# Computer Architecture

## Memory System

### Virendra Singh

Associate Professor

Computer Architecture and Dependable Systems Lab

Department of Electrical Engineering

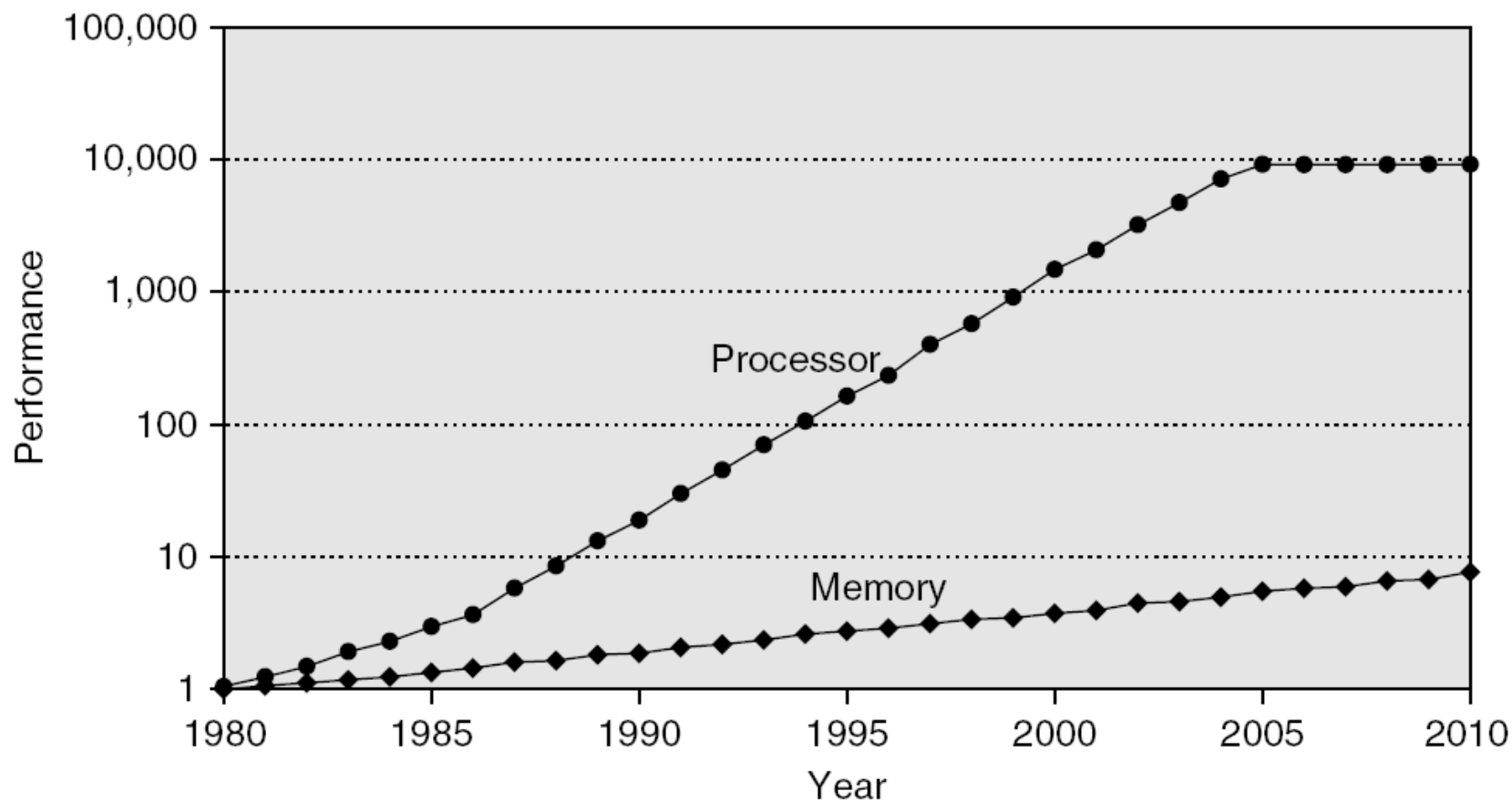Indian Institute of Technology Bombay

http://www.ee.iitb.ac.in/~viren/

E-mail: viren@ee.iitb.ac.in

*CS-683: Advanced Computer Architecture*

Lecture 6 (13 Aug 2013)

CADSL

# Memory Performance Gap

**CADSL**

# Why Memory Hierarchy?

- Need lots of bandwidth

$$BW = \frac{1.0\, inst}{cycle} \times \left[ \frac{1\, Ifetch}{inst} \times \frac{4B}{Ifetch} + \frac{0.4\, Dref}{inst} \times \frac{4B}{Dref} \right] \times \frac{1\, Gcycles}{\sec}$$

$$= \frac{5.6\, GB}{\sec}$$

- Need lots of storage
  - 64MB (minimum) to multiple TB
- Must be cheap per bit
  - (TB x anything) is a lot of money!
- These requirements seem incompatible

CADSL

# Memory Hierarchy Design

- Memory hierarchy design becomes more crucial with recent multi-core processors:
  - Aggregate peak bandwidth grows with # cores:
    - Intel Core i7 can generate two references per core per clock
    - Four cores and 3.2 GHz clock
      - 25.6 billion 64-bit data references/second +
      - 12.8 billion 128-bit instruction references
      - = 409.6 GB/s!
  - DRAM bandwidth is only 6% of this (25 GB/s)
  - Requires:
    - Multi-port, pipelined caches
    - Two levels of cache per core
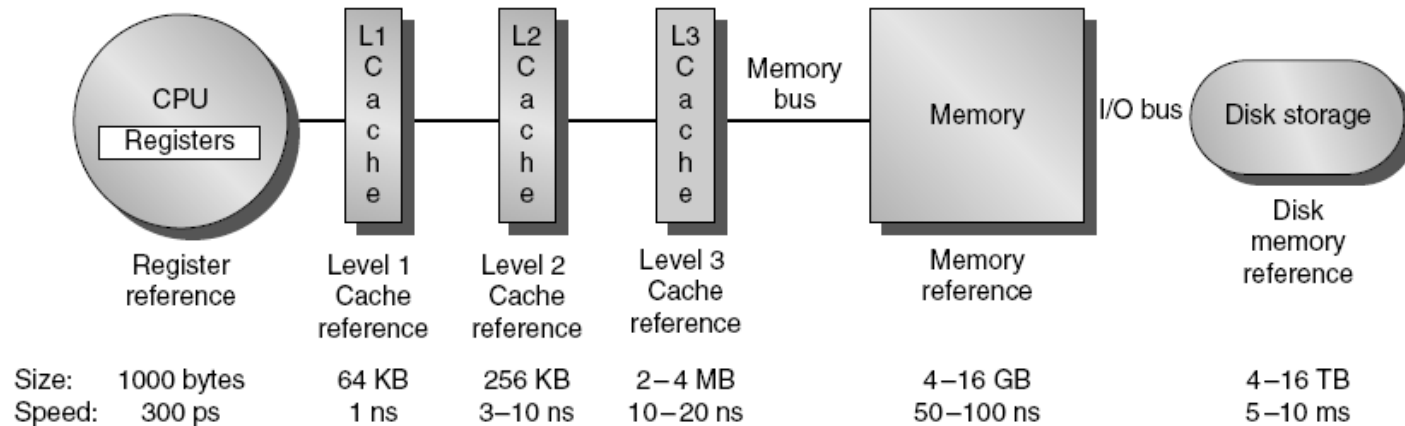    - Shared third-level cache on chip

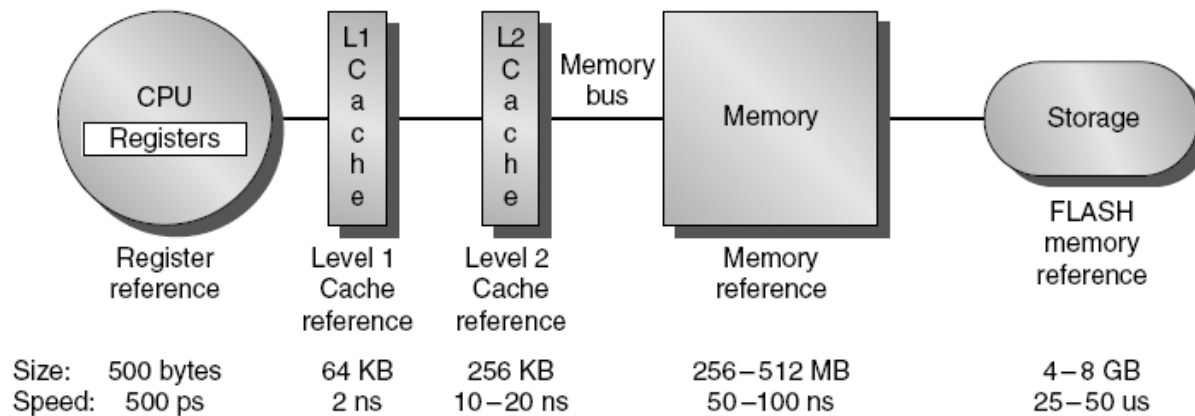**CADSL**

# Why Memory Hierarchy?

- Fast and small memories
  - Enable quick access (fast cycle time)
  - Enable lots of bandwidth (1+ L/S/I-fetch/cycle)
- Slower larger memories
  - Capture larger share of memory
  - Still relatively fast
- Slow huge memories
  - Hold rarely-needed state
  - Needed for correctness
- All together: provide appearance of large, fast memory with cost of cheap, slow memory

CADSL

# Memory Hierarchy



(a) Memory hierarchy for server

| Size: | 1000 bytes | 64 KB | 256 KB | 2–4 MB | 4–16 GB | 4–16 TB |
| Speed: | 300 ps | 1 ns | 3–10 ns | 10–20 ns | 50–100 ns | 5–10 ms |

(b) Memory hierarchy for a personal mobile device

| Size: | 500 bytes | 64 KB | 256 KB | 256–512 MB | 4–8 GB |
| Speed: | 500 ps | 2 ns | 10–20 ns | 50–100 ns | 25–50 us |

**CADSL**

# Why Does a Hierarchy Work?

- Locality of reference
  - Temporal locality
    - Reference same memory location repeatedly
  - Spatial locality
    - Reference near neighbors around the same time
- Empirically observed
  - Significant!
  - Even small local storage (8KB) often satisfies >90% of references to multi-MB data set

**CADSL**

# Why Locality?

- Analogy:
  - Library (Disk)
  - Bookshelf (Main memory)
  - Stack of books on desk (off-chip cache)
  - Opened book on desk (on-chip cache)
- Likelihood of:
  - Referring to same book or chapter again?
    - Probability decays over time
    - Book moves to bottom of stack, then bookshelf, then library
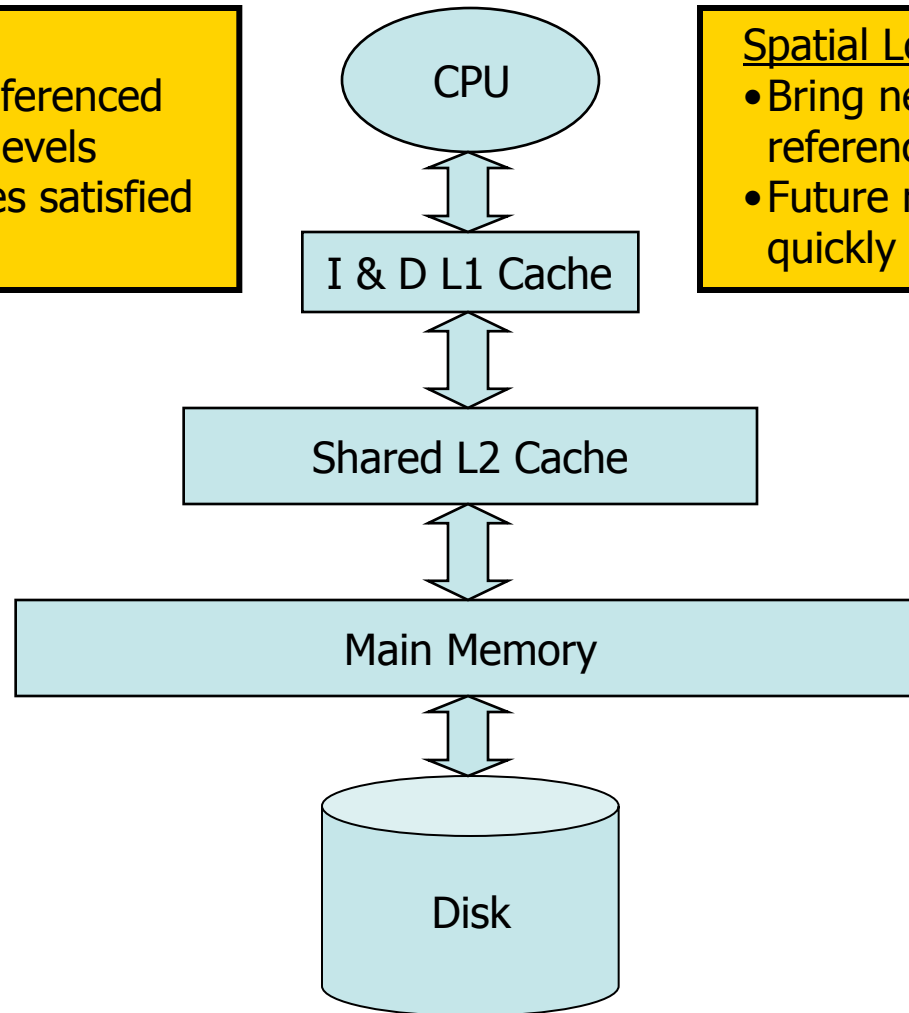  - Referring to chapter n+1 if looking at chapter n?

CADSL

# Memory Hierarchy

**Temporal Locality**
- Keep recently referenced items at higher levels
- Future references satisfied quickly

**Spatial Locality**
- Bring neighbors of recently referenced to higher levels
- Future references satisfied quickly

CPU

I & D L1 Cache

Shared L2 Cache

Main Memory

Disk

**CADSL**

# Performance

CPU execution time = (CPU clock cycles + memory stall cycles) x Clock Cycle time

Memory Stall cycles = Number of misses x miss penalty

= IC x misses/Instruction x miss penalty

=IC x memory access/instruction x miss rate x miss penalty

CADSL

# Four Burning Questions

- These are:
  - Placement
    - Where can a block of memory go?
  - Identification
    - How do I find a block of memory?
  - Replacement
    - How do I make space for new blocks?
  - Write Policy
    - How do I propagate changes?
- Consider these for caches
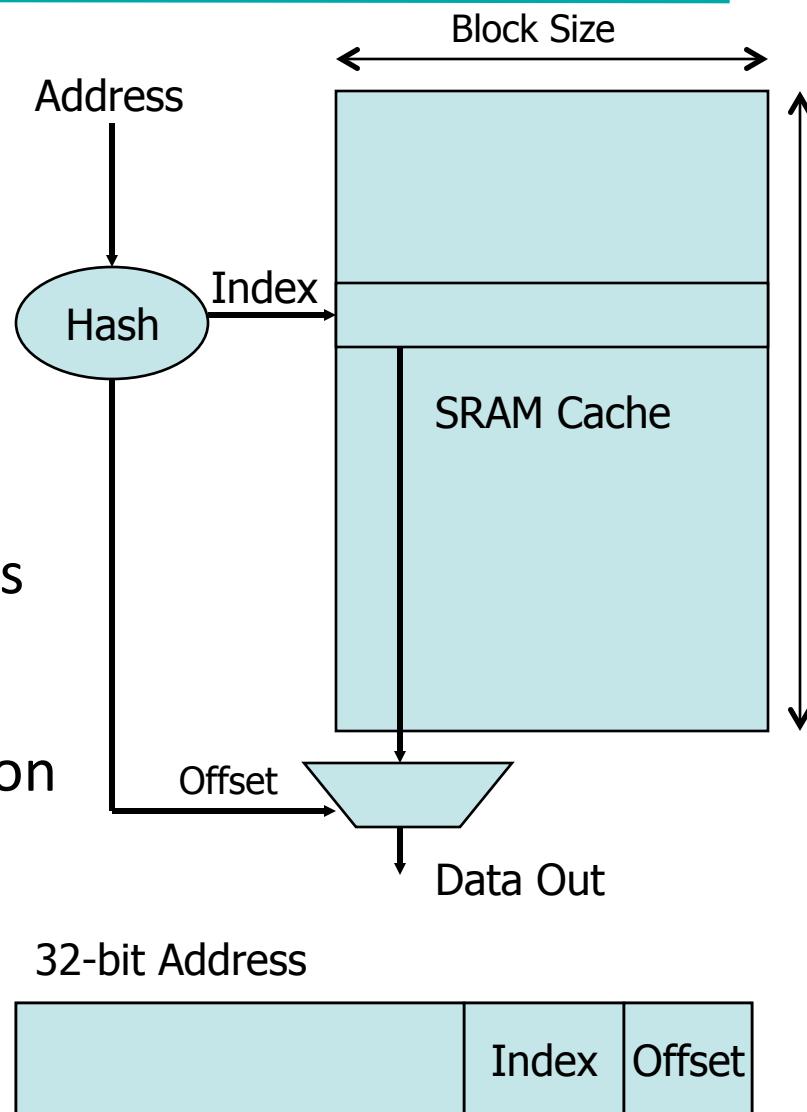  - Usually SRAM
- Will consider main memory, disks later

**CADSL**

# Placement

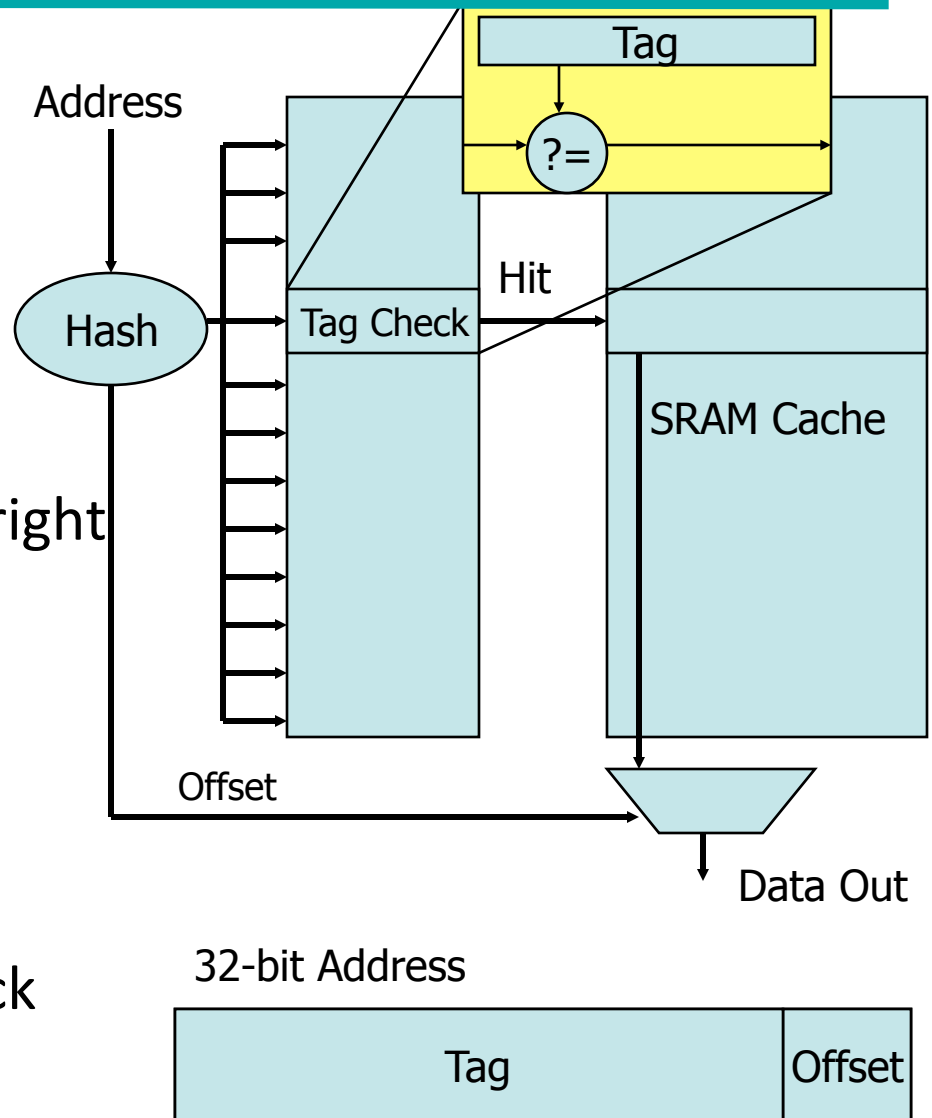| Memory Type | Placement | Comments |
|---|---|---|
| Registers | Anywhere; Int, FP, SPR | Compiler/programmer manages |
| Cache (SRAM) | Fixed in H/W | *Direct-mapped, set-associative, fully-associative* |
| DRAM | Anywhere | O/S manages |
| Disk | Anywhere | O/S manages |

CADSL

# Placement

- ## Address Range
  - Exceeds cache capacity
- ## Map address to finite capacity
  - Called a *hash*
  - Usually just masks high-order bits
- ## *Direct-mapped*
  - Block can only exist in one location
  - Hash collisions cause problems

Block Size

Address

Hash → Index

SRAM Cache

Offset

Data Out

32-bit Address

| | Index | Offset |
|---|---|---|

**CADSL**

# Placement

- *Fully-associative*
  - Block can exist anywhere
  - No more hash collisions
- *Identification*
  - How do I know I have the right block?
  - Called a *tag check*
    - Must store address tags
    - Compare against address
- Expensive!
  - Tag & comparator per block



Address

Hash

Tag

?=

Tag Check

Hit

SRAM Cache

Offset

Data Out

32-bit Address

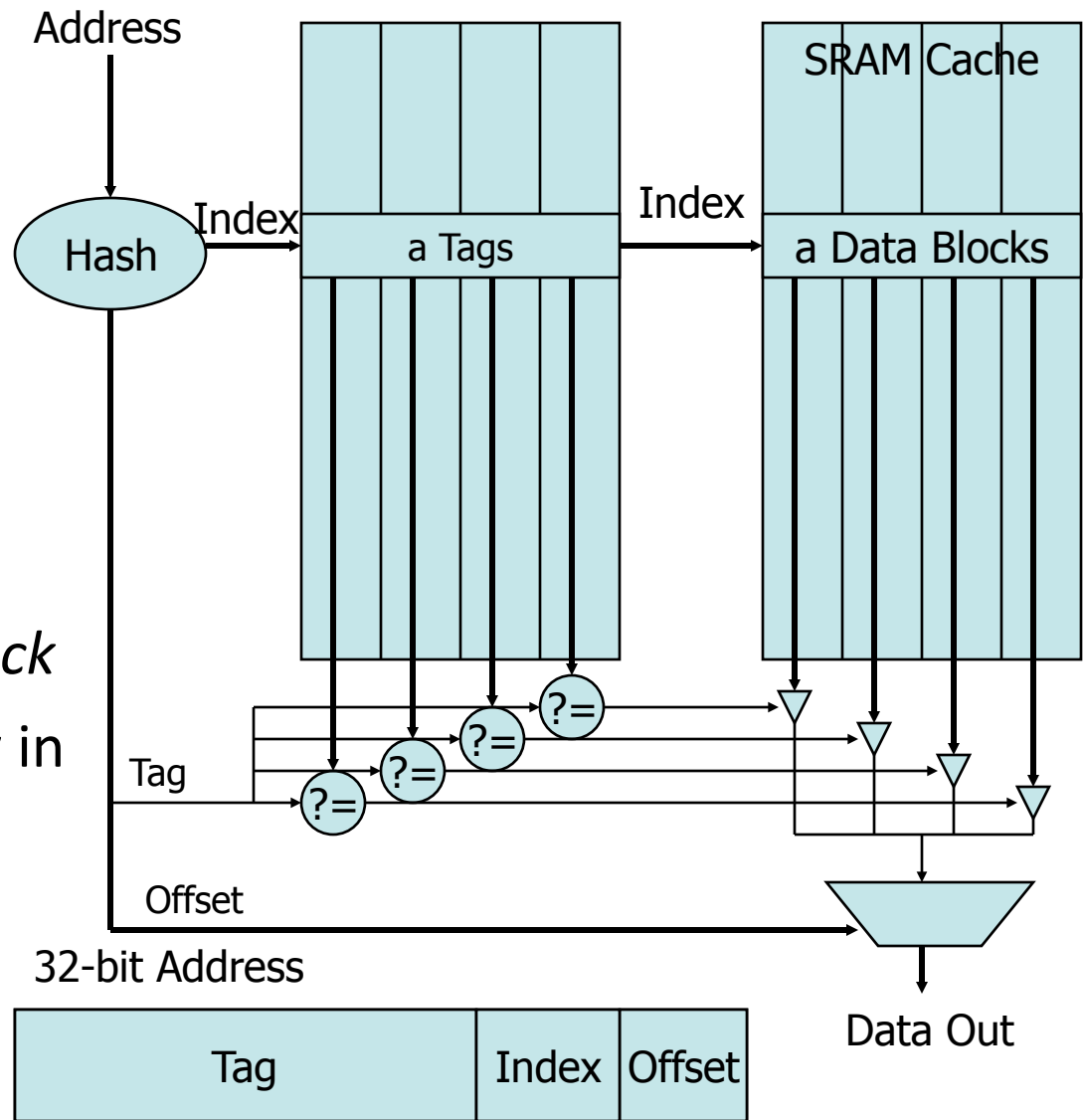| Tag | Offset |
|-----|--------|

**CADSL**

# Placement

- *Set-associative*
  - Block can be in a locations
  - Hash collisions:
    - a still OK
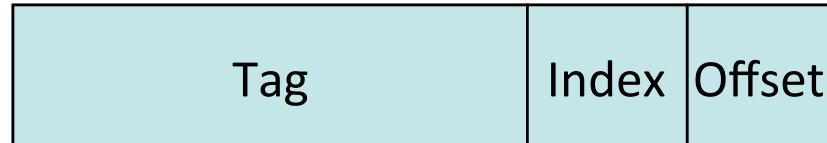- *Identification*
  - Still perform *tag check*
  - However, only a few in parallel

CADSL

# Placement and Identification

32-bit Address

| Tag | Index | Offset |
|-----|-------|--------|

| Portion | Length | Purpose |
|---------|--------|---------|
| Offset | $o = \log_2(\text{block size})$ | Select word within block |
| Index | $i = \log_2(\text{number of sets})$ | Select set of blocks |
| Tag | $t = 32 - o - i$ | ID block within set |

- Consider: <BS=block size, S=sets, B=blocks>
  - <64,64,64>: o=6, i=6, t=20: direct-mapped (S=B)
  - <64,16,64>: o=6, i=4, t=22: 4-way S-A (S = B / 4)
  - <64,1,64>: o=6, i=0, t=26: fully associative (S=1)
- Total size = BS x B = BS x S x (B/S)

CADSL

# Replacement

- Cache has finite size
  - What do we do when it is full?

- Analogy: desktop full?
  - Move books to bookshelf to make room

- Same idea:
  - Move blocks to next level of cache

CADSL

# Replacement

- How do we choose *victim*?
  - Verbs: *Victimize, evict, replace, cast out*
- Several policies are possible
  - FIFO (first-in-first-out)
  - LRU (least recently used)
  - NMRU (not most recently used)
  - Pseudo-random
- Pick victim within *set* where a = *associativity*
  - If a <= 2, LRU is cheap and easy (1 bit)
  - If a > 2, it gets harder
  - Pseudo-random works pretty well for caches

**CADSL**

# Write Policy

- Memory hierarchy
  - 2 or more copies of same block
    - Main memory and/or disk
    - Caches

- What to do on a write?
  - Eventually, all copies must be changed
  - Write must *propagate* to all levels

**CADSL**

# Write Policy

- Easiest policy: *write-through*
- Every write propagates directly through hierarchy
  - Write in L1, L2, memory, disk (?!?)
- Why is this a bad idea?
  - Very high bandwidth requirement
  - Remember, large memories are slow
- Popular in real systems only to the L2
  - Every write updates L1 and L2
  - Beyond L2, use *write-back* policy

**CADSL**

# Write Policy

- Most widely used: *write-back*
- Maintain *state* of each line in a cache
  - Invalid – not present in the cache
  - Clean – present, but not written (unmodified)
  - Dirty – present and written (modified)
- Store state in tag array, next to address tag
  - Mark dirty bit on a write
- On eviction, check dirty bit
  - If set, write back dirty line to next level
  - Called a *writeback* or *castout*

**CADSL**

# Write Policy

- Complications of write-back policy
  - Stale copies lower in the hierarchy
  - Must always check higher level for dirty copies before accessing copy in a lower level

- Not a big problem in uniprocessors
  - In multiprocessors: *the cache coherence problem*
- I/O devices that use DMA (direct memory access) can cause problems even in uniprocessors
  - Called coherent I/O
  - Must check caches for dirty copies before reading main memory

CADSL

# Thank You

**CADSL**