# Greetings

- Who is this...
  - David Thomas
  - *Currently*: Dept. of Electrical Engineering
  - *Previously*: MEng in Soft. Eng. in DoC
  - PhD student of Wayne Luk (*a long time ago...*)

- My interests
  - High-Performance Computing
  - FPGA and GPU acceleration
  - High-level Synthesis (C to gates)
  - Languages and parallelism
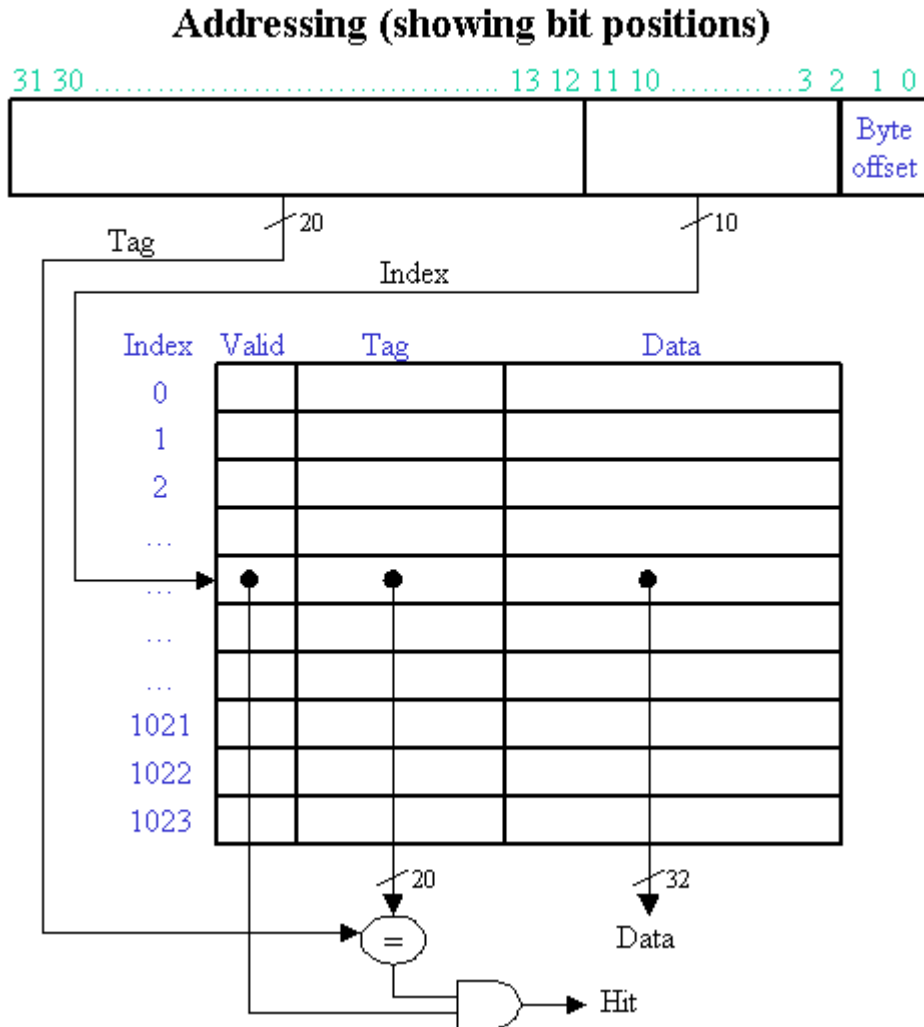
# Course administrivia

- Panopto will continue
  - Assuming the hardware and servers work

- If you post on piazza I'll eventually see it
  - Traditionally it is a ghost-town till exams

- I tend to do tutorials differently to Prof. Luk
  - Standard exercises + solutions still available

- Wed 1st of March: *No lecture*
  - Course is one lecture less than no. of lecture slots
  - Tutorial available as normal

# Cache features and performance

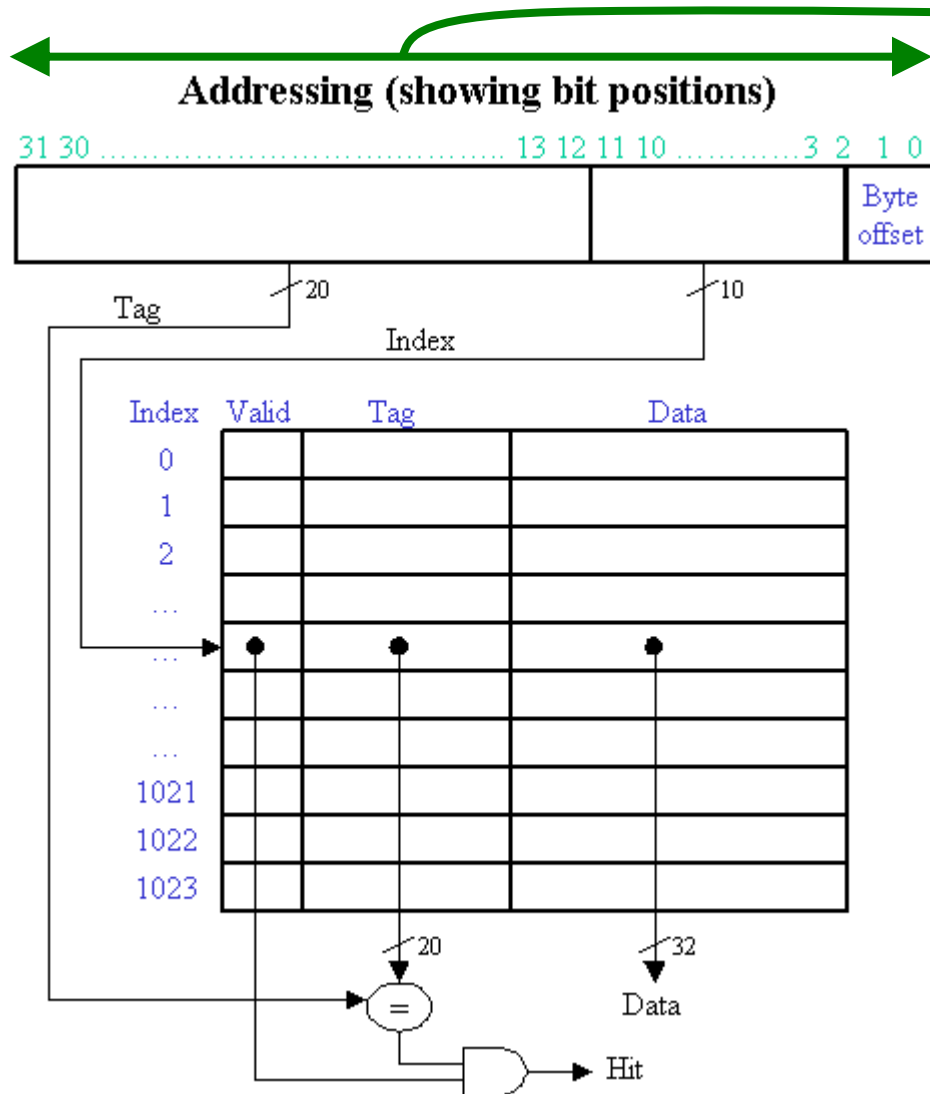(3rd Ed: p.492-496,505-511; 4th Ed: 475-479,487-492)

- so far:
  - locality principles
  - levels, blocks, hit, miss, miss penalty
  - direct-mapped cache, handling misses

- today:
  - review direct-mapped single- and multi-word cache
  - cache performance, read/write stall cycles
  - multi-level cache hierarchy

# Direct-mapped single-word cache

**Addressing (showing bit positions)**

31 30 .................................................... 13 12 11 10 ............ 3 2 1 0

Byte offset

Tag — 20
Index — 10

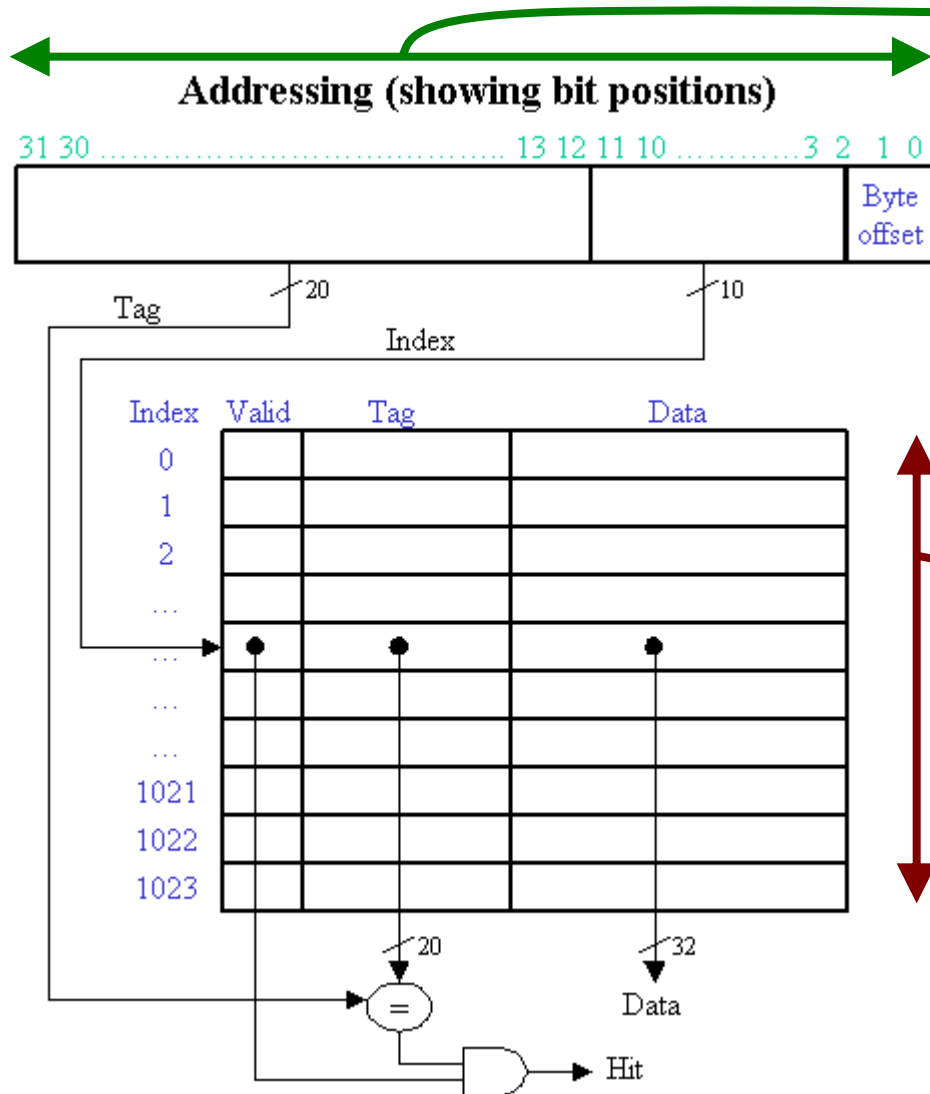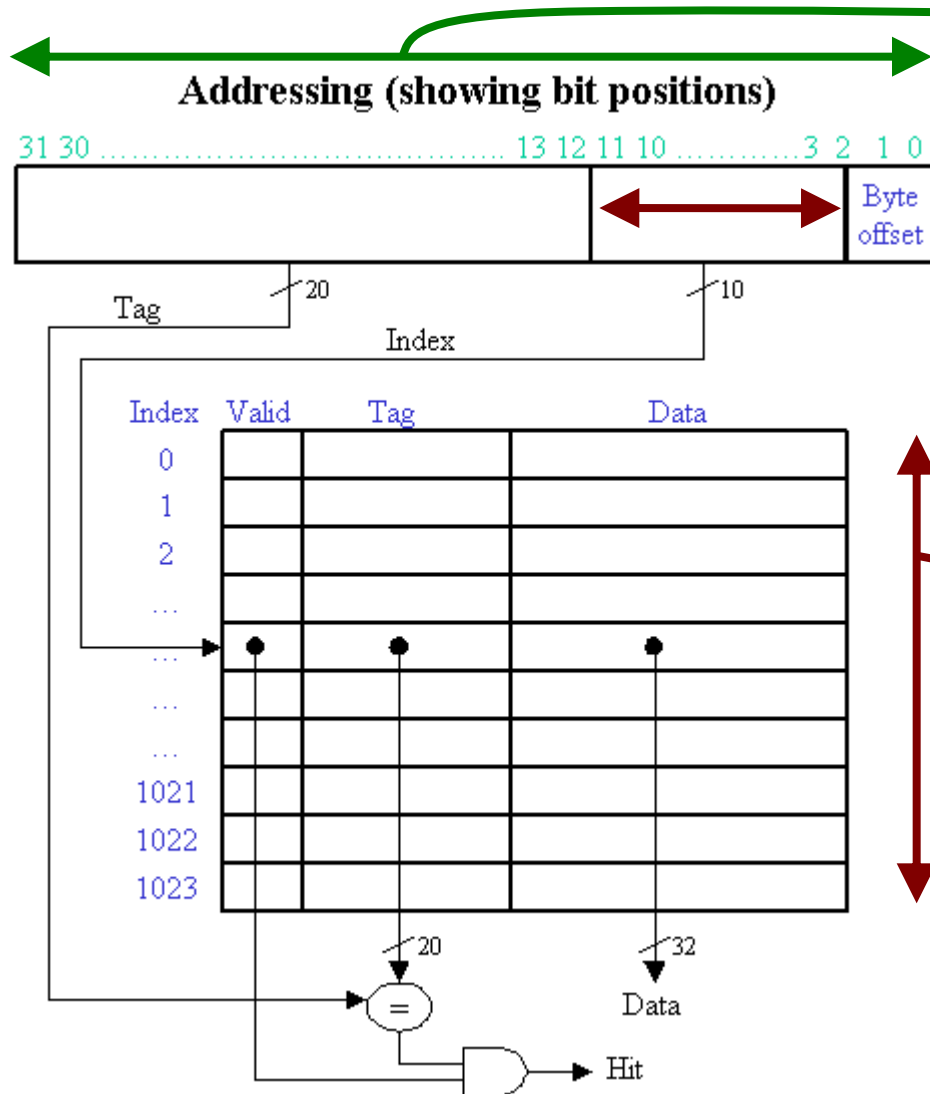| Index | Valid | Tag | Data |
|-------|-------|-----|------|
| 0 | | | |
| 1 | | | |
| 2 | | | |
| ... | | | |
| ... | • | • | • |
| ... | | | |
| ... | | | |
| 1021 | | | |
| 1022 | | | |
| 1023 | | | |

20

32

Data

=

Hit

- address becomes cache index + tag

- advantage?

- number of tag bits for a cache with n blocks, dealing with m-bit addresses?

- write misses: write to cache + memory

# Direct-mapped single-word cache

## Addressing (showing bit positions)

31 30 ............................................ 13 12 11 10 .......... 3 2 1 0

|  |  | Byte offset |

Tag — 20

Index — 10

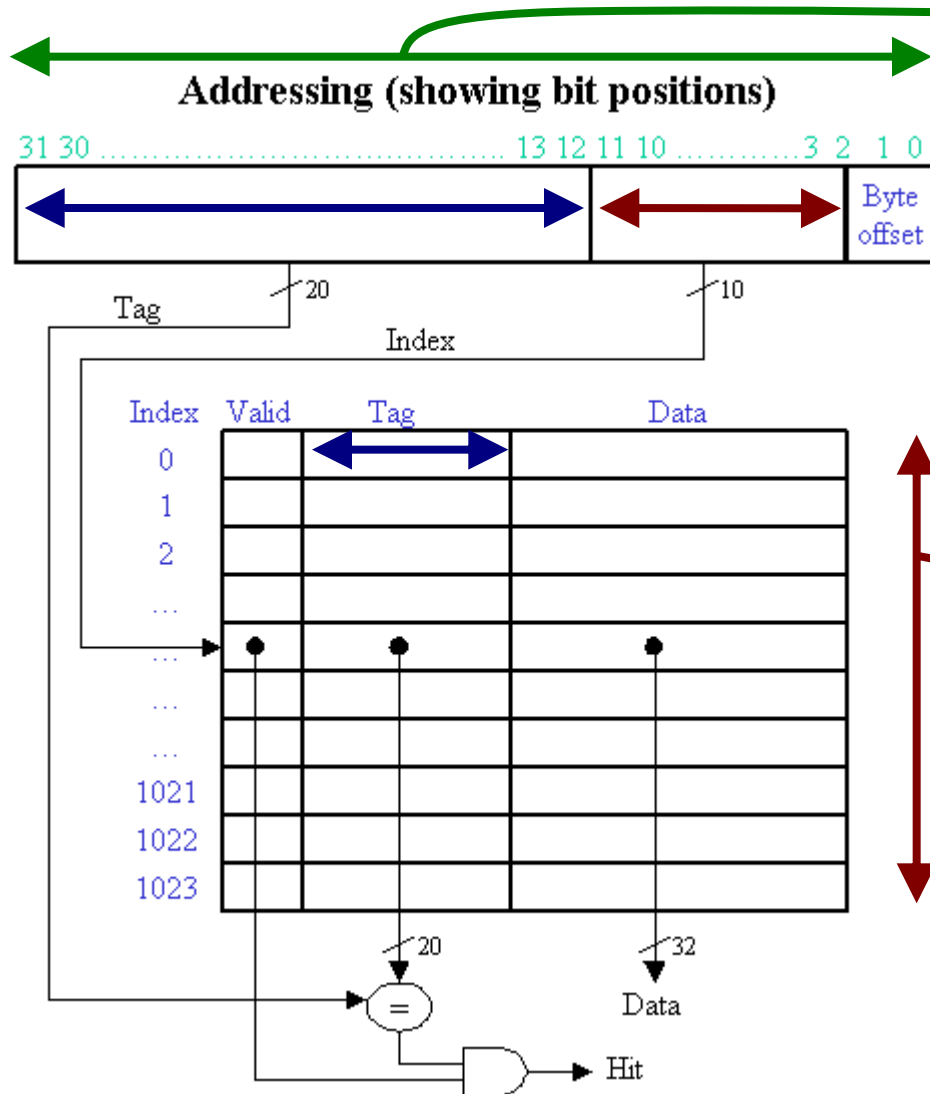| Index | Valid | Tag | Data |
|---|---|---|---|
| 0 |  |  |  |
| 1 |  |  |  |
| 2 |  |  |  |
| ... |  |  |  |
| ... | • | • | • |
| ... |  |  |  |
| ... |  |  |  |
| 1021 |  |  |  |
| 1022 |  |  |  |
| 1023 |  |  |  |

20

32

Data

( = )

Hit

- address becomes cache index + tag

- advantage?

- number of tag bits for a cache with n blocks, dealing with m-bit addresses?
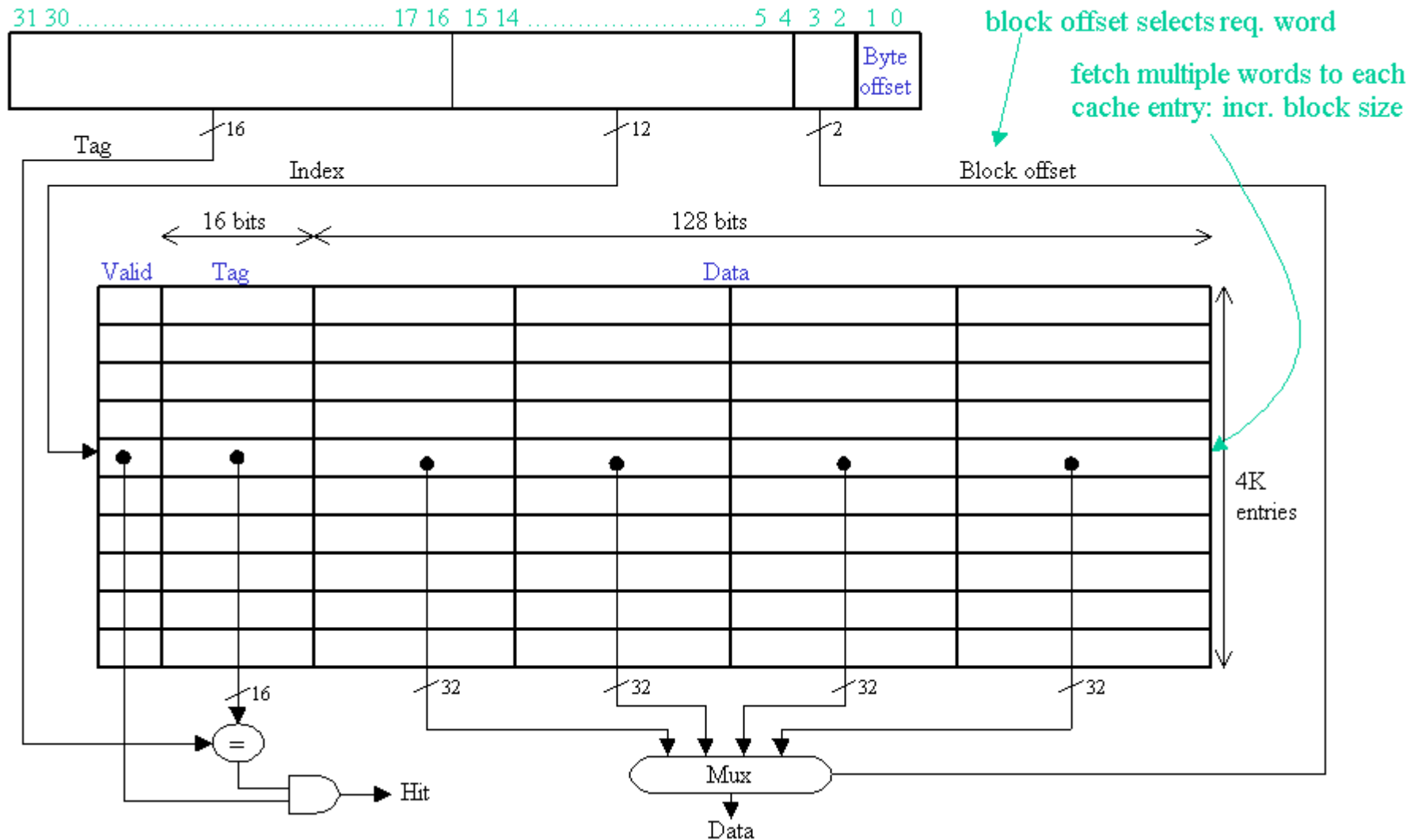
- write misses: write to cache + memory

# Direct-mapped single-word cache

**Addressing (showing bit positions)**

```
31 30 ........................................ 13 12 11 10 ............ 3 2  1 0
┌──────────────────────────────────────────┬────────────────────┬──────┐
│                                            │                    │ Byte │
│                                            │                    │offset│
└──────────────────────────────────────────┴────────────────────┴──────┘
              20                                    10
      Tag
                        Index
```

| Index | Valid | Tag | Data |
|-------|-------|-----|------|
| 0 | | | |
| 1 | | | |
| 2 | | | |
| ... | | | |
| ... | ● | ● | ● |
| ... | | | |
| ... | | | |
| 1021 | | | |
| 1022 | | | |
| 1023 | | | |

20                    32

= → Hit

Data

- address becomes cache index + tag

- advantage?

- number of tag bits for a cache with n blocks, dealing with m-bit addresses?
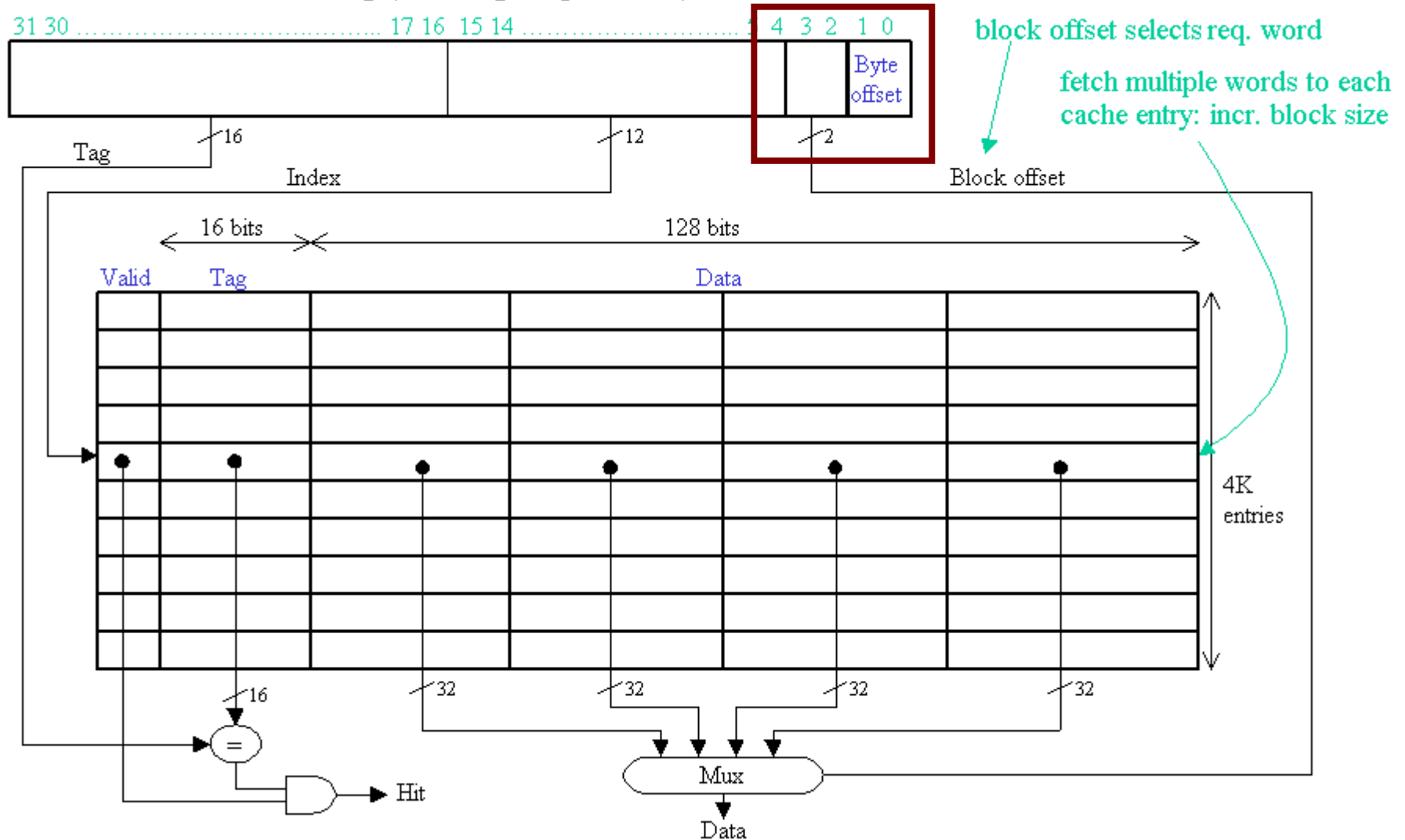
- write misses: write to cache + memory

# Direct-mapped single-word cache



Addressing (showing bit positions)

- address becomes cache index + tag

- advantage?

- number of tag bits for a cache with n blocks, dealing with m-bit addresses?

- write misses: write to cache + memory

# Direct-mapped single-word cache

## Addressing (showing bit positions)

31 30 ......................................... 13 12 11 10 ............. 3 2  1 0

Byte offset

Tag — 20

Index — 10

| Index | Valid | Tag | Data |
|-------|-------|-----|------|
| 0 | | | |
| 1 | | | |
| 2 | | | |
| ... | | | |
| ... | ● | ● | ● |
| ... | | | |
| ... | | | |
| 1021 | | | |
| 1022 | | | |
| 1023 | | | |

20          32

(=)          Data

Hit

- address becomes cache index + tag

- advantage?

- number of tag bits for a cache with n blocks, dealing with m-bit addresses?

- write misses: write to cache + memory

# Multi-word cache



**Addressing (showing bit positions)**

block offset selects req. word

fetch multiple words to each cache entry: incr. block size

# Multi-word cache

**Addressing (showing bit positions)**



block offset selects req. word

fetch multiple words to each cache entry: incr. block size

31 30 .................... 17 16 15 14 ....................... 1 4 3 2 1 0

Byte offset

Tag — 16

Index — 12 — 2

Block offset

16 bits — 128 bits

Valid   Tag   Data

16   32   32   32   32

4K entries

= 16

Mux

Hit

Data

# Direct-mapped caches

- byte: smallest addressable unit of memory

- word size : $w >= 1$, $w$ in bytes
  - natural granularity of CPU; size of registers

- block size : $k >= 1$, $k$ in words
  - granularity of cache to memory transfers

- cache size : $c >= k$, $c$ in words
  - total size of cache storage; number of blocks is ($c$ **div** $k$)

# Direct-mapped caches

- byte: smallest addressable unit of memory

- word size : $w >= 1$, $w$ in bytes
  - natural granularity of CPU; size of registers

- block size : $k >= 1$, $k$ in words
  - granularity of cache to memory transfers

- cache size : $c >= k$, $c$ in words
  - total size of cache storage; number of blocks is (`c` **div** `k`)

- Lookup `byte_address` in a direct-mapped cache
  ```
  1.  word_address  = (byte_address div w);
  2.  block_address = (word_address div k);
  3.  block_index   = (block_address mod (c div k) );
  ```

# Direct-mapped caches

Invariant : $0 <=$ `block_index` $<$ (`c` **div** `k`)

`block_index` specifies the ***only*** cache-block
where `byte_address` or `word_address` can be found

- cache size : $c >= k$, $c$ in words
  - total size of cache storage; number of blocks is (`c` **div** `k`)

- Lookup `byte_address` in a direct-mapped cache
```
1.  word_address  = (byte_address div w);
2.  block_address = (word_address div k);
3.  block_index   = (block_address mod (c div k) );
```

Words

1036

1040

1044

1048

1052

1056

1060

1064

1068

1072

| Blocks | Words | Bytes |
|--------|-------|-------|

1024          1036

                                    1040

                              1040

                                    1045

                              1044

1040

                              1048

                                    1052

                              1052          1053

                                    1054

                                    1055

                              1056

                              1060

1056

                              1064          1064

                              1068

                                    1070

1072          1072

| Blocks | Words | Bytes |
|--------|-------|-------|
| 1024 | 1036 | |
| | | ← 1040 |
| 1040 | 1040 | |
| | 1044 | ← 1045 |
| | 1048 | 1052 |
| | 1052 | 1053 |
| | 1056 | 1054 |
| 1056 | 1060 | 1055 |
| | 1064 | ← 1064 |
| | 1068 | ← 1070 |
| 1072 | 1072 | |

All byte addresses

```
addi $3, $0, 1044

lw   $4, 4($3)

addi $3, $3, 11

lb   $6, -4($3)

lw   $6, 0($3)
```

Blocks | Words | Bytes

1024 | 1036
| 1040 | 1040
1040 | 1044 | 1045
| 1048
| 1052 | 1052
| | 1053
| | 1054
1056 | 1056 | 1055
| 1060
| 1064 | 1064
| 1068
1072 | 1072 | 1070

```
addi $3, $0, 1044

lw   $4, 4($3)

addi $3, $3, 11

lb   $6, -4($3)

lw   $6, 0($3)
```

lw   : load word
lb   : load byte
addi : add immediate

lw $dest offset($base)

| Blocks | Words | Bytes |
|---|---|---|
| 1024 | 1036 | |
| | | 1040 |
| | 1040 | |
| | | 1045 |
| 1040 | 1044 | |
| | 1048 | 1052 |
| | | 1053 |
| | 1052 | 1054 |
| | | 1055 |
| | 1056 | |
| | 1060 | |
| 1056 | | 1064 |
| | 1064 | |
| | 1068 | 1070 |
| 1072 | 1072 | |

```
addi $3, $0, 1044

lw   $4, 4($3)

addi $3, $3, 11

lb   $6, -4($3)

lw   $6, 0($3)
```

$3

| Blocks | Words | Bytes |

1024

1036

1040

1040

1044

1040

1045

1048

1052

1052

1053

1054

1055

1056

1056

1060

1064

1064

1068

1070

1072

1072

```
addi $3, $0, 1044

lw   $4, 4($3)

addi $3, $3, 11

lb   $6, -4($3)

lw   $6, 0($3)
```

Blocks   Words   Bytes

$3

1024      1036
          1040        1040
1040      1044        1045
          1048        1052
          1052        1053
                      1054
          1056        1055
1056      1060
          1064        1064
          1068        1070
1072      1072

```
addi $3, $0, 1044

lw   $4, 4($3)

addi $3, $3, 11

lb   $6, -4($3)

lw   $6, 0($3)
```

Blocks | Words | Bytes

1024 | 1036

1040

1040 | 1040

1044

1048

1052

$3

1056

1056

1060

1064

1068

1072 | 1072

1040

1045

1052

1053

1054

1055

1064

1070

```
addi $3, $0, 1044

lw   $4, 4($3)

addi $3, $3, 1

lb   $6, -4($3)

lw   $6, 0($3)
```

Blocks    Words    Bytes

$3

1024
1040
1056
1072

1036
1040
1044
1048
1052
1056
1060
1064
1068
1072

1040
1045
1052
1053
1054
1055
1064
1070

```
addi $3, $0, 1044

lw    $4, 4($3)

addi $3, $3, 1

lb    $6, -4($3)

lw    $6, 0($3)
```

**$3**

| Blocks | Words | Bytes |
|---|---|---|

1024

1040

1056

1072

1036

1040

1044

1048

1052

1056

1060

1064

1068

1072

1040

1045

1052

1053

1054

1055

1064

1070

```
addi $3, $0, 1044

lw   $4, 4($3)

addi $3, $3, 1

lb   $6, -4($3)

lw   $6, 0($3)
```

Can't straddle words in one read
• Needs two accesses to RAM
• Would need *two* cycles

$3

| Blocks | Words | Bytes |
|--------|-------|-------|

1024

1036

1040

1040

1045

1044

1040

1048

1052

1052

1053

1056

1054

1055

1060

1056

1064

1064

1070

1068

1072

1072

# Impact of block size on miss rate



- ↑ block size, ↓ miss rate generally

- large block for small cache:
  ↑ miss rate - too few blocks

- ↑ block size: ↑ transfer time between cache and main memory

■ 1 KB
● 8 KB
● 16 KB
◆ 64 KB
◆ 256 KB

Total cache size (KBytes)

dt10 2017 10.26

# Multi-word cache: handling misses

- cache miss on read:
  - same way as single-word block
  - bring back the entire multi-word block from memory

- cache miss on write, given write-through cache:
  - single-word block: disregard hit or miss, just write to cache and write buffer / memory
  - do the same for multi-word block?

```
addi $3, $0, 1044
```

```
lw   $4, 12($3)
```

```
sw   $4, 0($3)
```

```
lw   $4, 4($3)
```

```
lw   $4, 24($3)
```

| $3 | 1044 |
|----|------|
| $4 | ?    |

Single Block
Cache

| ? |
|---|
| ? |
| ? |
| ? |

Blocks

1024

1040

1056

1072

Words

1036
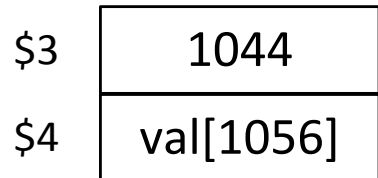1040
1044
1048
1052
1056
1060
1064
1068
1072

```
addi $3, $0, 1044
```
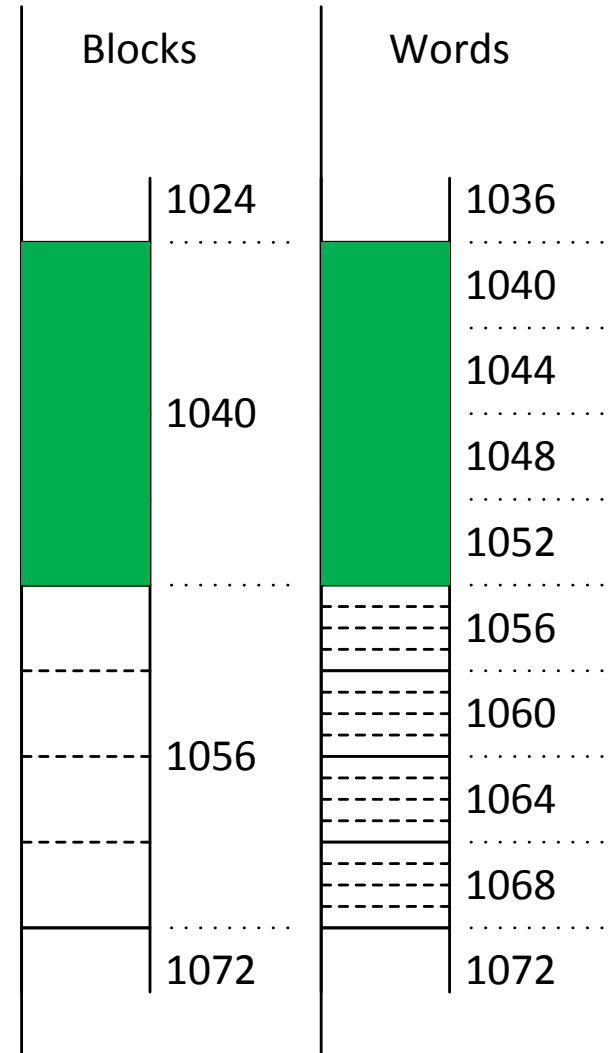
```
lw   $4, 12($3)
```

```
sw   $4, 0($3)
```

```
lw   $4, 4($3)
```

```
lw   $4, 24($3)
```

Single Block
Cache

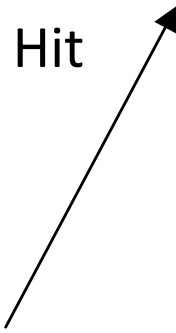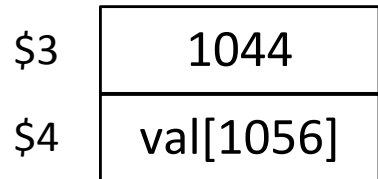Miss!

| ? |
| ? |
| ? |
| ? |

$3 | **1044**
$4 | ?

Blocks | Words

1024 | 1036

1040 | 1040

| 1044

| 1048

| 1052

1056 | 1056

| 1060

| 1064

| 1068

1072 | 1072

```
addi $3, $0, 1044
```

```
lw   $4, 12($3)
```

```
sw   $4, 0($3)
```

```
lw   $4, 4($3)
```

```
lw   $4, 24($3)
```

| $3 | 1044 |
|---|---|
| $4 | ? |

Single Block
Cache

**1056**

Blocks

1024

1040

1056

1072

Words

1036

1040

1044

1048

1052

1056

1060

1064

1068

1072

```
addi $3, $0, 1044

lw   $4, 12($3)

sw   $4, 0($3)

lw   $4, 4($3)

lw   $4, 24($3)
```

Single Block
Cache

Hit

$3     1044

$4     val[1056]

1056

Blocks          Words

1024            1036

1040            1040

1044

1040            1048

1052

1056            1056

1060

1056            1064

1072            1068

1072

```
addi $3, $0, 1044

lw   $4, 12($3)

sw   $4, 0($3)

lw   $4, 4($3)

lw   $4, 24($3)
```
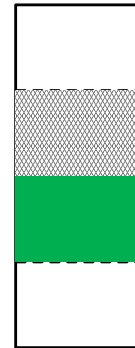
Single Block
Cache

Miss!

1056

$3   1044
$4   val[1056]

Blocks        Words

1024          1036

1040          1040

1044

1056          1048

1072          1052

1056

1060

1064

1068

1072

```
addi $3, $0, 1044

lw   $4, 12($3)

sw   $4, 0($3)

lw   $4, 4($3)

lw   $4, 24($3)
```

| $3 | 1044 |
|----|----------|
| $4 | val[1056] |

Single Block
Cache

**1040**

Blocks          Words

1024            1036

1040            1040

                1044

                1048

1056            1052

                1056

                1060

                1064

                1068

1072            1072

```
addi $3, $0, 1044

lw   $4, 12($3)

sw   $4, 0($3)

lw   $4, 4($3)

lw   $4, 24($3)
```

Single Block
Cache
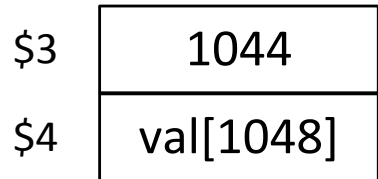
Hit

**1040**

$3 | 1044
$4 | val[1056]

Blocks | Words

1024 | 1036

1040 | 1040
     | 1044
     | 1048
     | 1052

1056 | 1056
     | 1060
     | 1064
     | 1068

1072 | 1072

```
addi $3, $0, 1044

lw   $4, 12($3)

sw   $4, 0($3)

lw   $4, 4($3)

lw   $4, 24($3)
```

Single Block
Cache

Hit

Blocks

Words

1024

1036

1040

1040

1044

1048

1052

1056

1056

1060

1064

1068

1072

1072

$3    1044

$4    val[1048]

**1040**

```
addi $3, $0, 1044

lw   $4, 12($3)

sw   $4, 0($3)

lw   $4, 4($3)

lw   $4, 24($3)
```

Miss!

Single Block
Cache

**1040**

$3  1044

$4  val[1048]

Blocks          Words

1024            1036

                1040

1040            1044

                1048

                1052

                1056

1056            1060

                1064

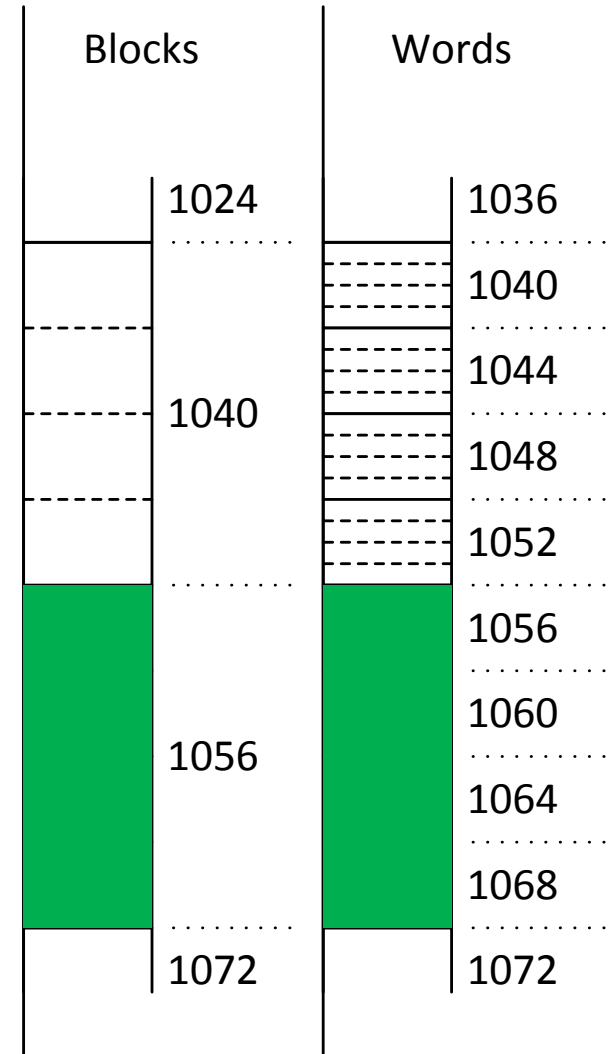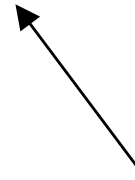                1068

1072            1072

```
addi $3, $0, 1044

lw   $4, 12($3)

sw   $4, 0($3)

lw   $4, 4($3)

lw   $4, 24($3)
```

| $3 | 1044 |
|----|------|
| $4 | val[1048] |

Blocks

Words

Single Block
Cache

**1040**

1024

1040

1056

1072

1036

1040

1044

1048

1052

1056

1060

1064

1068

1072

```
addi $3, $0, 1044

lw   $4, 12($3)

sw   $4, 0($3)

lw   $4, 4($3)

lw   $4, 24($3)
```

Single Block
Cache

Blocks

Words

1024

1036

1040

1040

1044

1048

1052

1056

1056

1060

1064

1068

1072

1072

1056

$3   1044

$4   val[1068]

# Cache performance

- CPU time = (execution cycles + memory stall cycles) × cycle time

# Cache performance

- CPU time = (execution cycles + memory stall cycles) × cycle time

- memory stall cycles = read stall cycles + write stall cycles

# Cache performance

- CPU time = (execution cycles + memory stall cycles)
  × cycle time

- memory stall cycles = read stall cycles + write stall cycles

- read stall cycles =

# Cache performance

- CPU time = (execution cycles + memory stall cycles)
  $$\times \text{ cycle time}$$

- memory stall cycles = read stall cycles + write stall cycles

- read stall cycles = $\dfrac{\text{reads}}{\text{program}} \times \dfrac{\text{read miss}}{\text{rate (\%)}} \times \dfrac{\text{read miss}}{\text{penalty (cycle)}}$

  $\underbrace{\qquad\qquad\qquad}_{\text{read miss per program}}$ $\underbrace{\qquad\qquad}_{\substack{\text{number of cycles} \\ \text{per read miss}}}$

dt10 2017  10.43

# Cache performance

- CPU time = (execution cycles + memory stall cycles)
  × cycle time

- memory stall cycles = read stall cycles + write stall cycles

- read stall cycles = $\underbrace{\dfrac{\text{reads}}{\text{program}} \times \begin{array}{c}\text{read miss}\\\text{rate (\%)}\end{array}}_{\text{read miss per program}} \times \underbrace{\begin{array}{c}\text{read miss}\\\text{penalty (cycle)}\end{array}}_{\substack{\text{number of cycles}\\\text{per read miss}}}$

- write stall cycles =

# Cache performance

- CPU time = (execution cycles + memory stall cycles) $\times$ cycle time

- memory stall cycles = read stall cycles + write stall cycles

- read stall cycles = $\underbrace{\dfrac{\text{reads}}{\text{program}} \times \begin{array}{l}\text{read miss}\\\text{rate (\%)}\end{array}}_{\text{read miss per program}} \times \underbrace{\begin{array}{l}\text{read miss}\\\text{penalty (cycle)}\end{array}}_{\substack{\text{number of cycles}\\\text{per read miss}}}$

- write stall cycles = $\dfrac{\text{writes}}{\text{program}} \times \begin{array}{l}\text{write miss}\\\text{rate (\%)}\end{array} \times \begin{array}{l}\text{write miss}\\\text{penalty (cycle)}\end{array}$

  + write buffer stalls (when full)

# Effect on CPU

- assume hit time insignificant (data transfer time dominated)

- let   c: CPI-no stall,          i: instruction miss rate
          p: miss penalty,        d: data miss rate
          n: instruction count,  f: load/store frequency

- total memory stall cycles: nip + nfdp

# Effect on CPU

- assume hit time insignificant (data transfer time dominated)

- let   $c$: CPI-no stall,         $i$: instruction miss rate
         $p$: miss penalty,       $d$: data miss rate
         $n$: instruction count,  $f$: load/store frequency

- total memory stall cycles: $nip + nfdp$

- total CPU cycles without stall: $nc$
  total CPU cycles with stall: $n(c + ip + fdp)$

# Effect on CPU

- assume hit time insignificant (data transfer time dominated)

- let   $c$: CPI-no stall,        $i$: instruction miss rate
       $p$: miss penalty,       $d$: data miss rate
       $n$: instruction count,  $f$: load/store frequency

- total memory stall cycles: $nip + nfdp$

- total CPU cycles without stall: $nc$
  total CPU cycles with stall: $n(c + ip + fdp)$

- % time on stall:
$$\frac{ip + fdp}{c + ip + fdp}$$

# Faster CPU

- same memory speed, halved CPI
  - % time on stall $= \dfrac{ip + fdp}{(\tfrac{1}{2})c + ip + fdp} > \dfrac{ip + fdp}{c + ip + fdp}$
  - lower CPI results in greater impact of stall cycles

# Faster CPU

- same memory speed, halved CPI
  - % time on stall = $\dfrac{ip + fdp}{(\frac{1}{2})c + ip + fdp} > \dfrac{ip + fdp}{c + ip + fdp}$

  - lower CPI results in greater impact of stall cycles

- same memory speed, halved clock cycle time t
  - miss penalty: 2p
  - total CPU time with new clock: n(c + 2ip + 2fdp)

# Faster CPU

- same memory speed, halved CPI
  - % time on stall $= \dfrac{ip + fdp}{(½)c + ip + fdp} > \dfrac{ip + fdp}{c + ip + fdp}$
  - lower CPI results in greater impact of stall cycles

- same memory speed, halved clock cycle time t
  - miss penalty: $2p$
  - total CPU time with new clock: $n(c + 2ip + 2fdp)$
  - performance improvement $= \dfrac{\text{exec. time with old clock}}{\text{exec. time with new clock}}$

    $= \dfrac{n(c + ip + fdp) \times t}{n(c + 2ip + 2fdp) \times t/2}$

    $= \dfrac{c + ip + fdp}{c + 2ip + 2fdp} \times 2 \; < \; 2$

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view

**Regs.**

**L1**

**L2**

**Mem**

**Disk**

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view
  - capacity view
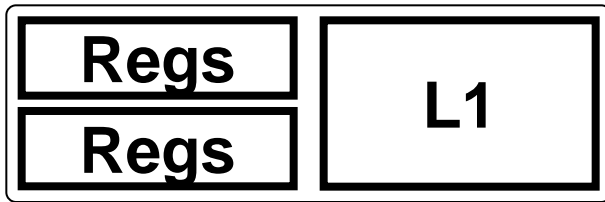
**Regs**  **L1**  **L2**  **Mem**  **Disk**

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view
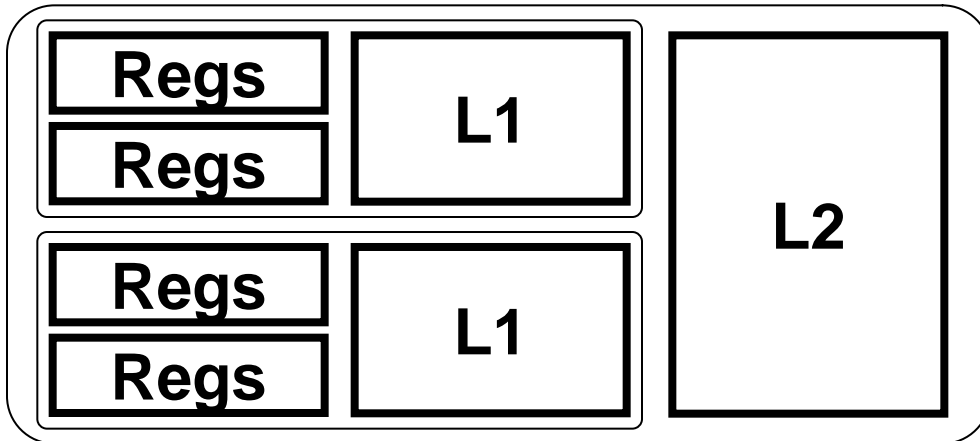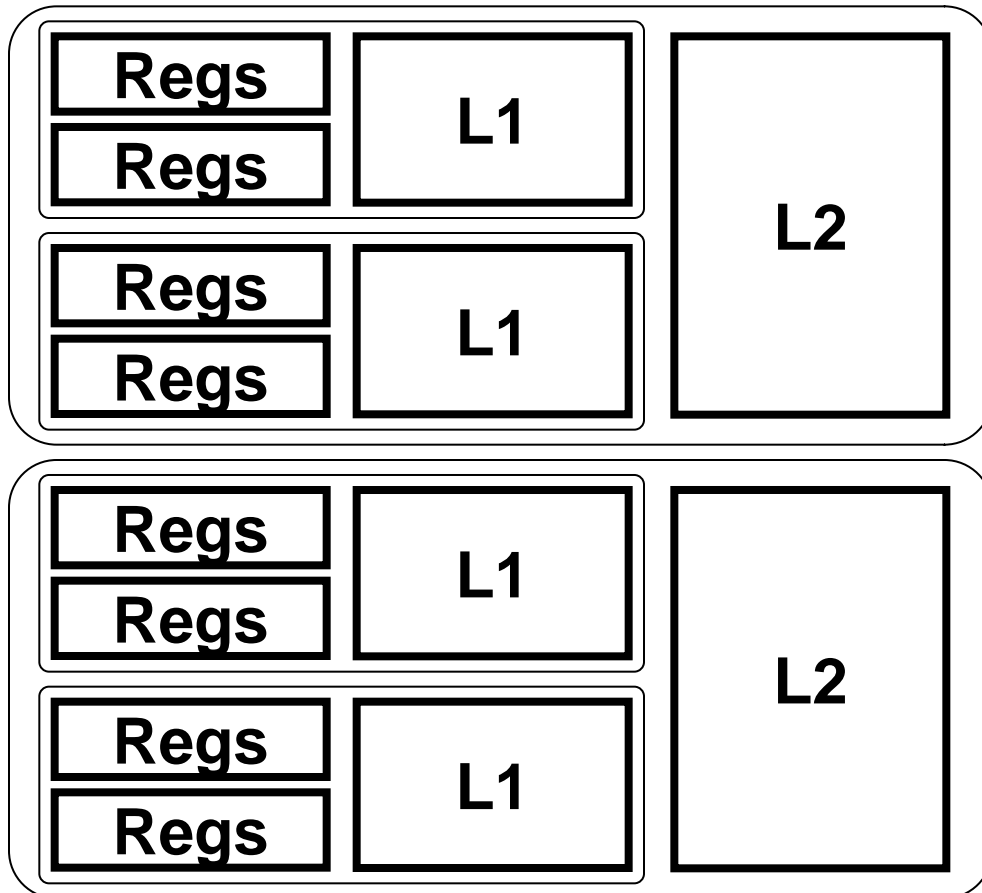  - capacity view
  - physical hierarchy

| **Regs** |

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
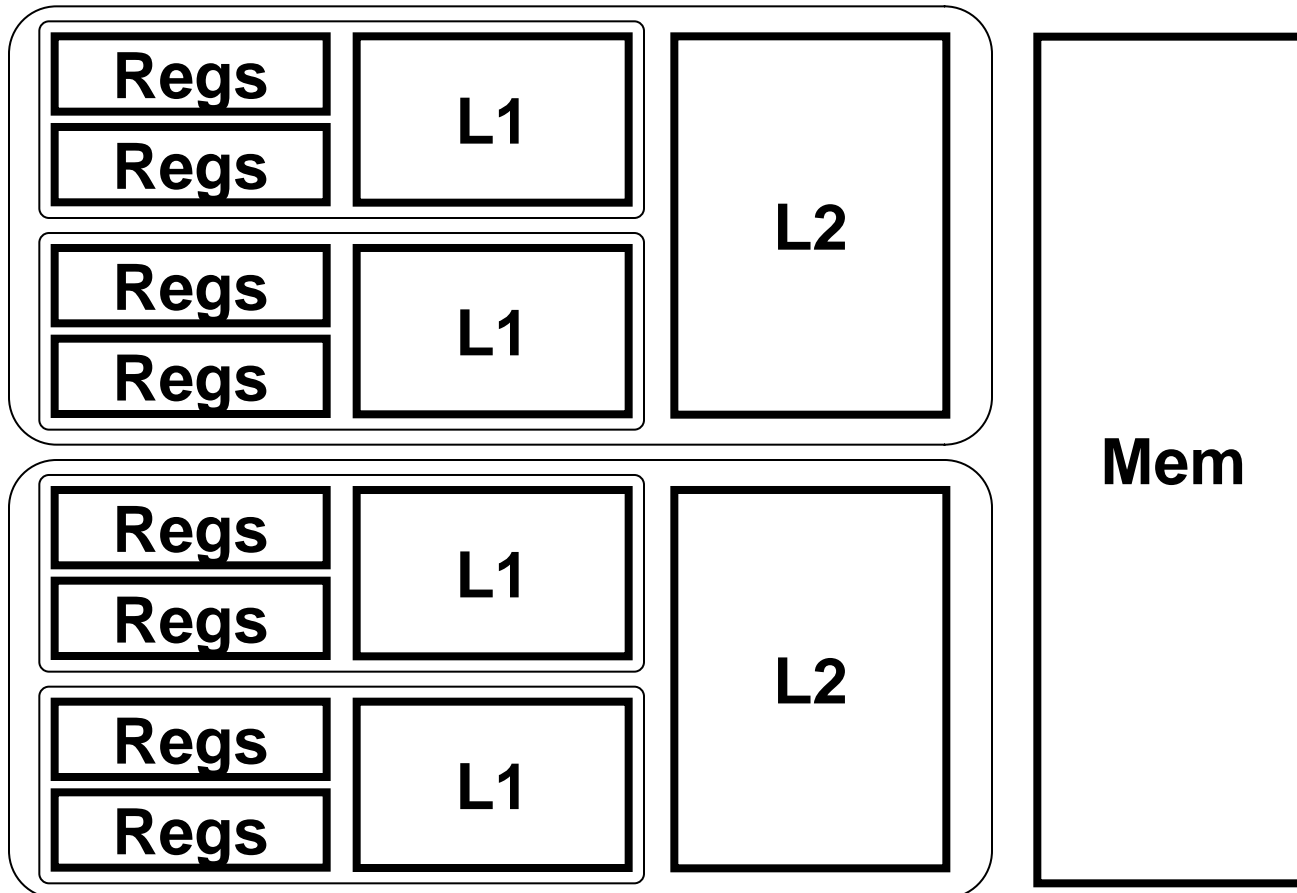  - performance view
  - capacity view
  - physical hierarchy

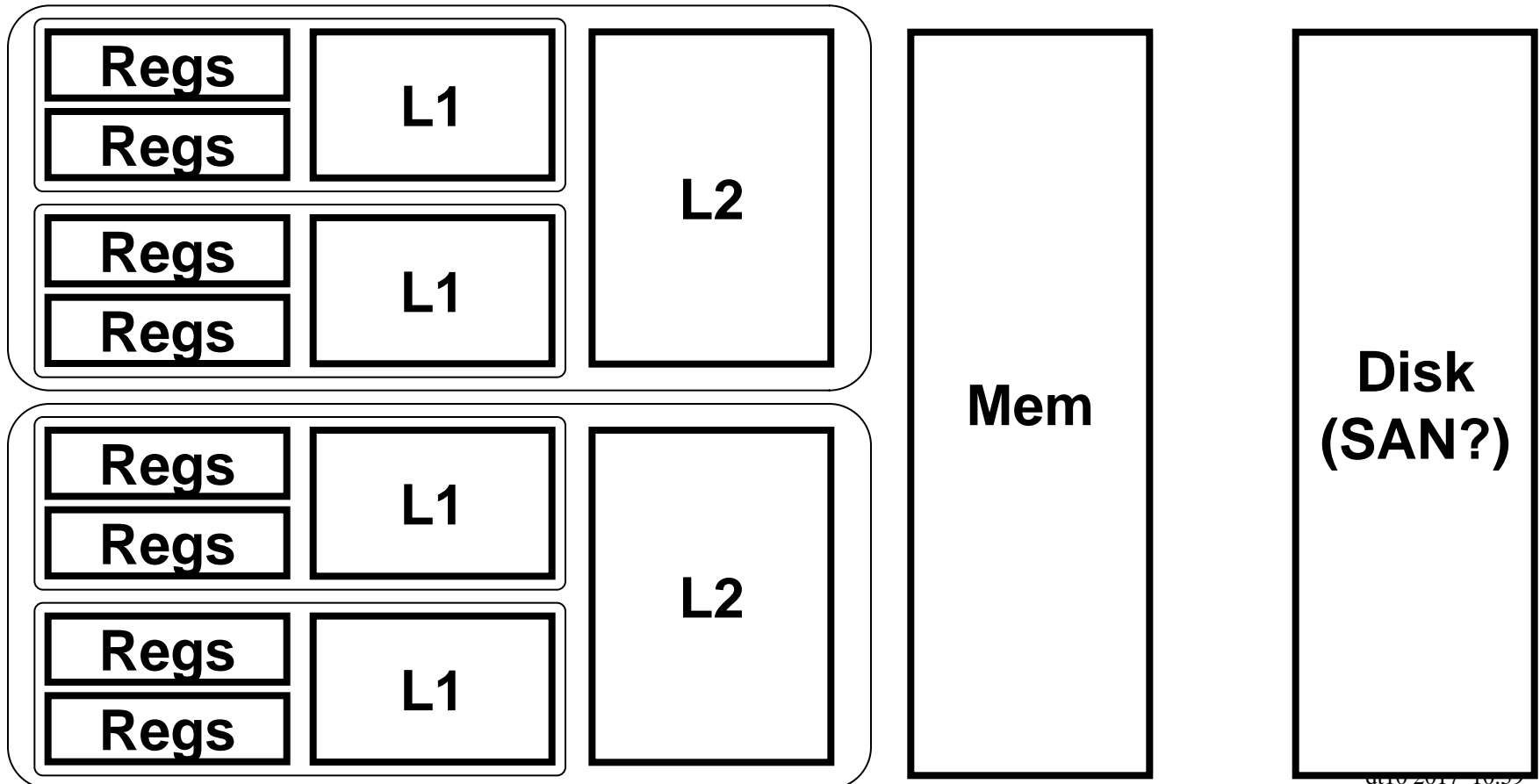| Regs | |
|------|----|
| Regs | L1 |

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view
  - capacity view
  - physical hierarchy

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view
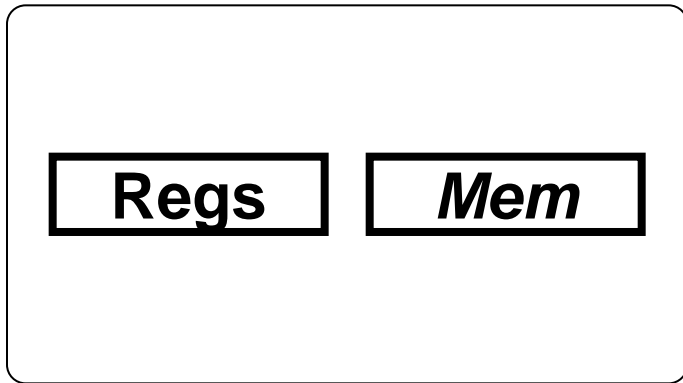  - capacity view
  - physical hierarchy

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view
  - capacity view
  - physical hierarchy

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view
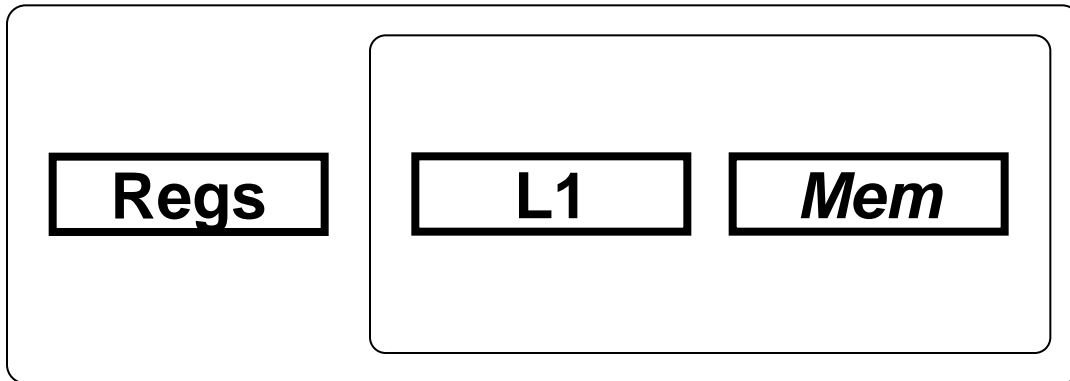  - capacity view
  - physical hierarchy

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view
  - capacity view
  - physical hierarchy
  - abstract hierarchy

| Regs | *Mem* |
|------|-------|

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view
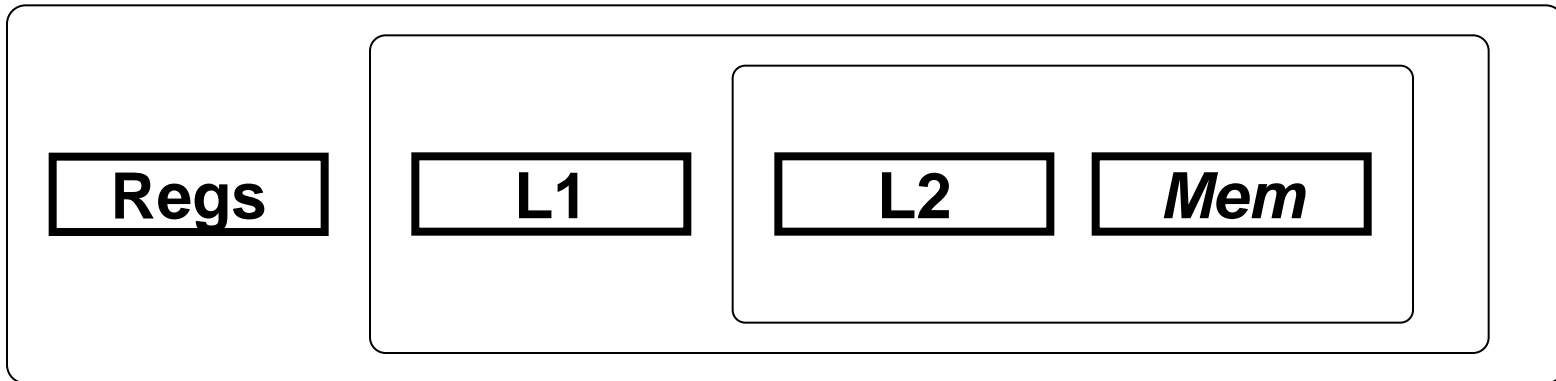  - capacity view
  - physical hierarchy
  - abstract hierarchy

| Regs | | L1 | *Mem* |

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view
  - capacity view
  - physical hierarchy
  - abstract hierarchy

| Regs | L1 | L2 | *Mem* |

# Multi-level cache hierarchy

- how can we look at the cache hierarchy?
  - performance view
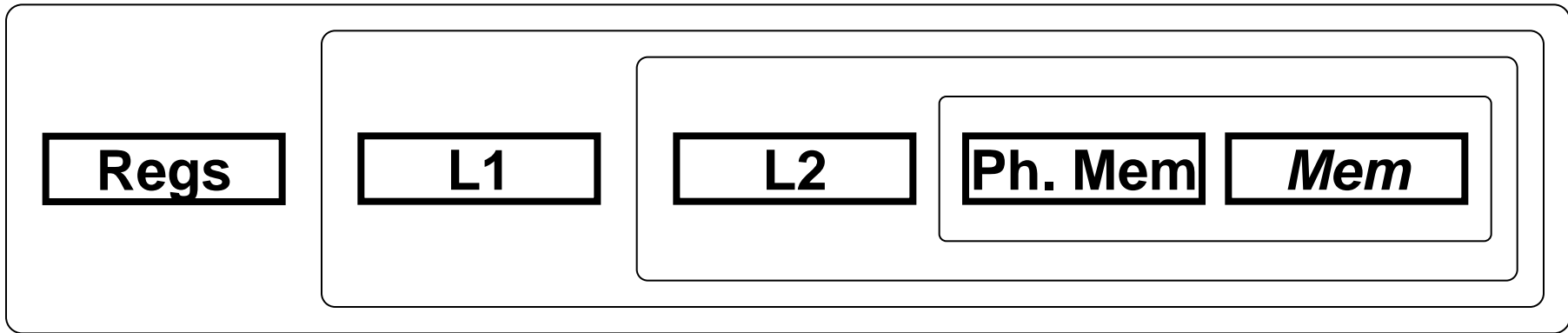  - capacity view
  - physical hierarchy
  - abstract hierarchy

| Regs | | L1 | | L2 | | Ph. Mem | *Mem* |

# Typical scale

- L1
  - size: tens of KB
  - hit time: complete in one clock cycle
  - miss rates: 1-5%

- L2:
  - size: hundreds of KB
  - hit time: a few clock cycles
  - miss rates: 10-20%

- L2 miss rate: fraction of L1 misses also miss L2
  - why so high?

- complex: different block size/placement for L1, L2