# cache accesses and C code (1)

```
int scaleFactor;

int scaleByFactor(int value) {
    return value * scaleFactor;
}
```

```
scaleByFactor:
    movl scaleFactor, %eax
    imull %edi, %eax
    ret
```

exericse: what data cache accesses does this function do?

# cache accesses and C code (1)

```
int scaleFactor;

int scaleByFactor(int value) {
    return value * scaleFactor;
}
```

```
scaleByFactor:
    movl scaleFactor, %eax
    imull %edi, %eax
    ret
```

exericse: what data cache accesses does this function do?
    4-byte read of scaleFactor
    8-byte read of return address

# possible scaleFactor use

```
for (int i = 0; i < size; ++i) {
    array[i] = scaleByFactor(array[i]);
}
```

# misses and code (2)

```
scaleByFactor:
    movl scaleFactor, %eax
    imull %edi, %eax
    ret
```

suppose each time this is called in the loop:

    return address located at address 0x7ffffffe43b8

    scaleFactor located at address 0x6bc3a0

with direct-mapped 32KB cache w/64 B blocks, what is their:

|        | return address | scaleFactor |
|--------|----------------|-------------|
| tag    |                |             |
| index  |                |             |
| offset |                |             |

# misses and code (2)

```
scaleByFactor:
    movl scaleFactor, %eax
    imull %edi, %eax
    ret
```

suppose each time this is called in the loop:
    return address located at address 0x7fffffffe43b8
    scaleFactor located at address 0x6bc3a0

with direct-mapped 32KB cache w/64 B blocks, what is their:

|        | return address | scaleFactor |
|--------|----------------|-------------|
| tag    | 0xfffffffc     | 0xd7        |
| index  | 0x10e          | 0x10e       |
| offset | 0x38           | 0x20        |

# misses and code (2)

```
scaleByFactor:
    movl scaleFactor, %eax
    imull %edi, %eax
    ret
```

suppose each time this is called in the loop:

return address located at address 0x7fffffe43b8
scaleFactor located at address 0x6bc3a0

with direct-mapped 32KB cache w/64 B blocks, what is their:

|        | return address | scaleFactor |
|--------|----------------|-------------|
| tag    | 0xfffffffc     | 0xd7        |
| index  | 0x10e          | 0x10e       |
| offset | 0x38           | 0x20        |

# conflict miss coincidences?

obviously I set that up to have the same index
    have to use exactly the right amount of stack space…

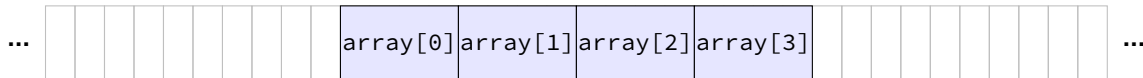but one of the reasons we'll want something better than direct-mapped cache

# C and cache misses (warmup 1)

```
int array[4];
...
int even_sum = 0, odd_sum = 0;
even_sum += array[0];
odd_sum += array[1];
even_sum += array[2];
odd_sum += array[3];
```

Assume everything but `array` is kept in registers (and the compiler does not do anything funny).

How many data cache misses on a 1-set direct-mapped cache with 8B blocks?

# some possiblities



... | | | | | | | | | | array[0] | array[1] | array[2] | array[3] | | | | | | | | | | | ...

Q1: how do cache blocks correspond to array elements?
not enough information provided!

# aside: alignment

compilers and malloc/new implementations usually try align values

align = make address be multiple of something

most important reason: don't cross cache block boundaries

# C and cache misses (warmup 2)

```
int array[4];
int even_sum = 0, odd_sum = 0;
even_sum += array[0];
even_sum += array[2];
odd_sum += array[1];
odd_sum += array[3];
```

Assume everything but array is kept in registers (and the compiler does not do anything funny).

Assume array[0] at beginning of cache block.

How many data cache misses on a 1-set direct-mapped cache with 8B blocks?

# C and cache misses (warmup 3)

```c
int array[8];
...
int even_sum = 0, odd_sum = 0;
even_sum += array[0];
odd_sum += array[1];
even_sum += array[2];
odd_sum += array[3];
even_sum += array[4];
odd_sum += array[5];
even_sum += array[6];
odd_sum += array[7];
```

Assume everything but `array` is kept in registers (and the compiler does not do anything funny), and array[0] at beginning of cache block.

How many data cache misses on a **2**-set direct-mapped cache with 8B blocks?

# C and cache misses (warmup 4)

```c
int array[8];
...
int even_sum = 0, odd_sum = 0;
even_sum += array[0];
even_sum += array[2];
even_sum += array[4];
even_sum += array[6];
odd_sum += array[1];
odd_sum += array[3];
odd_sum += array[5];
odd_sum += array[7];
```

Assume everything but `array` is kept in registers (and the compiler does not do anything funny).

How many data cache misses on a **2**-set direct-mapped cache with 8B blocks?

# arrays and cache misses (1)

```
int array[1024]; // 4KB array
int even_sum = 0, odd_sum = 0;
for (int i = 0; i < 1024; i += 2) {
    even_sum += array[i + 0];
    odd_sum +=  array[i + 1];
}
```

Assume everything but `array` is kept in registers (and the compiler does not do anything funny).

How many *data cache misses* on initially empty 2KB direct-mapped cache with 16B cache blocks?

# arrays and cache misses (2)

```c
int array[1024]; // 4KB array
int even_sum = 0, odd_sum = 0;
for (int i = 0; i < 1024; i += 2)
    even_sum += array[i + 0];
for (int i = 0; i < 1024; i += 2)
    odd_sum +=  array[i + 1];
```

Assume everything but `array` is kept in registers (and the compiler does not do anything funny).

How many *data cache misses* on initially empty 2KB direct-mapped cache with 16B cache blocks?
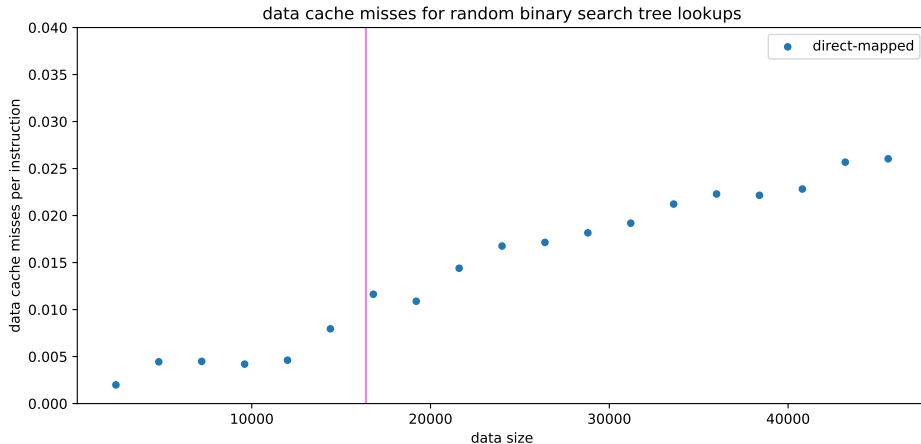
# arrays and cache misses (2b)

```
int array[1024]; // 4KB array
int even_sum = 0, odd_sum = 0;
for (int i = 0; i < 1024; i += 2)
    even_sum += array[i + 0];
for (int i = 0; i < 1024; i += 2)
    odd_sum +=  array[i + 1];
```

Assume everything but `array` is kept in registers (and the compiler does not do anything funny).

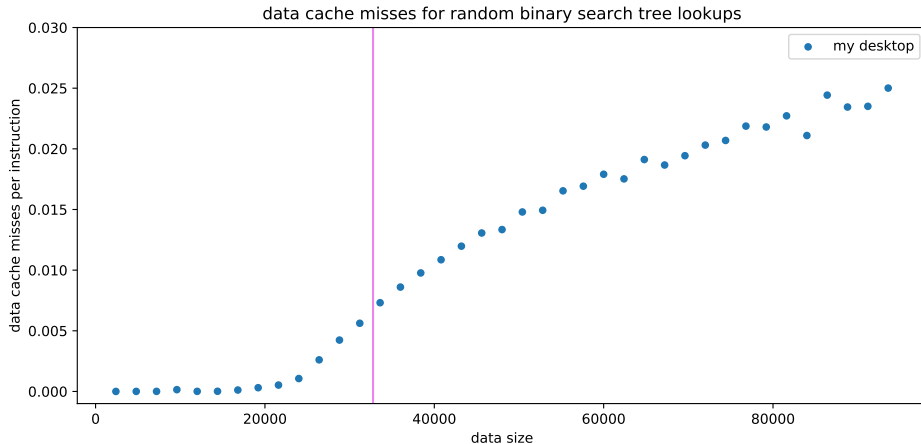How many *data cache misses* on initially empty 4KB direct-mapped cache with 16B cache blocks?

# simulated misses: BST lookups



data cache misses for random binary search tree lookups

(simulated 16KB direct-mapped data cache; excluding BST setup)

# actual misses: BST lookups



data cache misses for random binary search tree lookups
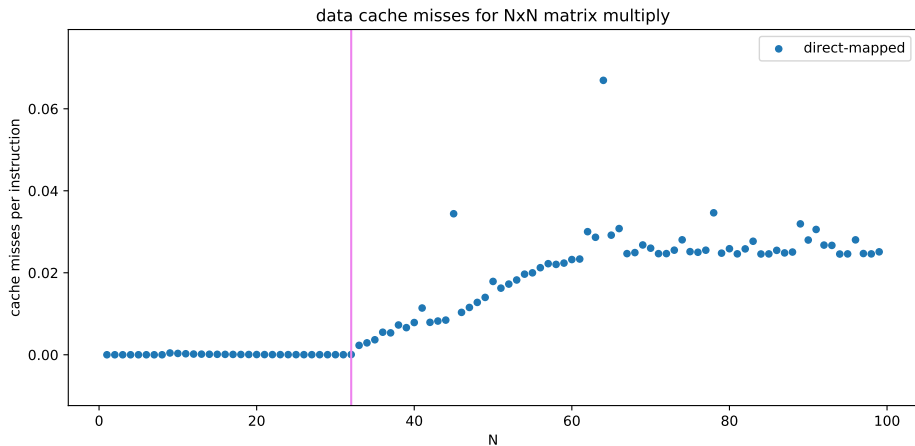
(actual 32KB more complex data cache)
(only one set of measurements + other things on machine + excluding initial load)

# simulated misses: matrix multiplies



data cache misses for NxN matrix multiply

(simulated 16KB direct-mapped data cache; excluding initial load)

# actual misses: matrix multiplies



cache misses for NxN matrix multiply

(actual 32KB more complex data cache; excluding matrix initial load)
(only one set of measurements + other things on machine)

# misses with skipping

```
int array1[512]; int array2[512];
...
for (int i = 0; i < 512; i += 1)
    sum += array1[i] * array2[i];
}
```

Assume everything but array1, array2 is kept in registers (and the compiler does not do anything funny).

About how many *data cache misses* on a 2KB direct-mapped cache with 16B cache blocks?
Hint: depends on relative placement of array1, array2

# best/worst case

array1[i] and array2[i] always different sets:
> = distance from array1 to array2 not multiple of # sets $\times$ bytes/set
> 2 misses every 4 i
> blocks of 4 array1[X] values loaded, then used 4 times before loading next block
> (and same for array2[X])

array1[i] and array2[i] same sets:
> = distance from array1 to array2 is multiple of # sets $\times$ bytes/set
> 2 misses every i
> block of 4 array1[X] values loaded, one value used from it,
> then, block of 4 array2[X] values replaces it, one value used from it, …

# worst case in practice?

two rows of matrix?

often sizeof(row) bytes apart

if the row size is multiple of number of sets $\times$ bytes per block, oops!

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 0 | | | 0 | | |
| 1 | 0 | | | 0 | | |

multiple places to put values with same index
avoid misses from two active values using same set
("conflict misses"))

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 0 | | **set 0** | 0 | | |
| 1 | 0 | | **set 1** | 0 | | |

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 0 | | | 0 | | |
| 1 | 0 | | | 0 | | |

$m = 8$ bit addresses
$S = 2 = 2^s$ sets
$s = 1$ (set) index bits

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$t = m - (s + b) = 6$ tag bits

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 0 | | |
| 1 | 0 | | | 0 | | |

| address (hex) | result |
|---------------|--------|
| 00000000 (00) | miss |
| 00000001 (01) | |
| 01100011 (63) | |
| 01100001 (61) | |
| 01100010 (62) | |
| 00000000 (00) | |
| 01100100 (64) | |

tag  index offset

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 0 | | |
| 1 | 0 | | | 0 | | |

| address (hex) | result |
|---------------|--------|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | |
| 01100001 (61) | |
| 01100010 (62) | |
| 00000000 (00) | |
| 01100100 (64) | |

tag   index offset

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 0 | | |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | | |

| address (hex) | result |
|---------------|--------|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | |
| 01100010 (62) | |
| 00000000 (00) | |
| 01100100 (64) | |

tag  index offset

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 1 | 011000 | mem[0x60]<br>mem[0x61] |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | | |

| address (hex) | result |
|---------------|--------|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | |
| 00000000 (00) | |
| 01100100 (64) | |

tag  index offset

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 1 | 000000 | mem[0x00] <br> mem[0x01] | 1 | 011000 | mem[0x60] <br> mem[0x61] |
| 1 | 1 | 011000 | mem[0x62] <br> mem[0x63] | 0 | | |

| address (hex) | result |
|---------------|--------|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | hit |
| 00000000 (00) | |
| 01100100 (64) | |

tag  index offset

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 1 | 000000 | mem[0x00] mem[0x01] | 1 | 011000 | mem[0x60] mem[0x61] |
| 1 | 1 | 011000 | mem[0x62] mem[0x63] | 0 | | |

| address (hex) | result |
|---------------|--------|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | hit |
| 00000000 (00) | hit |
| 01100100 (64) | |

tag  index offset

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 1 | 000000 | mem[0x00] mem[0x01] | 1 | 011000 | mem[0x60] mem[0x61] |
| 1 | 1 | 011000 | mem[0x62] mem[0x63] | 0 | | |

| address (hex) | result |
|---------------|--------|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | hit |
| 00000000 (00) | hit |
| 01100100 (64) | miss |

needs to replace block in set 0!

tag  index offset

# adding associativity

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|--------|----------------------|-------|--------|----------------------|
| 0 | 1 | 000000 | mem[0x00] mem[0x01] | 1 | 011000 | mem[0x60] mem[0x61] |
| 1 | 1 | 011000 | mem[0x62] mem[0x63] | 0 | | |

| address (hex) | result |
|---------------|--------|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | hit |
| 00000000 (00) | hit |
| 01100100 (64) | miss |

tag  indexoffset

# associative lookup possibilities

none of the blocks for the index are valid

none of the valid blocks for the index match the tag
 something else is stored there

one of the blocks for the index is valid and matches the tag

# replacement policies

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value |
|-------|-------|-----|-------|-------|-----|-------|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 1 | 011000 | mem[0x60]<br>mem[0x61] |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | | |

| address (hex) | result |
|---------------|--------|
| 000 | |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | hit |
| 00000000 (00) | hit |
| 01100100 (64) | miss |

how to decide where to insert 0x64?

# replacement policies

2-way set associative, 2 byte blocks, 2 sets

| index | valid | tag | value | valid | tag | value | LRU |
|-------|-------|--------|---------------------------|-------|--------|---------------------------|-----|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 1 | 011000 | mem[0x60]<br>mem[0x61] | 1 |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | | | 1 |

| address (hex) | result |
|--------------------|--------|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | hit |
| 00000000 (00) | hit |
| 01100100 (64) | miss |

track which block was read least recently
updated on every access

28

# example replacement policies

least recently used
> take advantage of temporal locality
> at least $\lceil \log_2(E!) \rceil$ bits per set for $E$-way cache
>> (need to store order of all blocks)

approximations of least recently used
> implementing least recently used is expensive
> really just need "avoid recently used" — much faster/simpler
> good approximations: $E$ to $2E$ bits

first-in, first-out
> counter per set — where to replace next

(pseudo-)random
> no extra information!
> actually works pretty well in practice

# associativity terminology

direct-mapped — one block per set

$E$-way set associative — $E$ blocks per set
$E$ ways in the cache

fully associative — one set total (everything in one set)

# Tag-Index-Offset formulas

| | |
|---|---|
| $m$ | memory addresses bits |
| $E$ | number of blocks per set ("ways") |
| $S = 2^s$ | number of sets |
| $s$ | (set) index bits |
| $B = 2^b$ | block size |
| $b$ | (block) offset bits |
| $t = m - (s + b)$ | tag bits |
| $C = B \times S \times E$ | cache size (excluding metadata) |

# arrays and cache misses (2)

```
int array[1024]; // 4KB array
int even_sum = 0, odd_sum = 0;
for (int i = 0; i < 1024; i += 2)
    even_sum += array[i + 0];
for (int i = 0; i < 1024; i += 2)
    odd_sum +=  array[i + 1];
```

Assume everything but `array` is kept in registers (and the compiler does not do anything funny).

How many *data cache misses* on initially empty 2KB direct-mapped cache with 16B cache blocks? Would a set-associative cache be better?

# explanation

2-way, 2KB set associative cache, 16B blocks

4 offset bits, 6 index bits

so addresses multiples $2^{10}$ bytes apart differ only in tag bits

example: array[0→3], array[256→259], array[512→515], array[768→771]
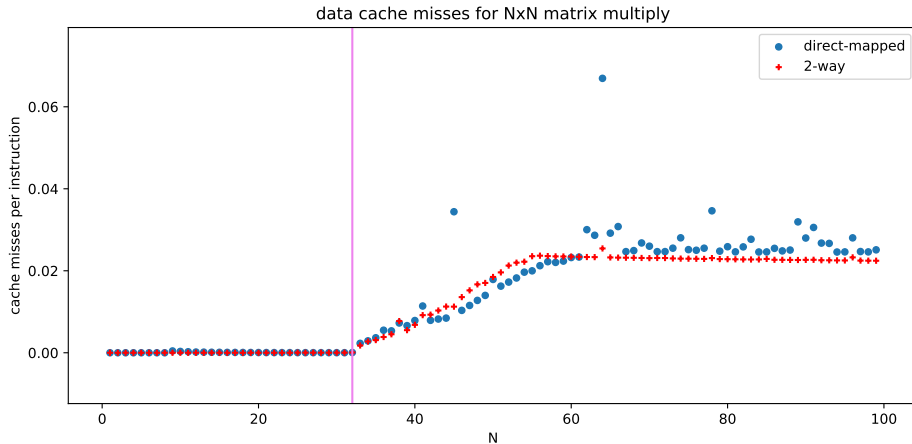
those all use the same set

but sets only holds 2 things

all misses

# simulated misses: BST lookups



data cache misses for random binary search tree lookups

# simulated misses: matrix multiplies



data cache misses for NxN matrix multiply

# handling writes

what about writing to the cache?

two decision points:

if the value is not in cache, do we add it?
    if yes: need to load rest of block
    if no: missing out on locality?

if value is in cache, when do we update next level?
    if immediately: extra writing
    if later: need to remember to do so

# allocate on write?

processor writes less than whole cache block

block not yet in cache
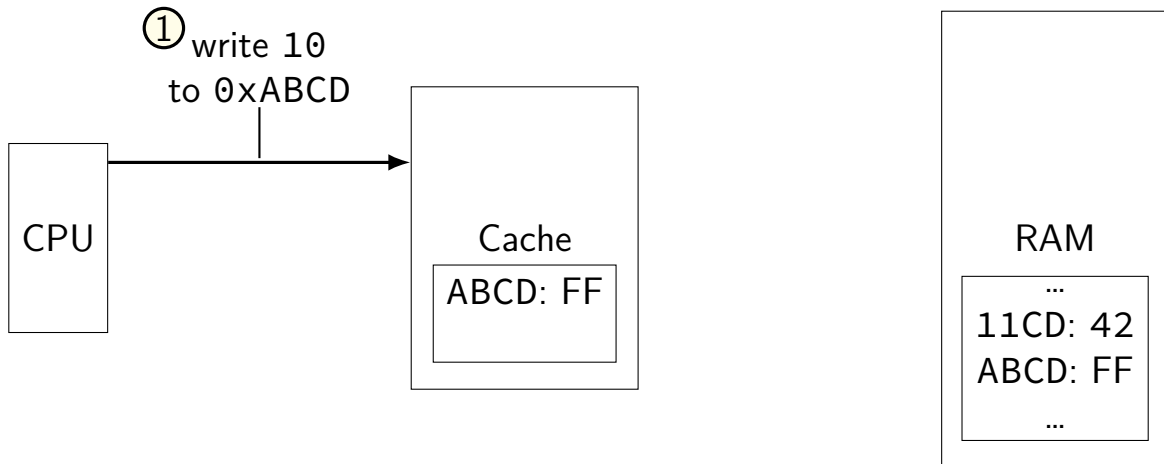
two options:

write-allocate
     fetch rest of cache block, replace written part
     (then follow write-through or write-back policy)

write-no-allocate
     don't use cache at all (send write to memory *instead*)
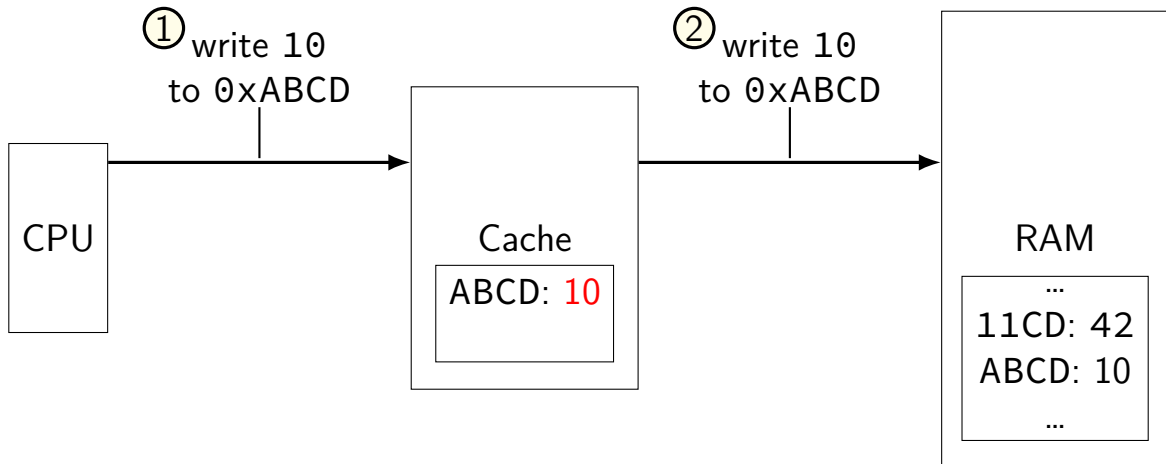     guess: not read soon?

# write-through v. write-back
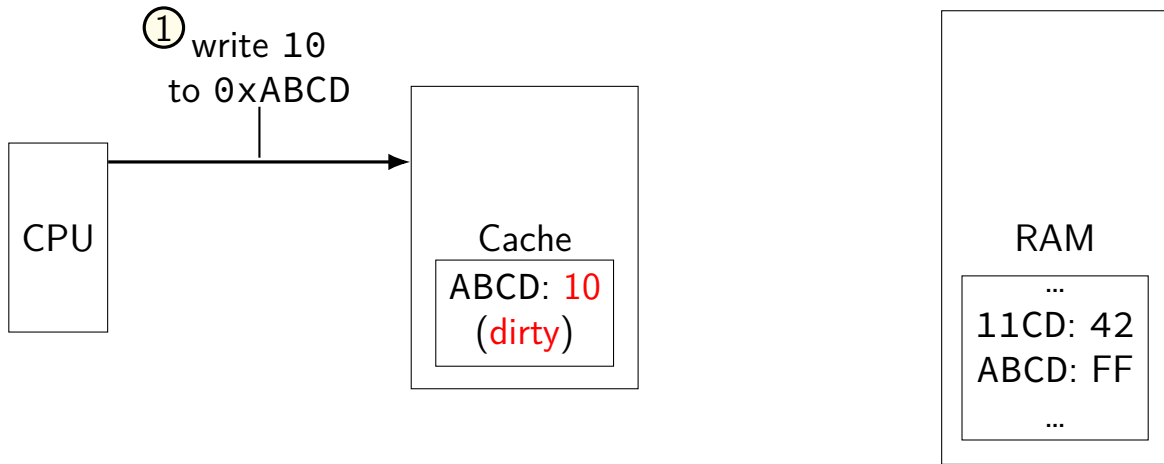
**option 1: write-through**
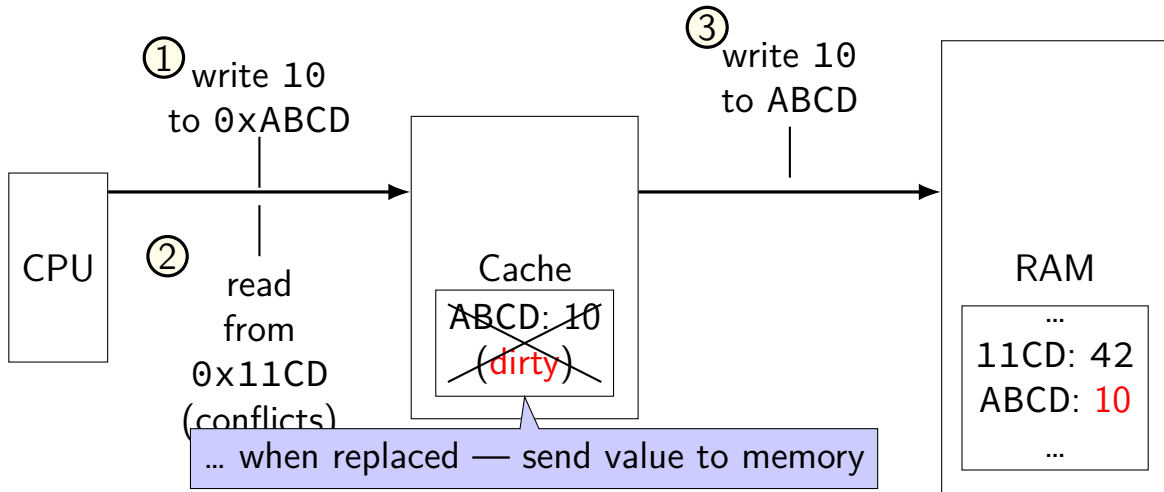
# write-through v. write-back

**option 1: write-through**

# write-through v. write-back

## option 2: write-back

# write-through v. write-back
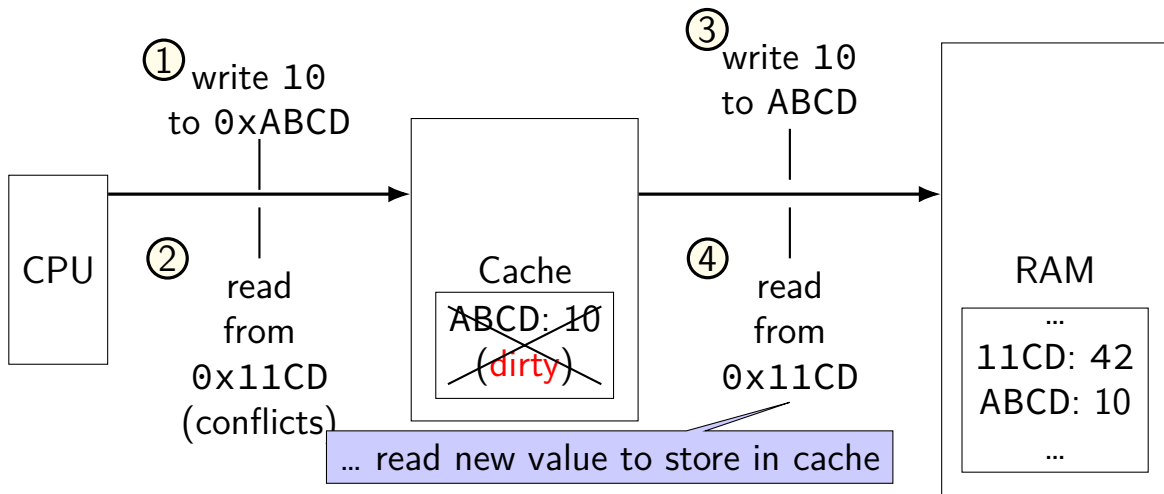
## option 2: write-back



CPU

① write 10 to 0xABCD

② read from 0x11CD (conflicts)

Cache

~~ABCD: 10~~
~~(dirty)~~

… when replaced — send value to memory

③ write 10 to ABCD

RAM

…
11CD: 42
ABCD: 10
…

# write-through v. write-back

# writeback policy

changed value!

2-way set associative, 4 byte blocks, 2 sets

| index | valid | tag | value | dirty | valid | tag | value | dirty | LRU |
|-------|-------|-----|-------|-------|-------|-----|-------|-------|-----|
| 0 | 1 | 000000 | mem[0x00] mem[0x01] | 0 | 1 | 011000 | mem[0x60]* mem[0x61]* | 1 | 1 |
| 1 | 1 | 011000 | mem[0x62] mem[0x63] | 0 | 0 | | | | 0 |

1 = dirty (different than memory)
needs to be written if evicted

# write-allocate + write-back

2-way set associative, LRU, writeback

| index | valid | tag | value | dirty | valid | tag | value | dirty | LRU |
|-------|-------|--------|----------------------------|-------|-------|--------|------------------------------|-------|-----|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 0 | 1 | 011000 | mem[0x60]*<br>mem[0x61]* | 1 | 1 |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | 0 | | | | 0 |

writing 0xFF into address 0x04?
index 0, tag 000001

# write-allocate + write-back

2-way set associative, LRU, writeback

| index | valid | tag | value | dirty | valid | tag | value | dirty | LRU |
|-------|-------|--------|------------------------------|-------|-------|--------|---------------------------------|-------|-----|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 0 | 1 | 011000 | mem[0x60]*<br>mem[0x61]* | 1 | 1 |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | 0 | | | | 0 |

writing 0xFF into address 0x04?
index 0, tag 000001
step 1: find least recently used block

40

# write-allocate + write-back

2-way set associative, LRU, writeback

| index | valid | tag | value | dirty | valid | tag | value | dirty | LRU |
|-------|-------|--------|---------------------------|-------|-------|--------|-----------------------------|-------|-----|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 0 | 1 | 011000 | mem[0x60]*<br>mem[0x61]* | ~~1~~ | 1 |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | 0 | | | | 0 |

writing 0xFF into address 0x04?
index 0, tag 000001
step 1: find least recently used block
step 2: possibly writeback old block

# write-allocate + write-back

2-way set associative, LRU, writeback

| index | valid | tag | value | dirty | valid | tag | value | dirty | LRU |
|-------|-------|-----|-------|-------|-------|-----|-------|-------|-----|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 0 | 1 | 000001 | 0xFF<br>mem[0x05] | 1 | 0 |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | 0 | | | | 0 |

writing 0xFF into address 0x04?
index 0, tag 000001
step 1: find least recently used block
step 2: possibly writeback old block
step 3a: read in new block – to get mem[0x05]
step 3b: update LRU information

# write-no-allocate + write-back

2-way set associative, LRU, writeback

| index | valid | tag | value | dirty | valid | tag | value | dirty | LRU |
|-------|-------|-----|-------|-------|-------|-----|-------|-------|-----|
| 0 | 1 | 000000 | mem[0x00]<br>mem[0x01] | 0 | 1 | 011000 | mem[0x60]*<br>mem[0x61]* | 1 | 1 |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | 0 | | | | 0 |

writing 0xFF into address 0x04?
step 1: is it in cache yet?
step 2: no, just send it to memory

# exercise (1)

2-way set associative, LRU, write-allocate, writeback

| index | valid | tag | value | dirty | valid | tag | value | dirty | LRU |
|-------|-------|-----|-------|-------|-------|-----|-------|-------|-----|
| 0 | 1 | 001100 | mem[0x30]<br>mem[0x31] | 0 | 1 | 010000 | mem[0x40]⋆<br>mem[0x41]⋆ | 1 | 0 |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | 1 | 001100 | mem[0x32]⋆<br>mem[0x33]⋆ | 1 | 1 |

for each of the following accesses, performed alone, would it
require (a) reading a value from memory (or next level of cache)
and (b) writing a value to the memory (or next level of cache)?

writing 1 byte to 0x33
reading 1 byte from 0x52
reading 1 byte from 0x50

# exercise (2)

2-way set associative, LRU, <span style="color:red">write-no-allocate, write-through</span>

| index | valid | tag | value | valid | tag | value | LRU |
|-------|-------|-----|-------|-------|-----|-------|-----|
| 0 | 1 | 001100 | mem[0x30] mem[0x31] | 1 | 010000 | mem[0x40] mem[0x41] | 0 |
| 1 | 1 | 011000 | mem[0x62] mem[0x63] | 1 | 001100 | mem[0x32] mem[0x33] | 1 |

for each of the following accesses, performed alone, would it require (a) reading a value from memory and (b) writing a value to the memory?

> writing 1 byte to 0x33
> reading 1 byte from 0x52
> reading 1 byte from 0x50

# fast writes



write appears to complete immediately when placed in buffer
memory can be much slower

# cache miss types

common to categorize misses:
    roughly "cause" of miss assuming cache block size fixed

*compulsory* (or *cold*) — first time accessing something
    adding more sets or blocks/set wouldn't change

*conflict* — sets aren't big/flexible enough
    a fully-associtive (1-set) cache of the same size would have done better

*capacity* — cache was not big enough

*coherence* — from sync'ing cache with other caches
    only issue with multiple cores

# making any cache look bad

1. access enough blocks, to fill the cache

2. access an additional block, replacing something

3. access last block replaced

4. access last block replaced

5. access last block replaced

…

but — typical real programs have locality

# cache optimizations

(assuming typical locality + keeping cache size constant if possible…)

|  | miss rate | hit time | miss penalty |
|---|---|---|---|
| increase cache size | better | worse | — |
| increase associativity | better | worse | worse? |
| increase block size | depends | worse | worse |
| add secondary cache | — | — | better |
| write-allocate | better | — | ? |
| writeback | — | — | ? |
| LRU replacement | better | ? | worse? |
| prefetching | better | — | — |

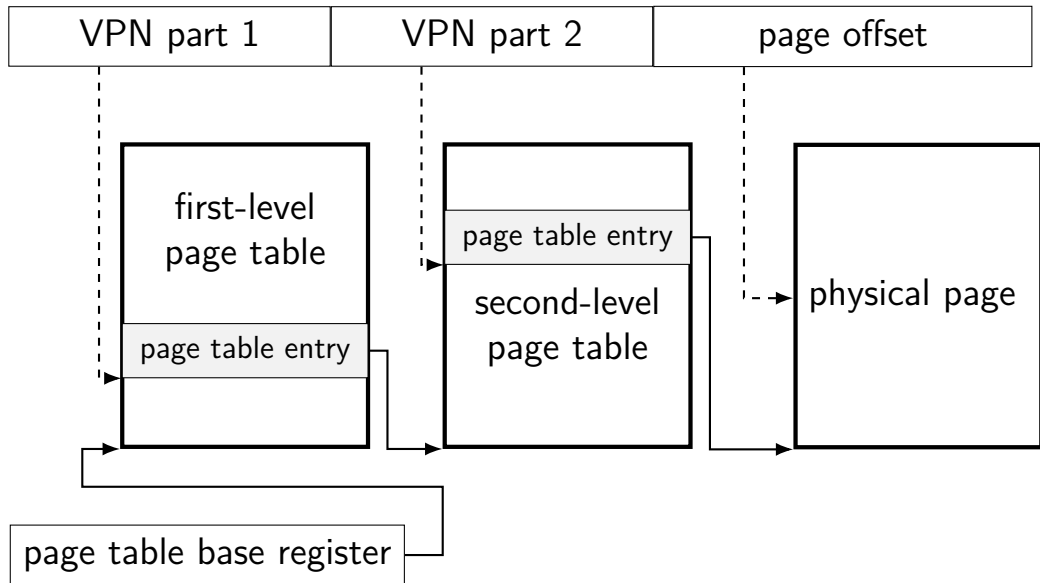prefetching = guess what program will use, access in advance

$$\text{average time} = \text{hit time} + \text{miss rate} \times \text{miss penalty}$$
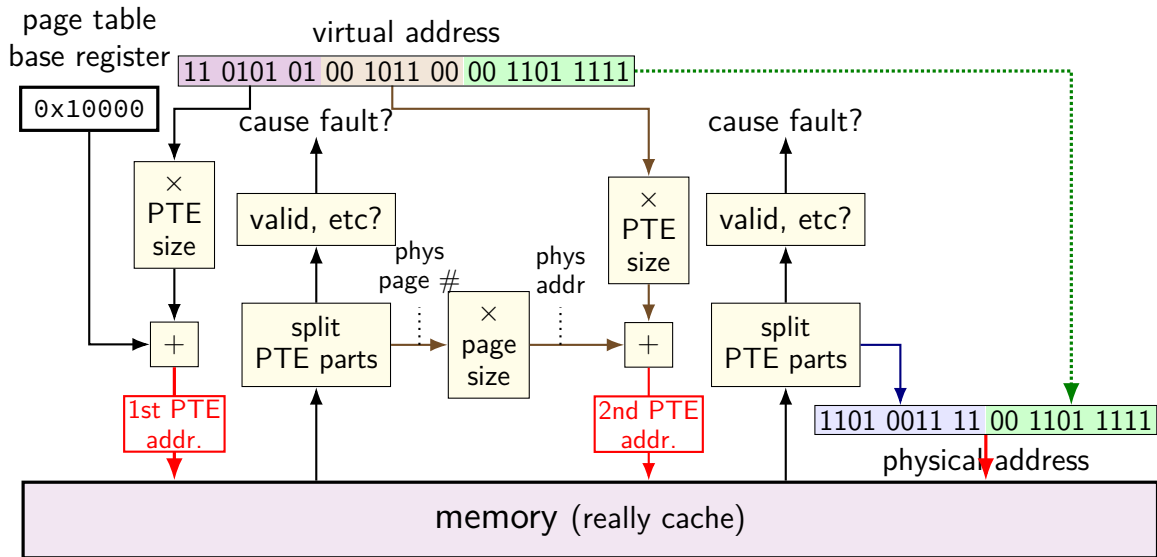
# cache optimizations by miss type

(assuming other listed parameters remain constant)

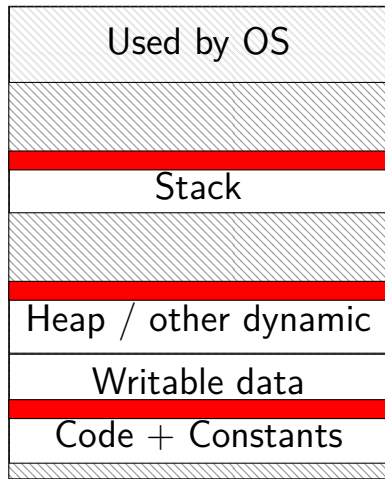|                        | capacity      | conflict      | compulsory   |
|------------------------|---------------|---------------|--------------|
| increase cache size    | fewer misses  | fewer misses  | —            |
| increase associativity | —             | fewer misses  | —            |
| increase block size    | more misses?  | more misses?  | fewer misses |
|                        |               |               |              |
| LRU replacement        | —             | fewer misses  | —            |
| prefetching            | —             | —             | fewer misses |

## another view

# two-level page table lookup



page table
base register

virtual address

`0x10000`

11 0101 01 00 1011 00 00 1101 1111

cause fault?

×
PTE
size

valid, etc?

+

1st PTE
addr.

split
PTE parts

phys
page #

×
page
size

phys
addr

cause fault?

×
PTE
size

valid, etc?

+

2nd PTE
addr.

split
PTE parts

1101 0011 11 00 1101 1111

physical address

memory (really cache)

# cache accesses and multi-level PTs

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

L1 cache hits — typically a couple cycles each?

so add 8 cycles to each program memory access?

not acceptable

# program memory active sets

| |
|---|
| Used by OS |
| |
| Stack |
| |
| Heap / other dynamic |
| Writable data |
| Code + Constants |

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 <span style="color:red">excellent temporal locality</span>

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 <span style="color:red">whole functions</span>, 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

caled 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

caled 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 tables

# page table entry cache

caled 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)
(and if spatial locality, maybe use larger page size?)

# page table entry cache

caled 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
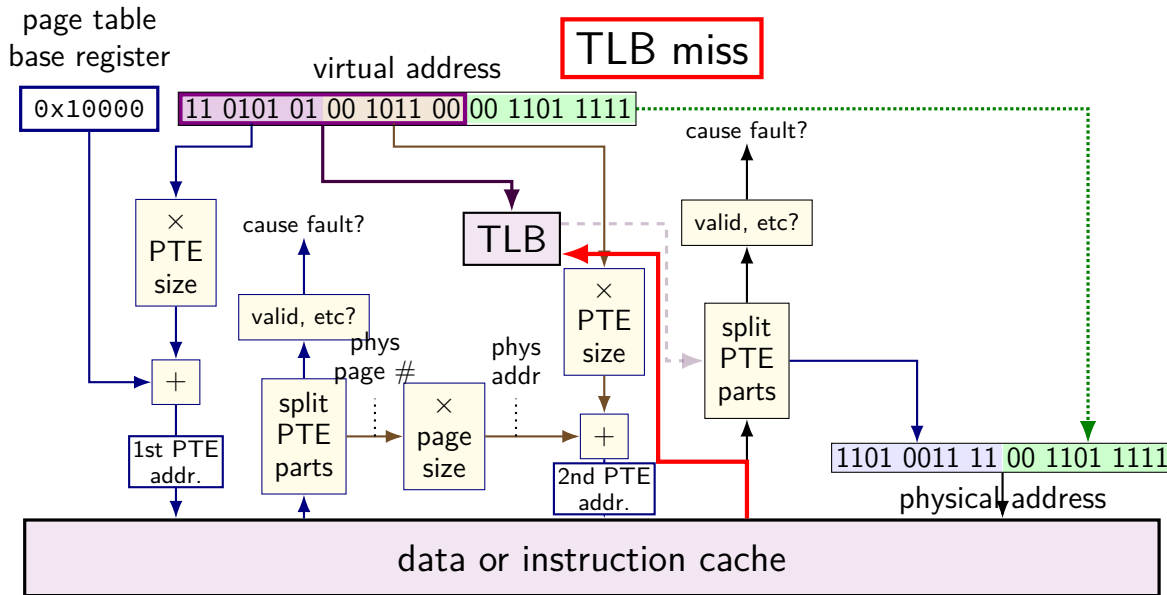enables highly associative cache designs

# TLB and multi-level page tables

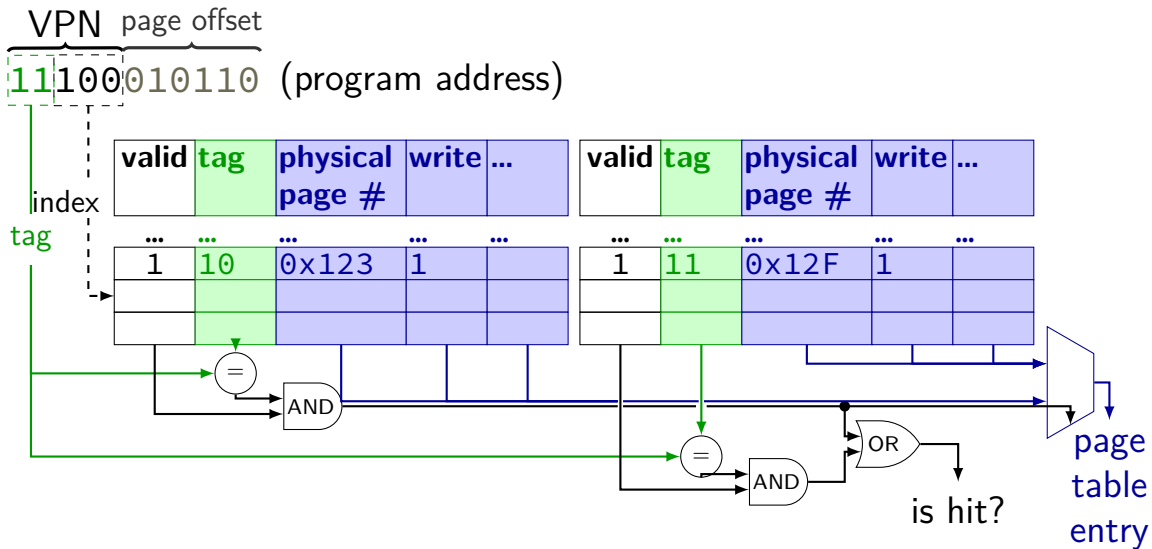TLB caches <span style="color:red">valid last-level page table entries</span>

doesn't matter which last-level page table
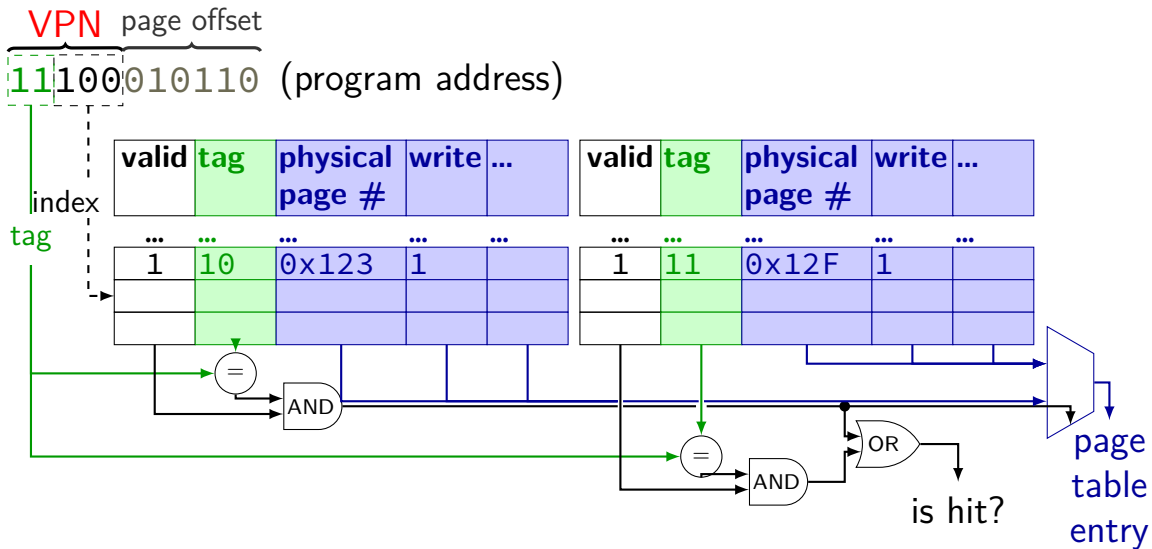
means TLB output can be used directly to form address
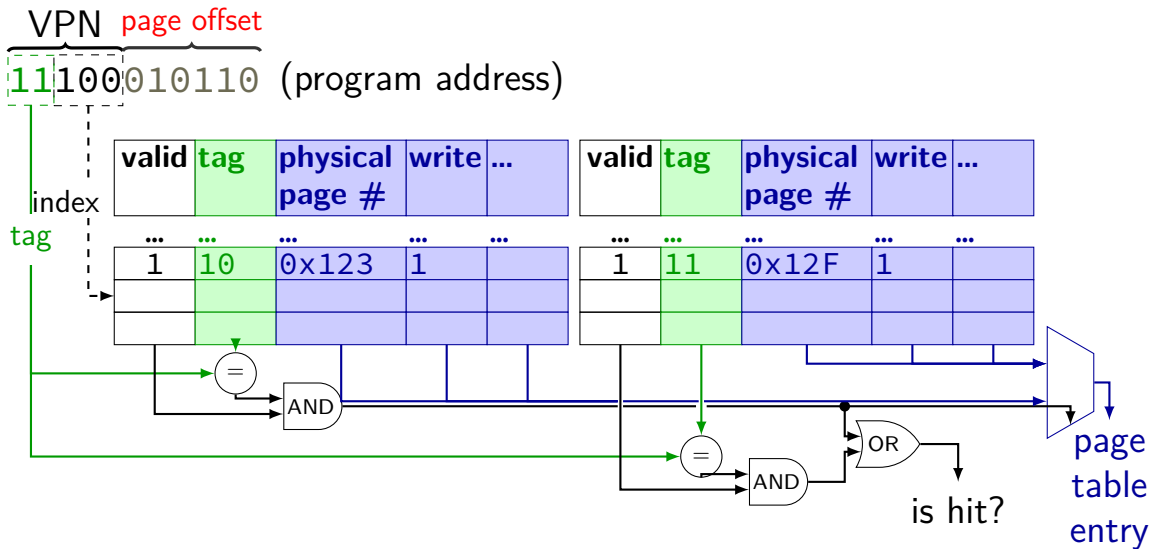
# TLB and two-level lookup

# TLB and two-level lookup



page table base register

0x10000

virtual address

11 0101 01 00 1011 00 00 1101 1111

TLB miss

cause fault?

× PTE size

cause fault?

valid, etc?

TLB

cause fault?

valid, etc?

× PTE size

split PTE parts

1st PTE addr.

split PTE parts

phys page #

× page size

phys addr

+

2nd PTE addr.

split PTE parts

1101 0011 11 00 1101 1111

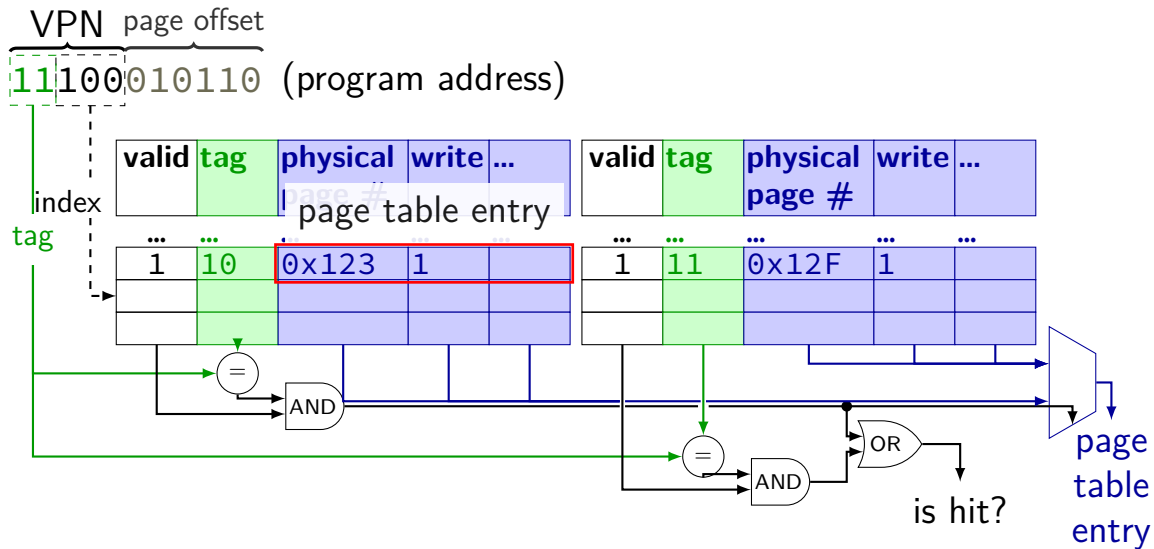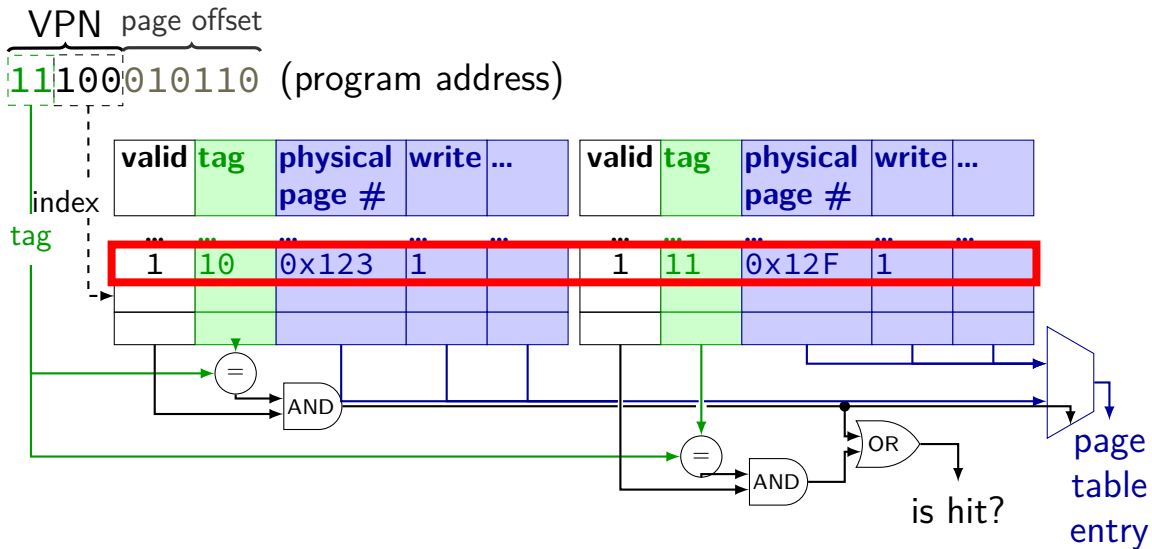physical address

data or instruction cache

# TLB organization (2-way set associative)

# TLB organization (2-way set associative)

# TLB organization (2-way set associative)

# TLB organization (2-way set associative)



VPN page offset

`11 100 010110` (program address)

page table entry

# TLB organization (2-way set associative)

# exercise: TLB access pattern (setup)

4-entry, 2-way TLB, LRU replacement policy, initially empty

4096 byte pages

how many index bits?

TLB index of virtual address 0x12345?

# exercise: TLB access pattern

4-entry, 2-way TLB, LRU replacement policy, initially empty

4096 byte pages

| type  | virtual    | physical  |
|-------|------------|-----------|
| read  | 0x440030   | 0x554030  |
| write | 0x440034   | 0x554034  |
| read  | 0x7FFFE008 | 0x556008  |
| read  | 0x7FFFE000 | 0x556000  |
| read  | 0x7FFFDFF8 | 0x5F8FF8  |
| read  | 0x664080   | 0x5F9080  |
| read  | 0x440038   | 0x554038  |
| write | 0x7FFFDFF0 | 0x5F8FF0  |

which are TLB hits? which are TLB misses? final contents of TLB?

**backup slides**

# arrays and cache misses (3)

```
int sum; int array[1024]; // 4KB array
for (int i = 8; i < 1016; i += 1) {
    int local_sum = 0;
    for (int j = i - 8; j < i + 8; j += 1) {
        local_sum += array[i] * (j - i);
    }
    sum += (local_sum - array[i]);
}
```

Assume everything but `array` is kept in registers (and the compiler does not do anything funny).

How many *data cache misses* on initially empty 2KB direct-mapped cache with 16B cache blocks?