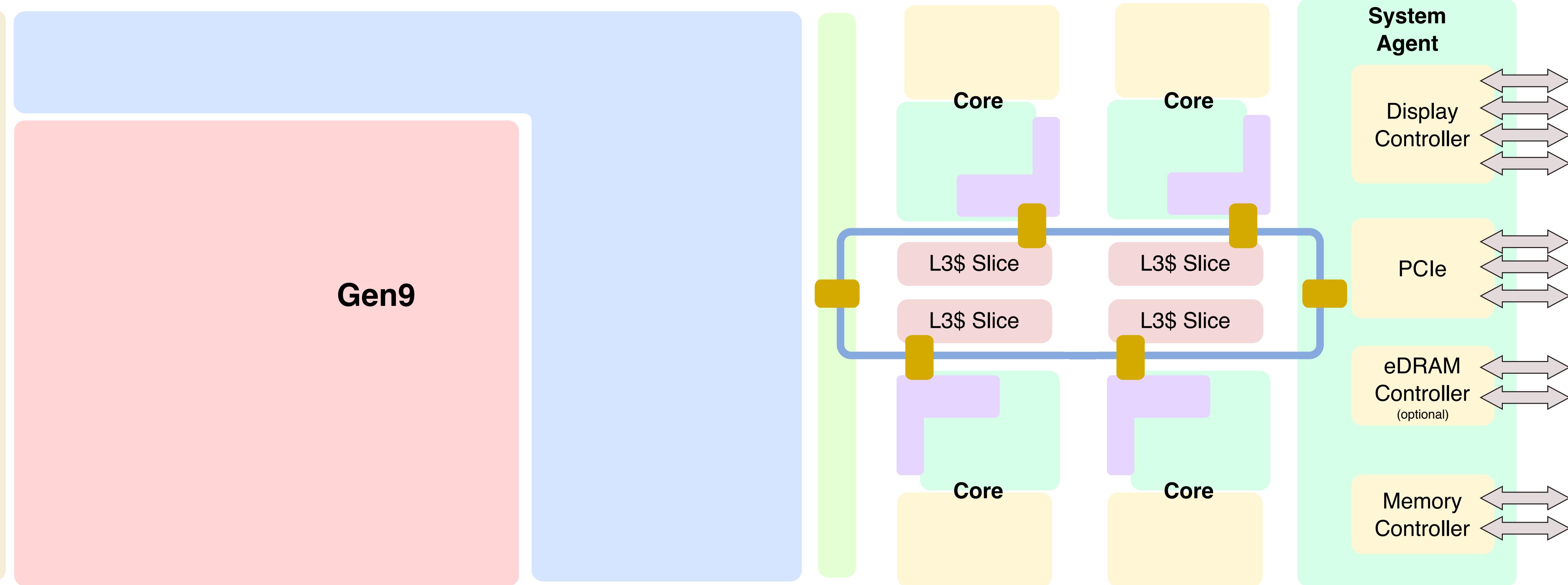
Recap of Key Concepts in Caching

- 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

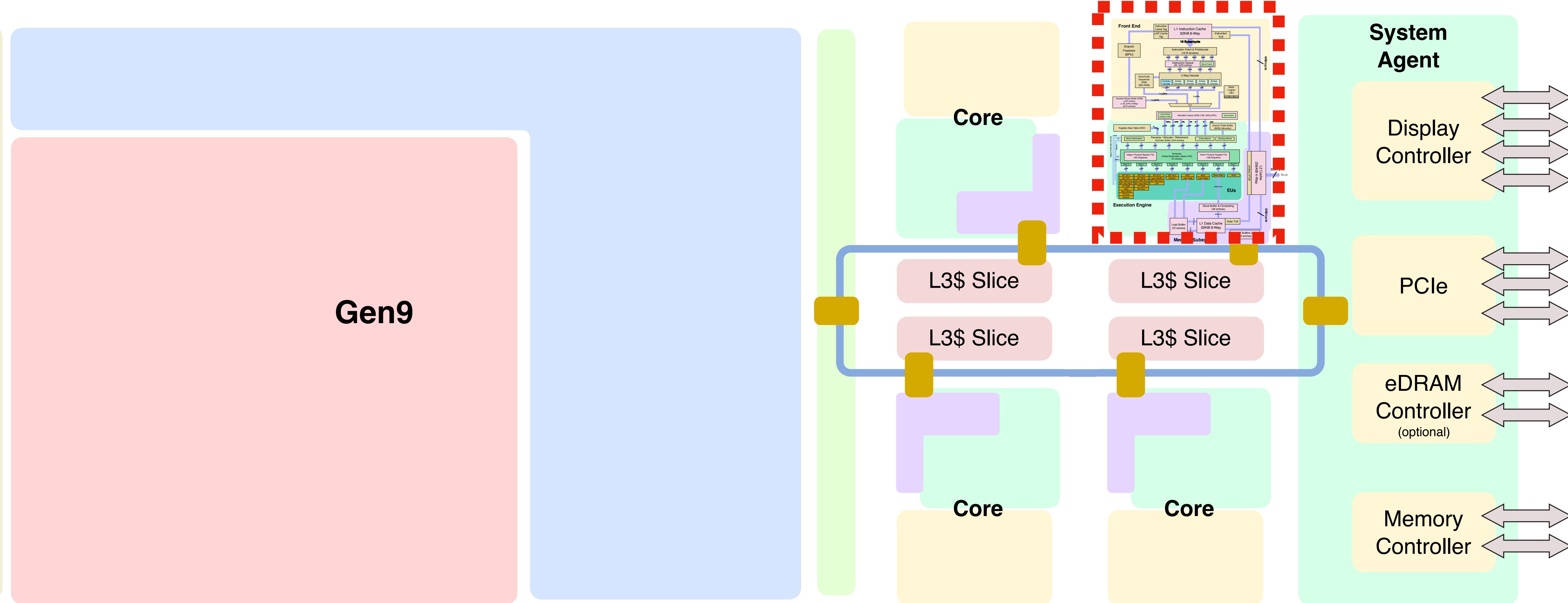


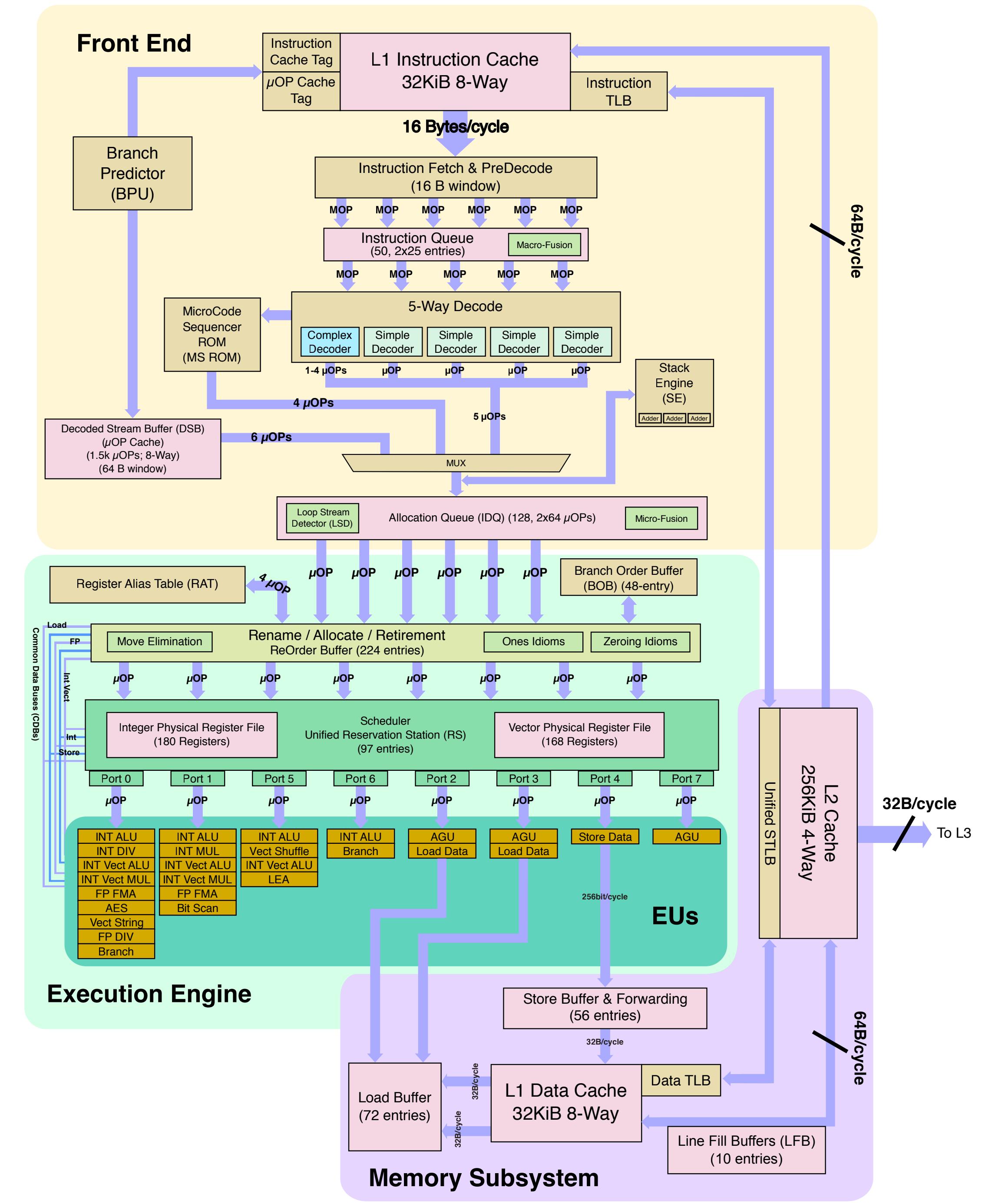
Caching & TLBs

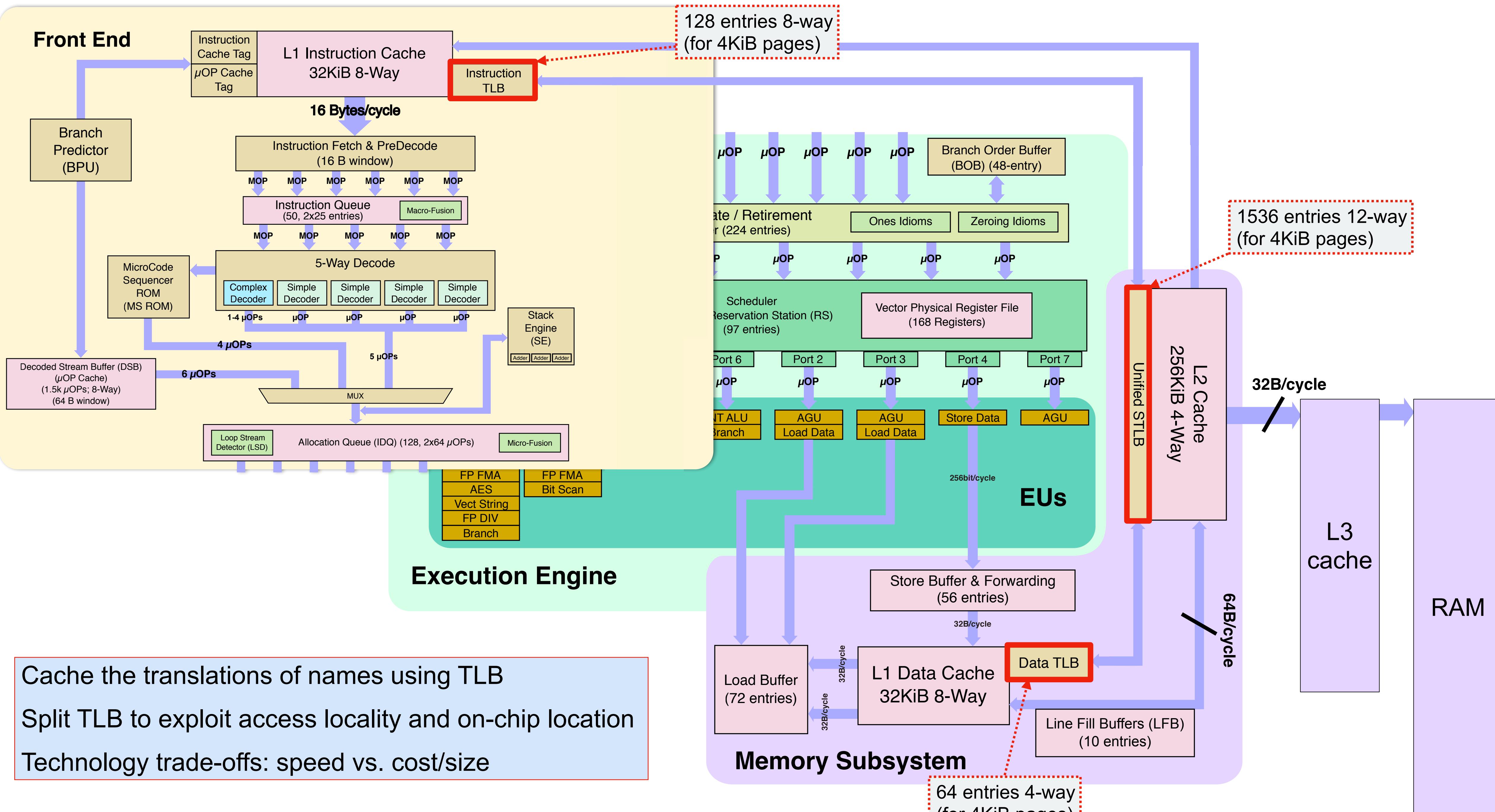
Intel Skylake Microarchitecture

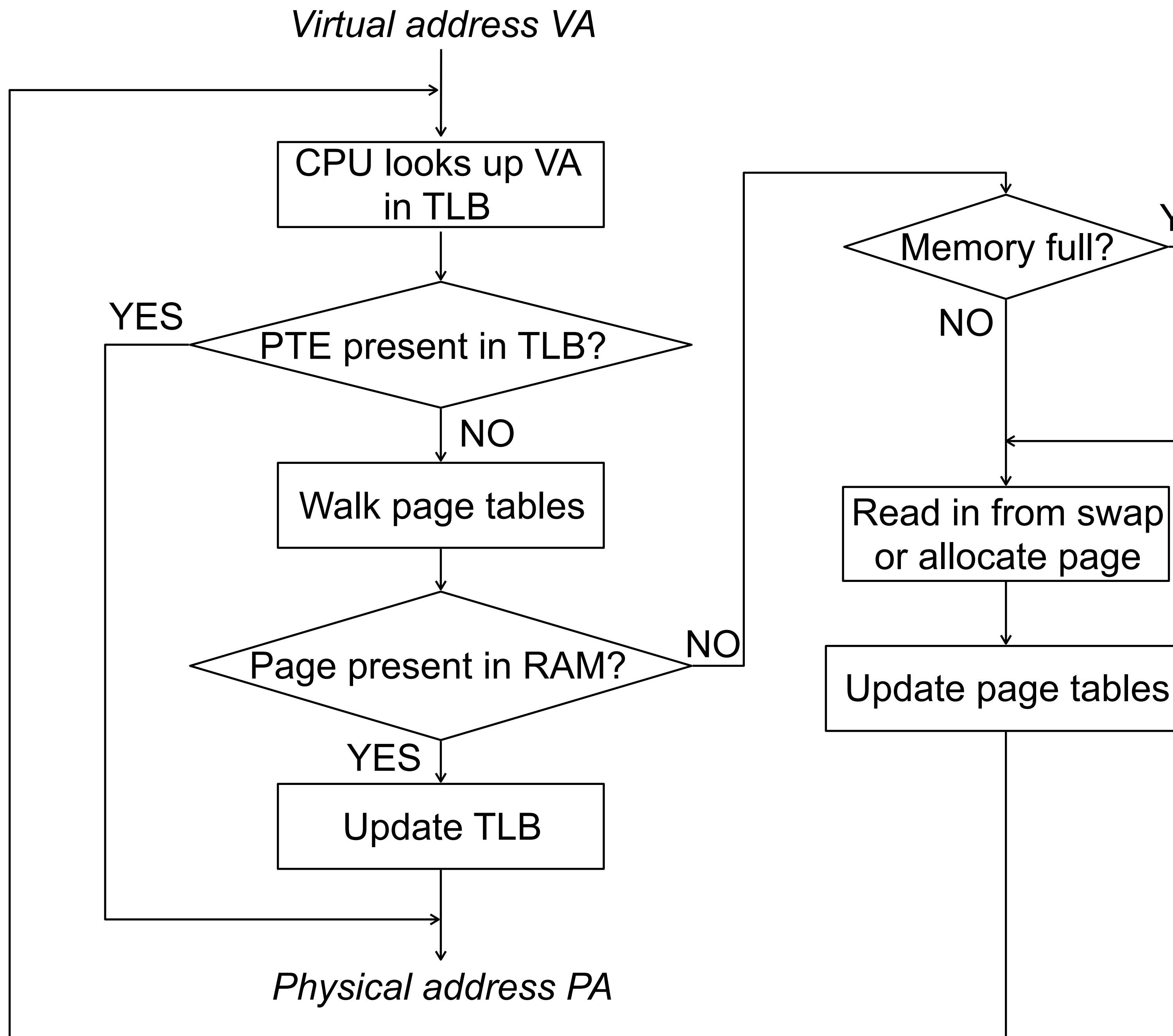


Intel Skylake Microarchitecture









This is HW-managed TLB !

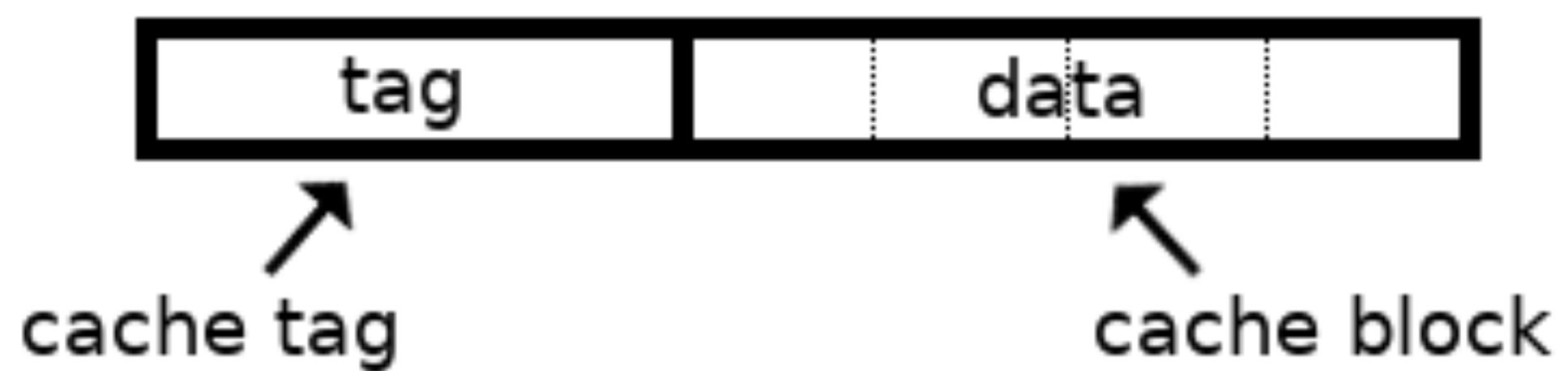
In SW-managed TLB (MIPS, SPARC, ...)

- 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)

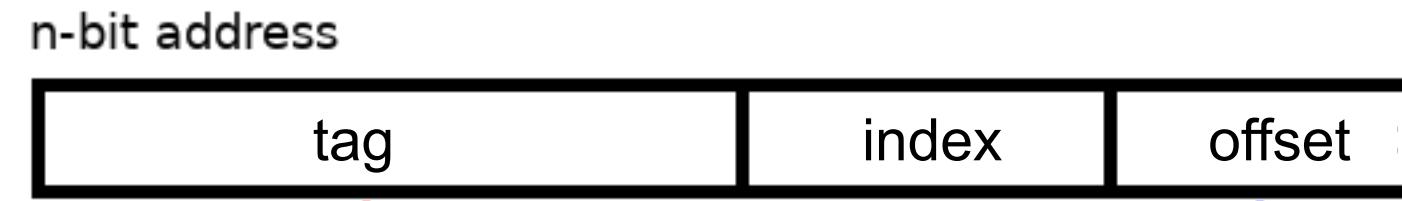


Address (cache lookup key)

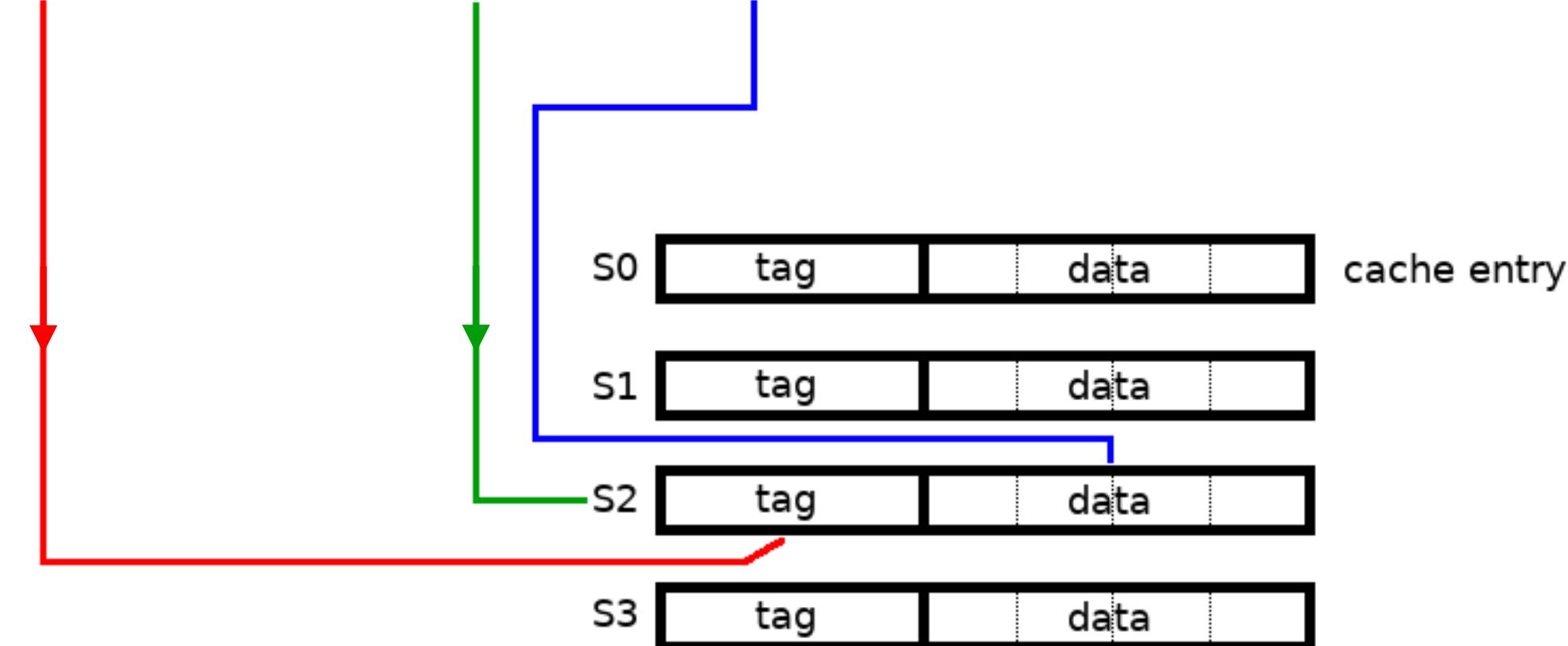
n-bit address



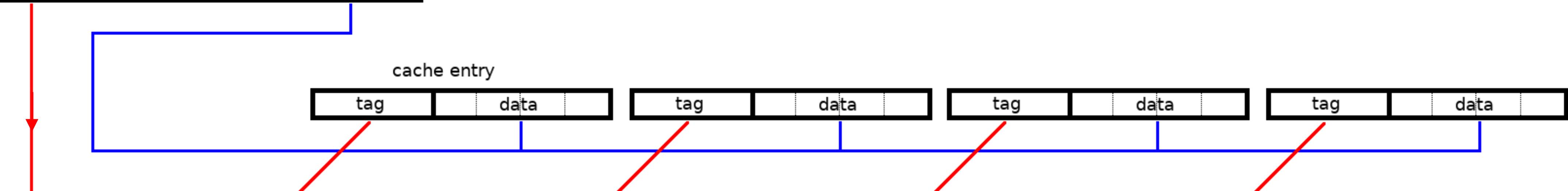
Cache structure



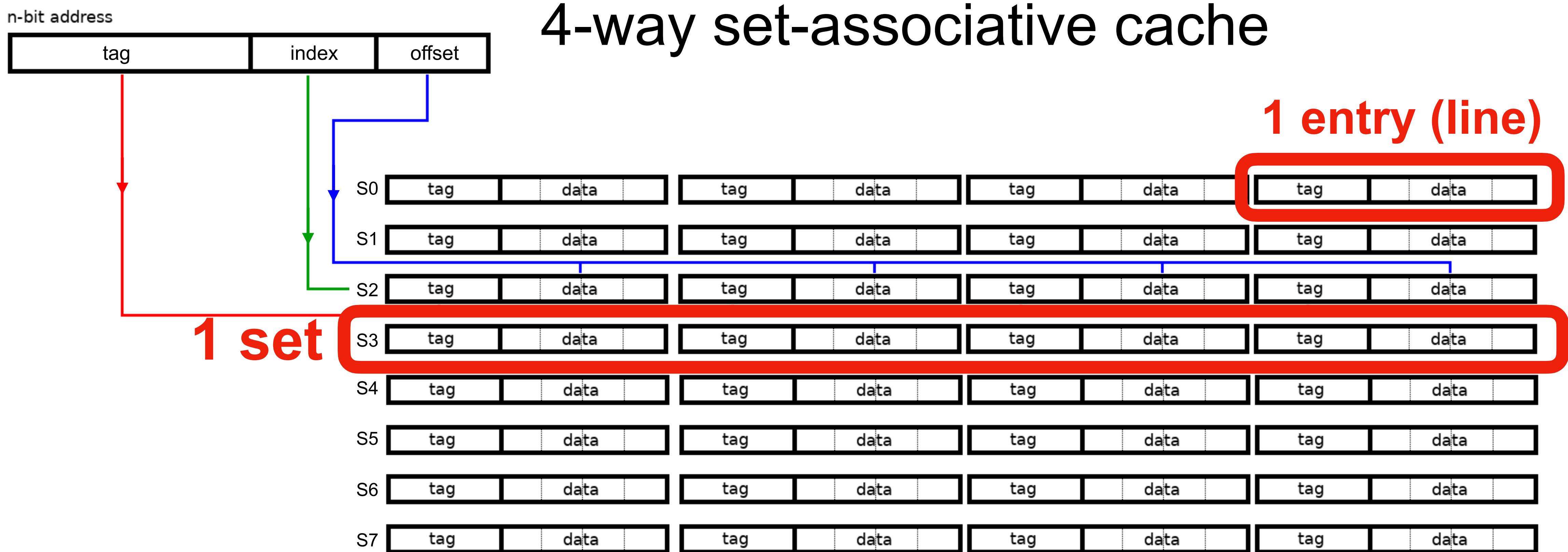
Direct-mapped cache



Fully associative cache

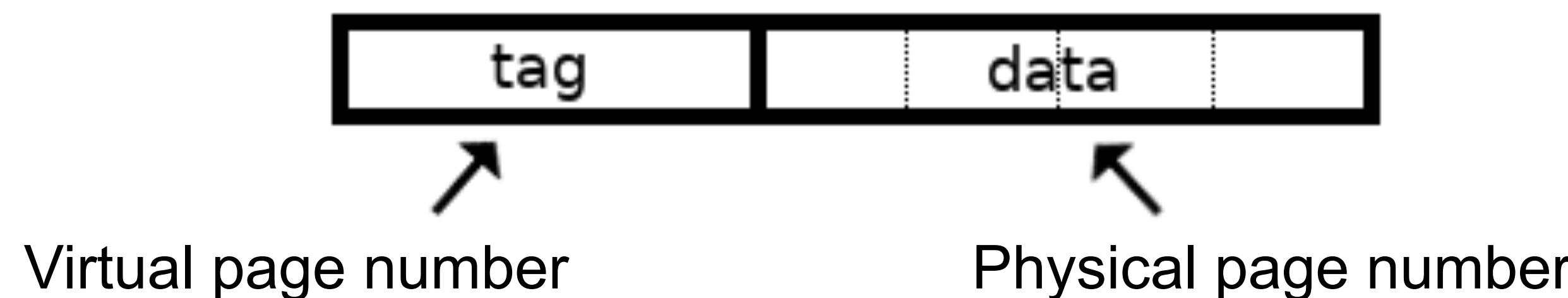


Cache structure



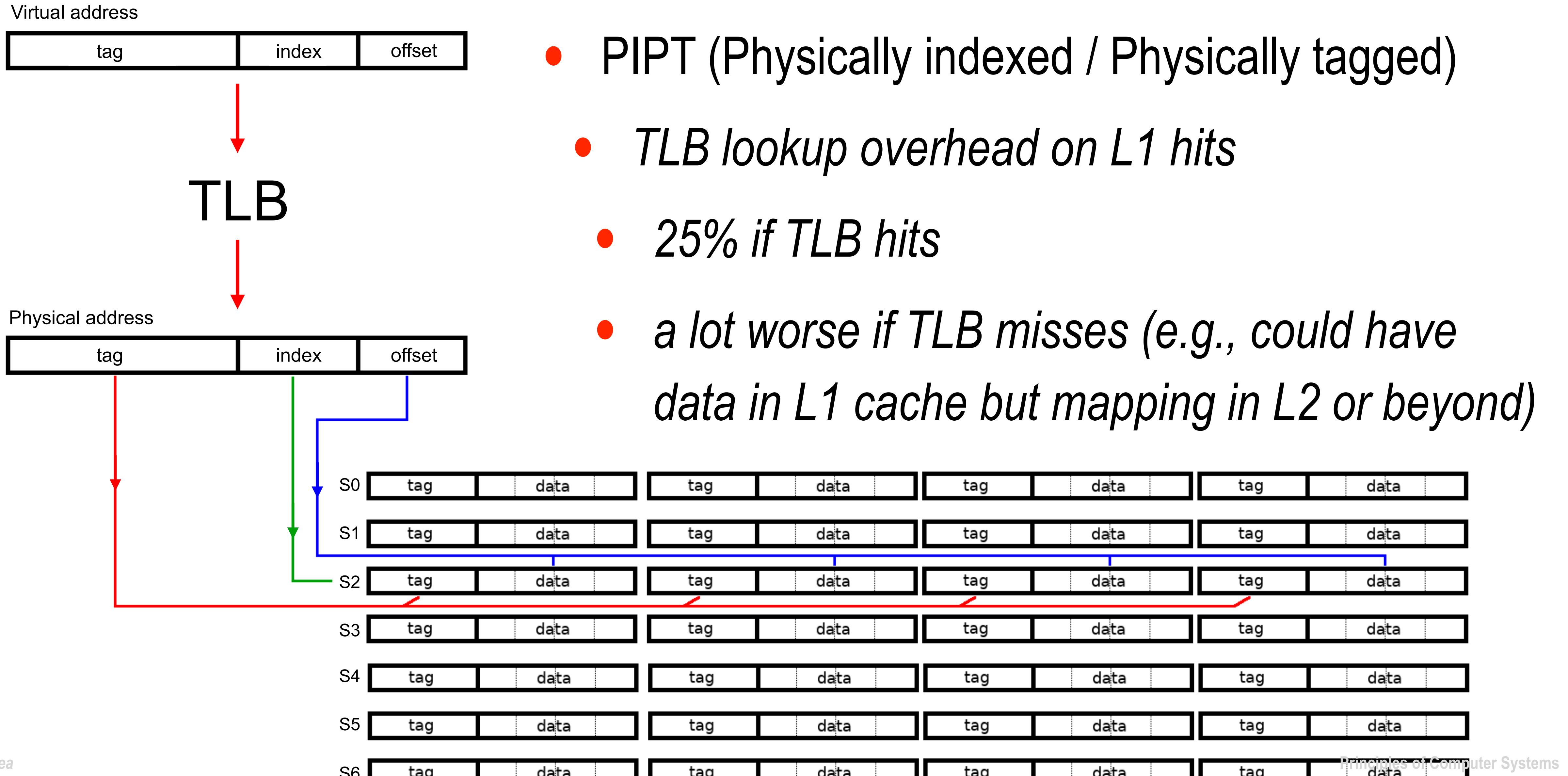
Overhead of memory-location name resolution

TLB cache line



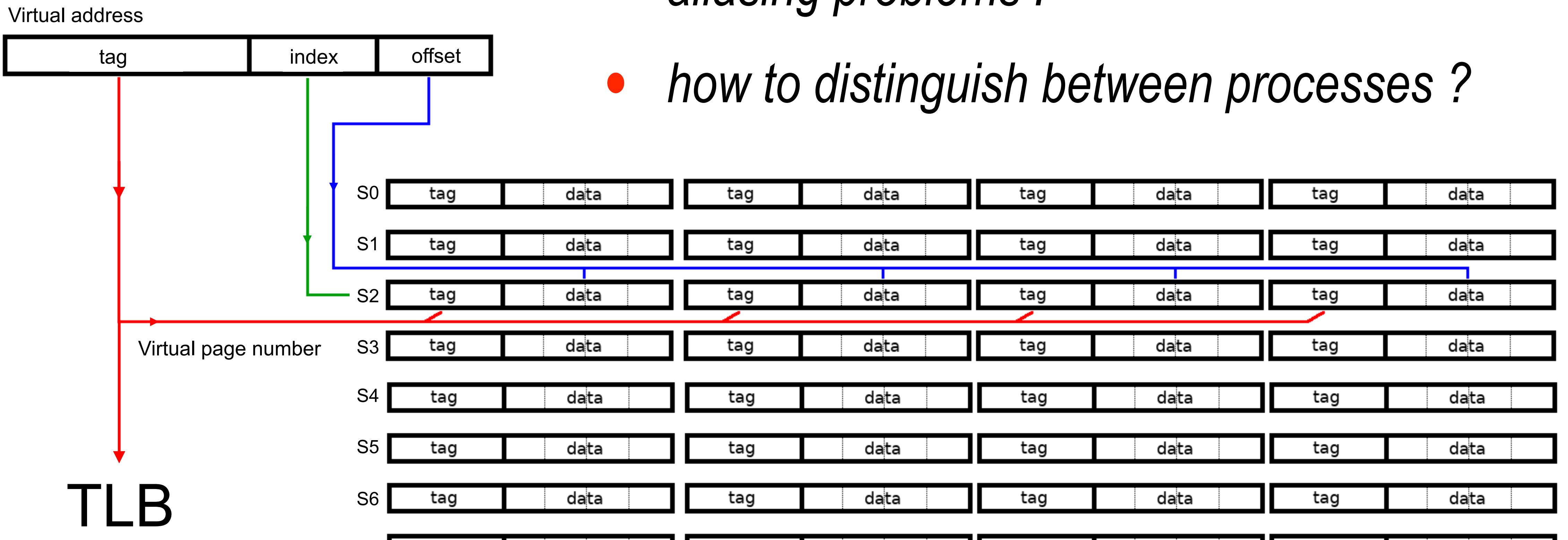
- Skylake (4 KiB pages)
 - $L1 i\text{-}TLB: 16 \text{ sets} * 8\text{-way set-associative} = 128 \text{ entries}$
 - $L1 d\text{-}TLB: 16 \text{ sets} * 4\text{-way set associative} = 64 \text{ entries}$
 - $L2 \text{ shared } TLB: 256 \text{ sets} * 12\text{-way set associative} = 1536 \text{ entries}$

PIPT L1 d/i-cache



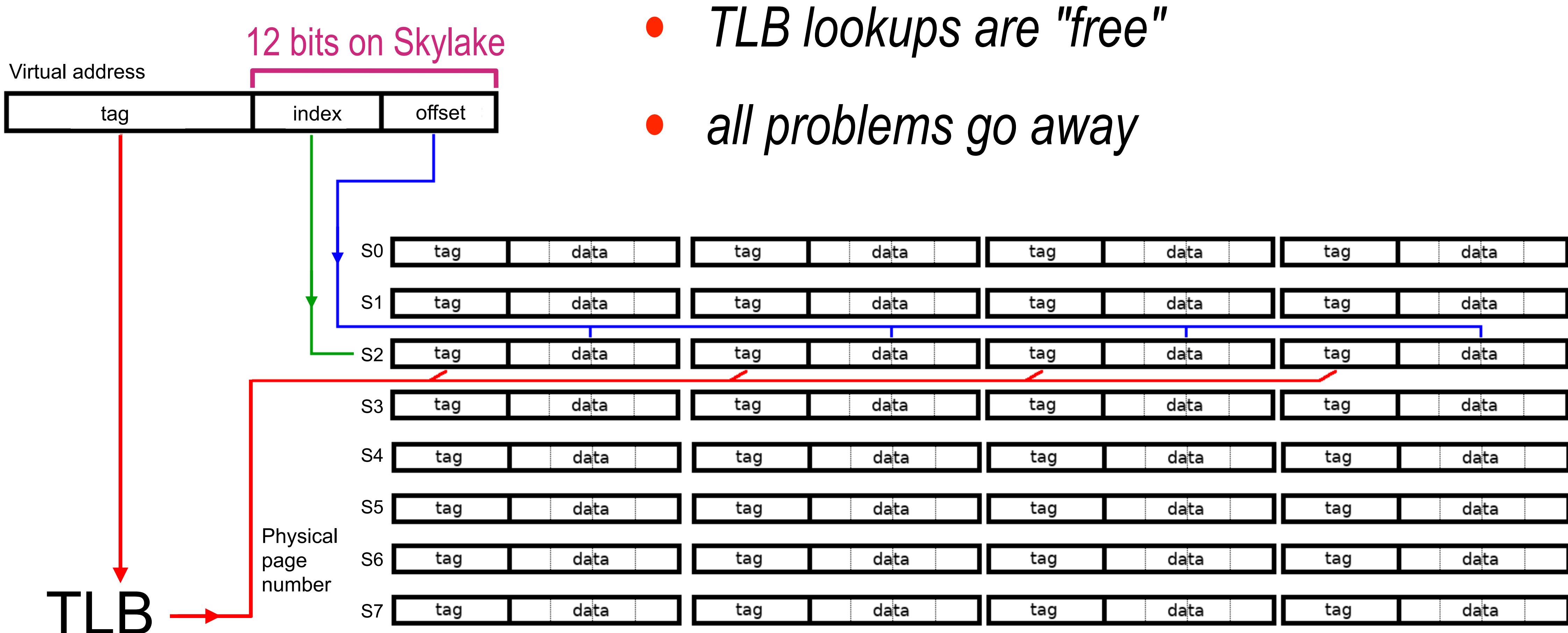
VIVT L1 d/i-cache

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



VIPT L1 d/i-cache

- VIPT (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 / ~2 cycles / ~30-40 cycles / ~100-120 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