# Advanced Caching Techniques

**Approaches to improving memory system performance**

- eliminate memory operations
- decrease the number of misses
- decrease the miss penalty
- decrease the cache/memory access times
- hide memory latencies
- increase cache throughput
- increase memory bandwidth

---

# Handling a Cache Miss the Old Way

(1) Send the address & read operation to the next level of the hierarchy

**(2) Wait for the data to arrive**

(3) Update the cache entry with data*, rewrite the tag, turn the valid bit on, clear the dirty bit (unless an instruction cache)

(4) Resend the memory address; this time there will be a hit.

* There are variations:

- send the requested word to the CPU as soon as it arrives at the cache (**early restart**)
- requested word is sent from memory first; then the rest of the block follows (**requested word first**)
- get data before replace the block

How do the variations improve memory system performance?

# Non-blocking Caches

**Non-blocking cache** (lockup-free cache)
- allows the CPU to continue executing instructions while a miss is handled
- some processors allow only 1 outstanding miss ("hit under miss")
- some processors allow multiple misses outstanding ("miss under miss")
- **miss status holding registers** (MSHR)
  - hardware structure for tracking outstanding misses
    - physical address of the block
    - which word in the block
    - destination register number (if data)
    - mechanism to merge requests to the same block
    - mechanism to insure accesses to the same location execute in program order

---

# Non-blocking Caches

**Non-blocking cache** (lockup-free cache)
- can be used with both in-order and out-of-order processors
  - **in-order processors** stall when an instruction that uses the load data is the next instruction to be executed (non-blocking loads)
  - **out-of-order processors** can execute instructions after the load consumer
- How do non-blocking caches improve memory system performance?

## Victim Cache

**Victim cache**

- small fully-associative cache
  - contains the most recently replaced blocks of a direct-mapped cache
  - alternative to 2-way set-associative cache
- check it on a cache miss
  - swap the direct-mapped block and victim cache block
- How do victim caches improve memory system performance?

- Why do victim caches work?

---

## Sub-block Placement

Divide a block into sub-blocks

| tag |  | I | data | V | data | V | data | I | data |
|-----|--|---|------|---|------|---|------|---|------|
| tag |  | V | data | V | data | V | data | V | data |
| tag |  | V | data | V | data | V | data | V | data |
| tag |  | I | data | I | data | I | data | I | data |

- **sub-block** = unit of transfer on a cache miss
- **valid bit**/sub-block
- Misses:
  - block-level miss: tags didn't match
  - sub-block-level miss: tags matched, valid bit was clear
- \+   the transfer time of a sub-block
- \+   fewer tags than if each sub-block were a block
- \-   less implicit prefetching
- how does sub-block placement improve memory system performance?

## Pseudo-set associative Cache

**Pseudo-set associative cache**

- access the cache
- if miss, invert the high-order index bit & access the cache again
- + miss rate of 2-way set associative cache
- + close to access time of direct-mapped cache
- - increase in hit time (relative to 2-way associative) if always try slow hit first
  - predict which is the fast-hit block
  - put the fast hit block in the same location; swap blocks if wrong
- How does pseudo-set associativity improve memory system performance?

## Pipelined Cache Access

**Pipelined cache access**

- simple 2-stage pipeline
  - access
  - data transfer
  - tag check & hit/miss logic with the shorter
- how pipelined caches improve memory system performance:

# Mechanisms for Prefetching

**Stream buffers**
- where prefetched instructions/data held
- if requested block in the stream buffer, then cancel the cache access

How save here?

# Trace Cache

**Trace cache contents**
- trace is analogous to a cache block
- contains instructions from the dynamic instruction stream
    + fetch statically noncontiguous instructions in a single cycle
    + a more efficient use of I-cache space
- low bits of next addresses (target & fall-through code) for last branch in a trace
- cache state is high branch address bits + predictions for all branches within the block

Effect on performance?

## Trace Cache

**Assessing a trace cache**

- assess trace cache + branch predictor, BTB, TLB, I-cache in parallel
- compare PC & prediction history of the current branch instruction to the trace cache tag
- hit: I-cache fetch ignored
- miss: use the I-cache
    - start constructing a new trace

Why does a trace cache work?

---

## Cache-friendly Compiler Optimizations

Exploit spatial locality
- **schedule for array misses**
  - hoist first load to a cache block

Improve spatial locality
- **group & transpose**
  - makes portions of vectors that are accessed together lie in memory together
- **loop interchange**
  - so inner loop follows memory layout

Improve temporal locality
- **loop fusion**
  - do multiple computations on the same portion of an array
- **tiling (also called blocking)**
  - do all computation on a small block of memory that will fit in the cache

## Tiling Example

```
/* before */
for (i=0; i<n; i=i+1)
     for (j=0; j<n; j=j+1){
          r = 0;
          for (k=0; k<n; k=k+1) {
               r = r + y[i,k] * z[k,j]; }
          x[i,j] = r;
          };
/* after */
for (jj=0; jj<n; jj=jj+B)
  for (kk=0; kk<n; kk=kk+B)
  for (i=0; i<n; i=i+1)
     for (j=jj; j<min(jj+B-1,n); j=j+1) {
          r = 0;
          for (k=kk; k<min(kk+B-1,n); k=k+1)
               {r = r + y[i,k] * z[k,j]; }
          x[i,j] = x[i,j] + r;
          };
```

## Memory Banks

**Interleaved memory:**
- multiple memory banks
  - word locations are assigned across banks
  - interleaving factor: number of banks
  - send a single address to all banks at once



- \+ get more data for one transfer: increases memory bandwidth
  - data is probably used *(why?)*
- \- larger DRAM chip capacity means fewer banks
- \- power issue

# Memory Banks

**Independent memory banks**

- different banks can be accessed at once, with different addresses
- allows parallel access, possibly parallel data transfer
- multiple memory controllers & separate address lines, one for each access
  - different controllers cannot access the same bank
- cheaper than dual porting

Effect on performance?

| | 21264 | R12000 | UltraSPARC-III | Pentium IV |
|---|---|---|---|---|
| **L1 I** onchip | 64KB | 32KB | 32KB | 12Kuop trace cache (~8-16KB) |
| | 2-way with set prediction | **2-way** | **4-way** | |
| | 64B block | **64B block** | 32B block | 6 uops/line |
| | **virtually indexed** | | virtually indexed, virtual tags | virtually indexed |
| | | 2-cycle access | pipelined 2-cycle access | |
| | | critical word first | | |
| **L1 D** onchip | 64KB **2-way** | 32KB 2-way, LRU replacement | **64KB 4-way** | 8KB 4-way |
| | 64B block | **32B block** | 32B block | 64B block |
| | write-back | | write-through **store compression** | write-through |
| | virtually indexed, **physical tags** | **physical tags** | virtually indexed | virtually indexed |
| | **TLB in parallel** | | TLB in parallel | |
| | 3 (int) or 4 (FP) cycle reads | 2-cycle access | | 2 cycle latency |
| | **phase-pipelined (read twice each cycle)** | | pipelined 2-cycle access | pipelined |
| | **miss under miss (32 loads or 8 blocks outstanding))** | nonblocking | nonblocking | nonblocking |
| | **victim cache** | critical word first | | requested word first |
| **L2** | external 1MB-16MB | external 1MB-16MB | external up to 8MB | onchip 256KB |
| | direct-mapped | 2-way pseudo, way prediction, LRU | direct-mapped | **8-way** |
| | 64B block | 128B blocks | 32B blocks | 128B block 64B "subblocks" |
| | write-back | write-back | write-back | write-back |
| | physical | | physical | physically indexed |
| | nonblocking | | | nonblocking |
| | 12 cycles | | 12 cycles | |
| | | | pipelined access | pipelined |
| **TLB** | 128 entries | 64 entries, each maps to 2 pages | | |
| | FA | FA | | |
| | dual-ported | | | |
| | multiple page sizes | **4KB - 16MB pages** | multiple page sizes | multiple page sizes |
| | **PAL code handling** | | software handling | hardware handling |

**Today's Memory Subsystems**

**Look for designs in common:**

---

**Advanced Caching Techniques**

**Approaches to improving memory system performance**

- eliminate memory operations

- decrease the number of misses

- decrease the miss penalty

- decrease the cache/memory access times

- hide memory latencies

- increase cache throughput

- increase memory bandwidth

## Wrap-up

Victim cache (reduce miss penalty)

TLB (reduce page fault time (penalty))

Hardware or compiler-based prefetching (reduce misses)

Cache-conscious compiler optimizations (reduce misses or hide miss penalty)

Coupling a write-through memory update policy with a write buffer (eliminate store ops/hide store latencies)

Handling the read miss before replacing a block with a write-back memory update policy (reduce miss penalty)

Sub-block placement (reduce miss penalty)

Non-blocking caches (hide miss penalty)

Merging requests to the same cache block in a non-blocking cache (hide miss penalty)

Requested word first or early restart (reduce miss penalty)

Cache hierarchies (reduce misses/reduce miss penalty)

Virtual caches (reduce miss penalty)

Pipelined cache accesses (increase cache throughput)

Pseudo-set associative cache (reduce misses)

Banked or interleaved memories (increase bandwidth)

Independent memory banks (hide latency)

Wider bus (increase bandwidth)