## Cache Performance

- CPUtime = Instruction count x CPI x Clock cycle time
- CPIexecution = CPI with ideal memory
- CPI = CPIexecution + Mem Stall cycles per instruction
- Mem Stall cycles per instruction =
- Mem accesses per instruction x Miss rate x Miss penalty
- CPUtime = Instruction Count x (CPIexecution + Mem Stall cycles per instruction) x Clock cycle time
- CPUtime = IC x (CPIexecution + Mem accesses per instruction x Miss rate x Miss penalty) x Clock cycle time
- Misses per instruction = Memory accesses per instruction x Miss rate
- CPUtime = IC x (CPIexecution + Misses per instruction x Miss penalty) x Clock cycle time

## Cache Performance

- Assuming the following execution and cache parameters:
    - Cache miss penalty = 50 cycles
    - Normal instruction execution CPI ignoring memory stalls = 2.0 cycles
    - Miss rate = 2%
    - Average memory references/instruction = 1.33
- CPU time = IC x [CPI execution + Memory accesses/instruction x Miss rate x Miss penalty ] x Clock cycle time
- CPUtime with cache = IC x (2.0 + (1.33 x 2% x 50)) x clock cycle time
- = IC x 3.33 x Clock cycle time

- *Lower CPI execution increases the impact of cache miss clock cycles*

## Cache Performance

- Suppose a CPU executes at Clock Rate = 200 MHz (5 ns per cycle) with a single level of cache.
- CPIexecution = 1.1
- Instruction mix: 50% arith/logic, 30% load/store, 20% control
- Assume a cache miss rate of 1.5% and a miss penalty of 50 cycles.
- CPI = $CPI_{execution}$ + mem stalls per instruction
- Mem Stalls per instruction =
- Mem accesses per instruction x Miss rate x Miss penalty
- Mem accesses per instruction = 1 + .3 = 1.3
- Mem Stalls per instruction = 1.3 x .015 x 50 = 0.975
- CPI = 1.1 + .975 = 2.075
- The ideal memory CPU with no misses is 2.075/1.1 = 1.88 times faster

## Cache Performance

- Suppose for the previous example we double the clock rate to 400 MHZ, how much faster is this machine, assuming similar miss rate, instruction mix?
- Since memory speed is not changed, the miss penalty takes more CPU cycles:
  - Miss penalty = 50 x 2 = 100 cycles.
  - CPI = 1.1 + 1.3 x .015 x **100** = 1.1 + 1.95 = 3.05
  - Speedup = (CPIold x Cold)/ (CPInew x Cnew)
  - = 2.075 x 2 / 3.05 = 1.36
- The new machine is only 1.36 times faster rather than 2

times faster due to the increased effect of cache misses.

- *CPUs with higher clock rate, have more cycles per cache miss and more memory impact on CPI.*

MK

22

---

## Cache Performance

- Suppose a CPU uses separate level one (L1) caches for instructions and data (Harvard memory architecture) with different miss rates for instruction and data access:
  - $CPI_{execution}$ = 1.1
  - Instruction mix: 50% arith/logic, 30% load/store, 20% control
  - Assume a cache miss rate of 0.5% for instruction fetch and a cache data miss rate of 6%.
  - A cache hit incurs no stall cycles while a cache miss incurs 200 stall cycles for both memory reads and writes. Find the resulting CPI using this cache? How much faster is the CPU with ideal memory?

  $CPI = CPI_{execution}$ + mem stalls per instruction

  Mem Stall cycles per instruction = Instruction Fetch Miss rate x Miss Penalty + Data Memory Accesses Per Instruction x Data Miss Rate x Miss Penalty

  Mem Stall cycles per instruction = 1 x 0.5/100 x 200 + 0.3 x 6/100 x 200 = 1 + 3.6 = 4.6

  $CPI = CPI_{execution}$ + mem stalls per instruction = 1.1 + 4.6 = 5.7

The CPU with ideal cache (no misses) is 5.7/1.1 = 5.18 times faster
With no cache the CPI would have been = 1.1 + 1.3 X 200 = 261.1

MK

23

---

## Cache Performance

| Size | Instruction cache | Data cache | Unified cache |
|---|---|---|---|
| 1 KB | 3.06% | 24.61% | 13.34% |
| 2 KB | 2.26% | 20.57% | 9.78% |
| 4 KB | 1.78% | 15.94% | 7.24% |
| 8 KB | 1.10% | 10.19% | 4.57% |
| 16 KB | 0.64% | 6.47% | 2.87% |
| 32 KB | 0.39% | 4.82% | 1.99% |
| 64 KB | 0.15% | 3.77% | 1.35% |
| 128 KB | 0.02% | 2.88% | 0.95% |

MK

24

## Write Policy

1 <u>Write Though</u>:  Data is written to both the cache block and to a block of main memory.
   - The lower level always has the most updated data; an important feature for I/O and multiprocessing.
   - Easier to implement than write back.
   - <u>A write buffer</u> is often used to reduce CPU write stall while data is written to memory.

2 <u>Write back</u>:  Data is written or updated only to the cache block.  The modified or dirty cache block is written to main memory when it's being replaced from cache.
   - Writes occur at the speed of cache
   - A status bit called a dirty or modified bit, is used to indicate whether the block was modified while in cache; if not the block is not written back to main memory when replaced.
   - Uses less memory bandwidth than write through.

## Write Policy

<u>Write Allocate:</u>

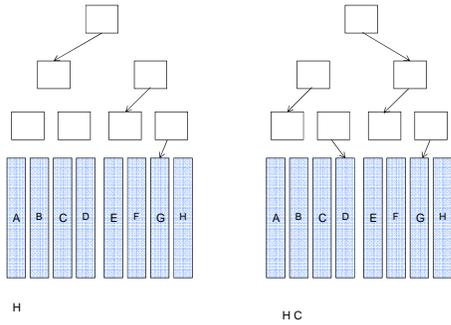 The cache block is loaded on a write miss followed by write hit actions.

<u>No-Write Allocate:</u>

The block is modified in the lower level (lower cache level, or main
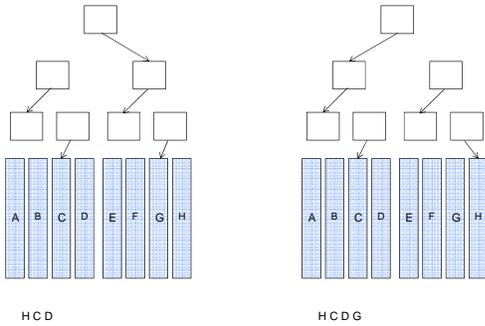
memory) and not loaded into cache.

## LRU

- A list to keep track of the order of access to every block in the set.
- The least recently used block is replaced (if needed).
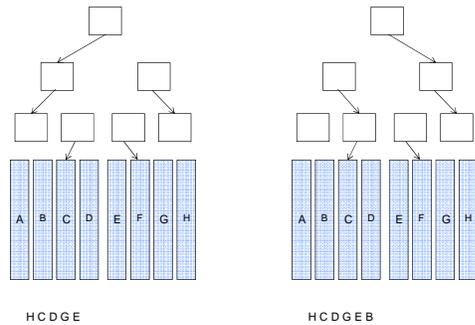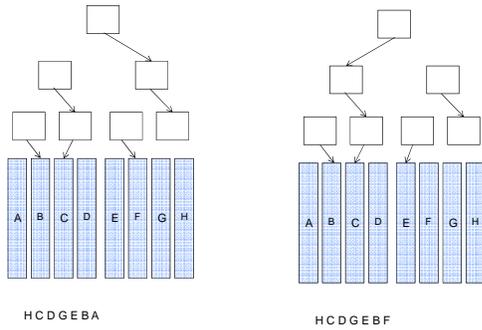- How many bits we need for that?

# Pseudo LRU

A B C D E F G H     A B C D E F G H

H          H C

28

# Psuedo LRU

A B C D E F G H     A B C D E F G H

H C D          H C D G

29

# Psuedo LRU

A B C D E F G H     A B C D E F G H

H C D G E          H C D G E B

30

## Psuedo LRU



H C D G E B A

H C D G E B F

31

---

## Example

- Which has a lower miss rate 16KB cache for both instruction or data, or a combined 32KB cache? (0.64%, 6.47%, 1.99%).
- Assume hit=1cycle and miss =50 cycles. 75% of memory references are instruction fetch.  *reads*
- Miss rate of split cache=0.75*0.64%+0.25*6.47%=2.1%
- Slightly worse than 1.99% for combined cache. But, what about average memory access time?
- Split cache: 75%(1+0.64%*50)+25%(1+6.47%*50) = 2.05 cycles.
- Combined cache:    Extra cycle for load/store
  75%(1+1.99%*50)+25%(1+1+1.99%*50) = 2.24

32

---

## Example

- A CPU with  $CPI_{execution}$ = 1.1  Mem accesses per instruction = 1.3
- Uses a unified L1 Write Through, No Write Allocate,  with:
  - No write buffer.
  - Perfect Write buffer
  - A realistic write buffer that eliminates 85% of write stalls
- Instruction mix:   50% arith/logic,  15% load, 15% store, 20% control
- Assume a cache miss rate of 1.5% and a miss penalty of 50 cycles.
  CPI =    $CPI_{execution}$  +   mem stalls per instruction
  % reads =  1.15/1.3 =   88.5%        % writes  =  .15/1.3 =  11.5%

33

## Example

- A CPU with $CPI_{execution}$ = 1.1 uses a unified L1 with with <u>write back</u>, with <u>write allocate</u>, and the <u>probability a cache block is dirty = 10%</u>
- Instruction mix: 50% arith/logic, 15% load, 15% store, 20% control
- Assume a cache miss rate of 1.5% and a miss penalty of 50 cycles.

*[handwritten:]* 1.3 mem ref/inst    50+50

$$\frac{1.5}{100}\left(1.3 \times 50 \times 0.9 + 1.3 \times 0.1 \times 100\right)$$

## Example

- CPU with $CPI_{execution}$ = 1.1 running at clock rate = 500 MHz
- 1.3 memory accesses per instruction.
- $L_1$ cache operates at 500 MHz with a miss rate of 5%
- $L_2$ cache operates at 250 MHz with local miss rate 40%, ($T_2$ = 2 cycles)
- Memory access penalty, M = 100 cycles. Find CPI.

*[handwritten:]*
$$1.3\left(\frac{5}{100} \times 0.6 \times 2 + \frac{5}{100} \times 0.4 \times 100\right)$$

## Example

- CPU with $CPI_{execution}$ = 1.1 running at clock rate = 500 MHz
- 1.3 memory accesses per instruction.
- For $L_1$ :
  - Cache operates at 500 MHz with a miss rate of 1-H1 = 5%
  - Write though to $L_2$ with perfect write buffer with write allocate
- For $L_2$:
  - Cache operates at 250 MHz with local miss rate 1- H2 = 40%, ($T_2$ = 2 cycles)
  - Write back to main memory with write allocate
  - Probability a cache block is dirty = 10%
- Memory access penalty, M = 100 cycles. Find CPI.

*[handwritten:]*
$$0.05\left(0.6 \times 2 + 0.4 \times 0.9 \times 100 \\ + 0.4 \times 0.1 \times 700\right)$$

## Example

- CPU with $CPI_{execution}$ = 1.1 running at clock rate = 500 MHz
- 1.3 memory accesses per instruction.
- $L_1$ cache operates at 500 MHz with a miss rate of 5%
- $L_2$ cache operates at 250 MHz with a local miss rate 40%, ($T_2$ = 2 cycles)
- $L_3$ cache operates at 100 MHz with a local miss rate 50%, ($T_3$ = 5 cycles)
- Memory access penalty, M= 100 cycles.  Find CPI.

HW