

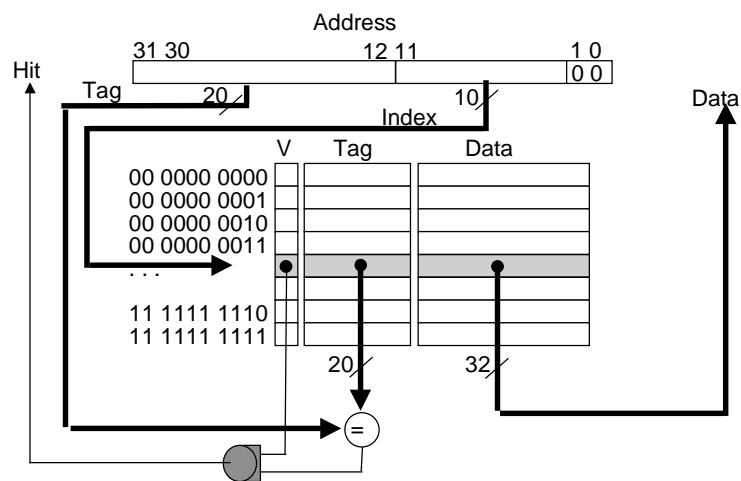


Cache Behavior

Constructing an effective cache requires the balancing of many properties. Bigger is always better, but how it is arranged is also important.

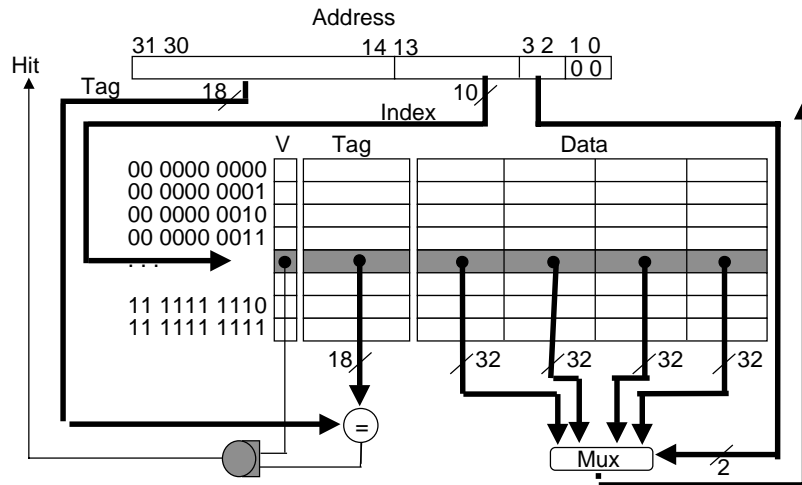
© Larry Snyder, 2000 All rights reserved

Recall Direct Mapped Cache



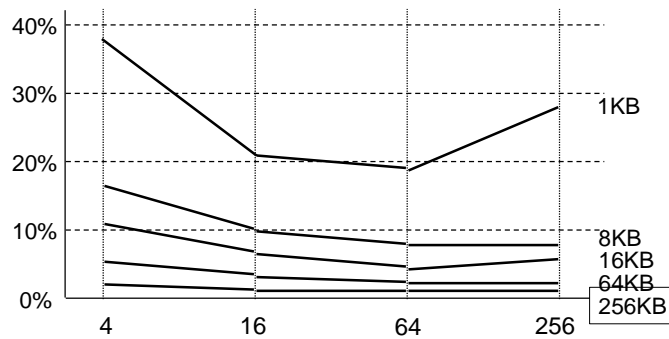
© Larry Snyder, 2000 All rights reserved

Direct Mapped Cache (4 word blks)



© Larry Snyder, 2000 All rights reserved

Miss rate vs block size



VAX traces (Agarwal, 1987)

© Larry Snyder, 2000 All rights reserved

Benefits of Multiword Blocks

Increasing block size improves performance, to a point

Larger blocks increase benefits of spatial locality

Larger blocks = fewer blocks for a given cache size = greater likelihood a useful block is flushed when another block is brought in (conflict misses)

Memory request techniques --

Early restart

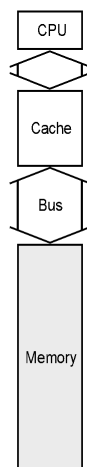
Requested word first

Pgm	Wd	Inst Miss	Data Miss	Effective Miss Ratio
GCC	1	6.1%	2.1%	5.4%
	4	2.0%	1.7%	1.9%
Spice	1	1.2%	1.3%	1.2%
	4	0.3%	0.6%	0.4%

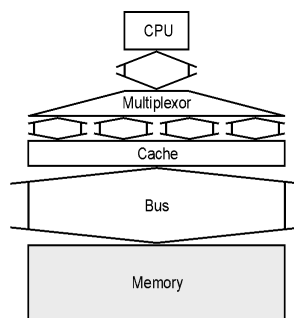
© Larry Snyder, 2000 All rights reserved

Alternative Organizations ...

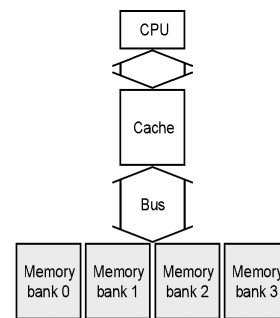
Clks	Operation
1	Send Address
10	Access DRAM
1	Send Data Wd



a. One-word-wide memory organization



b. Wide memory organization



c. Interleaved memory organization

Miss Penalty: $1+4 \times 10+4 \times 1=45$ $1+1 \times 10+1=12$ $1+1 \times 10+4 \times 1=15$

Bytes per Cycle: $4 \times 4 / 45 = 0.35$ $4 \times 4 / 12 = 1.33$ $4 \times 4 / 15 = 1.0$

(a)

(b)

(c)

© Larry Snyder, 2000 All rights reserved

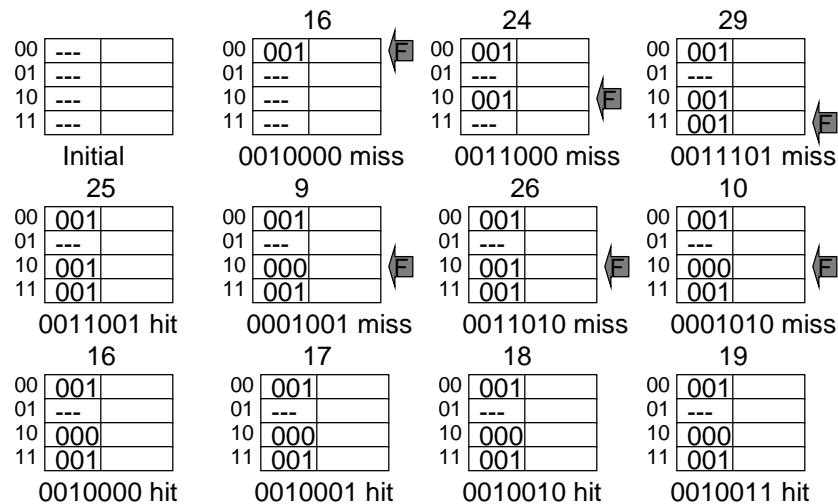
Writing vs Reading Cache

- Writing has two basic forms
 - Write through
 - Write back
- Since writing not on critical path, write buffers
- When an element is written, need it be kept in cache?
 - Load on write ... especially of block > word

© Larry Snyder, 2000 All rights reserved

Reference Sequence

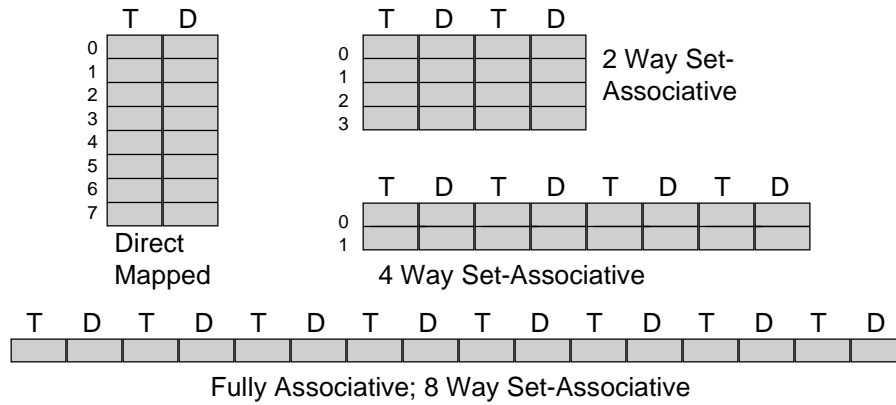
Memory Refs: 16, 24, 29, 25, 9, 26, 10, 16, 17, 18, 19



© Larry Snyder, 2000 All rights reserved

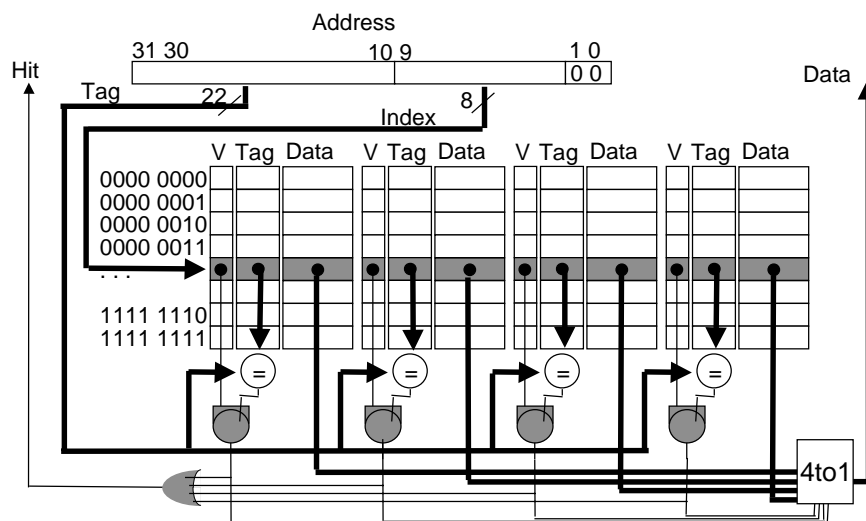
Alternative Designs

Possible arrangements for cache elements



© Larry Snyder, 2000 All rights reserved

Set Associativity



© Larry Snyder, 2000 All rights reserved

Associativity

- 8-way is pretty much the upper limit except for TLB and memory
- Replacement policy
 - Optimal -- a concept that is not realizable
 - Least recently used (LRU)
 - Random
 - Pgm control