# Exam Review

# SEQ without stages

register file
PC+9
srcA    R[srcA]
rB  %rsp   srcB    R[srcB]
0xF   dstM
0xF   dstE
%rsp
ALU
aluA
valE
aluB
8
0
Data in
Data out
Addr in
write?
Instr.
Mem.
PC
next R[dstM]
next R[dstE]
instr.
length

# SEQ with stages

decode
execute
memory
fetch
register file
PC+9
rA   srcA   R[srcA]
rB  %rsp   srcB   R[srcB]
0xF   dstM
0xF   dstE
%rsp
ALU
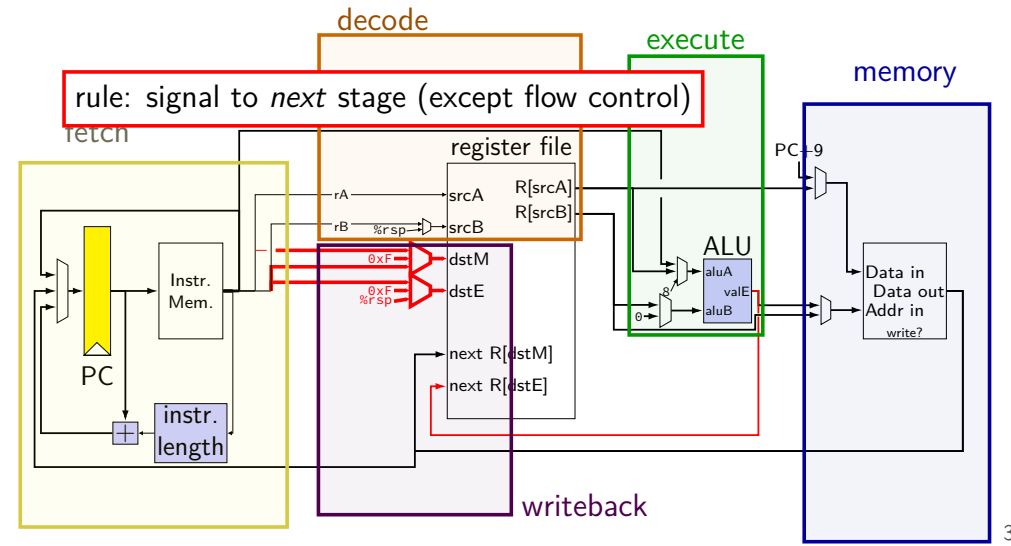aluA
valE
aluB
8
0
Data in
Data out
Addr in
write?
Instr.
Mem.
PC
next R[dstM]
next R[dstE]
instr.
length
writeback

# SEQ with stages

decode
execute
memory
fetch
register file
PC+9
rA   srcA   R[srcA]
rB  %rsp   srcB   R[srcB]
0xF   dstM
0xF   dstE
%rsp
ALU
aluA
valE
aluB
8
0
Data in
Data out
Addr in
write?
Instr.
Mem.
PC
next R[dstM]
next R[dstE]
instr.
length
writeback

## SEQ with stages

decode · execute · memory · fetch

rule: signal to *next* stage (except flow control)

register file

rA → srcA · R[srcA]
rB → %rsp → srcB · R[srcB]

PC+9

0xF → dstM
%rsp → dstE

ALU
aluA
valE
aluB
0 →

Data in
Data out
Addr in
write?

next R[dstM]
next R[dstE]

Instr. Mem.

PC

instr. length

writeback

3

## SEQ with stages (actually sequential)

decode · execute · memory · fetch

0xF
%rsp

register file

rA → srcA · R[srcA]
rB → %rsp → srcB · R[srcB]

PC+9

dstM
dstE

ALU
aluA
valE
aluB
0 →

Data in
Data out
Addr in
write?

next R[dstM]
next R[dstE]

Instr. Mem.

PC

instr. length

writeback

4

## adding pipeline registers

decode · execute · memory · fetch

0xF
%rsp

register file

rA → srcA · R[srcA]
rB → %rsp → srcB · R[srcB]

9

dstM
dstE

ALU
aluA
valE
aluB
0 →

Data in
Data out
Addr in
write?

next R[dstM]
next R[dstE]

Instr. Mem.

PC

instr. length

writeback

5

## adding pipeline registers

decode · execute · memory · fetch

0xF
%rsp

not shown — control logic

register file

rA → srcA · R[srcA]
rB → %rsp → srcB · R[srcB]

9

dstM
dstE

ALU
aluA
valE
aluB
0 →

Data in
Data out
Addr in
write?

next R[dstM]
next R[dstE]

Instr. Mem.

PC

instr. length

writeback

5

## pipeline with different hazards

example: 4-stage pipeline:
fetch/decode/execute+memory/writeback

```
                    // 4 stage   // 5 stage
addq %rax, %r8   //            // W
subq %rax, %r9   // W          // M
xorq %rax, %r10  // EM         // E
andq %r8,  %r11  // D          // D
```

## pipeline with different hazards

example: 4-stage pipeline:
fetch/decode/execute+memory/writeback

```
                    // 4 stage   // 5 stage
addq %rax, %r8   //            // W
subq %rax, %r9   // W          // M
xorq %rax, %r10  // EM         // E
andq %r8,  %r11  // D          // D
```

addq/andq is hazard with 5-stage pipeline

addq/andq is **not** a hazard with 4-stage pipeline

## exercise: different pipeline

split execute into two stages: F/D/E1/E2/M/W

result only available after second execute stage

where does forwarding, stalls occur?

| cycle # | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| addq %rcx, %r9 | | F | D | E1 | E2 | M | W | | |
| addq %r9, %rbx | | | | | | | | | |
| addq %rax, %r9 | | | | | | | | | |
| rmmovq %r9, (%rbx) | | | | | | | | | |

## exercise: different pipeline

split execute into two stages: F/D/E1/E2/M/W

| cycle # | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| addq %rcx, %r9 | | F | D | E1 | E2 | M | W | | |
| addq %r9, %rbx | | | | | | | | | |
| addq %rax, %r9 | | | | | | | | | |
| rmmovq %r9, (%rbx) | | | | | | | | | |

# exercise: different pipeline

split execute into two stages: F/D/E1/E2/M/W

| cycle # | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| addq %rcx, %r9 | F | D | E1 | E2 | M | W | | | |
| addq %r9, %rbx | | F | D | E1 | E2 | M | W | | |
| addq %rax, %r9 | | | F | D | E1 | E2 | M | W | |
| rmmovq %r9, (%rbx) | | | | F | D | E1 | E2 | M | W |

---

# exercise: different pipeline

split execute into two stages: F/D/E1/E2/M/W

| cycle # | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| addq %rcx, %r9 | F | D | E1 | E2 | M | W | | | |
| addq %r9, %rbx | | F | D | E1 | E2 | M | W | | |
| addq %r9, %rbx | | F | D | D | E1 | E2 | M | W | |
| addq %rax, %r9 | | | F | D | E1 | E2 | M | W | |
| addq %rax, %r9 | | | F | F | D | E1 | E2 | M | W |
| rmmovq %r9, (%rbx) | | | | F | D | E1 | E2 | M | W |
| rmmovq %r9, (%rbx) | | | | | F | D | E1 | E2 | M | W |

---

# exercise: different pipeline

split execute into two stages: F/D/E1/E2/M/W

| cycle # | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| addq %rcx, %r9 | F | D | E1 | E2 | M | W | | | |
| addq %r9, %rbx | | F | D | E1 | E2 | M | W | | |
| addq %r9, %rbx | | F | D | D | E1 | E2 | M | W | |
| addq %rax, %r9 | | | F | D | E1 | E2 | M | W | |
| addq %rax, %r9 | | | F | F | D | E1 | E2 | M | W |
| rmmovq %r9, (%rbx) | | | | F | D | E1 | E2 | M | W |
| rmmovq %r9, (%rbx) | | | | | F | D | E1 | E2 | M | W |

---

# exercise: different pipeline

split execute into two stages: F/D/E1/E2/M/W

| cycle # | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| addq %rcx, %r9 | F | D | E1 | E2 | M | W | | | |
| addq %r9, %rbx | | F | D | E1 | E2 | M | W | | |
| addq %r9, %rbx | | F | D | D | E1 | E2 | M | W | |
| addq %rax, %r9 | | | F | D | E1 | E2 | M | W | |
| addq %rax, %r9 | | | F | F | D | E1 | E2 | M | W |
| rmmovq %r9, (%rbx) | | | | F | D | E1 | E2 | M | W |
| rmmovq %r9, (%rbx) | | | | | F | D | E1 | E2 | M | W |

## ex.: dependencies and hazards (1)
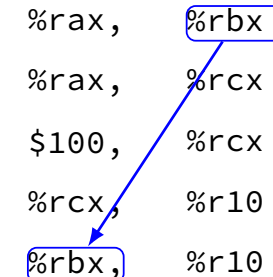
```
addq      %rax,    %rbx

subq      %rax,    %rcx

irmovq    $100,    %rcx

addq      %rcx,    %r10

addq      %rbx,    %r10
```
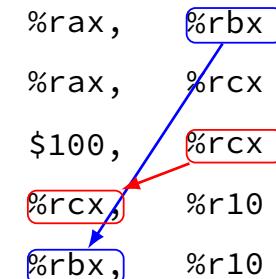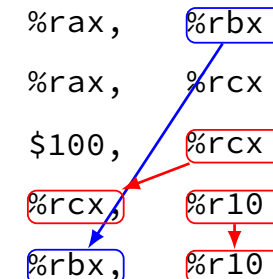
where are dependencies?
which are hazards in our pipeline?
which are resolved with forwarding?

# ex.: dependencies and hazards (2)

```
        mrmovq    0(%rax)  %rbx

        addq      %rbx     %rcx

        jne       foo

foo:    addq      %rcx     %rdx

        mrmovq    (%rdx)   %rcx
```

where are dependencies?
which are hazards in our pipeline?
which are resolved with forwarding?

# 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 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

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

## 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] 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] 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] mem[0x01] | 0 | | |
| 1 | 1 | 011000 | mem[0x62] 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] 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) | |
| 00000000 (00) | |
| 01100100 (64) | |

tag index offset

11

---

## 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) | |
| 01100100 (64) | |

tag index offset

11

---

## 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

11

---

## 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
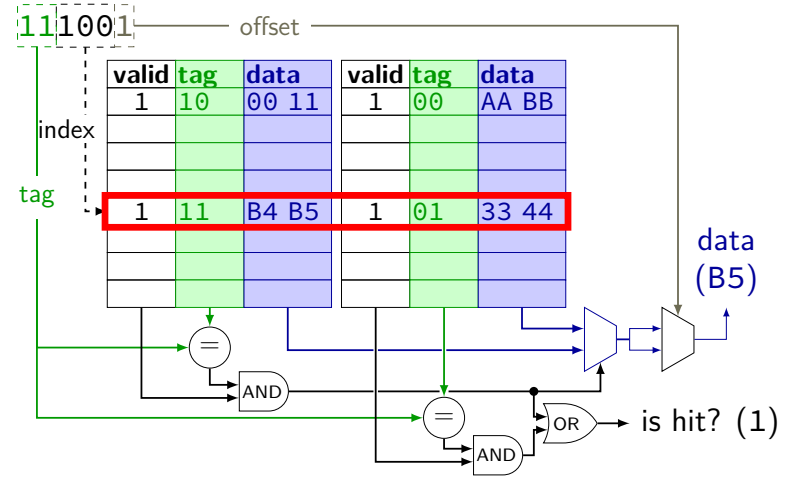
11

## 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) | hit |
| 01100100 (64) | miss |

tag   index offset

11

---

## cache operation (associative)



12

---

## cache operation (associative)



12

---

## cache operation (associative)



12

## 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

## Tag-Index-Offset formulas (complete)

| | |
|---|---|
| $m$ | memory addreses bits (Y86-64: 64) |
| $E$ | 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) |

## replacement policies

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 |
|---|---|
| 000 | how to decide where to insert 0x64? |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | hit |
| 00000000 (00) | hit |
| 01100100 (64) | miss |

## replacement policies

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

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

## example replacement policies

least recently used and approximations
    take advantage of temporal locality
    exact: $\lceil \log_2(E!) \rceil$ bits per set for $E$-way cache
    good approximations: $E$ to $2E$ bits

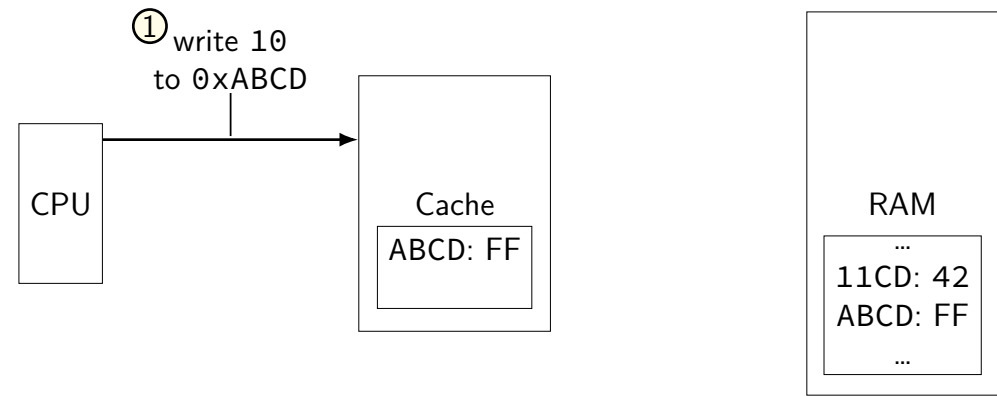first-in, first-out
    counter per set — where to replace next

(pseudo-)random
    no extra information!
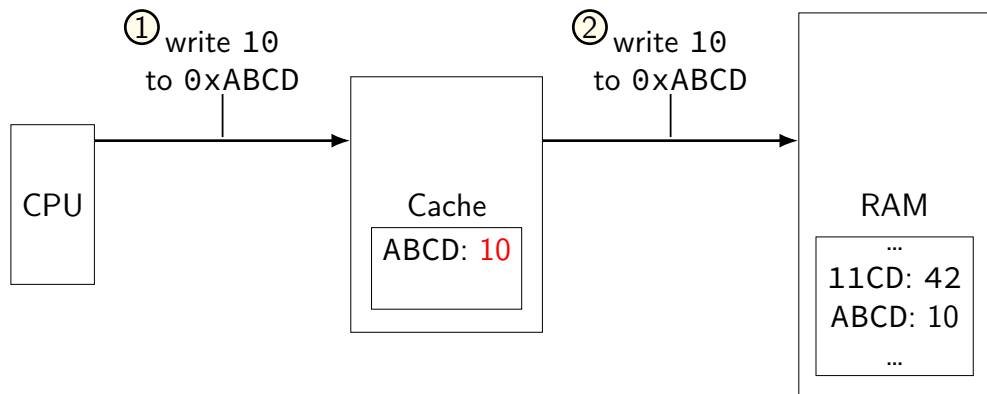
---

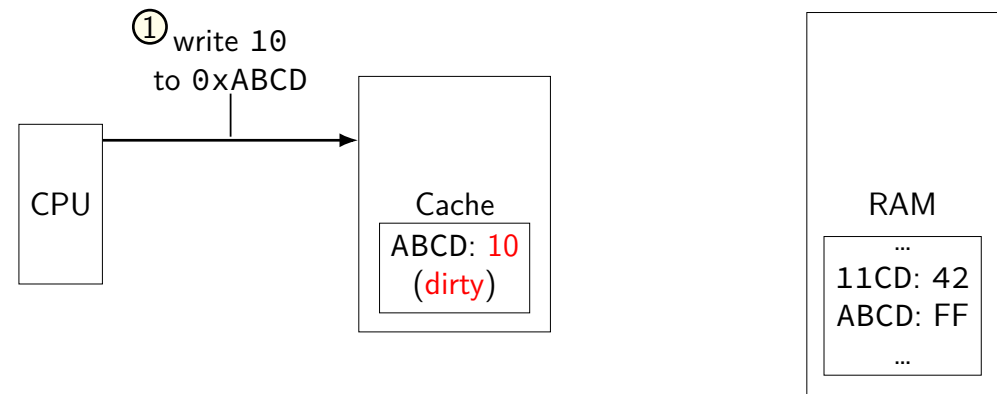## write-through v. write-back

**option 1: write-through**



①write 10 to 0xABCD

CPU

Cache
ABCD: FF

RAM
...
11CD: 42
ABCD: FF
...

---

## write-through v. write-back

**option 1: write-through**



①write 10 to 0xABCD

②write 10 to 0xABCD

CPU

Cache
ABCD: 10

RAM
...
11CD: 42
ABCD: 10
...

---

## write-through v. write-back

**option 2: write-back**



①write 10 to 0xABCD

CPU

Cache
ABCD: 10
(dirty)

RAM
...
11CD: 42
ABCD: FF
...

## write-through v. write-back

**option 2: write-back**

① write 10 to 0xABCD

② read from 0x11CD (conflicts)

③ write 10 to ABCD

CPU

Cache
ABCD: 10 ~~(dirty)~~

RAM
...
11CD: 42
ABCD: 10
...

... when replaced — send value to memory

---

## write-through v. write-back

① write 10 to 0xABCD

② read from 0x11CD (conflicts)

③ write 10 to ABCD

④ read from 0x11CD

CPU

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

RAM
...
11CD: 42
ABCD: 10
...

... read new value to store in cache

---

## 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]<br>mem[0x01] | 0 | 1 | 011000 | mem[0x60]⋆<br>mem[0x61]⋆ | 1 | 1 |
| 1 | 1 | 011000 | mem[0x62]<br>mem[0x63] | 0 | 0 | | | | 0 |

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

---

## 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

**write-no-allocate**
   send write through to memory
   guess: not read soon?

# write-allocate

2-way set associative, LRU, writeback

| 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 |

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

20

# write-allocate

2-way set associative, LRU, writeback

| 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 |

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

20

# write-allocate

2-way set associative, LRU, writeback

| 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 |

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

20

# write-allocate

2-way set associative, LRU, writeback

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

20

## write-no-allocate

2-way set associative, LRU, writeback

| 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 |

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

## fast writes



write 10 to 0xABCD

write 20 to 0x1234

CPU

**write buffer**
0xABCD: 10
0x1234: 20

Cache

RAM

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

## exercise (1)

initial cache: 64-byte blocks, 64 sets, 8 ways/set

If we leave the other parameters listed above unchanged, which will probably reduce the number of capacity misses in a typical program? (Multiple may be correct.)
  A.  quadrupling the block size (256-byte blocks, 64 sets, 8 ways/set)
  B.  quadrupling the number of sets
  C.  quadrupling the number of ways/set

## exercise (2)

initial cache: 64-byte blocks, 8 ways/set, 64KB cache

If we leave the other parameters listed above unchanged, which will probably reduce the number of capacity misses in a typical program? (Multiple may be correct.)
  A.  quadrupling the block size (256-byte block, 8 ways/set, 64KB cache)
  B.  quadrupling the number of ways/set
  C.  quadrupling the cache size

# exercise (3)

initial cache: 64-byte blocks, 8 ways/set, 64KB cache

If we leave the other parameters listed above unchanged, which will probably reduce the number of conflict misses in a typical program? (Multiple may be correct.)
  A.  quadrupling the block size (256-byte block, 8 ways/set, 64KB cache
  B.  quadrupling the number of ways/set
  C.  quadrupling the cache size

# C 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 a 2KB direct-mapped cache with 16B cache blocks?

# C 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 = 1; 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 a 2KB direct-mapped cache with 16B cache blocks? Would a set-associtiave cache be better?

# thinking about cache storage (1)

2KB direct-mapped cache with 16B blocks —

set 0: address 0 to 15, (0 to 15) + 2KB, (0 to 15) + 4KB, …

set 1: address 16 to 31, (16 to 31) + 2KB, (16 to 31) + 4KB, …

…

set 127: address 2032 to 2047, (2032 to 2047) + 2KB, …

# thinking about cache storage (1)

2KB direct-mapped cache with 16B blocks —

set 0: address 0 to 15, (0 to 15) + 2KB, (0 to 15) + 4KB, …

set 1: address 16 to 31, (16 to 31) + 2KB, (16 to 31) + 4KB, …

…

set 127: address 2032 to 2047, (2032 to 2047) + 2KB, …

# thinking about cache storage (1)

2KB direct-mapped cache with 16B blocks —

set 0: address 0 to 15, (0 to 15) + 2KB, (0 to 15) + 4KB, …
    block at 0: array[0] through array[3]

set 1: address 16 to 31, (16 to 31) + 2KB, (16 to 31) + 4KB, …
    block at 16: array[4] through array[7]

…

set 127: address 2032 to 2047, (2032 to 2047) + 2KB, …
    block at 2032: array[508] through array[511]

# thinking about cache storage (1)

2KB direct-mapped cache with 16B blocks —

set 0: address 0 to 15, (0 to 15) + 2KB, (0 to 15) + 4KB, …
    block at 0: array[0] through array[3]
    block at 0+2KB: array[512] through array[515]

set 1: address 16 to 31, (16 to 31) + 2KB, (16 to 31) + 4KB, …
    block at 16: array[4] through array[7]
    block at 16+2KB: array[516] through array[519]

…

set 127: address 2032 to 2047, (2032 to 2047) + 2KB, …
    block at 2032: array[508] through array[511]
    block at 2032+2KB: array[1020] through array[1023]

# thinking about cache storage (2)

2KB 2-way set associative cache with 16B blocks: block addresses —

set 0: address 0, 0 + 2KB, 0 + 4KB, …

set 1: address 16, 16 + 2KB, 16 + 4KB, …

…

set 63: address 1008, 2032 + 2KB, 2032 + 4KB …

# thinking about cache storage (2)

2KB 2-way set associative cache with 16B blocks: block addresses
—

set 0: address 0, 0 + 2KB, 0 + 4KB, …
    block at 0: array[0] through array[3]


set 1: address 16, 16 + 2KB, 16 + 4KB, …
    address 16: array[4] through array[7]

…

set 63: address 1008, 2032 + 2KB, 2032 + 4KB …
    address 1008: array[252] through array[255]

---

---

# thinking about cache storage (2)

2KB 2-way set associative cache with 16B blocks: block addresses
—

set 0: address 0, 0 + 2KB, 0 + 4KB, …
    block at 0: array[0] through array[3]
    block at 0+1KB: array[256] through array[259]
    block at 0+2KB: array[512] through array[515]
    …

set 1: address 16, 16 + 2KB, 16 + 4KB, …
    address 16: array[4] through array[7]

…

set 63: address 1008, 2032 + 2KB, 2032 + 4KB …
    address 1008: array[252] through array[255]

---

# C and cache misses (3)

```
typedef struct {
    int a_value, b_value;
    int boring_values[126];
} item;
item items[8]; // 4 KB array
int a_sum = 0, b_sum = 0;
for (int i = 0; i < 8; ++i)
    a_sum += items[i].a_value;
for (int i = 0; i < 8; ++i)
    b_sum += items[i].b_value;
```

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

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

## C and cache misses (3, rewritten?)

```
item array[1024]; // 4 KB array
int a_sum = 0, b_sum = 0;
for (int i = 0; i < 1024; i += 128)
    a_sum += array[i];
for (int i = 1; i < 1024; i += 128)
    b_sum += array[i];
```

## C and cache misses (4)

```
typedef struct {
    int a_value, b_value;
    int boring_values[126];
} item;
item items[8]; // 4 KB array
int a_sum = 0, b_sum = 0;
for (int i = 0; i < 8; ++i)
    a_sum += items[i].a_value;
for (int i = 0; i < 8; ++i)
    b_sum += items[i].b_value;
```

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

How many *data cache misses* on a 4-way set associative 2KB direct-mapped cache with 16B cache blocks?

## a note on matrix storage

$A$ — $N \times N$ matrix

represent as array

makes dynamic sizes easier:

```
float A_2d_array[N][N];
float *A_flat = malloc(N * N);
```

`A_flat[i * N + j]` === `A_2d_array[i][j]`

## locality exercise (1)
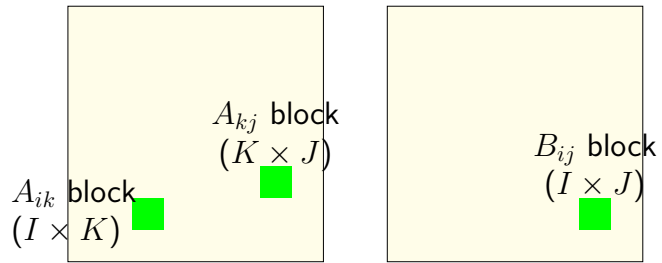
```
/* version 1 */
for (int i = 0; i < N; ++i)
    for (int j = 0; j < N; ++j)
        A[i] += B[j] * C[i * N + j]

/* version 2 */
for (int j = 0; j < N; ++j)
    for (int i = 0; i < N; ++i)
        A[i] += B[j] * C[i * N + j];
```

exercise: which has better temporal locality in A? in B? in C?
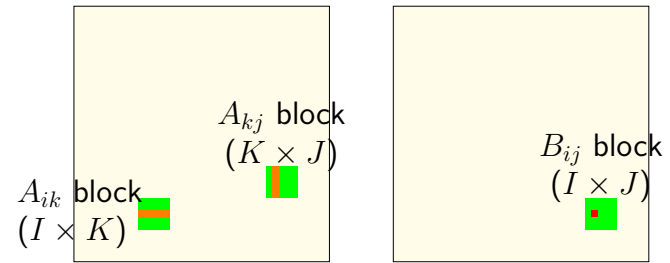how about spatial locality?

## array usage: block



$A_{ik}$ block $(I \times K)$
$A_{kj}$ block $(K \times J)$
$B_{ij}$ block $(I \times J)$

inner loop keeps "blocks" from $A$, $B$ in cache

## array usage: block



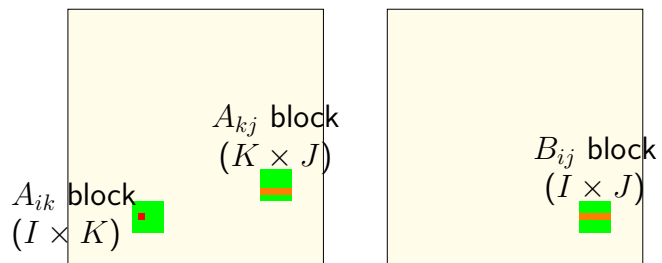$A_{ik}$ block $(I \times K)$
$A_{kj}$ block $(K \times J)$
$B_{ij}$ block $(I \times J)$

$B_{ij}$ calculation uses strips from $A$
$K$ calculations for one load (cache miss)

## array usage: block



$A_{ik}$ block $(I \times K)$
$A_{kj}$ block $(K \times J)$
$B_{ij}$ block $(I \times J)$
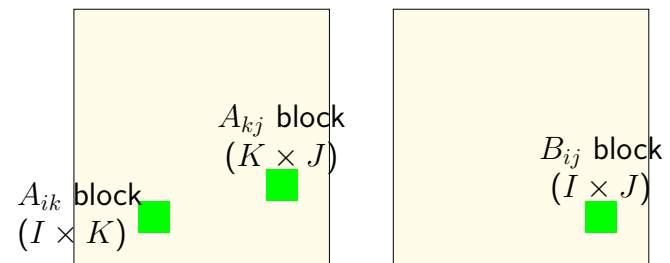
$A_{ik}$ calculation uses strips from $A$, $B$
$J$ calculations for one load (cache miss)

## array usage: block



$A_{ik}$ block $(I \times K)$
$A_{kj}$ block $(K \times J)$
$B_{ij}$ block $(I \times J)$

(approx.) $KIJ$ fully cached calculations
for $KI + IJ + KJ$ loads
(assuming everything stays in cache)

# cache blocking efficiency

load $I \times K$ elements of $A_{ik}$:
    do $> J$ multiplies with each

load $K \times J$ elements of $A_{kj}$:
    do $I$ multiplies with each

load $I \times J$ elements of $B_{ij}$:
    do $K$ adds with each

bigger blocks — more work per load!

catch: $IK + KJ + IJ$ elements must fit in cache

# cache blocking rule of thumb

fill the most of the cache with useful data
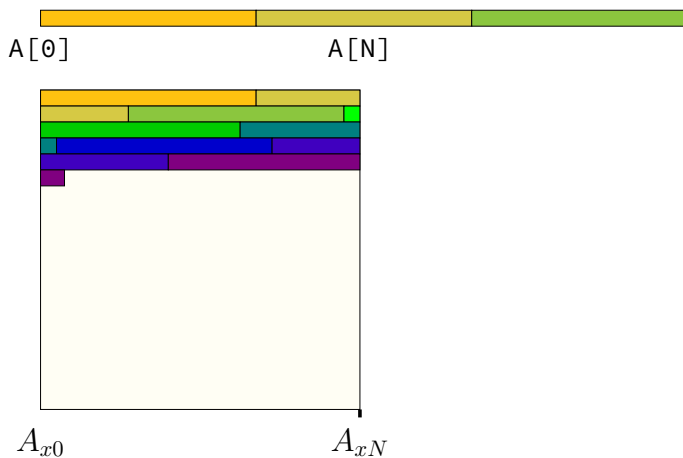
and do as much work as possible from that

example: my desktop 32KB L1 cache

$I = J = K = 48$ uses $48^2 \times 3$ elements, or $27$KB.

assumption: conflict misses aren't important

# 'flat' 2D arrays and cache blocks



`A[0]`        `A[N]`

$A_{x0}$        $A_{xN}$