

计算机组成与系统结构

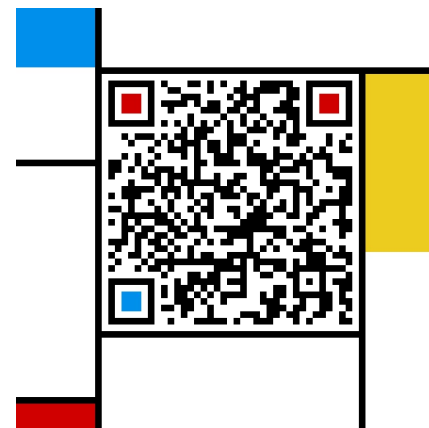
Computer Organization & System Architecture

Huang Kejie (黄科杰) 百人计划研究员

Office: 玉泉校区老生仪楼 304

Email address: huangkejie@zju.edu.cn

HP: 17706443800



Handling Stores with Write-Through

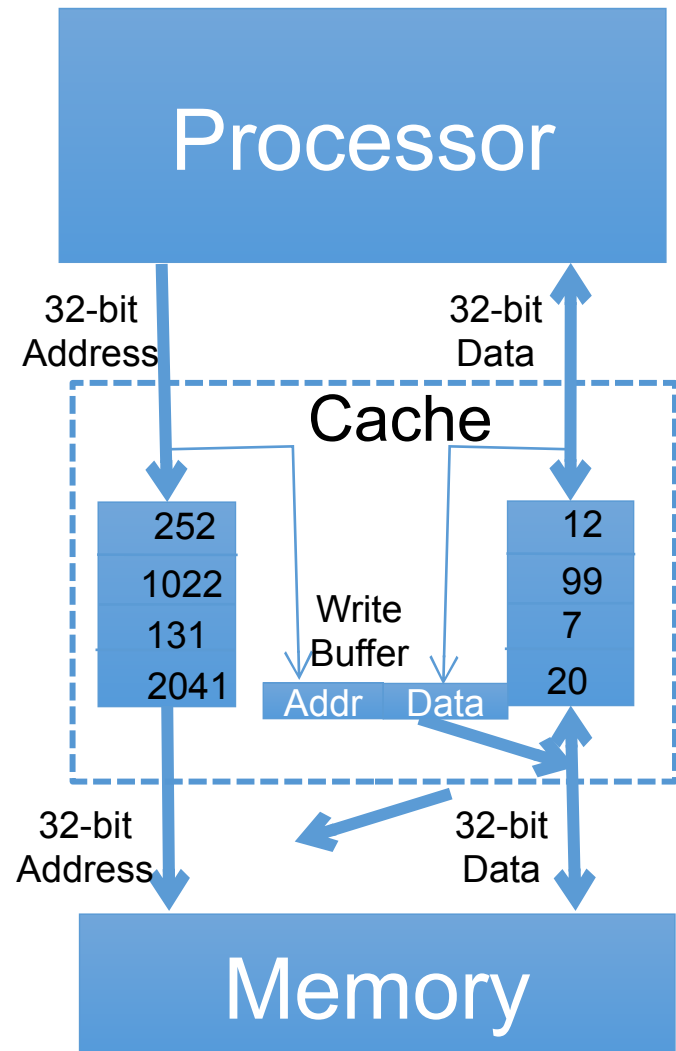
- Store instructions write to memory, changing values
- Need to make sure cache and memory have same values on writes: two policies

1) **Write-Through Policy:** write cache and write through the cache to memory

- Every write eventually gets to memory
- Too slow, so include Write Buffer to allow processor to continue once data in Buffer
- Buffer updates memory in parallel to processor

Write-Through Cache

- Write both values in cache and in memory
- Write buffer stops CPU from stalling if memory cannot keep up
- Write buffer may have multiple entries to absorb bursts of writes
- What if store misses in cache?

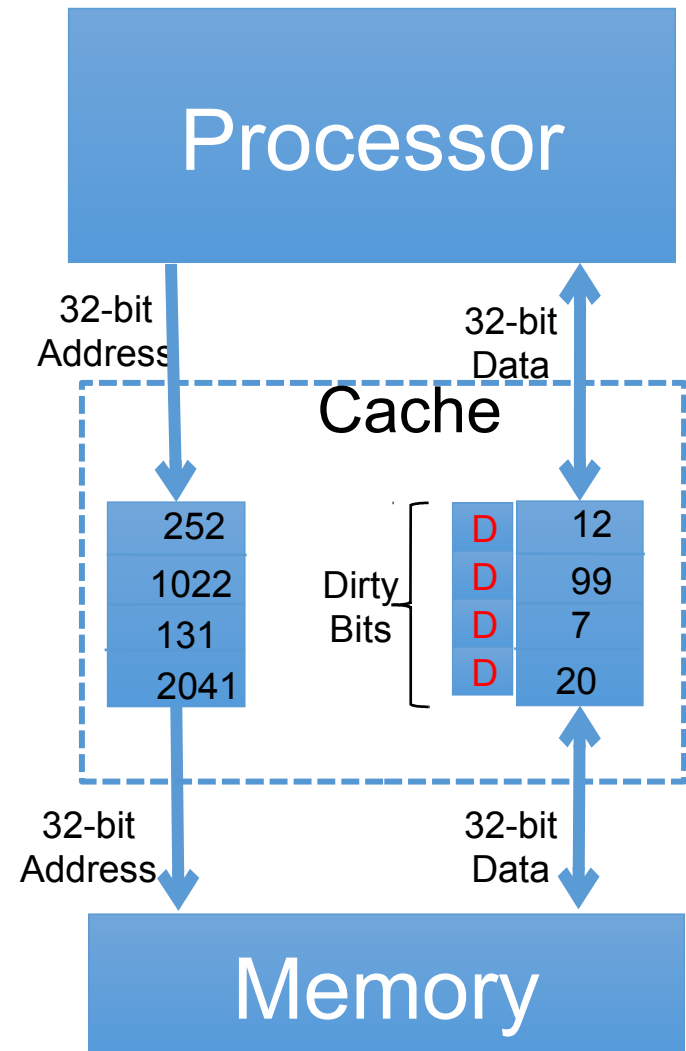


Handling Stores with Write-Back

- 2) **Write-Back Policy**: write only to cache and then write cache block back to memory when evict block from cache
- Writes collected in cache, only single write to memory per block
 - Include bit to see if wrote to block or not, and then only write back if bit is set
 - Called “**Dirty**” bit (writing makes it “dirty”)

Write-Back Cache

- Store/cache hit, write data in cache only and set dirty bit
 - Memory has stale value
- Store/cache miss, read data from memory, then update and set dirty bit
 - “Write-allocate” policy
- Load/cache hit, use value from cache
- On any miss, write back evicted block, only if dirty. Update cache with new block and clear dirty bit



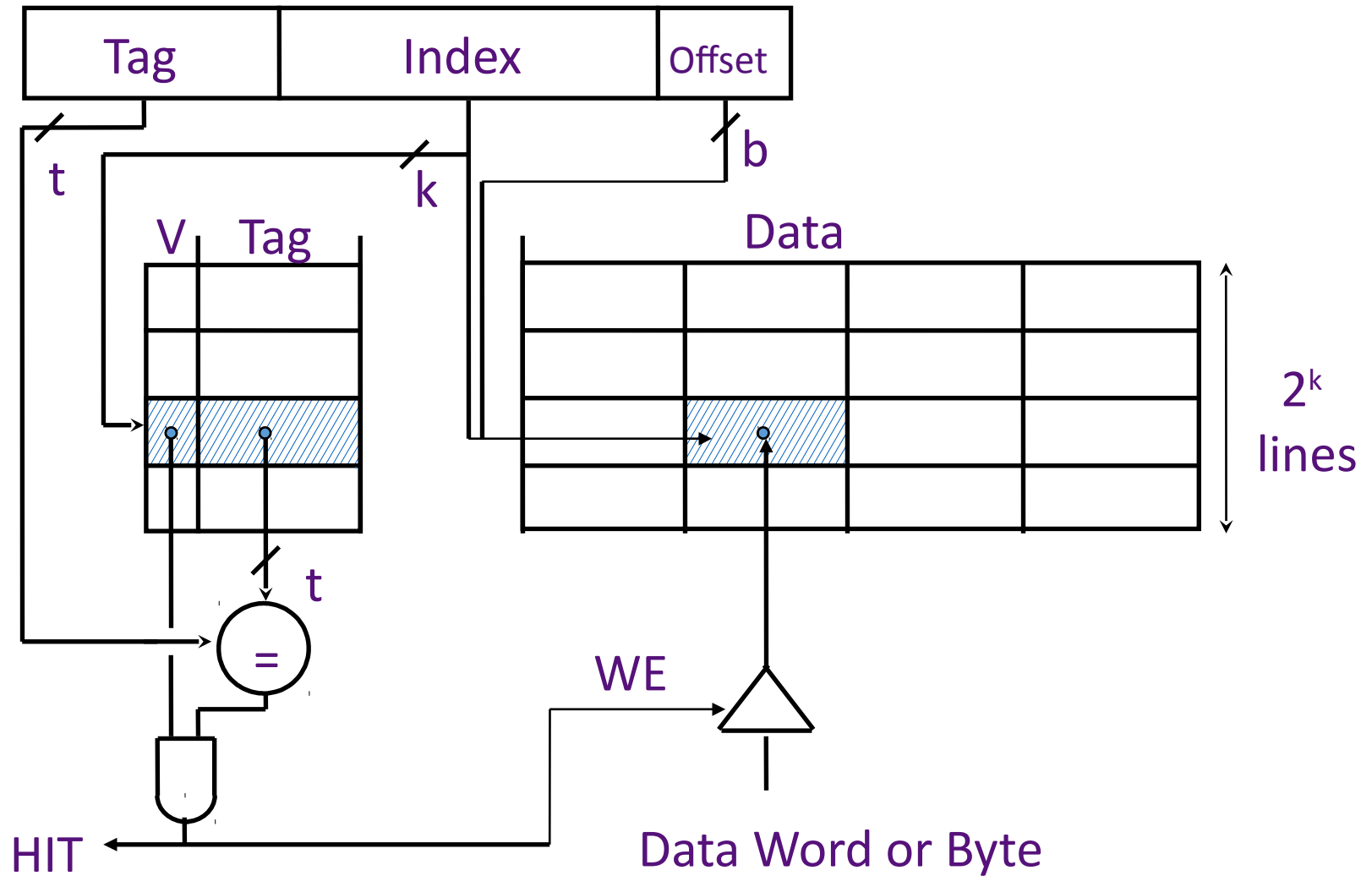
Write-Through vs. Write-Back

- Write-Through:
 - Simpler control logic
 - More predictable timing
simplifies processor control logic
 - Easier to make reliable,
since memory always has
copy of data (big idea:
Redundancy!)
- Write-Back
 - More complex control logic
 - More variable timing (0,1,2
memory accesses per
cache access)
 - Usually reduces write traffic
 - Harder to make reliable,
sometimes cache has only
copy of data

Write Policy Choices

- Cache Hit:
 - **Write through:** writes both cache & memory on every access
 - Generally higher memory traffic but simpler pipeline & cache design
 - **Write back:** writes cache only, memory written only when dirty entry evicted
 - A dirty bit per line reduces write-back traffic
 - Must handle 0, 1, or 2 accesses to memory for each load/store
- Cache Miss:
 - **No write allocate:** only write to main memory
 - **Write allocate** (aka fetch on write): fetch into cache
- Common combinations:
 - Write through and no write allocate
 - Write back with write allocate

Write Performance



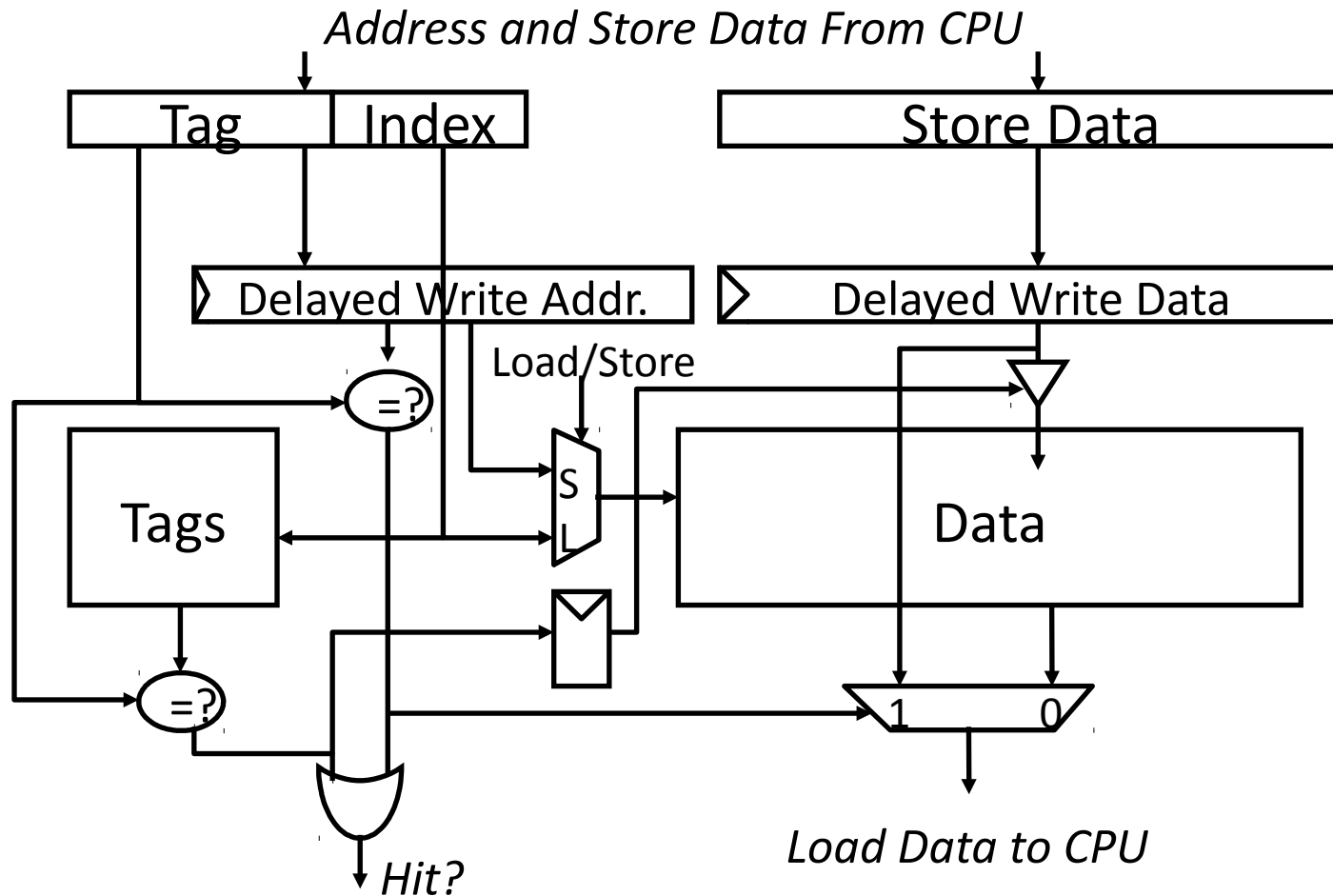
Reducing Write Hit Time

Problem: Writes take two cycles in memory stage, one cycle for tag check plus one cycle for data write if hit

Solutions:

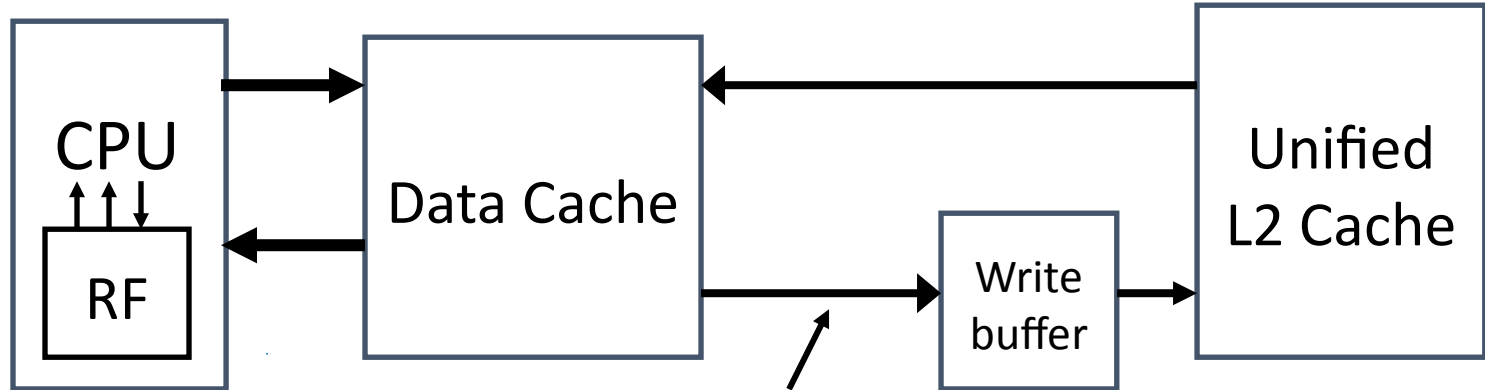
- Design data RAM that can perform read and write in one cycle, restore old value after tag miss
- Fully-associative (CAM Tag) caches: Word line only enabled if hit
- Pipelined writes: Hold write data for store in single buffer ahead of cache, write cache data during next store's tag check

Pipelining Cache Writes



Data from a store hit is written into data portion of cache during tag access of subsequent store

Write Buffer to Reduce Read Miss Penalty



Evicted dirty lines for writeback cache
OR

All writes in writethrough cache

Processor is not stalled on writes, and read misses can go ahead of write to main memory

Problem: Write buffer may hold updated value of location needed by a read miss

Simple solution: on a read miss, wait for the write buffer to go empty

Faster solution: Check write buffer addresses against read miss addresses, if no match, allow read miss to go ahead of writes, else, return value in write buffer

Reducing Tag Overhead with Sub-Blocks

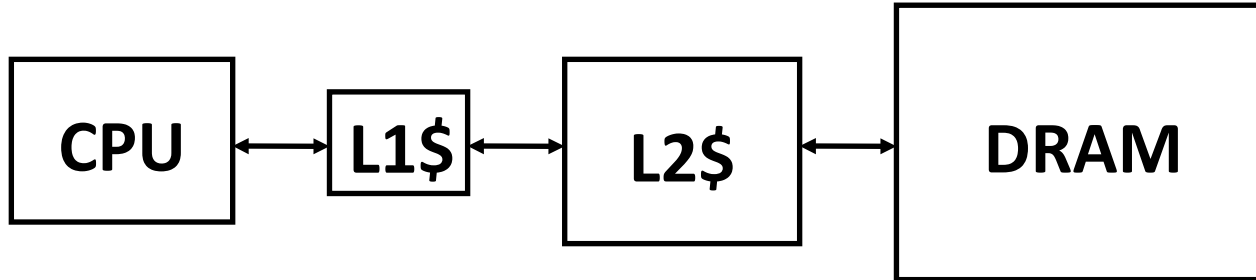
- **Problem:** Tags are too large, i.e., too much overhead
 - Simple solution: Larger lines, but miss penalty could be large.
- **Solution:** Sub-block placement (aka sector cache)
 - A valid bit added to units smaller than full line, called sub-blocks
 - Only read a sub-block on a miss
 - *If a tag matches, is the word in the cache?*

100
300
204

1		1		1		1	
1		1		0		0	
0		1		0		1	

Reducing Tag Overhead with Sub-Blocks

- **Problem:** A memory cannot be large and fast
- **Solution:** Increasing sizes of cache at each level



Local miss rate = misses in cache / accesses to cache

Global miss rate = misses in cache / CPU memory accesses

Misses per instruction = misses in cache / number of instructions

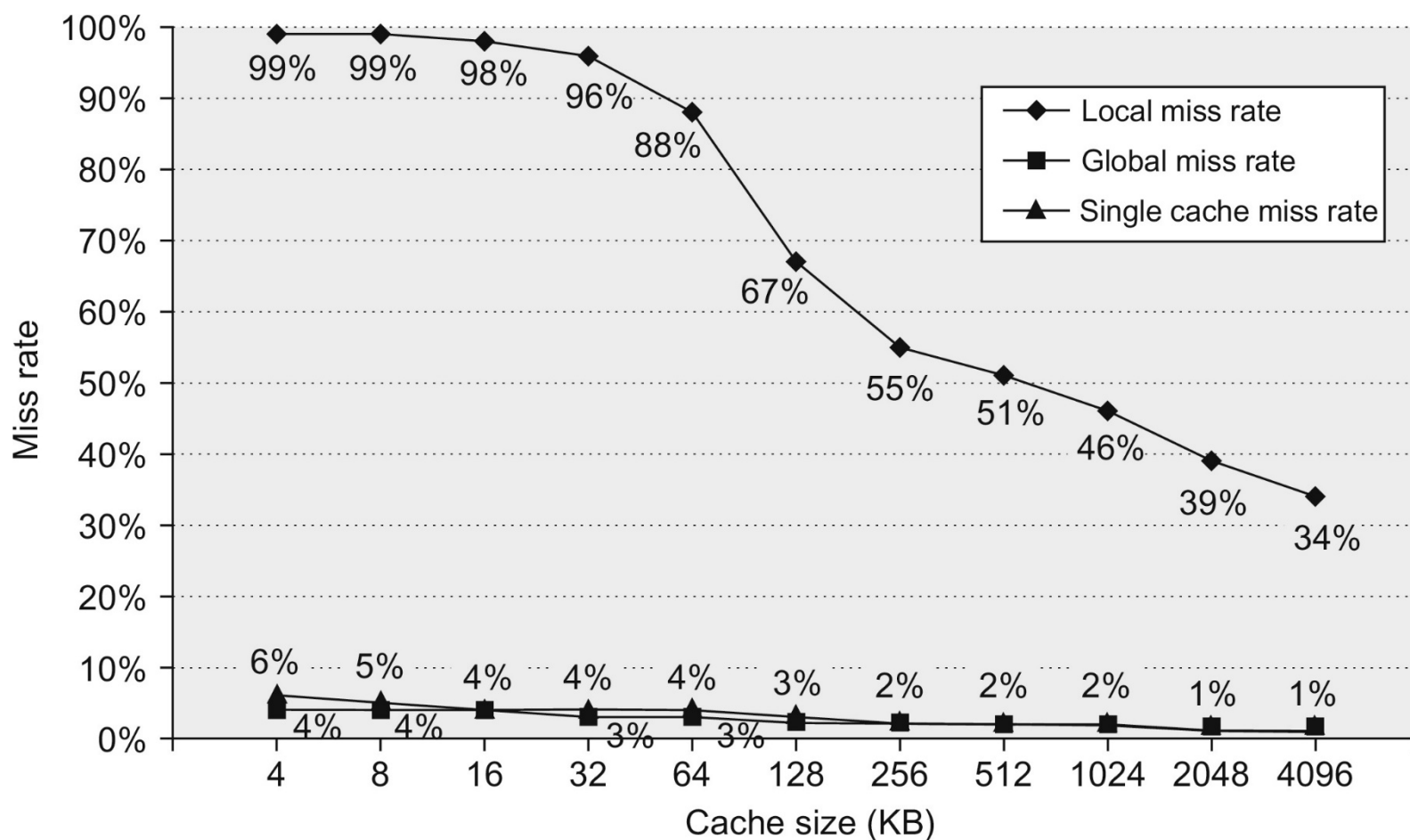


Figure B.14 Miss rates versus cache size for multilevel caches. Second-level caches *smaller* than the sum of the two 64 KiB first-level caches make little sense, as reflected in the high miss rates. After 256 KiB the single cache is within 10% of the global miss rates. The miss rate of a single-level cache versus size is plotted against the local miss rate and global miss rate of a second-level cache using a 32 KiB first-level cache. The L2 caches (unified) were two-way set associative with replacement. Each had split L1 instruction and data caches that were 64 KiB two-way set associative with LRU replacement. The block size for both L1 and L2 caches was 64 bytes. Data were collected as in Figure B.4.

Presence of L2 influences L1 design

- Use smaller L1 if there is also L2
 - Trade increased L1 miss rate for reduced L1 hit time
 - Backup L2 reduces L1 miss penalty
 - Reduces average access energy
- Use simpler write-through L1 with on-chip L2
 - Write-back L2 cache absorbs write traffic, doesn't go off-chip
 - At most one L1 miss request per L1 access (no dirty victim write back) simplifies pipeline control
 - Simplifies coherence issues
 - Simplifies error recovery in L1 (can use just parity bits in L1 and reload from L2 when parity error detected on L1 read)

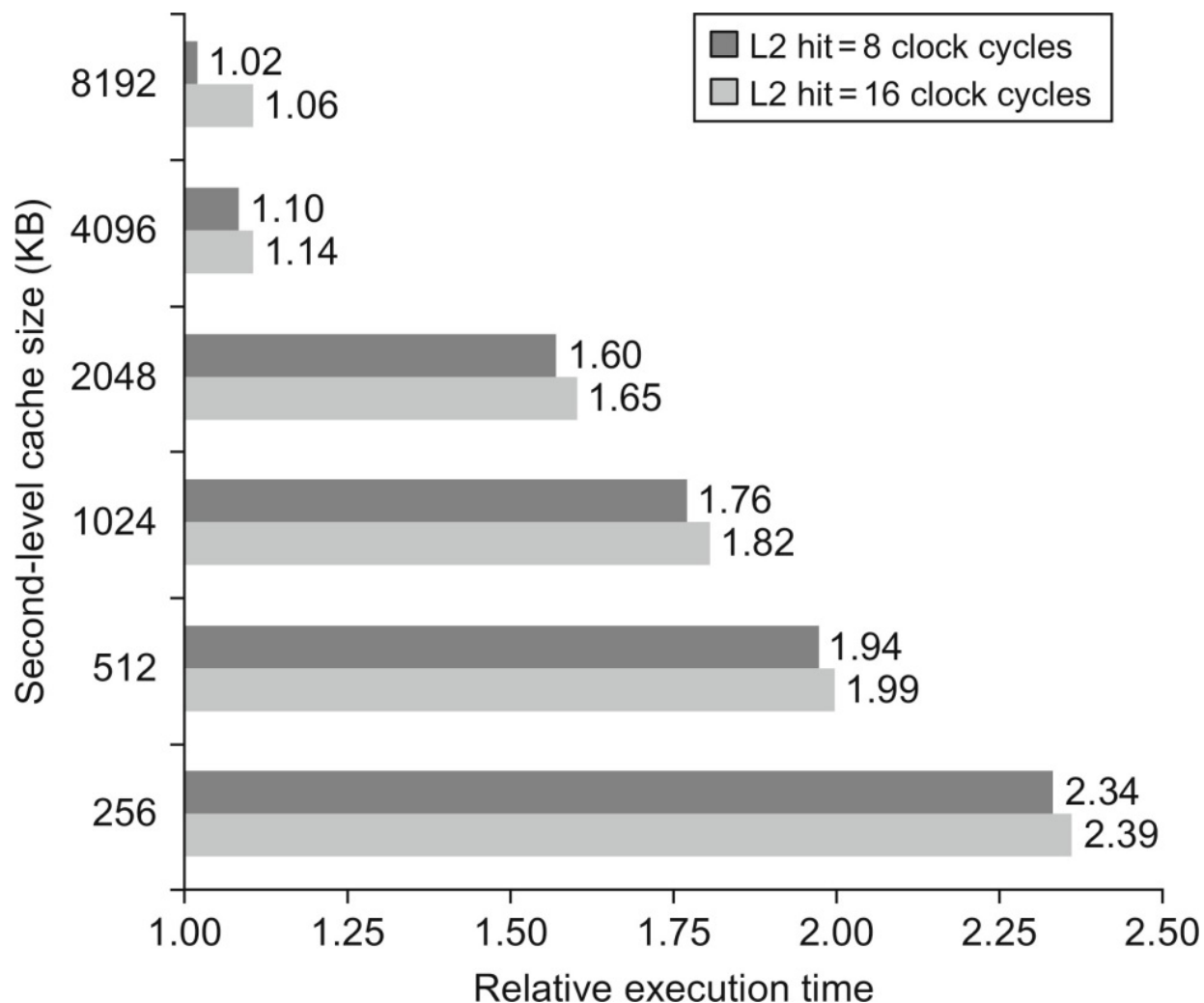
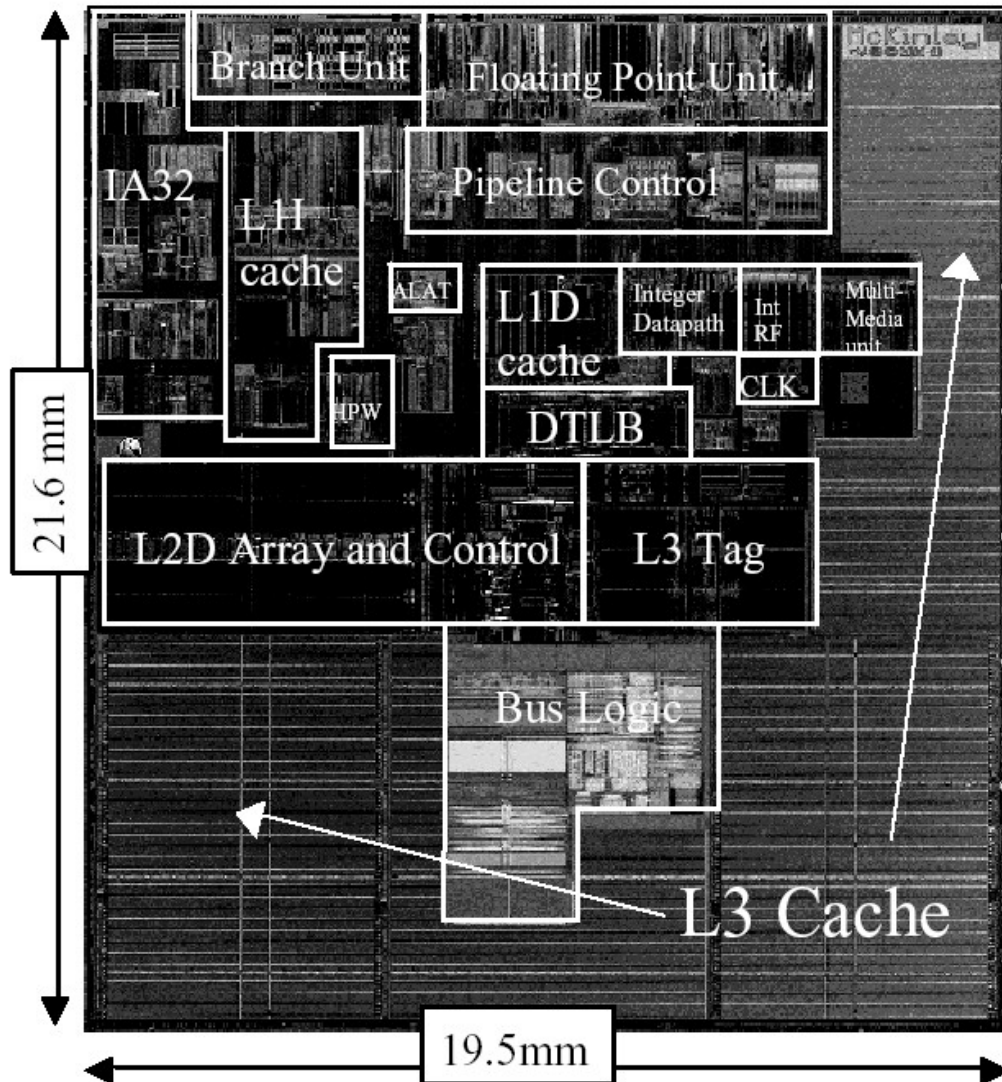


Figure B.15 Relative execution time by second-level cache size. The two bars are for different clock cycles for an L2 cache hit. The reference execution time of 1.00 is for an 8192 KiB second-level cache with a 1-clock-cycle latency on a second-level hit. These data were collected the same way as in Figure B.14, using a simulator to imitate the Alpha 21264.

Inclusion Policy

- Inclusive multilevel cache:
 - Inner cache can only hold lines also present in outer cache
 - External coherence snoop access need only check outer cache
- Exclusive multilevel caches:
 - Inner cache may hold lines not in outer cache
 - Swap lines between inner/outer caches on miss
 - Used in AMD Athlon with 64KB primary and 256KB secondary cache
- Why choose one type or the other?

Itanium-2 On-Chip Caches (Intel/HP, 2002)



Level 1: 16KB, 4-way s.a., 64B line, quad-port (2 load+2 store), single cycle latency

Level 2: 256KB, 4-way s.a., 128B line, quad-port (4 load or 4 store), five cycle latency

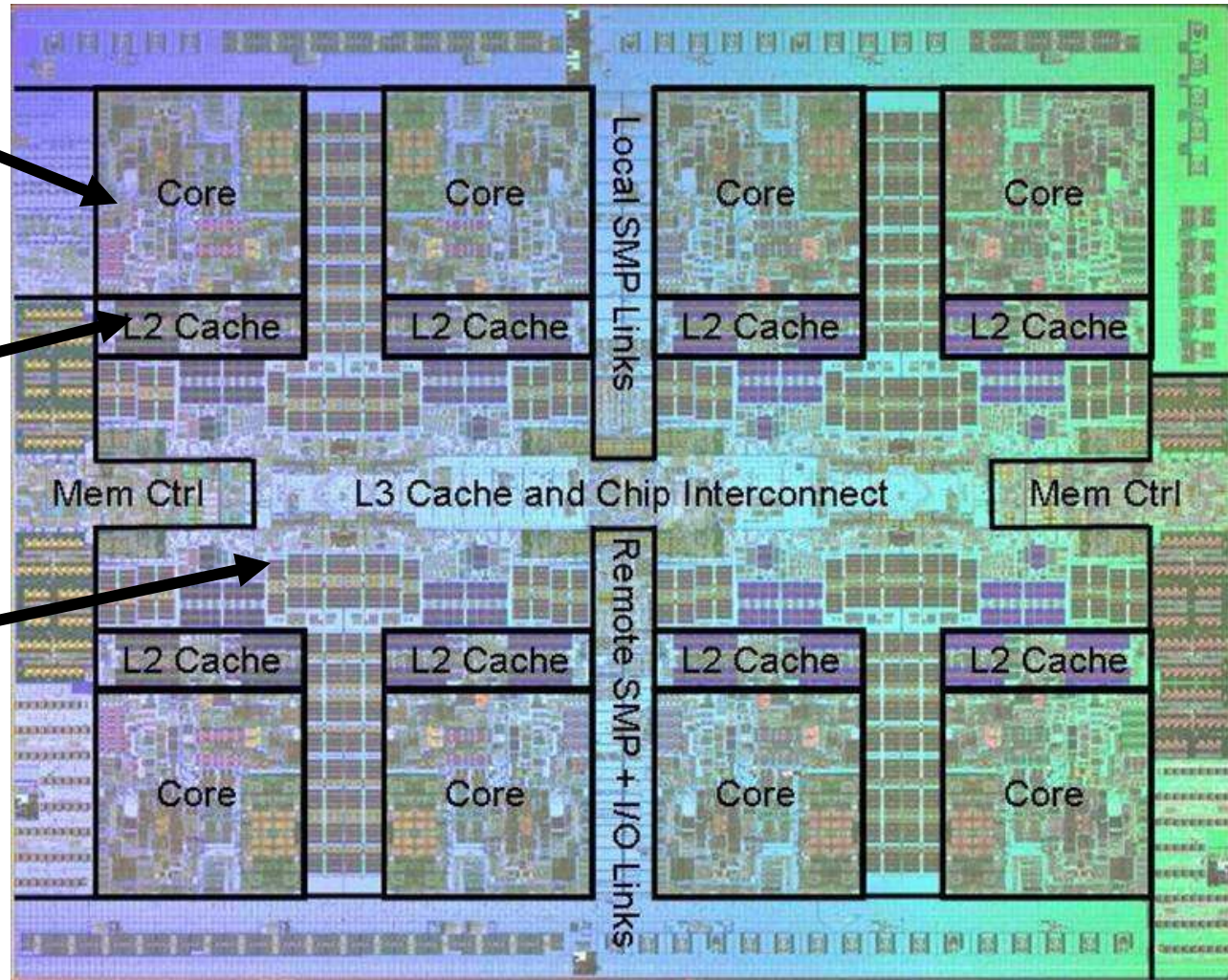
Level 3: 3MB, 12-way s.a., 128B line, single 32B port, twelve cycle latency

Power 7 On-Chip Caches [IBM 2009]

32KB L1 I\$/core
32KB L1 D\$/core
3-cycle latency

256KB Unified L2\$/core
8-cycle latency

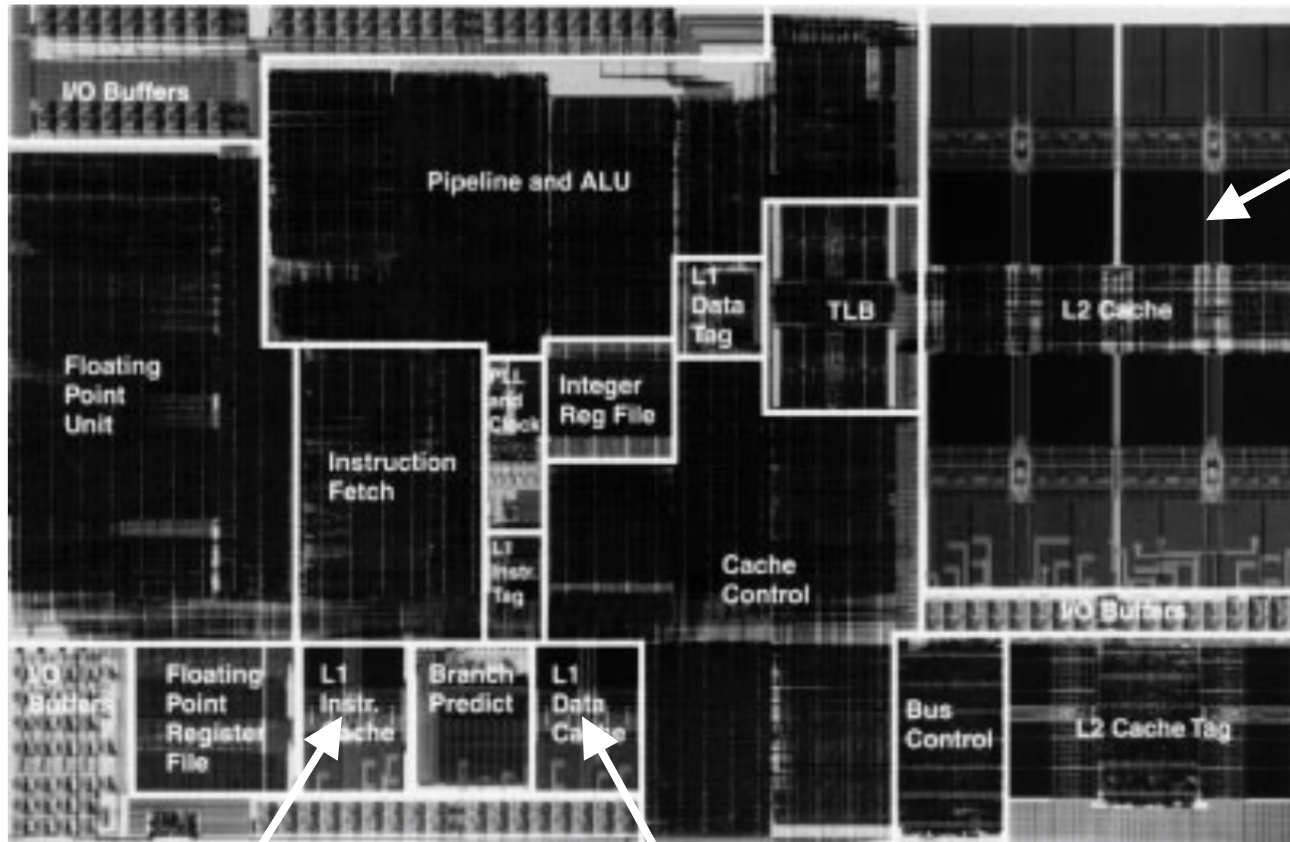
32MB Unified Shared L3\$
Embedded DRAM (eDRAM)
25-cycle latency to local
slice



IBM z196 Mainframe Caches 2010

- 96 cores (4 cores/chip, 24 chips/system)
 - Out-of-order, 3-way superscalar @ 5.2GHz
- L1: 64KB I-\$/core + 128KB D-\$/core
- L2: 1.5MB private/core (144MB total)
- L3: 24MB shared/chip (eDRAM) (576MB total)
- L4: 768MB shared/system (eDRAM)

Exponential X704 PowerPC Processor(1997)



32KB L2 8-way
Set-Associative
Write-Back
Unified Cache

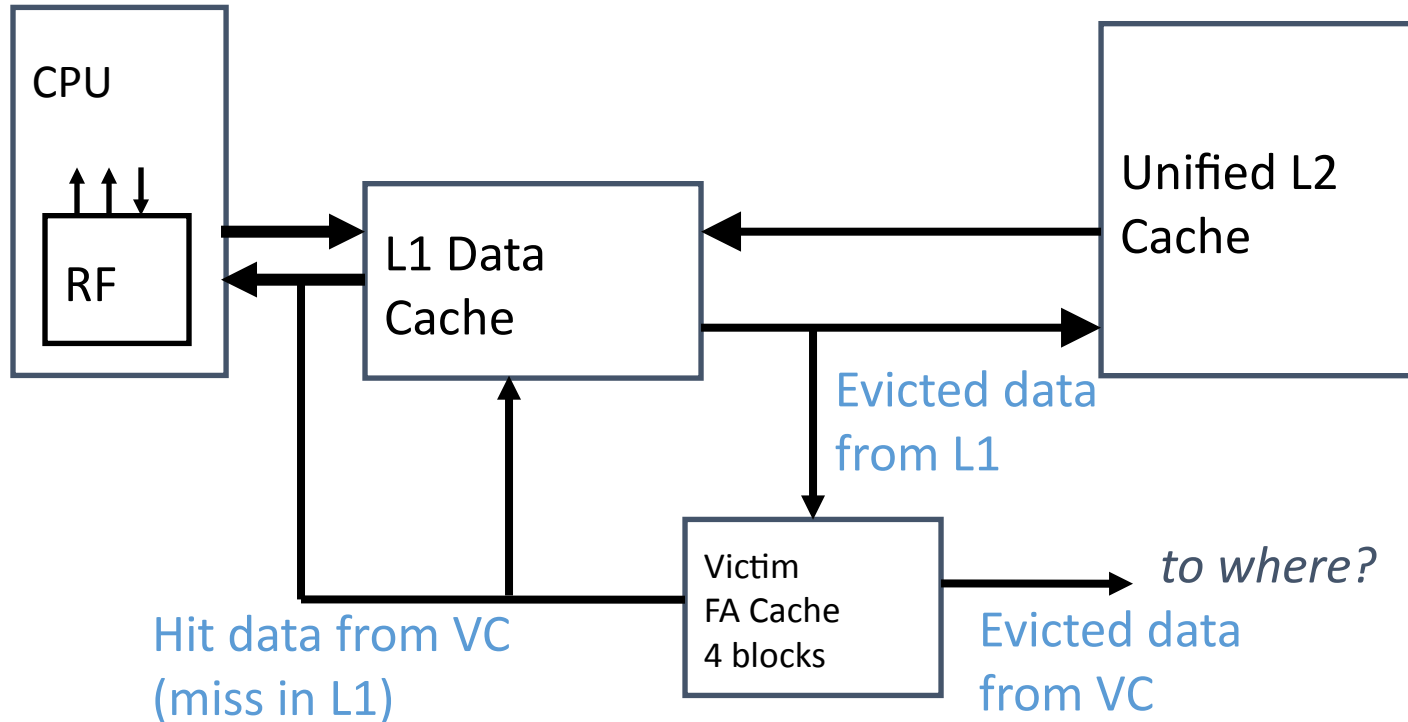
0.5 μ m BiCMOS
Ran at 410-533MHz
when other PC
processors were
much lower clock
rate

Project delayed –
missed market
window for Apple

2KB L1 Direct-Mapped
Instruction Cache

2KB L1 Direct-Mapped
Write-Through Data
Cache

Victim Caches (HP 7200)



Victim cache is a small associative backup cache, added to a direct-mapped cache, which holds recently evicted lines

- First look up in direct-mapped cache
- If miss, look in victim cache
- If hit in victim cache, swap hit line with line now evicted from L1
- If miss in victim cache, L1 victim -> VC, VC victim->?
- Fast hit time of direct mapped but with reduced conflict misses

MIPS R10000 Off-Chip L2 Cache (Yeager, IEEE Micro 1996)

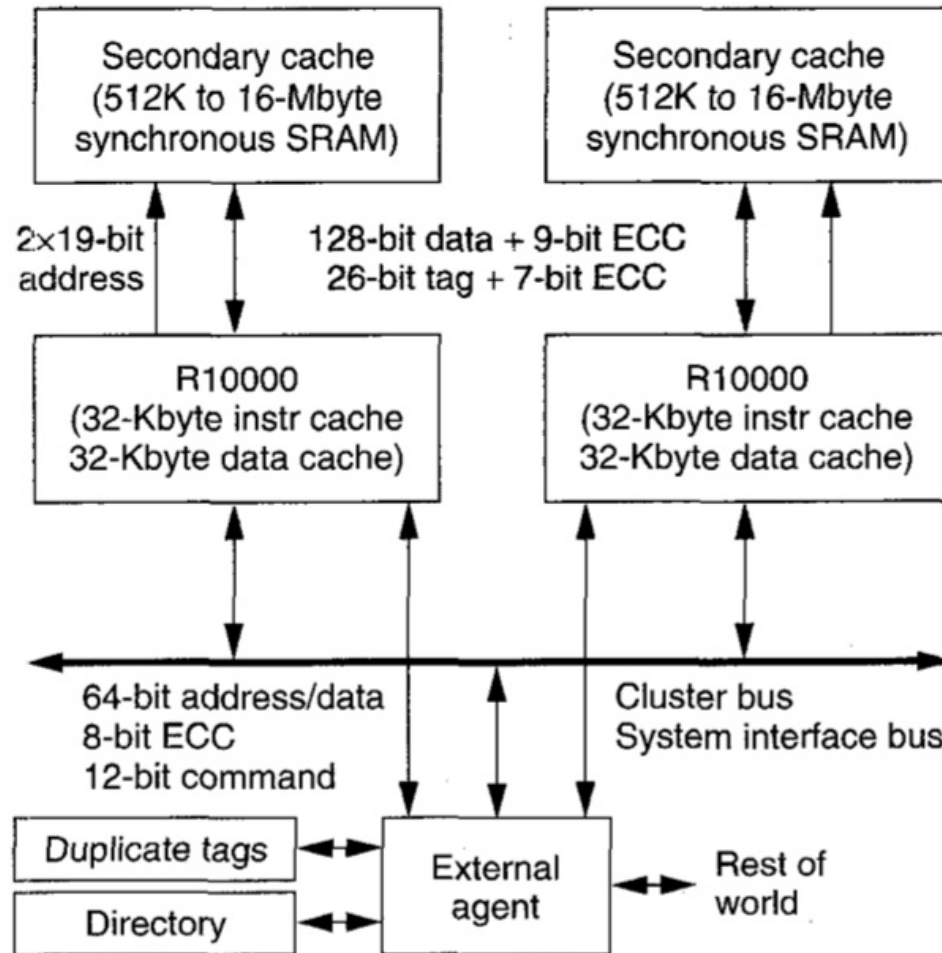
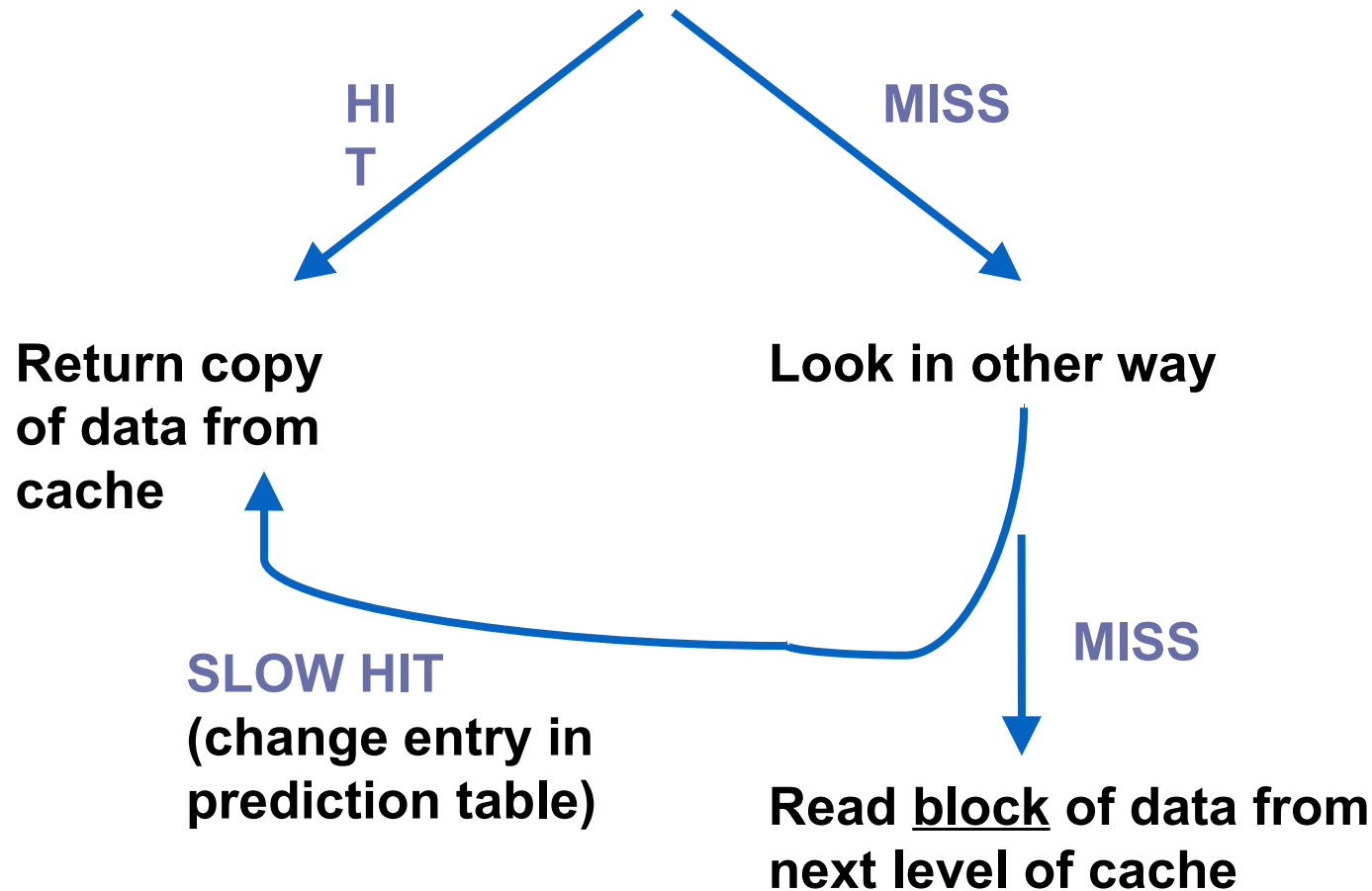


Figure 1. System configuration. The cluster bus directly connects as many as four chips.

Way-Predicting Caches (MIPS R10000 L2 cache)

- Use processor address to index into way-prediction table
- Look in predicted way at given index, then:



R10000 L2 Cache Timing Diagram

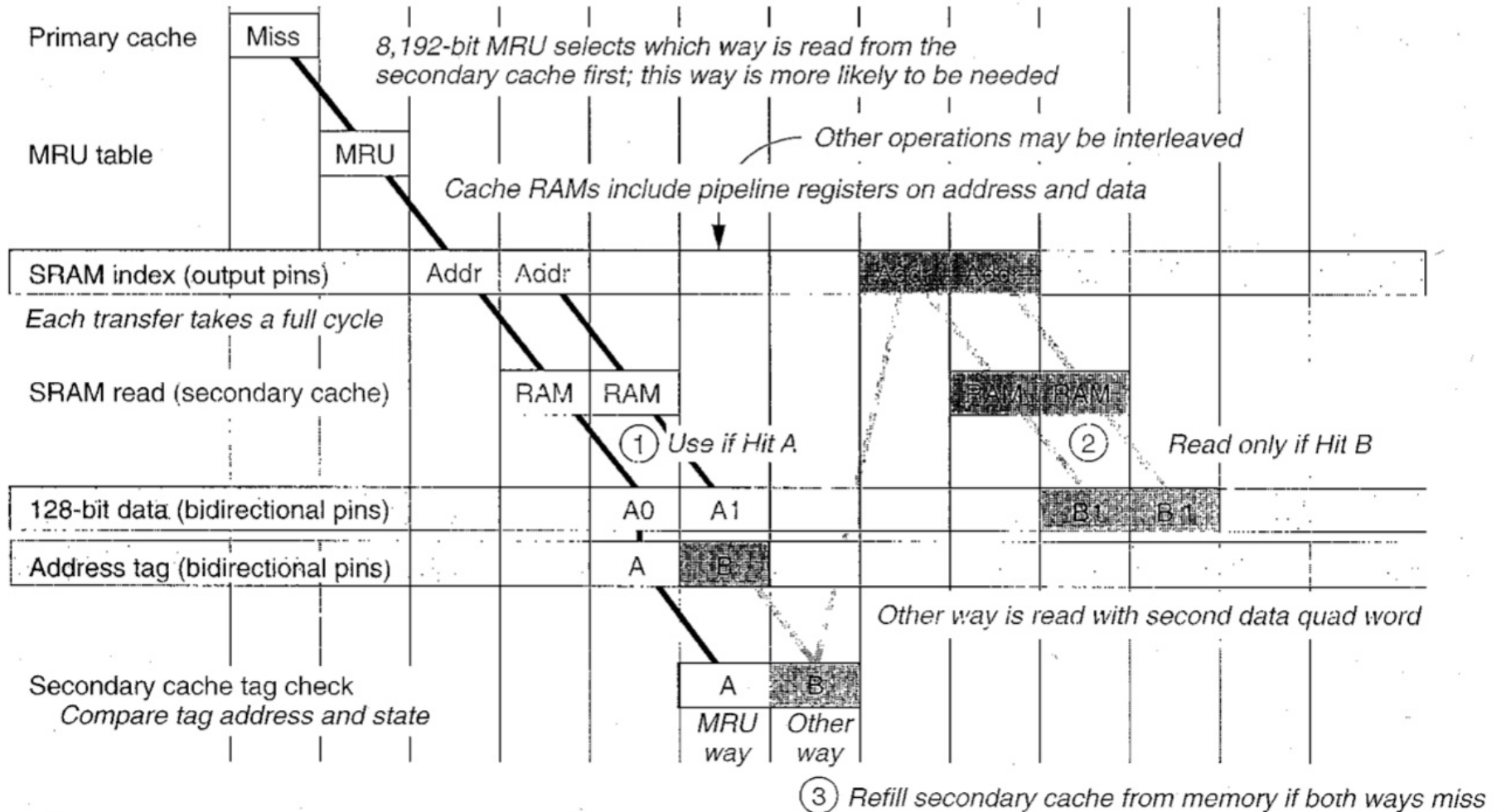
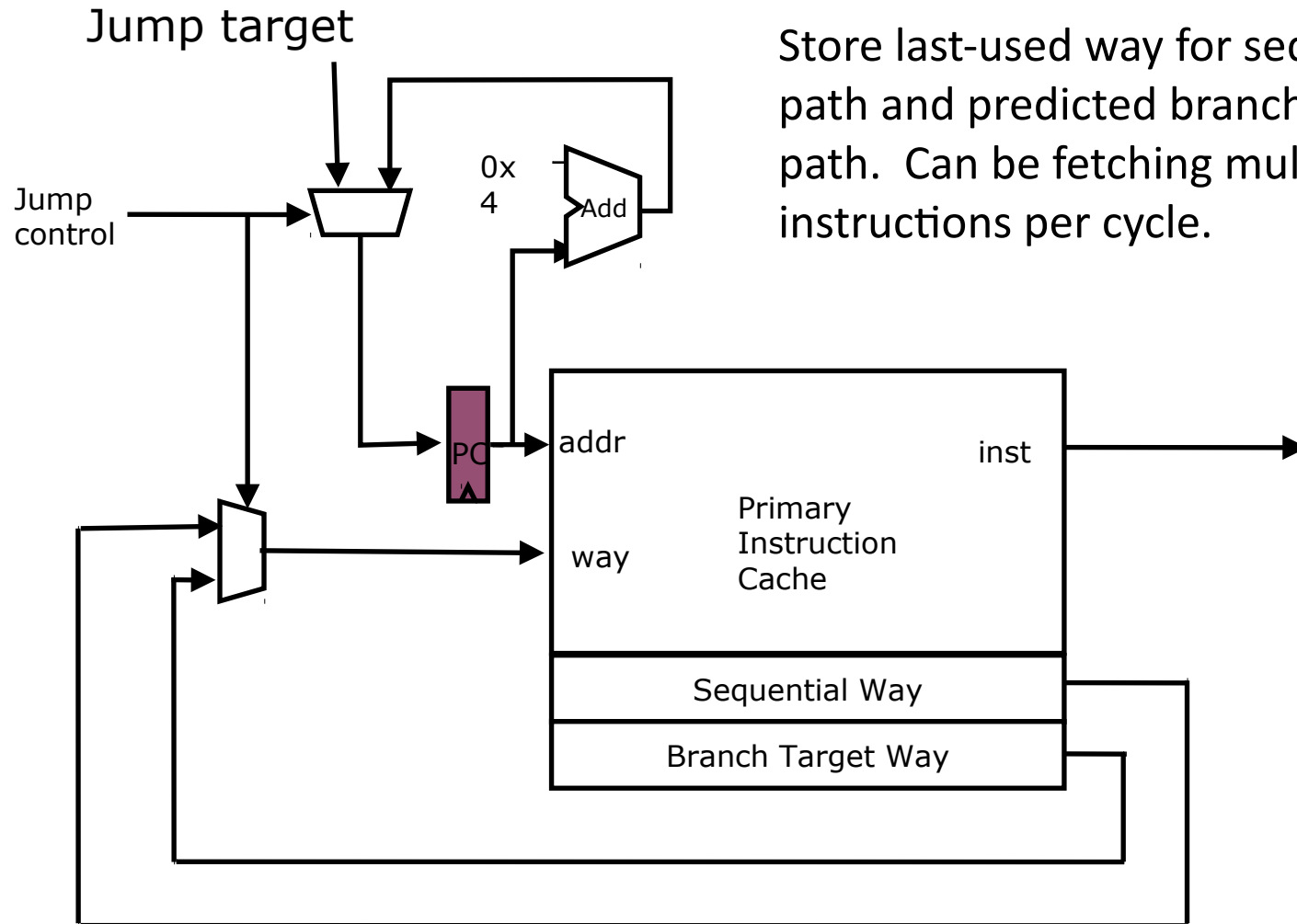


Figure 12. Refill from the set-associative secondary cache. In this example, the secondary clock equals the processor's internal pipeline clock. It may be slower.

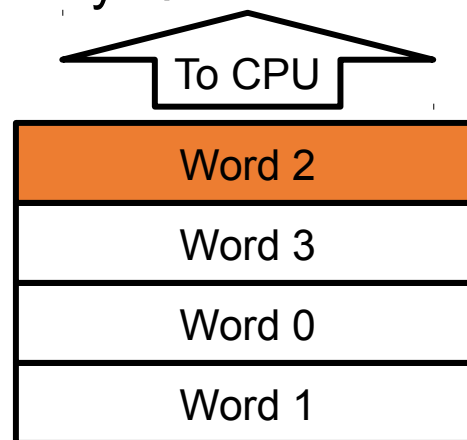
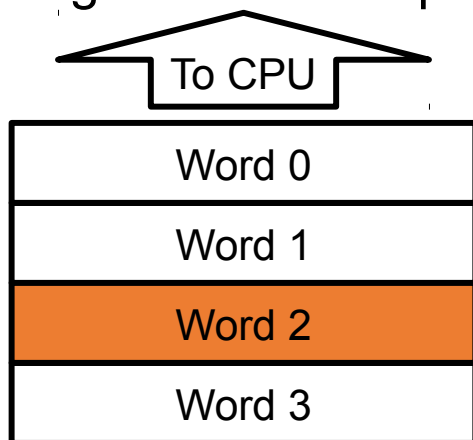
Way-Predicting Instruction Cache (Alpha 21264-like)



Store last-used way for sequential path and predicted branch taken path. Can be fetching multiple instructions per cycle.

Reduce Miss Penalty of Long Blocks: Early Restart and Critical Word First

- Don't wait for full block before restarting CPU
- *Early restart*—As soon as the requested word of the block arrives, send it to the CPU and let the CPU continue execution
- *Critical Word First*—Request the missed word first from memory and send it to the CPU as soon as it arrives; let the CPU continue execution while filling the rest of the words in the block
 - Long blocks more popular today \Rightarrow Critical Word 1st Widely used

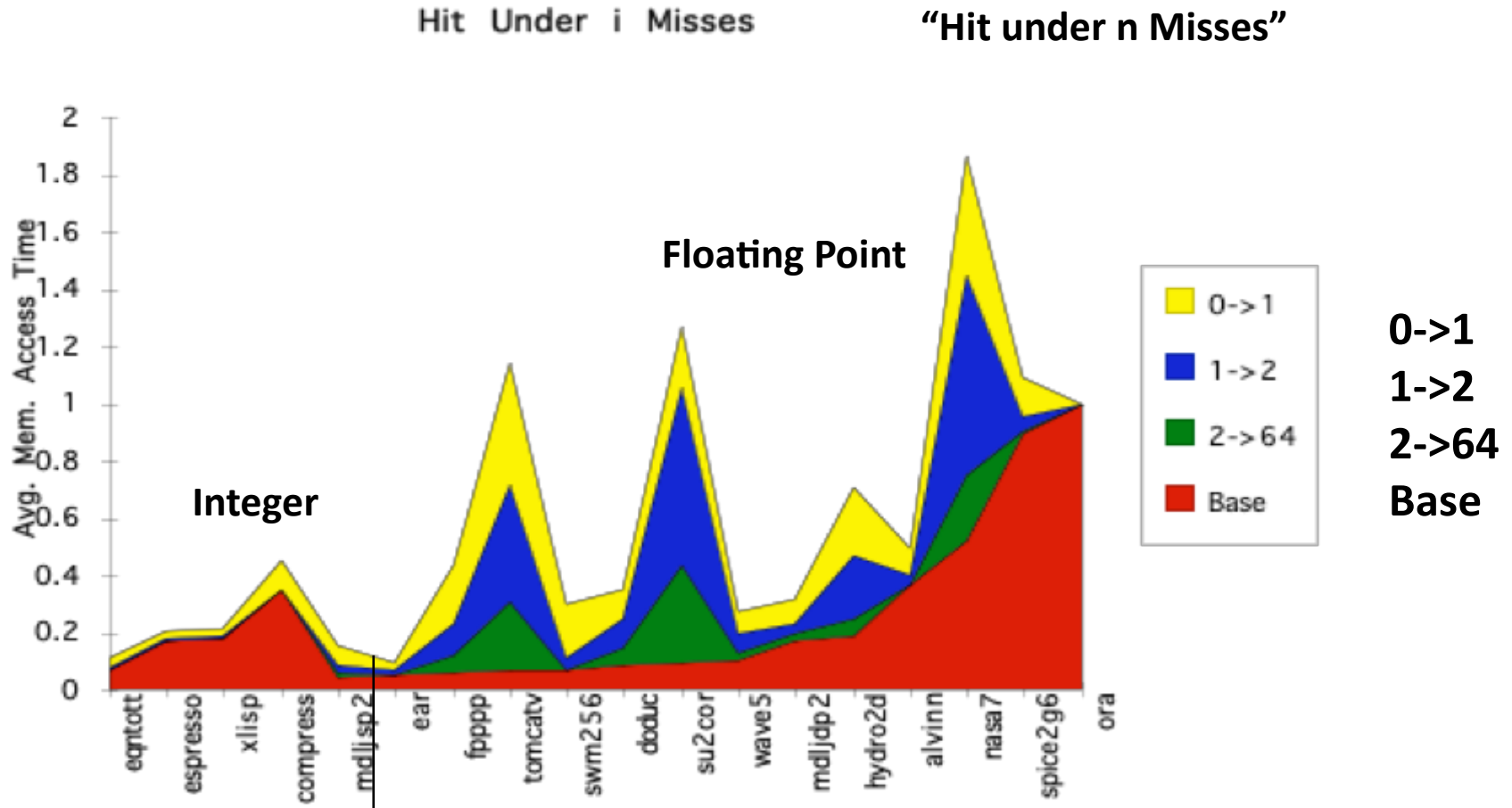


Rest of line filled in
with wrap-around
on cache line

Increasing Cache Bandwidth with Non-Blocking Caches

- **Non-blocking cache** or **lockup-free cache** allow data cache to continue to supply cache hits during a miss
 - requires Full/Empty bits on registers or out-of-order execution
- **“hit under miss”** reduces the effective miss penalty by working during miss vs. ignoring CPU requests
- **“hit under multiple miss”** or **“miss under miss”** may further lower the effective miss penalty by overlapping multiple misses
 - Significantly increases the complexity of the cache controller as there can be multiple outstanding memory accesses, and can get miss to line with outstanding miss (secondary miss)
 - Requires pipelined or banked memory system (otherwise cannot support multiple misses)
 - Pentium Pro allows 4 outstanding memory misses
 - Cray X1E vector supercomputer allows 2,048 outstanding memory misses

Value of Hit Under Miss for SPEC (old data)



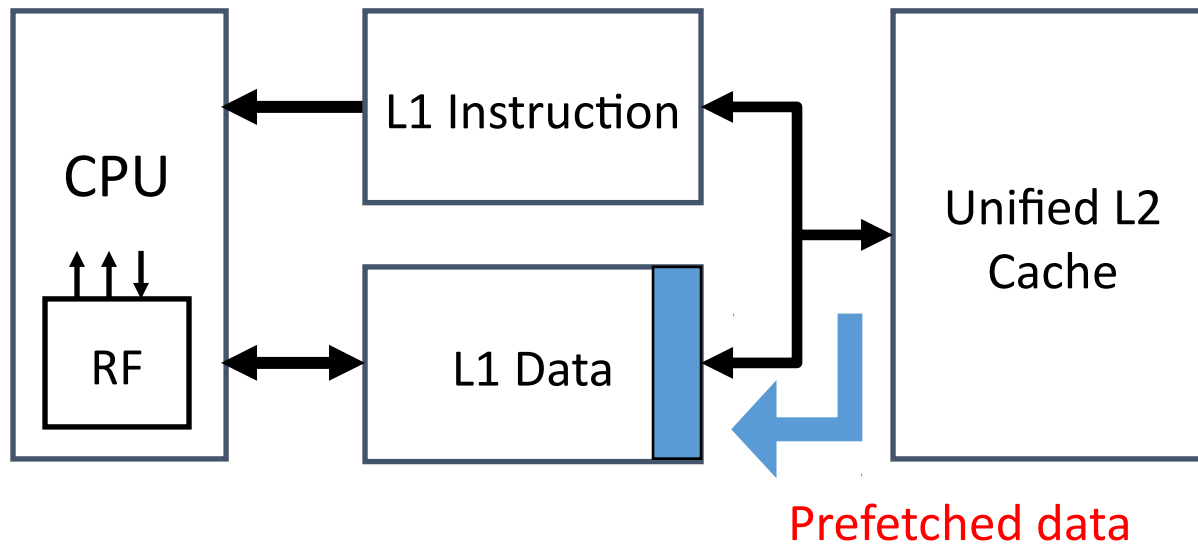
- FP programs on average: AMAT= 0.68 -> 0.52 -> 0.34 -> 0.26
- Int programs on average: AMAT= 0.24 -> 0.20 -> 0.19 -> 0.19
- 8 KB Data Cache, Direct Mapped, 32B block, 16 cycle miss, SPEC 92

Prefetching

- Speculate on future instruction and data accesses and fetch them into cache(s)
 - Instruction accesses easier to predict than data accesses
- Varieties of prefetching
 - Hardware prefetching
 - Software prefetching
 - Mixed schemes
- What types of misses does prefetching affect?

Issues in Prefetching

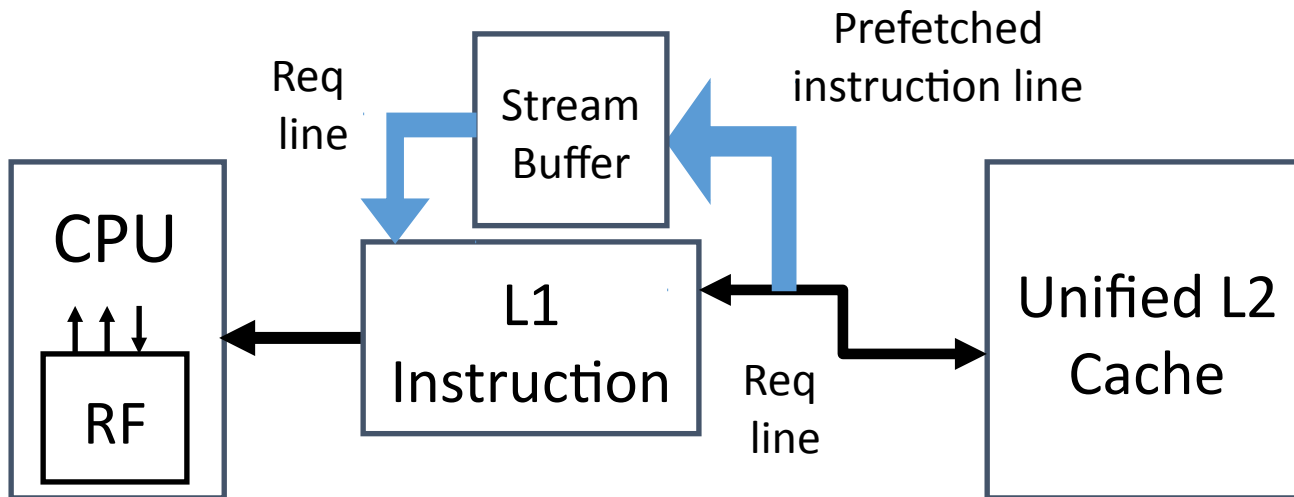
- Usefulness – should produce hits
- Timeliness – not late and not too early
- Cache and bandwidth pollution



Hardware Instruction Prefetching

Instruction prefetch in Alpha AXP 21064

- Fetch two lines on a miss; the requested line (i) and the next consecutive line (i+1)
- Requested line placed in cache, and next line in instruction stream buffer
- If miss in cache but hit in stream buffer, move stream buffer line into cache and prefetch next line (i+2)



Hardware Data Prefetching

- Prefetch-on-miss:
 - Prefetch $b + 1$ upon miss on b
- One-Block Lookahead (OBL) scheme
 - Initiate prefetch for block $b + 1$ when block b is accessed
 - Why is this different from doubling block size?
 - Can extend to N -block lookahead
- Strided prefetch
 - If observe sequence of accesses to line b , $b+N$, $b+2N$, then prefetch $b+3N$ etc.
- Example: IBM Power 5 [2003] supports eight independent streams of strided prefetch per processor, prefetching 12 lines ahead of current access

Software Prefetching

```
for(i=0; i < N; i++) {  
    prefetch( &a[i + 1] );  
    prefetch( &b[i + 1] );  
    SUM = SUM + a[i] * b[i];  
}
```

Software Prefetching Issues

- Timing is the biggest issue, not predictability
 - If you prefetch very close to when the data is required, you might be too late
 - Prefetch too early, cause pollution
 - Estimate how long it will take for the data to come into L1, so we can set P appropriately
 - *Why is this hard to do?*

```
for(i=0; i < N; i++) {  
    prefetch( &a[i + P] );  
    prefetch( &b[i + P] );  
    SUM = SUM + a[i] * b[i];  
}
```

Must consider cost of prefetch instructions

Software Prefetching Example

[“Data prefetching on the HP PA8000”, Santhanam et al., 1997]

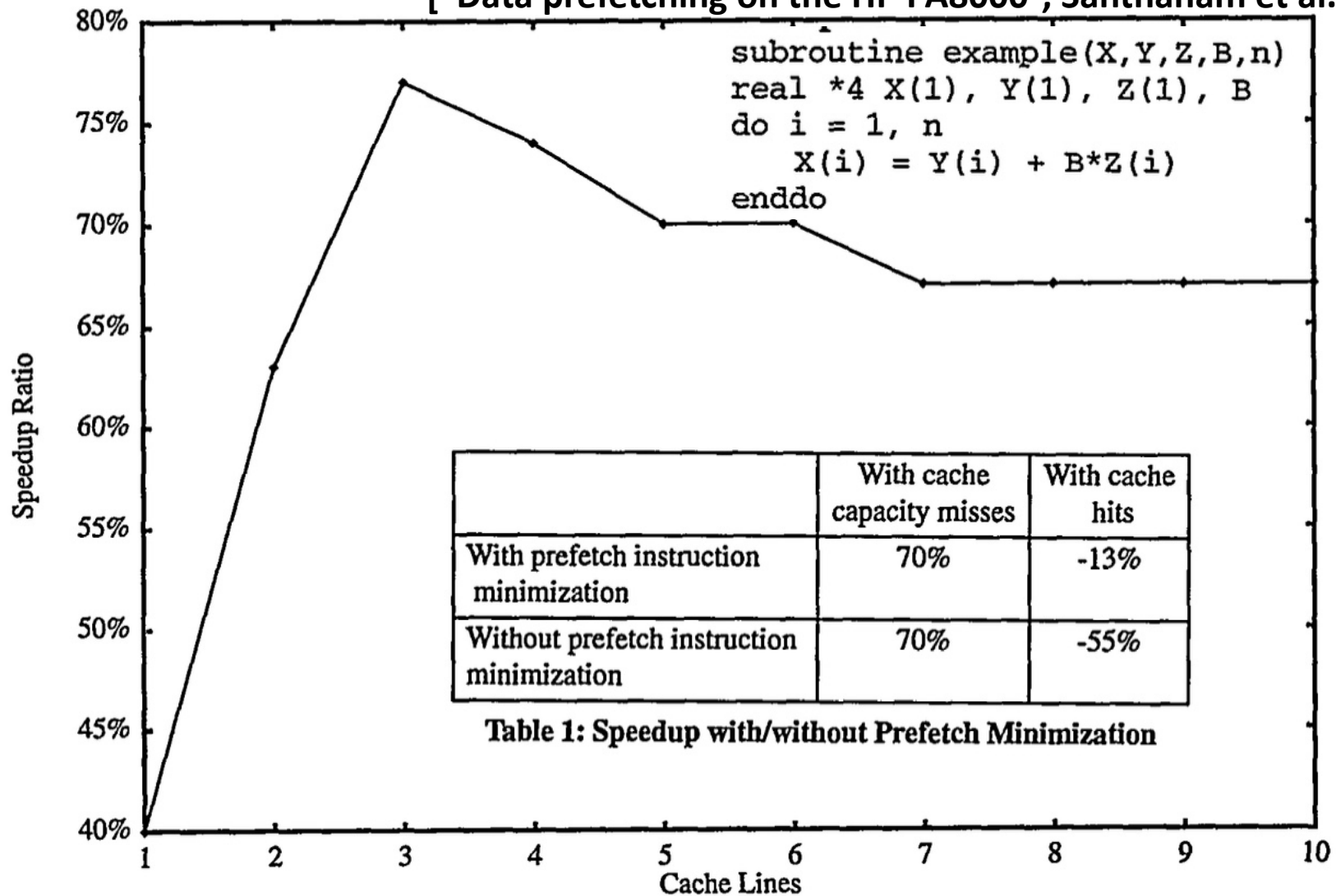


Figure 2: Speedup Ratio for Different Prefetch Distances

Compiler Optimizations

- Restructuring code affects the data access sequence
 - Group data accesses together to improve spatial locality
 - Re-order data accesses to improve temporal locality
- Prevent data from entering the cache
 - Useful for variables that will only be accessed once before being replaced
 - Needs mechanism for software to tell hardware not to cache data (“no-allocate” instruction hints or page table bits)
- Kill data that will never be used again
 - Streaming data exploits spatial locality but not temporal locality
 - Replace into dead cache locations

Loop Interchange

```
for (j=0; j < N; j++) {  
    for (i=0; i < M; i++) {  
        x[i][j] = 2 * x[i][j];  
    }  
}
```



```
for (i=0; i < M; i++) {  
    for (j=0; j < N; j++) {  
        x[i][j] = 2 * x[i][j];  
    }  
}
```

What type of locality does this improve?

Loop Fusion

```
for(i=0; i < N; i++)  
    a[i] = b[i] * c[i];
```

```
for(i=0; i < N; i++)  
    d[i] = a[i] * c[i];
```

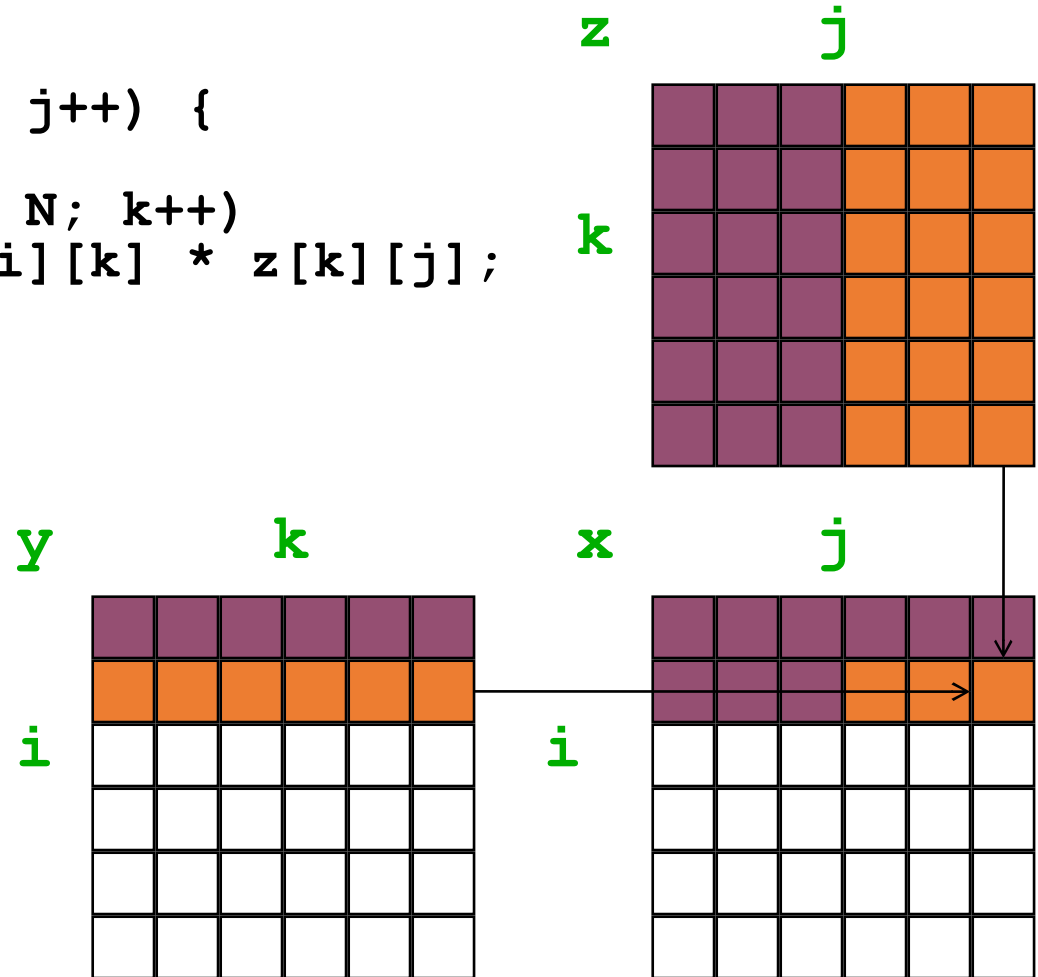


```
for(i=0; i < N; i++)  
{  
    a[i] = b[i] * c[i];  
    d[i] = a[i] * c[i];  
}
```

What type of locality does this improve?

Matrix Multiply, Naïve Code

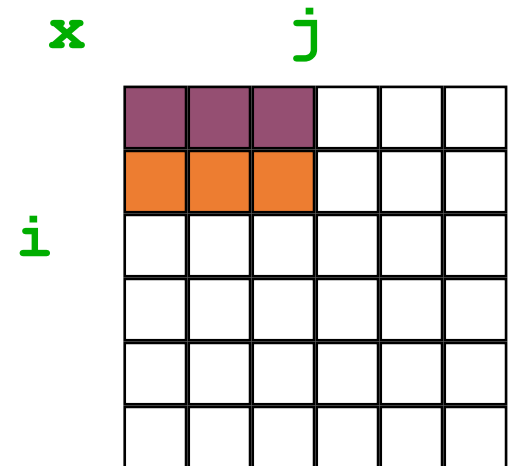
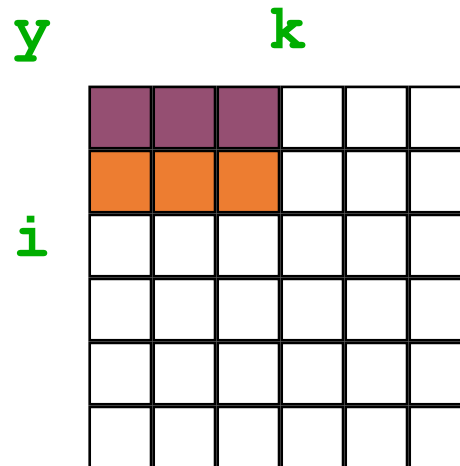
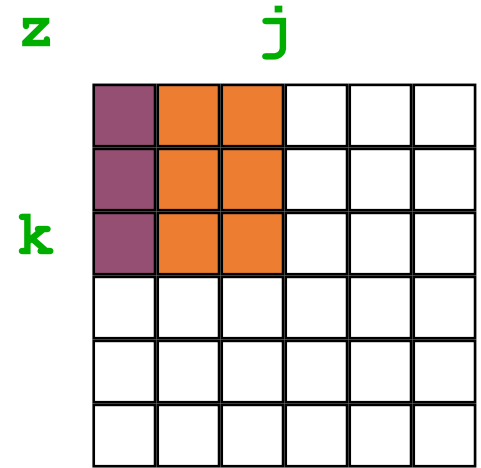
```
for(i=0; i < N; i++)
  for(j=0; j < N; j++) {
    r = 0;
    for(k=0; k < N; k++)
      r = r + y[i][k] * z[k][j];
    x[i][j] = r;
  }
```



☐ Not touched ☒ Old access ☒ New access

Matrix Multiply with Cache Tiling

```
for(jj=0; jj < N; jj=jj+B)
  for(kk=0; kk < N; kk=kk+B)
    for(i=0; i < N; i++)
      for(j=jj; j < min(jj+B,N); j++) {
        r = 0;
        for(k=kk; k < min(kk+B,N); k++)
          r = r + y[i][k] * z[k][j];
        x[i][j] = x[i][j] + r;
      }
}
```



What type of locality does this improve?