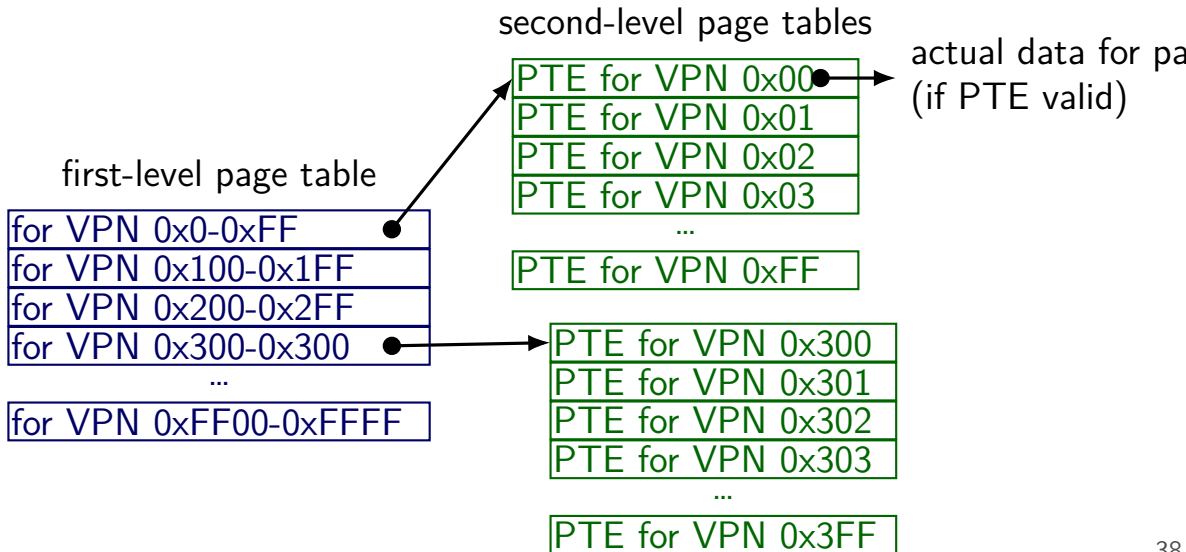


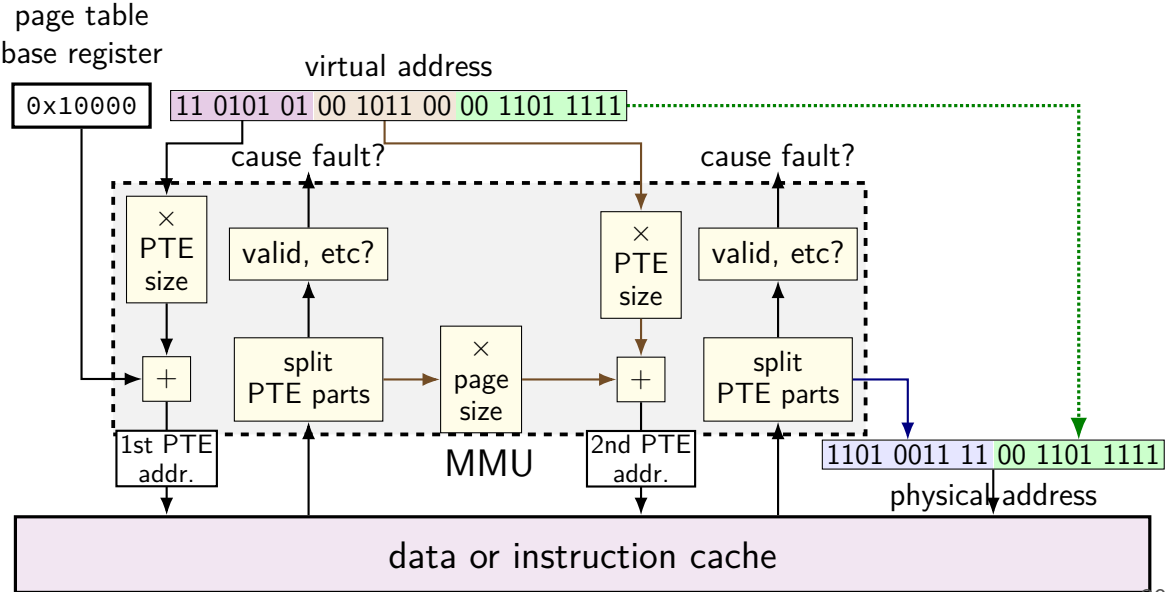
Virtual Memory (2)

two-level page tables

two-level page table for 65536 pages (16-bit VPN)



two-level page table lookup



current x86-64 page tables

4-level page table

512 PTEs of 8 bytes each for each page table

choice: exactly one page per page table

allows OS to allocate new page table space in one page units

(just like program memory)

page table space exercise (1)

4-level page table

512 PTEs of 8 bytes each for each page table

suppose a process has exactly one page allocated

how much space for page tables?

page table space exercise (1)

4-level page table

512 PTEs of 8 bytes each for each page table

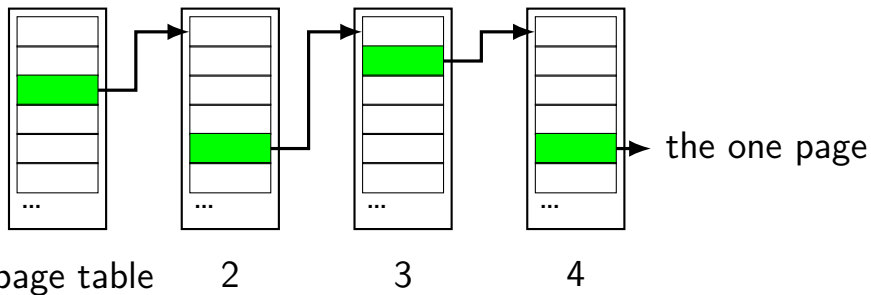
suppose a process has exactly one page allocated

how much space for page tables?

1 page at each level (4KB each)

exactly one valid entry in each of them

page table space exercise (1)



4 page tables at 1 page/page table
plus 1 page of data
5 pages total

page table space exercise (2)

4-level page table

512 PTEs of 8 bytes each for each page table

suppose a process has exactly two pages allocated:

one at address 0x0, one at address 0x200000000000

how much space for page tables?

page table space exercise (2)

4-level page table

512 PTEs of 8 bytes each for each page table

suppose a process has exactly two pages allocated:

one at address 0x0, one at address 0x200000000000

how much space for page tables?

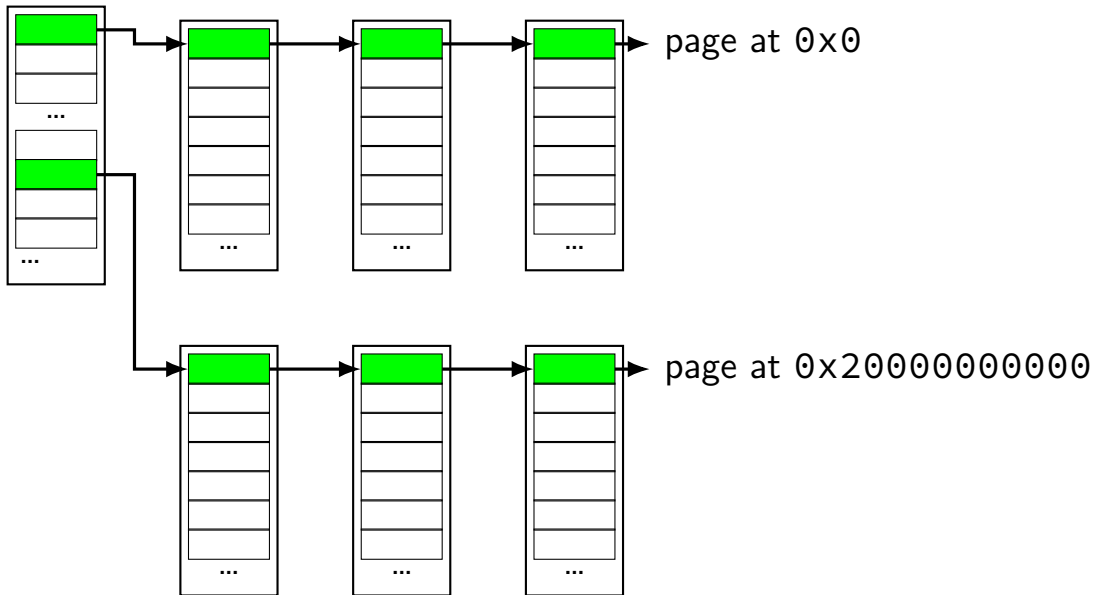
1 shared first-level PT, with two valid entries

two second-level PTs, each with one valid entry

two third-level PTs, each with one valid entry

two fourth-level PTs, each with one valid entry

page table space exercise (2)



page table space exercise (3)

4-level page table; each PT: 512 PTEs of 8 bytes

suppose a process has 100 pages of stack, 100 pages of code+constants (contiguous)

stack and code+constants far apart

how much space for page tables?

page table space exercise (3)

4-level page table; each PT: 512 PTEs of 8 bytes

suppose a process has 100 pages of stack, 100 pages of code+constants (contiguous)

stack and code+constants far apart

how much space for page tables? — *minimum*:

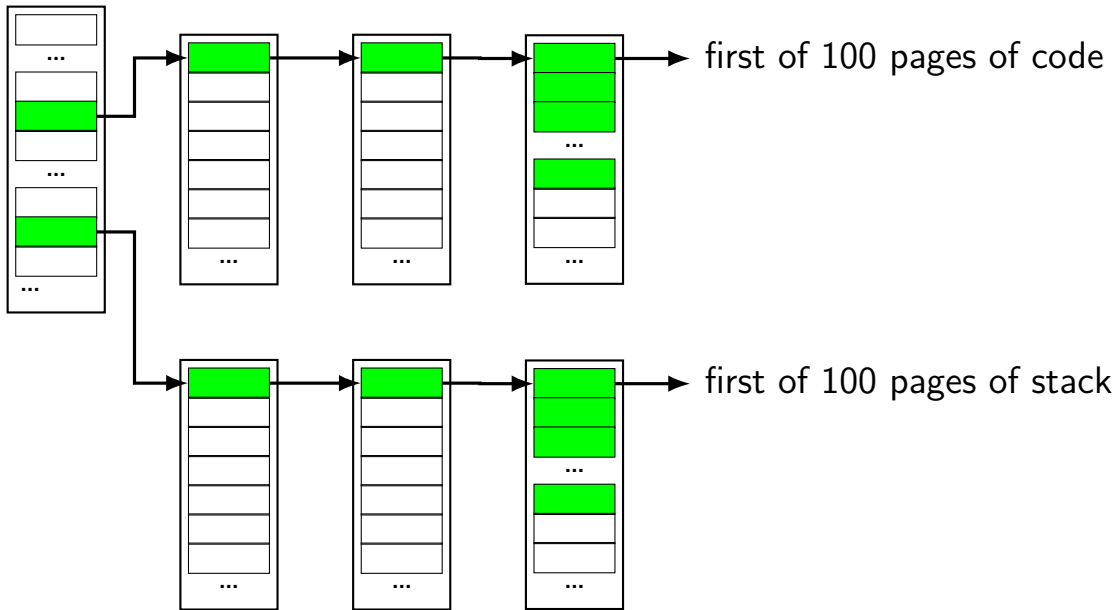
1 shared first-level PT, with two valid entries

two second-level PT, each with one valid entry

two third-level PT, each with one valid entry

two fourth-level PT, each with 100 valid entries

page table space exercise (3)



page table space exercise (3)

4-level page table; each PT: 512 PTEs of 8 bytes

suppose a process has 100 pages of stack, 100 pages of code+constants (contiguous)

how much space for page tables?

page table space exercise (3)

4-level page table; each PT: 512 PTEs of 8 bytes

suppose a process has 100 pages of stack, 100 pages of code+constants (contiguous)

how much space for page tables? — *maximum*:

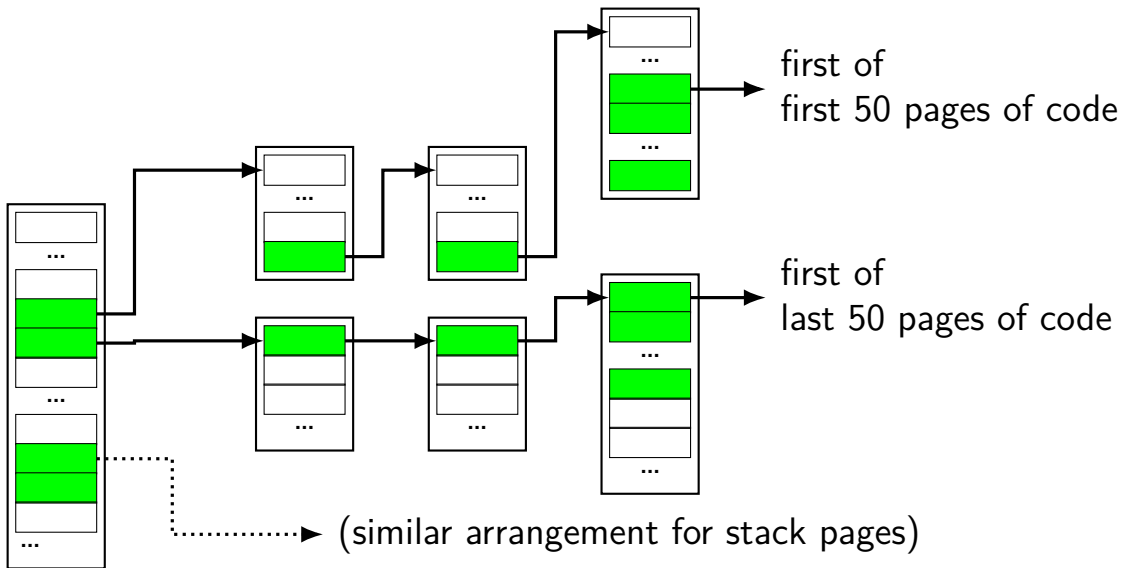
1 shared first-level PT, with four valid entries

four second-level PT, each with one valid entry
two for stack, two for code+constants

four third-level PT, each with one valid entry

four fourth-level PT, each with 50 valid entries

page table space exercise (3)



page table space exercise (4)

4-level page table; each PT: 512 PTEs of 8 bytes

suppose a process has 200 pages, randomly distributed in PT

about how much space for page tables?

page table space exercise (4)

4-level page table; each PT: 512 PTEs of 8 bytes

suppose a process has 200 pages, randomly distributed in PT

about how much space for page tables?

about 165 ($\pm \sim 8$) entries in first-level PT

(some pages randomly share first-level PT entries)

about 165 second-level PTs, 200 third-level, 200 fourth-level

a bit less than 600 page tables — almost 2400 KB

cache accesses and multi-level PTs

four-level page tables — four cache accesses per memory access

L1 cache hits — typically a couple cycles each?

so add 8 cycles to each memory access?

not acceptable

TLBs

cache accesses and multi-level PTs

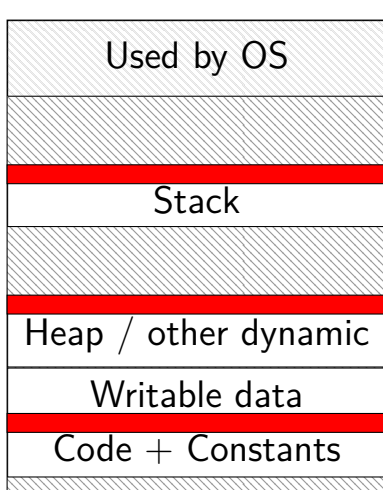
four-level page tables — four cache accesses per memory access

L1 cache hits — typically a couple cycles each?

so add 8 cycles to each memory access?

not acceptable

program memory active sets



0xFFFF FFFF FFFF FFFF

0xFFFF 8000 0000 0000

0x7F...

small areas of memory active at a time
one or two pages in each area?

0x0000 0000 0040 0000

page table entries and locality

page table entries have **excellent temporal locality**

typically one or two pages of the stack active

typically one or two pages of code active

typically one or two pages of heap/globals active

each page contains **whole functions**, arrays, stack frames, etc.

page table entries and locality

page table entries have **excellent temporal locality**

typically one or two pages of the stack active

typically one or two pages of code active

typically one or two pages of heap/globals active

each page contains **whole functions**, arrays, stack frames, etc.

needed page table entries are **very small**

page table entry cache

called a **TLB** (translation lookaside buffer)

very small cache of page table entries

L1 cache	TLB
physical addresses	virtual page numbers
bytes from memory	page table entries
tens of bytes per block	one page table entry per block
usually thousands of blocks	usually tens of entries

page table entry cache

called a **TLB** (translation lookaside buffer)

very small cache of page table entries

L1 cache	TLB
physical addresses	virtual page numbers
bytes from memory	page table entries
tens of bytes per block	one page table entry per block
usually thousands of blocks	usually tens of entries
only caches the page table lookup itself (generally) just entries from the last-level page table	

page table entry cache

called a **TLB** (translation lookaside buffer)

very small cache of page table entries

L1 cache	TLB
physical addresses	virtual page numbers
bytes from memory	page table entries
tens of bytes per block	one page table entry per block
usually thousands of blocks	usually tens of entries

not much spatial locality between page table entries
(they're used for kilobytes of data already)

page table entry cache

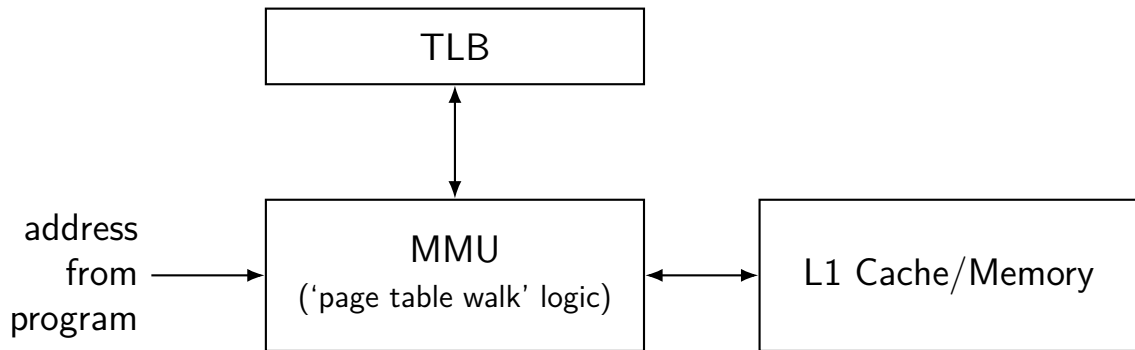
called a **TLB** (translation lookaside buffer)

very small cache of page table entries

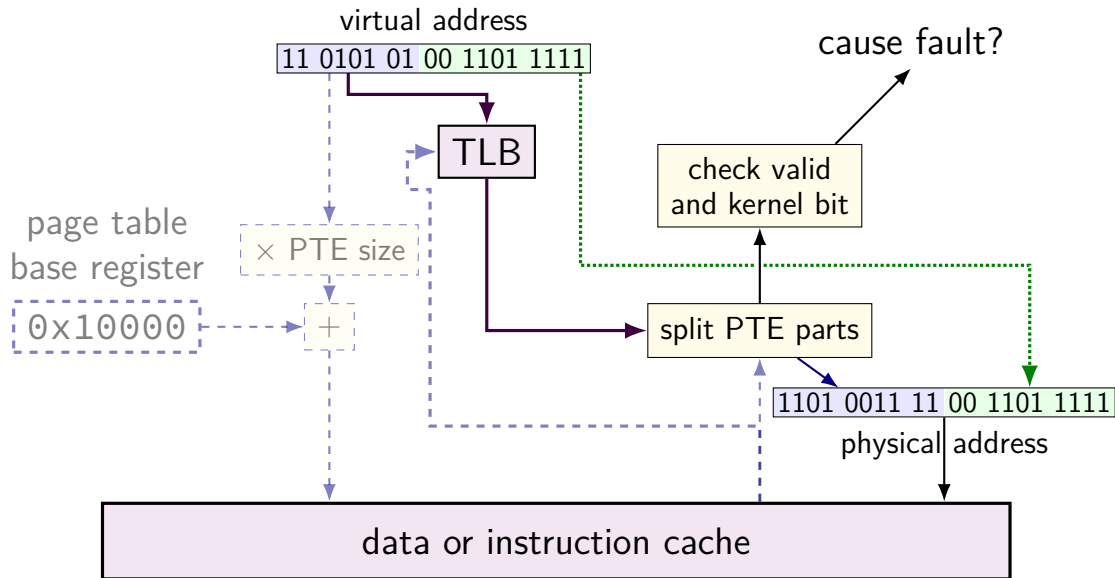
L1 cache	TLB
physical addresses	virtual page numbers
bytes from memory	page table entries
tens of bytes per block	one page table entry per block
usually thousands of blocks	usually tens of entries

few active page table entries at a time
enables highly associative cache designs

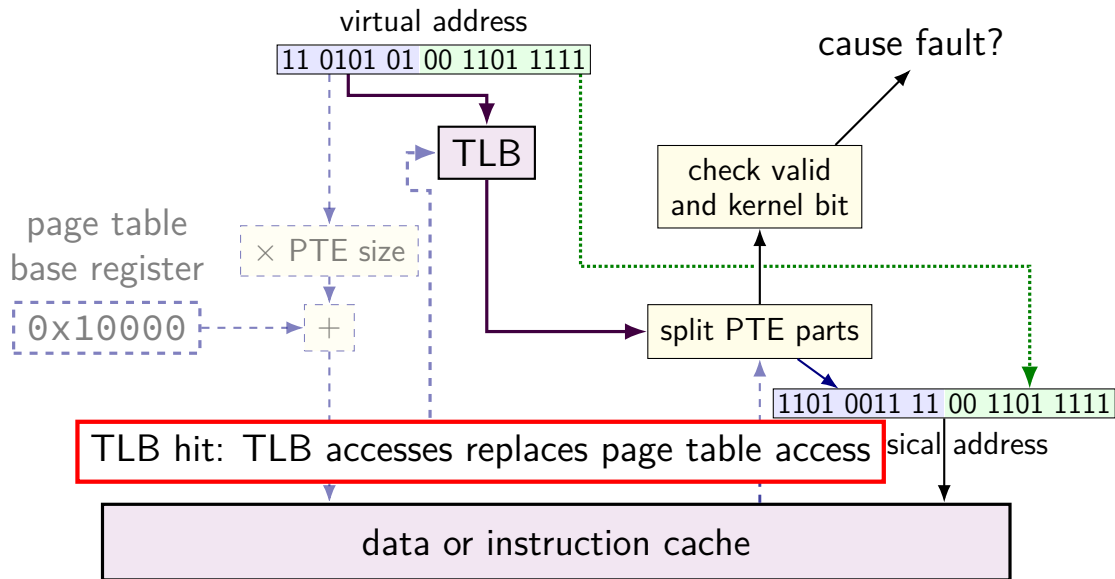
TLB and the MMU (1)



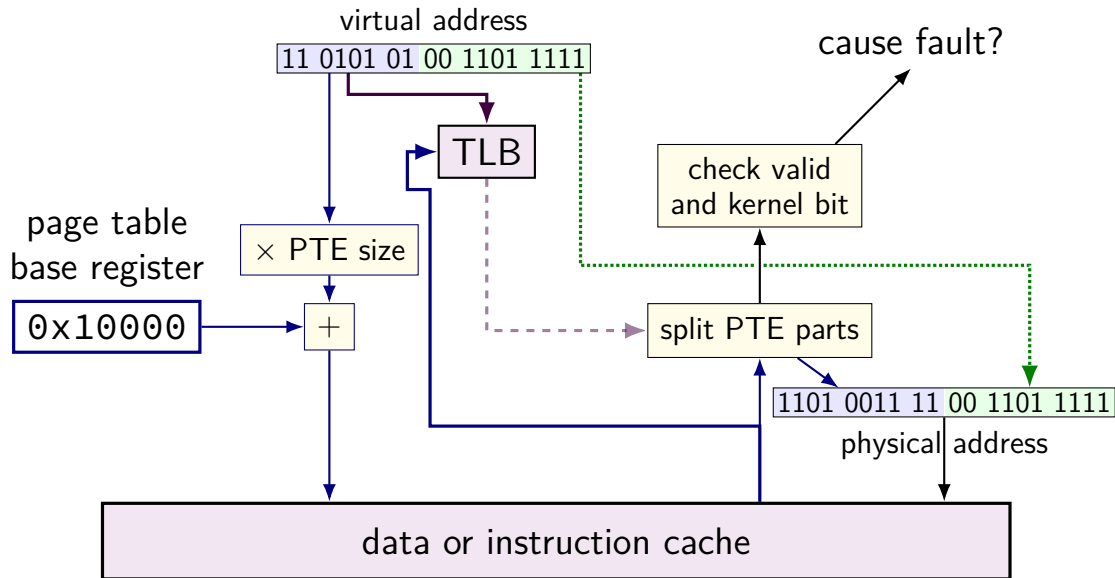
TLB and the MMU (2)



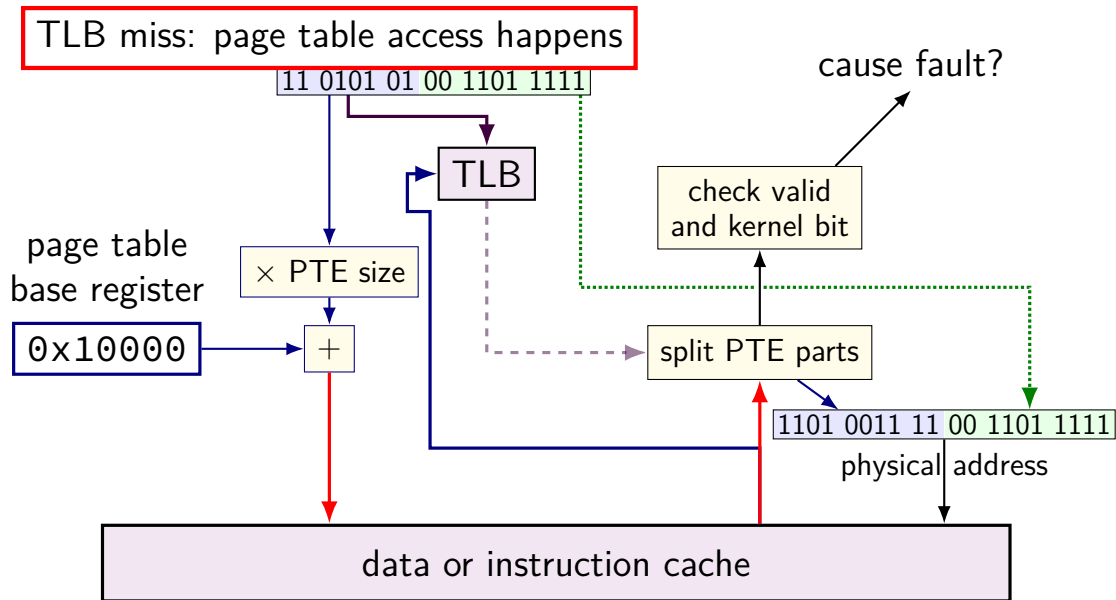
TLB and the MMU (2)



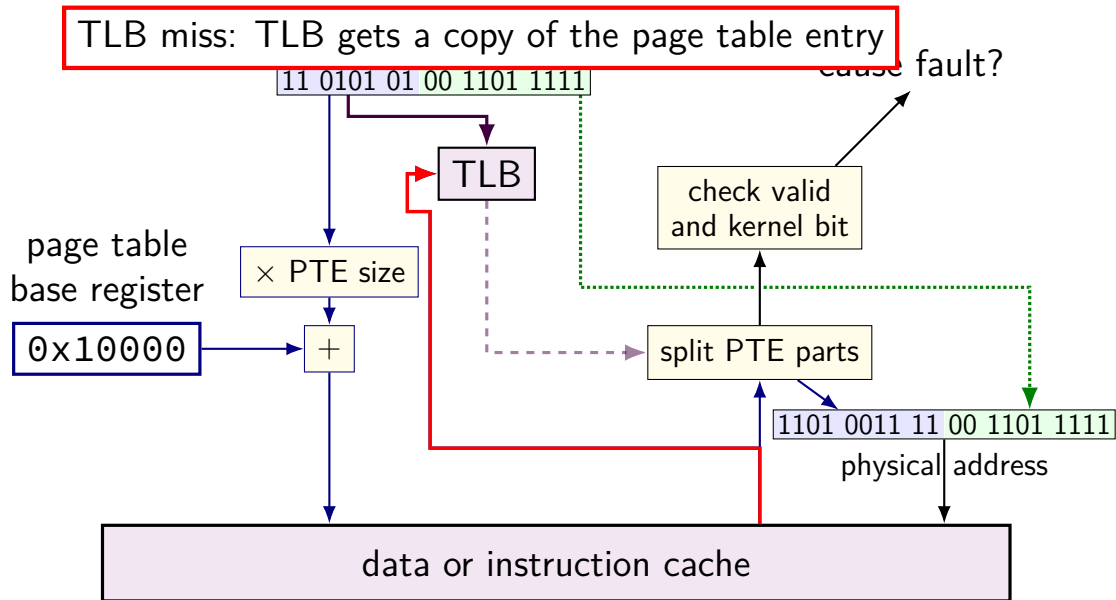
TLB and the MMU (2)



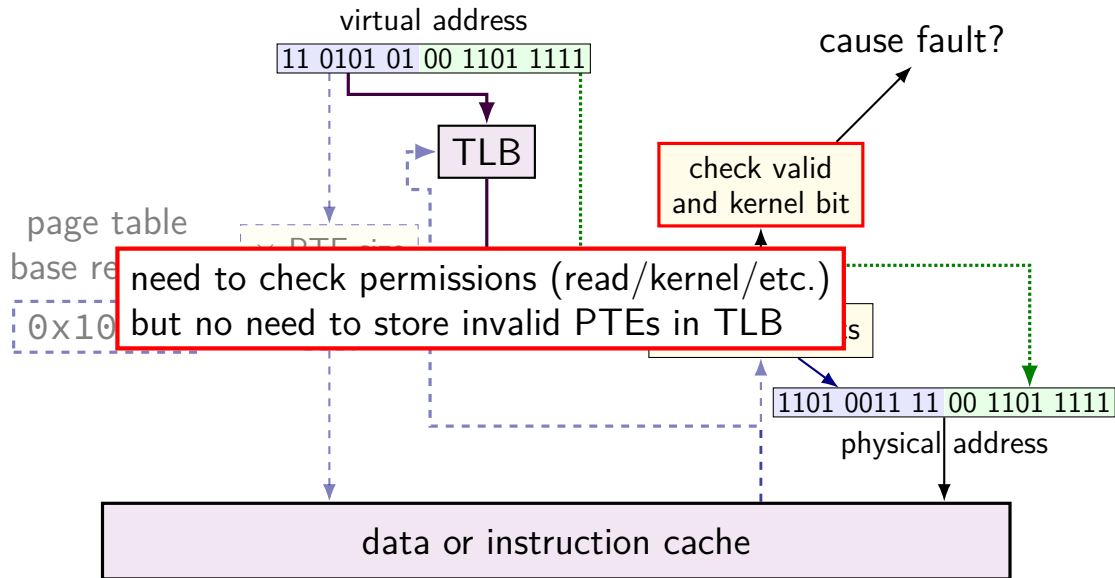
TLB and the MMU (2)



TLB and the MMU (2)



TLB and the MMU (2)



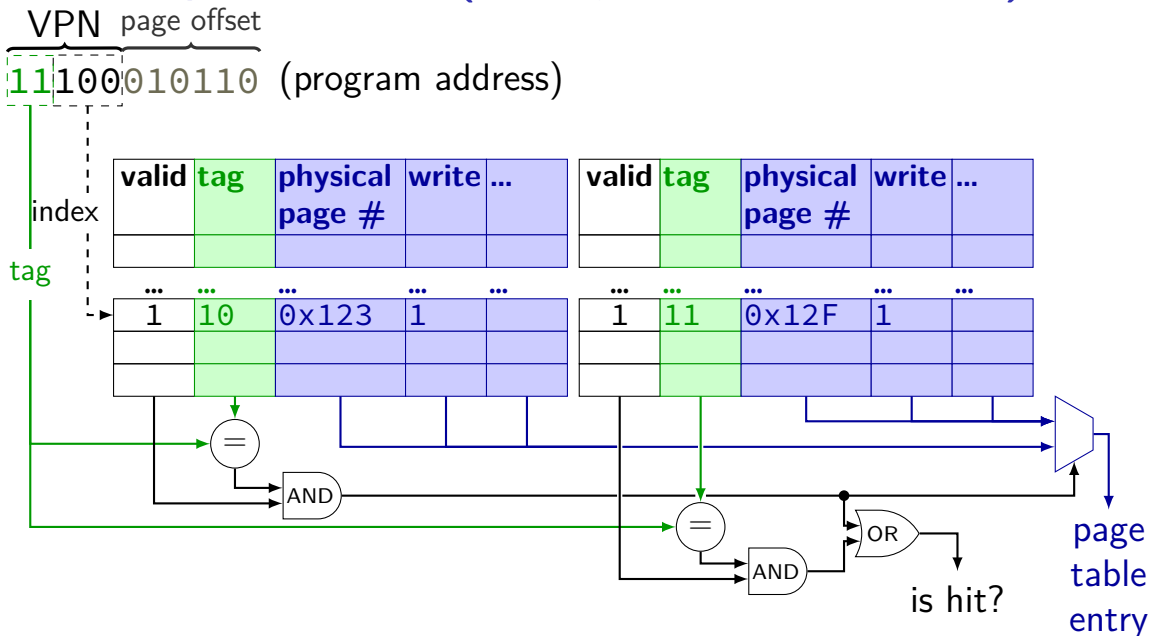
TLB and multi-level page tables

TLB caches **valid last-level page table entries**

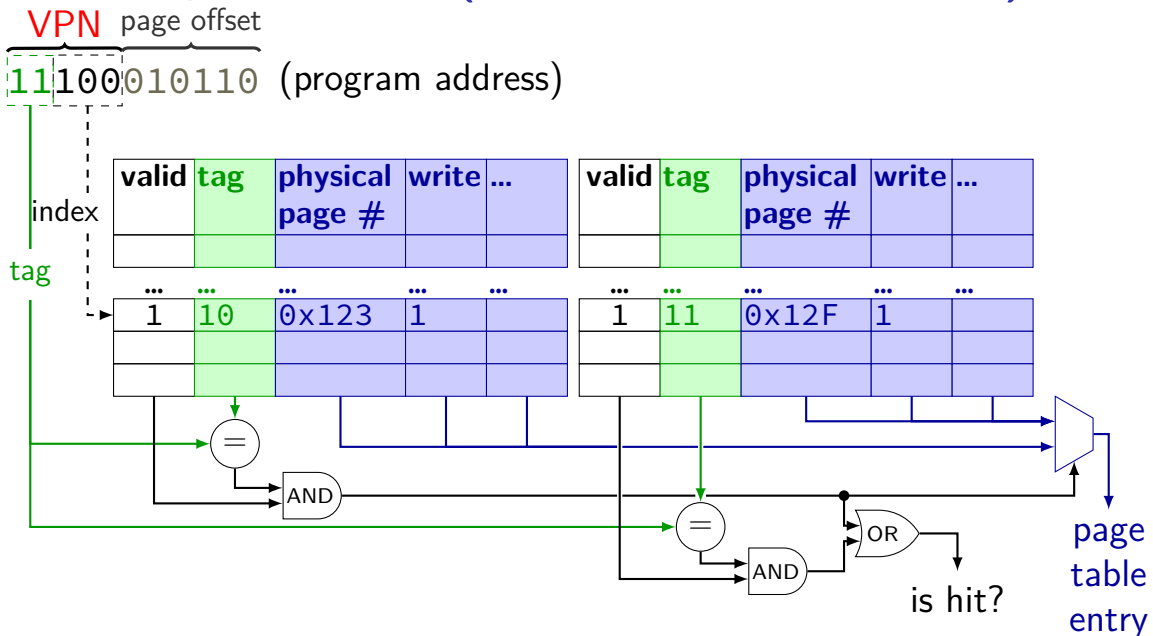
doesn't matter which last-level page table

means TLB output can be used directly to form address

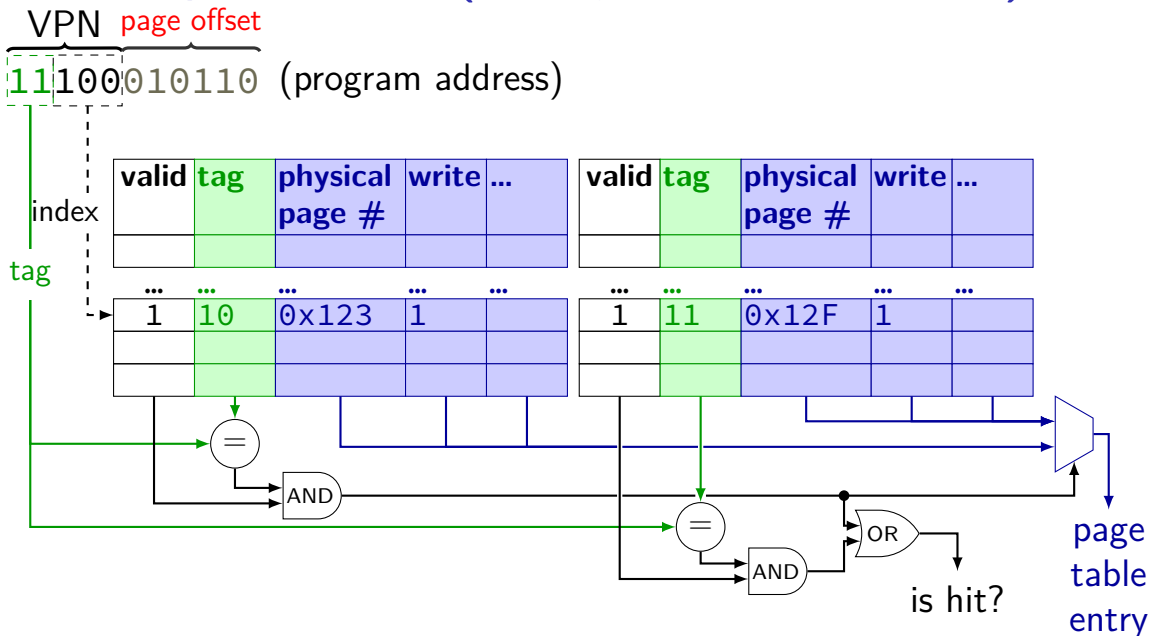
TLB organization (2-way set associative)



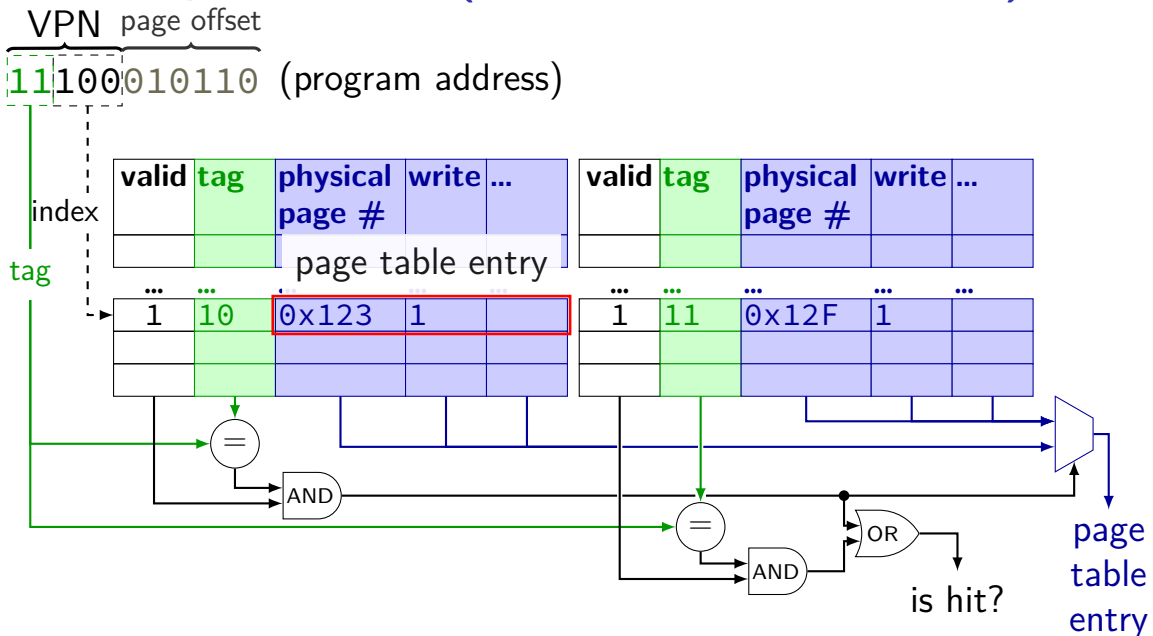
TLB organization (2-way set associative)



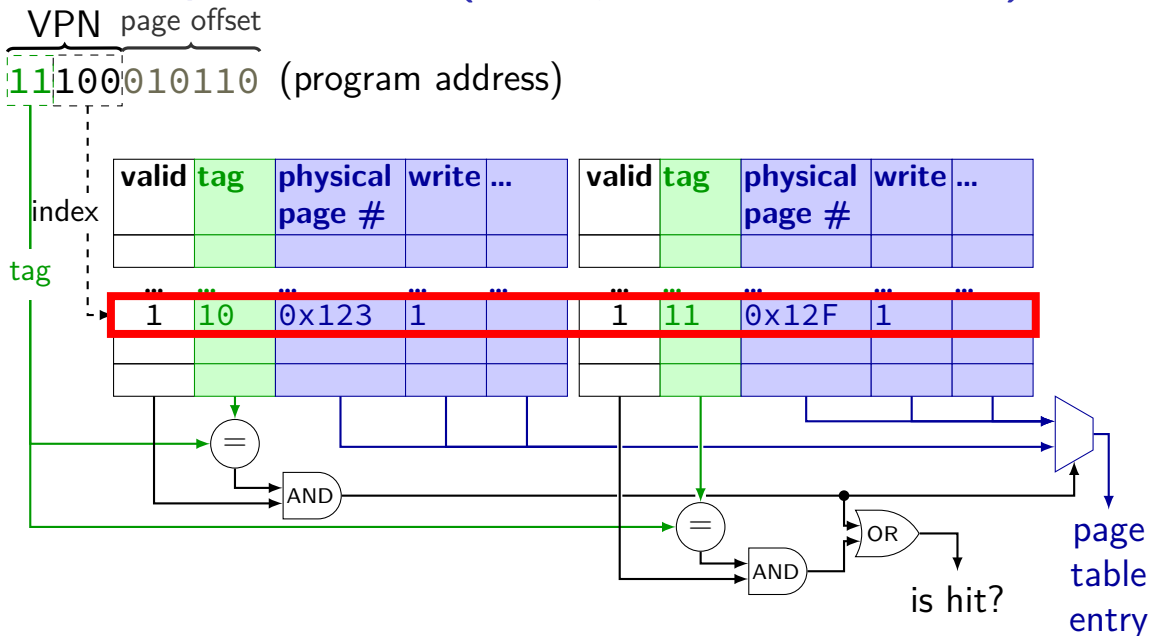
TLB organization (2-way set associative)



TLB organization (2-way set associative)



TLB organization (2-way set associative)



address splitting for TLBs (1)

my desktop:

4KB (2^{12} byte) pages; 48-bit virtual address

64-entry, 4-way L1 data TLB

TLB index bits?

TLB tag bits?

address splitting for TLBs (1)

my desktop:

4KB (2^{12} byte) pages; 48-bit virtual address

64-entry, 4-way L1 data TLB

TLB index bits?

$$64/4 = 16 \text{ sets} \text{ — } 4 \text{ bits}$$

TLB tag bits?

$$48 - 12 = 36 \text{ bit virtual address} \text{ — } 36 - 4 = 32 \text{ bit TLB tag}$$

address splitting for TLBs (2)

my desktop:

4KB (2^{12} byte) pages; 48-bit virtual address

1536-entry ($3 \cdot 2^9$), 12-way L2 TLB

TLB index bits?

TLB tag bits?

address splitting for TLBs (2)

my desktop:

4KB (2^{12} byte) pages; 48-bit virtual address

1536-entry ($3 \cdot 2^9$), 12-way L2 TLB

TLB index bits?

$$1536/12 = 128 \text{ sets} \text{ — } 7 \text{ bits}$$

TLB tag bits?

$$48 - 12 = 36 \text{ bit virtual address} \text{ — } 36 - 7 = 29 \text{ bit TLB tag}$$

address splitting exercise (3)

384-entry, 3-way set-associative TLB

32-bit virtual address; 8KB pages

2-level page table; 4 byte PTEs

256 entries in first level; 2048 in second

split the address `0x12345678`

address splitting exercise (3)

384-entry, 3-way set-associative TLB

32-bit virtual address; 8KB pages

2-level page table; 4 byte PTEs

256 entries in first level; 2048 in second

split the address `0x12345678`

0001 0010 0011 0100 0101 0110 0111 1000

address splitting exercise (3)

384-entry, 3-way set-associative TLB

32-bit virtual address; 8KB pages

2-level page table; 4 byte PTEs

256 entries in first level; 2048 in second

split the address 0x12345678

0001 0010 0011 0100 0101 0110 0111 1000

13-bit page offset 1 0110 0111 1000

address splitting exercise (3)

384-entry, 3-way set-associative TLB

32-bit virtual address; 8KB pages

2-level page table; 4 byte PTEs

256 entries in first level; 2048 in second

split the address 0x12345678

0001 0010 0011 0100 0101 0110 0111 1000

13-bit page offset 1 0110 0111 1000

32 - 13 = 19-bit VPN 0001 0010 0011 0100 010

address splitting exercise (3)

384-entry, 3-way set-associative TLB

32-bit virtual address; 8KB pages

2-level page table; 4 byte PTEs

256 entries in first level; 2048 in second

split the address 0x12345678

0001 0010 0011 0100 0101 0110 0111 1000

13-bit page offset 1 0110 0111 1000

32 - 13 = 19-bit VPN 0001 0010 0011 0100 010

8-bit first part of VPN 0001 0010

address splitting exercise (3)

384-entry, 3-way set-associative TLB

32-bit virtual address; 8KB pages

2-level page table; 4 byte PTEs

256 entries in first level; 2048 in second

split the address 0x12345678

0001 0010 0011 0100 0101 0110 0111 1000

13-bit page offset 1 0110 0111 1000

32 - 13 = 19-bit VPN 0001 0010 0011 0100 010

8-bit first part of VPN 0001 0010

11-bit second part of VPN 0011 0100 010

address splitting exercise (3)

384-entry, 3-way set-associative TLB

32-bit virtual address; 8KB pages

2-level page table; 4 byte PTEs

256 entries in first level; 2048 in second

split the address 0x12345678

0001 0010 0011 0100 0101 0110 0111 1000

13-bit page offset 1 0110 0111 1000

32 - 13 = 19-bit VPN 0001 0010 0011 0100 010

8-bit first part of VPN 0001 0010

11-bit second part of VPN 0011 0100 010

7-bit TLB index 0100 010

address splitting exercise (3)

384-entry, 3-way set-associative TLB

32-bit virtual address; 8KB pages

2-level page table; 4 byte PTEs

256 entries in first level; 2048 in second

split the address 0x12345678

0001 0010 0011 0100 0101 0110 0111 1000

13-bit page offset 1 0110 0111 1000

32 - 13 = 19-bit VPN 0001 0010 0011 0100 010

8-bit first part of VPN 0001 0010

11-bit second part of VPN 0011 0100 010

7-bit TLB index 0100 010

19 - 7 = 12-bit TLB tag 0001 0010 0011

TLB example: address splitting

16-bit virtual addresses

64-byte pages

8-entry, 2-way TLB

TLB access pattern

64-byte pages

8 entries, 4 sets

index	V	tag	PTE
00	0	????	—
01	0	????	—
10	0	????	—
11	0	????	—

V	tag	PTE	LRU
0	????	—	?
0	????	—	?
0	????	—	?
0	????	—	?

address (hex)	hit?
000100000000 (100)	
110100000001 (D01)	
000100001010 (10A)	
110100100001 (D21)	
000011111100 (0FC)	
110011111000 (CF8)	
111100101000 (F23)	

TLB access pattern

64-byte pages

8 entries, 4 sets

index	V	tag	PTE
00	0	????	—
01	0	????	—
10	0	????	—
11	0	????	—

V	tag	PTE	LRU
0	????	—	?
0	????	—	?
0	????	—	?
0	????	—	?

address (hex)	hit?
000100000000 (100)	
110100000001 (D01)	
000100001010 (10A)	
110100100001 (D21)	
000011111100 (0FC)	
110011111000 (CF8)	
111100101000 (F23)	

VPN

TLB access pattern

64-byte pages

8 entries, 4 sets

index	V	tag	PTE
00	0	????	—
01	0	????	—
10	0	????	—
11	0	????	—

V	tag	PTE	LRU
0	????	—	?
0	????	—	?
0	????	—	?
0	????	—	?

address (hex)	hit?
000100000000 (100)	
110100000001 (D01)	
000100001010 (10A)	
110100100001 (D21)	
000011111100 (0FC)	
110011111000 (CF8)	
111100101000 (F23)	

VPN

tag index

TLB access pattern

64-byte pages

8 entries, 4 sets

index	V	tag	PTE	V	tag	PTE	LRU
00	1	0001	for VPN 000100	0	????		way 1
01	0	????	—	0	????	—	?
10	0	????	—	0	????	—	?
11	0	????	—	0	????	—	?

address (hex)	hit?
000100 000000 (100)	miss
110100 000001 (D01)	
000100 001010 (10A)	
110100 100001 (D21)	
000011 111100 (0FC)	
110011 111000 (CF8)	
111100 101000 (F23)	

VPN

tag index

TLB access pattern

64-byte pages

8 entries, 4 sets

index	V	tag	PTE
00	1	0001	for VPN 000100
01	0	????	—
10	0	????	—
11	0	????	—

V	tag	PTE
1	1101	for VPN 110100
0	????	—
0	????	—
0	????	—

LRU
way 0
?
?
?

address (hex)	hit?
000100 (000000 (100))	miss
110100 (000001 (D01))	miss
000100 (001010 (10A))	
110100 (100001 (D21))	
000011 (111100 (0FC))	
110011 (111000 (CF8))	
111100 (101000 (F23))	

VPN

tag index

TLB access pattern

64-byte pages

8 entries, 4 sets

index	V	tag	PTE
00	1	0001	for VPN 000100
01	0	????	—
10	0	????	—
11	0	????	—

V	tag	PTE	LRU
1	1101	for VPN 110100	way 1
0	????	—	?
0	????	—	?
0	????	—	?

address (hex)	hit?
000100 (000000 (100))	miss
110100 (000001 (D01))	miss
000100 (001010 (10A))	hit
110100 (100001 (D21))	
000011 (111100 (0FC))	
110011 (111000 (CF8))	
111100 (101000 (F23))	

VPN

tag index

TLB access pattern

64-byte pages

8 entries, 4 sets

index	V	tag	PTE
00	1	0001	for VPN 000100
01	0	????	—
10	0	????	—
11	0	????	—

V	tag	PTE	LRU
1	1101	for VPN 110100	way 0
0	????	—	?
0	????	—	?
0	????	—	?

address (hex)	hit?
000100 (000000 (100))	miss
110100 (000001 (D01))	miss
000100 (001010 (10A))	hit
110100 (100001 (D21))	hit
000011 (111100 (0FC))	
110011 (111000 (CF8))	
111100 (101000 (F23))	

VPN

tag index

TLB access pattern

64-byte pages

8 entries, 4 sets

index	V	tag	PTE
00	1	0001	for VPN 000100
01	0	????	—
10	0	????	—
11	1	0000	for VPN 000011

V	tag	PTE
1	1101	for VPN 110100
0	????	—
0	????	—
	????	—

LRU
way 0
?
?
way 1

address (hex)	hit?
000100 (000000 (100))	miss
110100 (000001 (D01))	miss
000100 (001010 (10A))	hit
110100 (100001 (D21))	hit
000011 (111100 (0FC))	miss
110011 (111000 (CF8))	
111100 (101000 (F23))	

VPN

tag index

TLB access pattern

64-byte pages

8 entries, 4 sets

index	V	tag	PTE
00	1	0001	for VPN 000100
01	0	????	—
10	0	????	—
11	1	0000	for VPN 000011

V	tag	PTE
1	1101	for VPN 110100
0	????	—
0	????	—
1	1100	for VPN 110011

LRU
way 0
?
?
way 0

address (hex)	hit?
000100 (100)	miss
110100 (D01)	miss
000100001010 (10A)	hit
110100100001 (D21)	hit
000011111100 (0FC)	miss
110011111000 (CF8)	miss
111100101000 (F23)	

VPN

tag index

TLB access pattern

64-byte pages

8 entries, 4 sets

index	V	tag	PTE
00	1	1111	for VPN 111100
01	0	????	—
10	0	????	—
11	1	0000	for VPN 000011

V	tag	PTE
1	1101	for VPN 110100
0	????	—
0	????	—
1	1100	for VPN 110011

LRU
way 1
?
?
way 0

address (hex)	hit?
000100 (100)	miss
110100 (D01)	miss
000100001010 (10A)	hit
110100100001 (D21)	hit
000011111100 (0FC)	miss
110011111000 (CF8)	miss
111100101000 (F23)	miss

VPN

tag index

changing page tables

what happens to TLB when page table base pointer is changed?

e.g. context switch

most entries in TLB refer to things from **wrong process**

oops — read from the wrong process's stack?

changing page tables

what happens to TLB when page table base pointer is changed?

e.g. context switch

most entries in TLB refer to things from **wrong process**

oops — read from the wrong process's stack?

option 1: **invalidate** all TLB entries

side effect on “change page table base register” instruction

changing page tables

what happens to TLB when page table base pointer is changed?

e.g. context switch

most entries in TLB refer to things from **wrong process**

oops — read from the wrong process's stack?

option 1: **invalidate** all TLB entries

side effect on “change page table base register” instruction

option 2: TLB entries contain process ID

set by OS (special register)

checked by TLB in addition to TLB tag, valid bit

editing page tables

what happens to TLB when OS changes a page table entry?

invalid to valid — nothing needed

TLB doesn't contain invalid entries

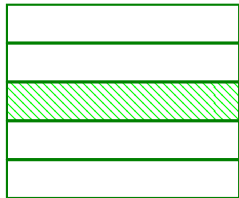
MMU will check memory again

valid to invalid — OS needs to tell processor to invalidate it
special instruction (x86: `invlpg`)

valid to other valid — OS needs to tell processor to invalidate it

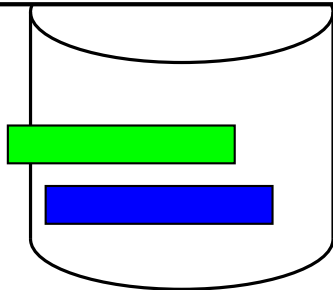
TLB shutdown

program A pages

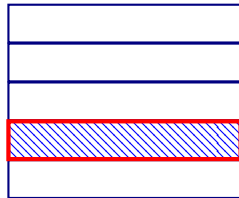


...

mark evicted page invalid in page table



program B pages



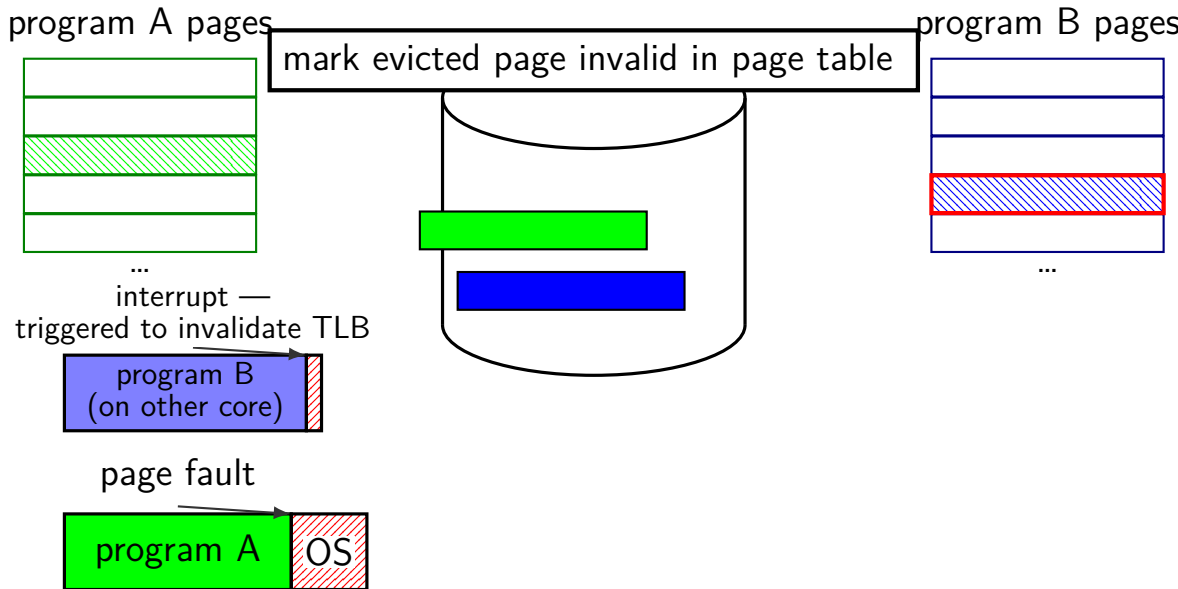
...

program B
(on other core)

page fault

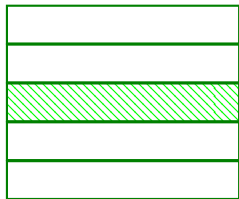


TLB shutdown



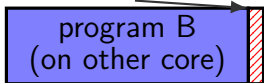
TLB shutdown

program A pages



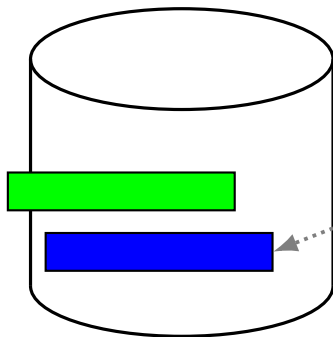
...

interrupt —
triggered to invalidate TLB



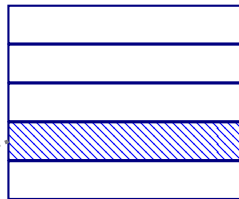
page fault

start read



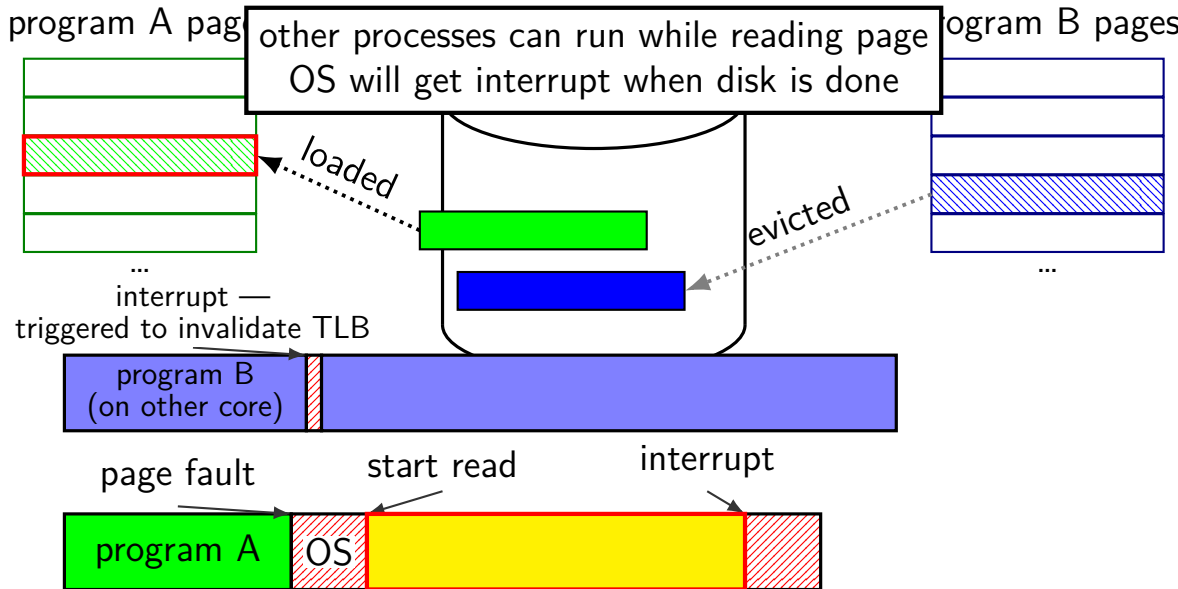
evicted

program B pages



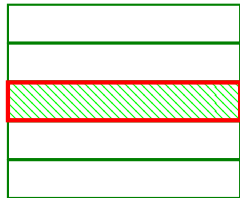
...

TLB shutdown



TLB shutdown

program A pages



...

process A's page table updated and restarted from point of fault

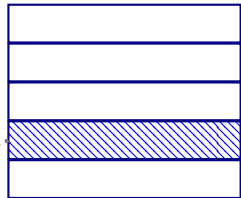
loaded



evicted



program B pages



...

interrupt —
triggered to invalidate TLB



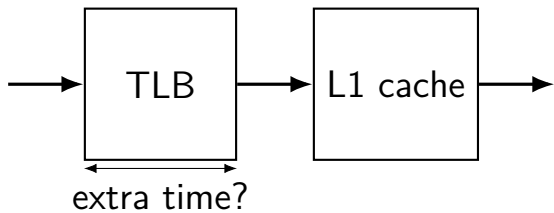
page fault

start read

interrupt



TLBs and performance



L1 caches and page numbers (Intel Skylake)

physical address (48 bits)		
PPN (36 bit)	page offset (12 bit)	
L1 cache tag (36 bit)	L1 index (6 bit)	L1 offset (6 bit)

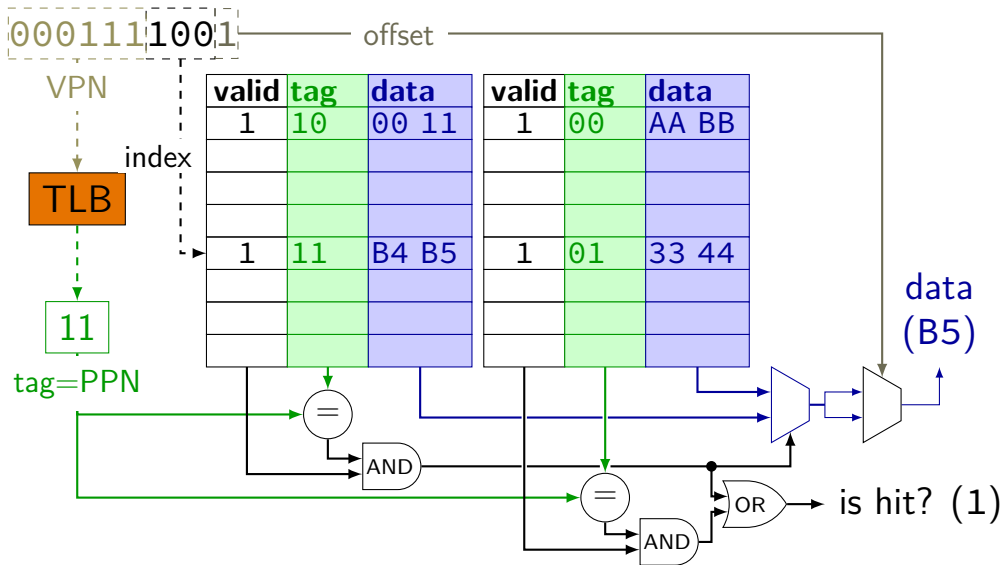
L1 caches and page numbers (Intel Skylake)

physical address (48 bits)		
PPN (36 bit)	page offset (12 bit)	
L1 cache tag (36 bit)	L1 index (6 bit)	L1 offset (6 bit)

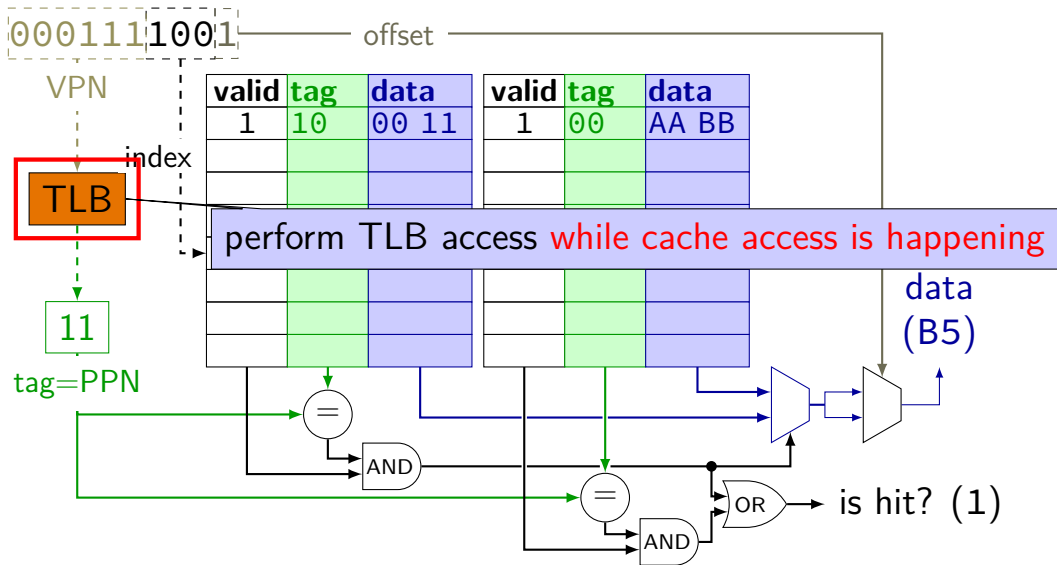
not a coincidence

why did Intel make this decision?

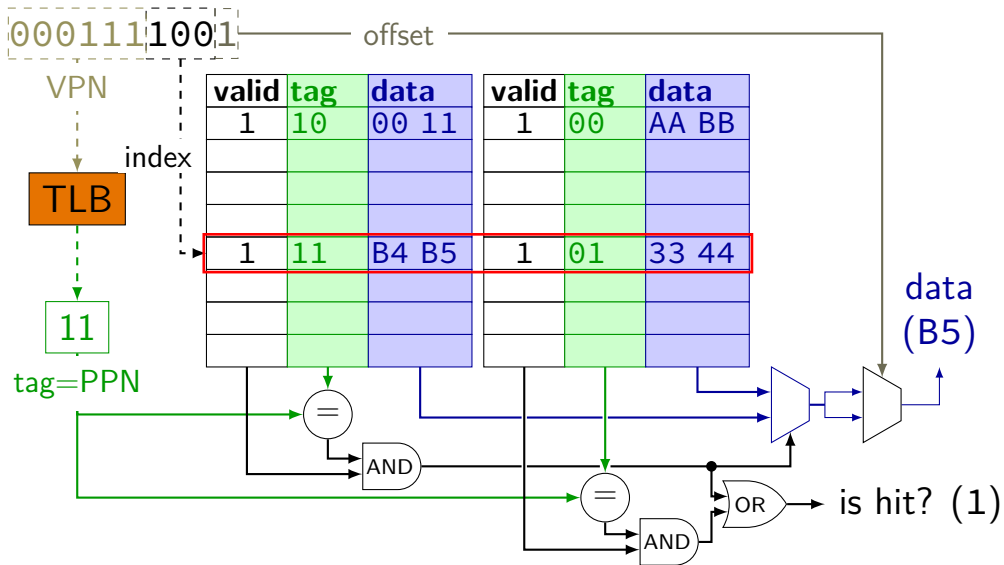
overlapping TLB and cache access



overlapping TLB and cache access



overlapping TLB and cache access



virtually-indexed, physically-tagged

called virtually-indexed, physically-tagged cache

requirement: **index contained entirely in page offset**

do not need to do translation to start cache access

tag overlaps with PPN

example: tag=PPN

(but tag could include part of page offset, too)

do TLB access **while retrieving cache set**

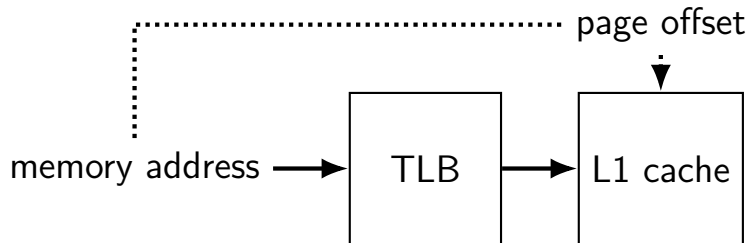
most common design in current processors

reason for highly associative (e.g. 8-way) L1 caches

physical caches

so far: caches use **physical addresses**:

means cache lookup can't complete without TLB
(and can't start without index from physical address)



address splitting

16-bit virtual addresses

64-byte pages

256B, 8-way L1 cache with 16B blocks

can TLB and cache access overlap?

x86-64 page table entries (1)

6	6	6	6	5	5	5	5	5	5	5	5		M ¹	M-1			3	3	3	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1	1	0									
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								
X	Prot. Key ⁴	Ignored	Rsvd.	Address of 4KB page frame														Ign.	G	P	A	D	A	P	C	W	T	U	/	R	/	1	PTE: 4KB page														
D				Ignored																																										0	PTE: not present

present = valid

R/W = writes allowed?

U/S = kernel-only? ("user/supervisor")

XD = execute-disable?

A = accessed? (MMU sets to 1 on page read/write)

D = helps support replacement policies for swapping

x86-64 page table entries (1)

6 6 6 6 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 1 1 1 1 1 1 1															3 2 1 0 9 8 7 6 5 4 3 2 1											0
X D	Prot. Key ⁴	Ignored	Rsvd.	Address of 4KB page frame															Ign.	G	P A T	D	A	P C D	P W T	U / S	R / W	1	PTE: 4KB page										
																												0	PTE: not present										

present = valid

R/W = writes allowed?

U/S = kernel-only? (“user/supervisor”)

XD = execute-disable?

A = accessed? (MMU sets to 1 on page read/write)

D = dirty? (MMU sets to 1 on page write)

helps support writeback policy for swapping

