



ECE/CS 472/572
Computer Architecture:
Cache Advanced

Prof. Lizhong Chen

Spring 2019

Example: Intrinsity FastMATH

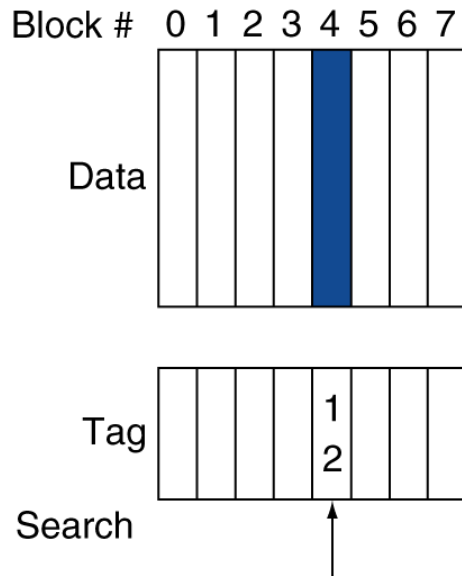
- Embedded MIPS processor
 - 12-stage pipeline
 - Instruction and data access on each cycle
- Split cache: separate I-cache and D-cache
 - Each having 16KB:
 - 256 blocks \times 16 words/block, direct mapped
 - D-cache: write-through or write-back
- SPEC2000 miss rates
 - I-cache: 0.4%, D-cache: 11.4%
 - Weighted average: 3.2%

Associative Caches

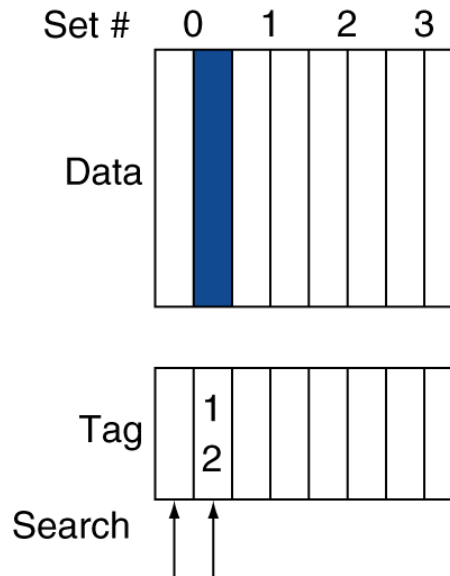
- Fully associative
 - Allow a given block to go in any cache entry
 - Requires all entries to be searched at once
 - Comparator per entry (expensive)
- n -way set associative
 - Each set contains n entries
 - Block number determines which set
 - (Block number) modulo (#Sets in cache)
 - Search all entries in a given set at once
 - n comparators

Associative Cache Example

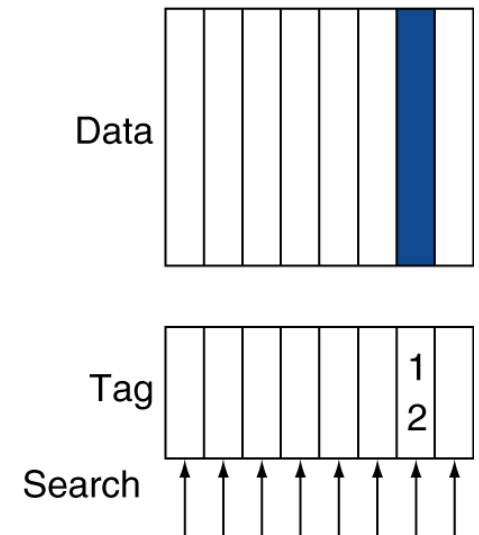
Direct mapped



Set associative



Fully associative



Spectrum of Associativity

- For a cache with 8 entries

**One-way set associative
(direct mapped)**

Block	Tag	Data
0		
1		
2		
3		
4		
5		
6		
7		

Two-way set associative

Set	Tag	Data	Tag	Data
0				
1				
2				
3				

Four-way set associative

Set	Tag	Data	Tag	Data	Tag	Data	Tag	Data
0								
1								

Eight-way set associative (fully associative)

Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data

Associativity Example

- Compare 4-block caches
 - Direct mapped, 2-way set associative, fully associative
 - Block access sequence: 0, 8, 0, 6, 8
- Direct mapped

Block address	Cache index	Hit/miss	Cache content after access			
			0	1	2	3
0	0	miss	Mem[0]			
8	0	miss	Mem[8]			
0	0	miss	Mem[0]			
6	2	miss	Mem[0]		Mem[6]	
8	0	miss	Mem[8]		Mem[6]	

Associativity Example

■ 2-way set associative (LRU)

Block address	Cache index	Hit/miss	Cache content after access			
			Set 0		Set 1	
0	0	miss	Mem[0]			
8	0	miss	Mem[0]	Mem[8]		
0	0	hit	Mem[0]	Mem[8]		
6	0	miss	Mem[0]	Mem[6]		
8	0	miss	Mem[8]	Mem[6]		

■ Fully associative (LRU)

Block address		Hit/miss	Cache content after access			
0		miss	Mem[0]			
8		miss	Mem[0]	Mem[8]		
0		hit	Mem[0]	Mem[8]		
6		miss	Mem[0]	Mem[8]	Mem[6]	
8		hit	Mem[0]	Mem[8]	Mem[6]	

How Much Associativity

- Increased associativity decreases miss rate
 - But with diminishing returns
- Simulation of a system with 64KB D-cache, 16-word blocks, SPEC2000
 - 1-way: 10.3%
 - 2-way: 8.6%
 - 4-way: 8.3%
 - 8-way: 8.1%
- Overhead of increasing associativity?

- 4KB cache



Replacement Policy

- Direct mapped: no choice
- Set associative
 - Prefer non-valid entry, if there is one
 - Otherwise, choose among entries in the set
- Optimal (OPT): oracle & mirror future accesses
- Least-recently used (LRU)
 - Choose the one unused for the longest time
 - Simple for 2-way, manageable for 4-way and 8-way, a bit too hard beyond that
- Random
 - Gives approximately the same performance as LRU for high associativity

Multilevel Caches

- Primary cache attached to CPU (\$L1)
 - Small, but fast
- Level-2 cache services misses from primary cache (\$L2)
 - Larger, slower, but still faster than main memory
- Main memory services L2 cache misses
- Some high-end systems include L3 cache

Multilevel Cache Example

- Given
 - CPU base CPI = 1, clock rate = 4GHz
 - Miss rate/instruction = 2%
 - Main memory access time = 100ns
- Actual CPI with only primary cache:
 - Miss penalty = $100\text{ns} / 0.25\text{ns} = 400$ cycles
 - Effective CPI = $1 + 0.02 \times 400 = 9$

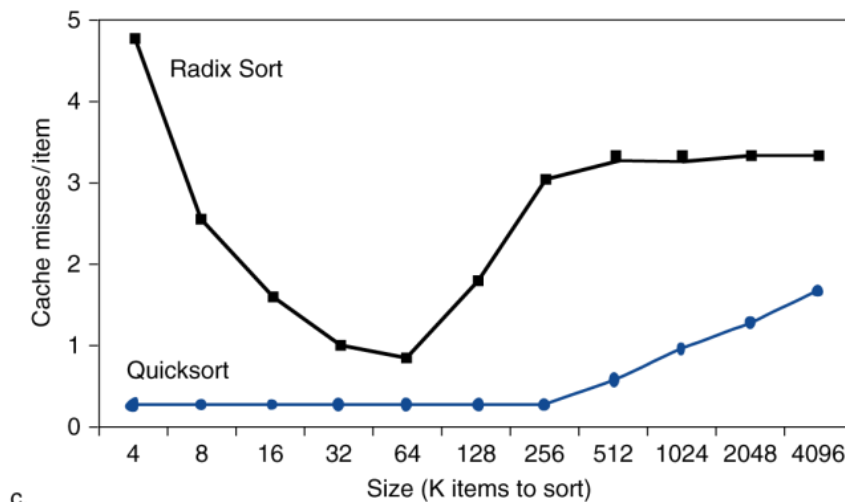
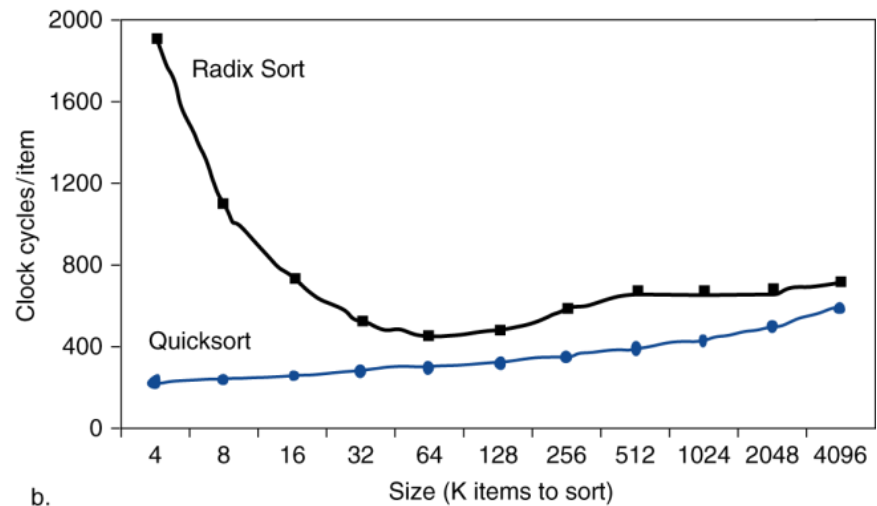
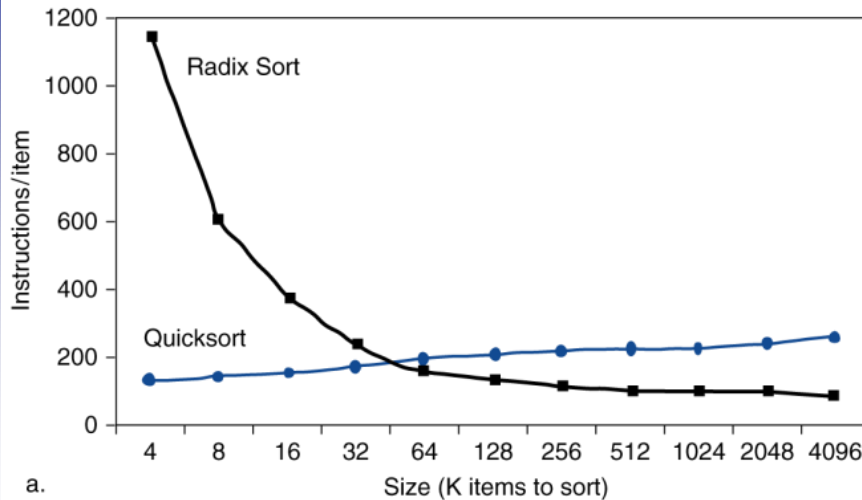
Multilevel Cache Example (cont.)

- Now add L-2 cache
 - Access time = 5ns
 - **Global miss rate** to main memory = 0.5%
- Primary miss with L-2 hit
 - Penalty = $5\text{ns} / 0.25\text{ns} = 20$ cycles
- Primary miss with L-2 miss
 - Extra penalty = 400 cycles
- $\text{CPI} = 1 + 0.02 \times 20 + 0.005 \times 400 = 3.4$
- $\text{Speedup} = 9 / 3.4 = 2.6$

Multilevel Cache Considerations

- Primary cache
 - Focus on minimizing hit time
- L-2 cache
 - Focus on low miss rate to avoid DRAM access
 - Hit time has less overall impact
- Results
 - L-1: small and fast
 - L-2: large and slow

Interactions with Software



- Standard algorithm analysis often ignores the impact of memory hierarchy

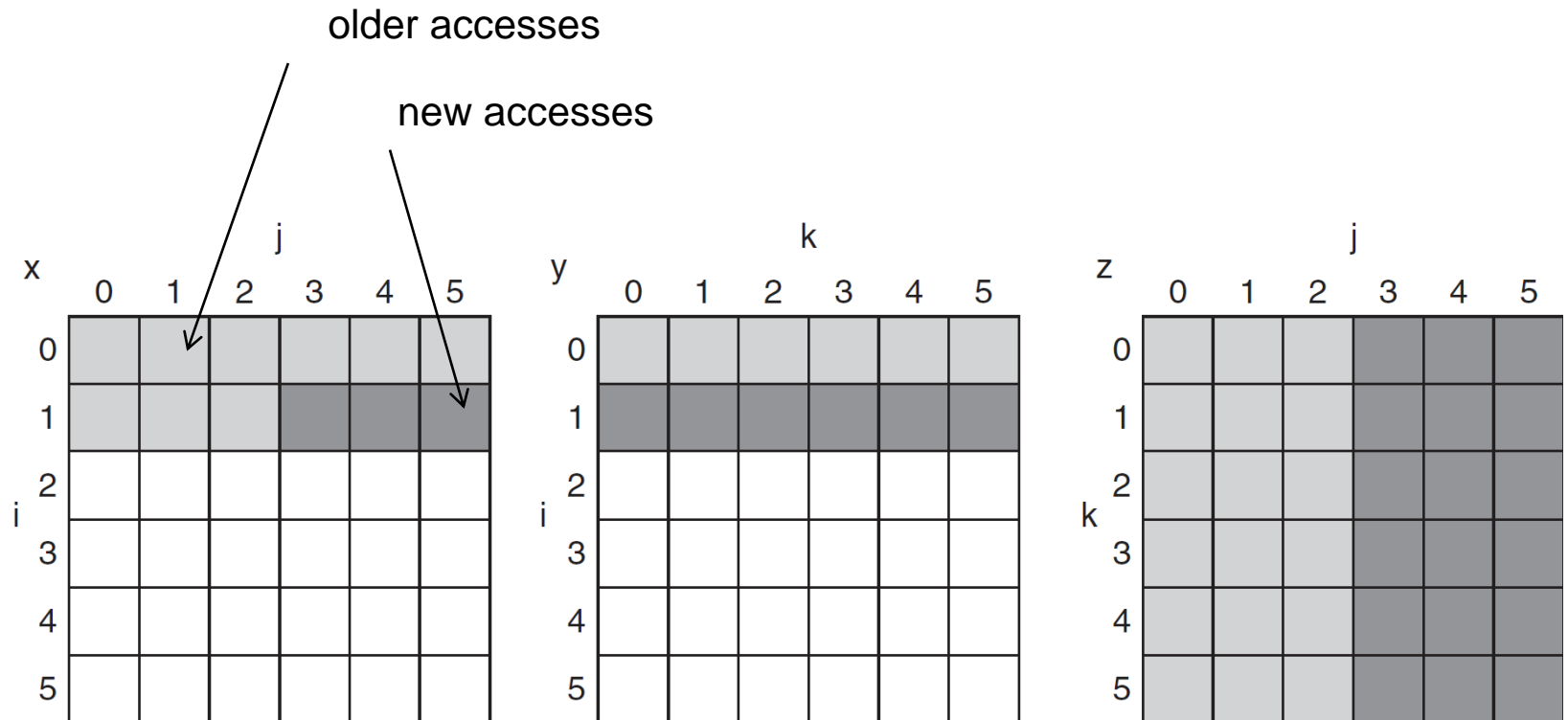
Software Optimization via Blocking

- Goal: maximize accesses to data before it is replaced
- Consider inner loops of DGEMM:
 - Double-precision GEneral Matrix Multiply (DGEMM)

```
for (int j = 0; j < n; ++j)
{
    double cij = C[i+j*n]; // cij = C[i][j]
    for( int k = 0; k < n; k++ )
        cij += A[i+k*n] * B[k+j*n];
        // cij += A[i][k]*B[k][j];
    C[i+j*n] = cij;
}
```


DGEMM Access Pattern

- C, A, and B arrays



Cache Blocked DGEMM

```
1 #define BLOCKSIZE 32
2 void do_block (int n, int si, int sj, int sk, double *A, double
3 *B, double *C)
4 {
5   for (int i = si; i < si+BLOCKSIZE; ++i)
6     for (int j = sj; j < sj+BLOCKSIZE; ++j)
7       {
8         double cij = C[i+j*n];/* cij = C[i][j] */
9         for( int k = sk; k < sk+BLOCKSIZE; k++ )
10          cij += A[i+k*n] * B[k+j*n];/* cij+=A[i][k]*B[k][j] */
11        C[i+j*n] = cij;/* C[i][j] = cij */
12      }
13 }
14 void dgemm (int n, double* A, double* B, double* C)
15 {
16   for ( int sj = 0; sj < n; sj += BLOCKSIZE )
17     for ( int si = 0; si < n; si += BLOCKSIZE )
18       for ( int sk = 0; sk < n; sk += BLOCKSIZE )
19         do_block(n, si, sj, sk, A, B, C);
20 }
```

Blocked DGEMM Access Pattern

Diagram illustrating the access pattern for matrix **x** in a blocked DGEMM operation. The matrix is 6x6, with row index **i** and column index **j**. The access pattern shows a 3x3 block of elements (rows 0-2, columns 0-2) being accessed, indicated by gray shading.

	0	1	2	3	4	5
0	Gray	Gray	Gray			
1	Dark Gray	Dark Gray	Dark Gray			
2						
3						
4						
5						

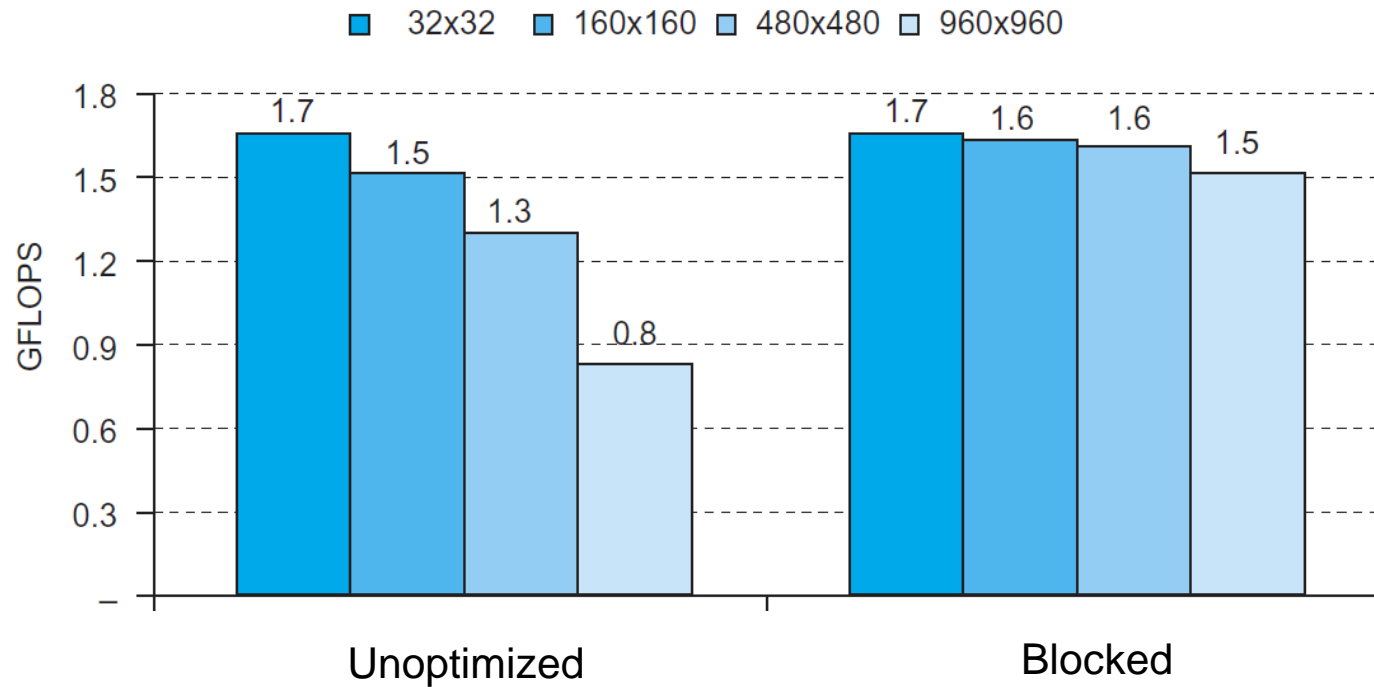
Diagram illustrating the access pattern for matrix **y** in a blocked DGEMM operation. The matrix is 6x6, with row index **i** and column index **k**. The access pattern shows a 3x3 block of elements (rows 0-2, columns 0-2) being accessed, indicated by gray shading.

	0	1	2	3	4	5
0	Gray	Gray	Gray			
1	Dark Gray	Dark Gray	Dark Gray			
2						
3						
4						
5						

Diagram illustrating the access pattern for matrix **z** in a blocked DGEMM operation. The matrix is 6x6, with row index **k** and column index **j**. The access pattern shows a 3x3 block of elements (rows 0-2, columns 0-2) being accessed, indicated by gray shading.

	0	1	2	3	4	5
0	Gray	Dark Gray	Dark Gray			
1	Gray	Dark Gray	Dark Gray			
2	Gray	Dark Gray	Dark Gray			
3						
4						
5						

Blocked DGEMM Access Pattern



- Optimized (Blocked): less than 10% slowdown even with 900 times larger

Sources of Misses: 3Cs

- Cold misses
 - First access to a block
- Capacity misses
 - Due to finite cache size
 - A replaced block is later accessed again
- Conflict misses
 - In a non-fully associative cache
 - Due to competition for entries in a set

Cache Design Trade-offs

Design change	Effect on miss rate	Negative performance effect
Increase cache size	Decrease capacity misses	May increase access time
Increase associativity	Decrease conflict misses	May increase access time
Increase block size	Decrease cold misses	Increases miss penalty. For very large block size, may increase miss rate due to pollution.

Multilevel On-Chip Caches

Characteristic	ARM Cortex-A8	Intel Nehalem
L1 cache organization	Split instruction and data caches	Split instruction and data caches
L1 cache size	32 KiB each for instructions/data	32 KiB each for instructions/data per core
L1 cache associativity	4-way (I), 4-way (D) set associative	4-way (I), 8-way (D) set associative
L1 replacement	Random	Approximated LRU
L1 block size	64 bytes	64 bytes
L1 write policy	Write-back, Write-allocate(?)	Write-back, No-write-allocate
L1 hit time (load-use)	1 clock cycle	4 clock cycles, pipelined
L2 cache organization	Unified (instruction and data)	Unified (instruction and data) per core
L2 cache size	128 KiB to 1 MiB	256 KiB (0.25 MiB)
L2 cache associativity	8-way set associative	8-way set associative
L2 replacement	Random(?)	Approximated LRU
L2 block size	64 bytes	64 bytes
L2 write policy	Write-back, Write-allocate (?)	Write-back, Write-allocate
L2 hit time	11 clock cycles	10 clock cycles
L3 cache organization	-	Unified (instruction and data)
L3 cache size	-	8 MiB, shared
L3 cache associativity	-	16-way set associative
L3 replacement	-	Approximated LRU
L3 block size	-	64 bytes
L3 write policy	-	Write-back, Write-allocate
L3 hit time	-	35 clock cycles