# CPE 431/531

# Chapter 5 – Large and Fast: Exploiting Memory Hierarchy

# Dr. Rhonda Kay Gaede

THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE

# 5.1 Introduction

- Programmers always want _____ amounts of _____memory. Caches give that _____

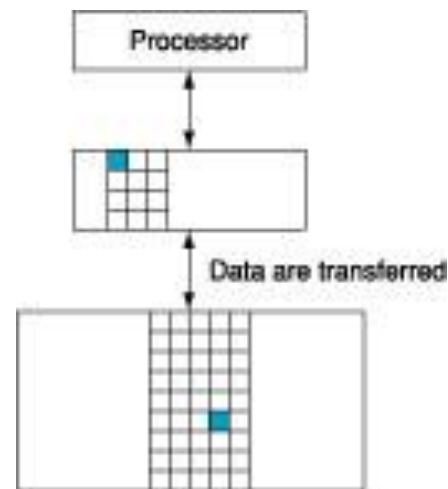- Principle of Locality
  - Temporal Locality -

    _____

    _____

  - Spatial Locality -    _____

    _____

- Build a memory _____.

# 5.1 Cache Terminology

- Data is copied between only ___ levels at a time.

- The minimum data unit is a _____.

- If the data appears  in the upper level, this situation is called a ___. The data not appearing in the upper level is called a ____.
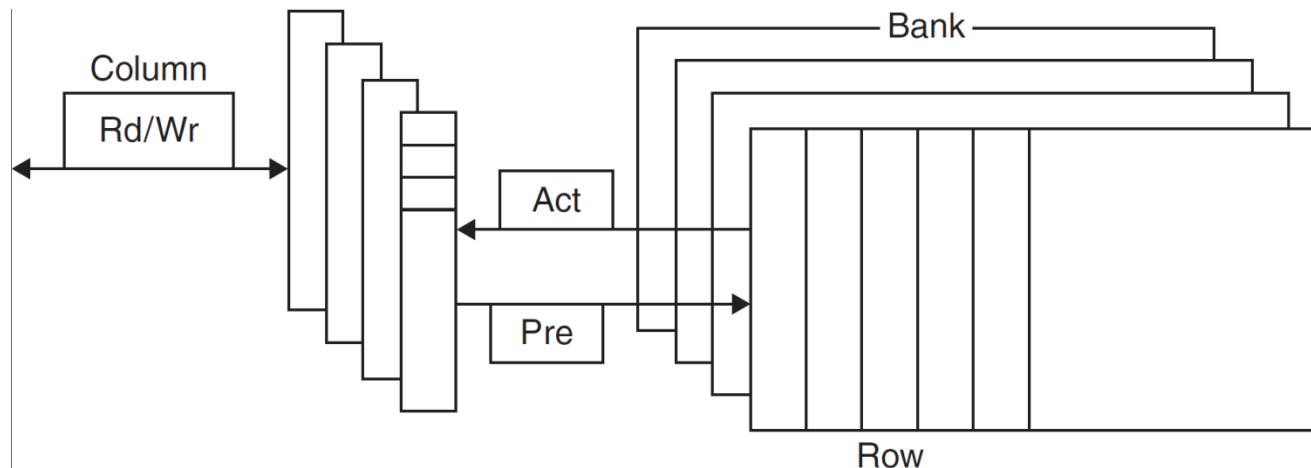
# 5.1 More Terminology

- The _____ is the fraction of memory accesses found in the upper level.

- The _____ is the fraction of memory accesses not found in the upper level.

- The _____ is the time to access the upper level of the memory hierarchy.

- The _____ is the time to replace a block in the upper level with the corresponding block from the lower level, plus the time to deliver this block to the processor.

# 5.2 Memory Technologies

- _____ RAM (____)
  - 0.5ns – 2.5ns, $2000 – $5000 per GB

- _____RAM (____)
  - 50ns – 70ns, $20 – $75 per GB

- _____ disk
  - 5ms – 20ms, $0.20 – $2 per GB

- Ideal memory
  - Access time of _____
  - Capacity and cost/GB of ____

# 5.2 DRAM Technology

- Data stored as a _____ in a _____
  - Single _____ used to access the _____
  - Must periodically be _____
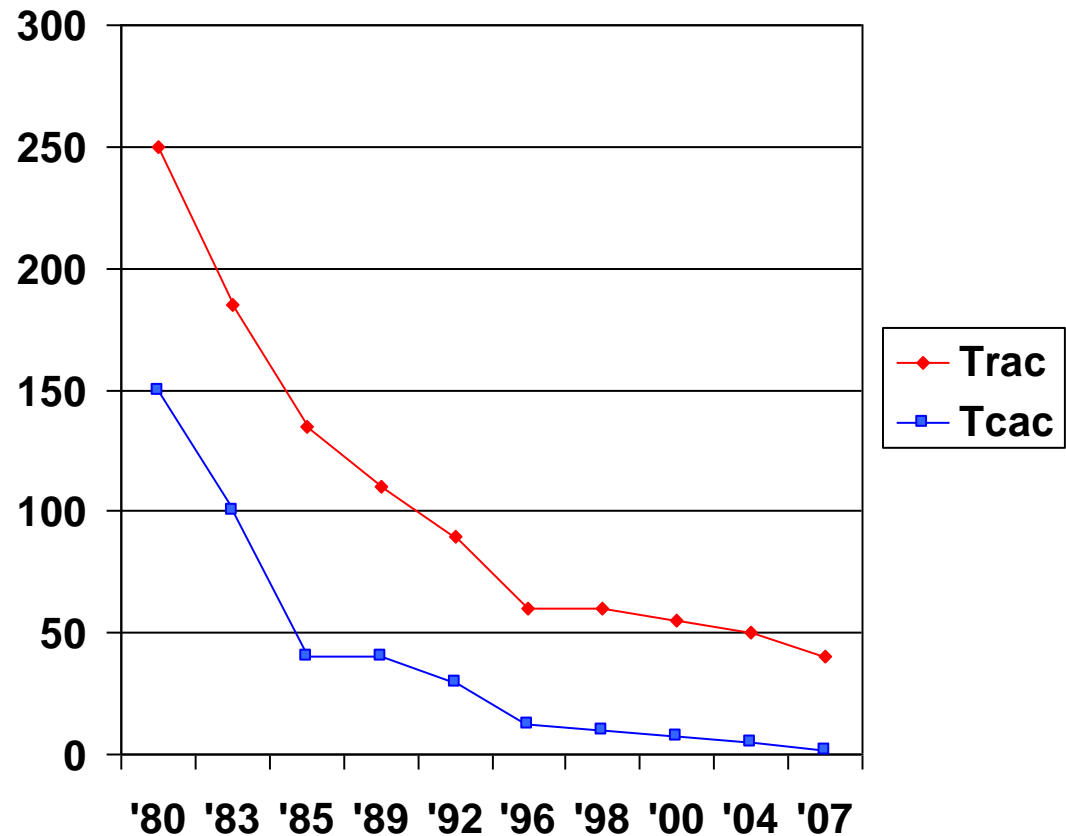    - _____ contents and _____ back
    - Performed on a DRAM "____"

# 5.2 Advanced DRAM Organization

- Bits in a DRAM are organized as a _____ _____
  - DRAM accesses an _____ _____
  - _____ mode: supply _____ words from a ____ with _____ _____
- _____ data rate (DDR) DRAM
  - Transfer on _____ and _____ clock edges
- _____ data rate (QDR) DRAM
  - Separate DDR _____ and _____

# 5.2 DRAM Generations

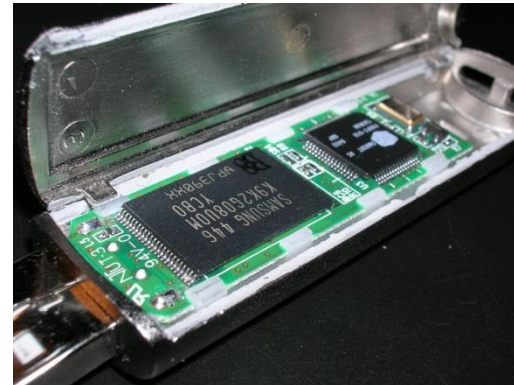| Year | Capacity | $/GB |
|------|----------|------|
| 1980 | 64Kbit | $1500000 |
| 1983 | 256Kbit | $500000 |
| 1985 | 1Mbit | $200000 |
| 1989 | 4Mbit | $50000 |
| 1992 | 16Mbit | $15000 |
| 1996 | 64Mbit | $10000 |
| 1998 | 128Mbit | $4000 |
| 2000 | 256Mbit | $1000 |
| 2004 | 512Mbit | $250 |
| 2007 | 1Gbit | $50 |

# 5.2 DRAM Performance Factors

- _____ _____

  - Allows _____ words to be read and refreshed in _____

- _____ _____

  - Allows for _____ accesses in bursts without needing to send _____ _____

  - Improves _____

- _____ _____

  - Allows _____ access to _____ DRAMs

  - Improves _____

# 5.2 Flash Storage

- _____ semiconductor storage
  - ____× – _____× faster than _____
  - _____, _____ power, more _____
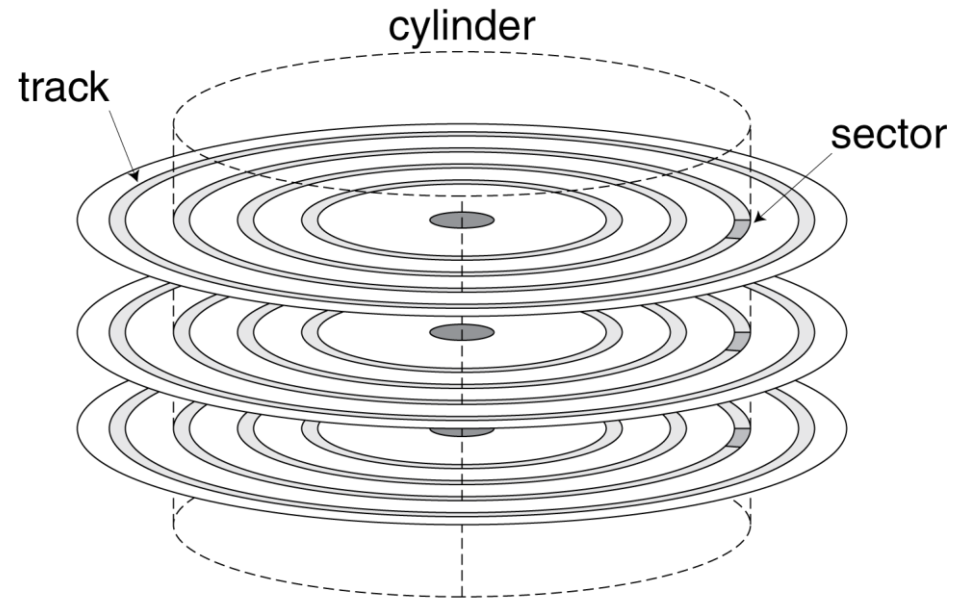  - But more $/GB (between _____ and _____)

# 5.2 Flash Types

- _____ flash: bit cell like a _____ gate
  - _____ read/write access
  - Used for _____ memory in _____ systems
- _____ flash: bit cell like a _____ gate
  - _____ (bits/area), but _____ access
  - _____ per GB
  - Used for _____ _____, _____ _____, …
- Flash bits wears out after _____ of accesses
  - Not suitable for direct _____ or _____ replacement
  - _____ _____: _____ data to less used blocks

# 5.2 Disk Storage

- _____, _____, _____ storage

# 5.2 Disk Sectors and Access

- Each _____ records
  - Sector \_\_\_
  - Data (\_\_\_\_ bytes, _____ bytes proposed)
  - _____ correcting code (ECC)
    - Used to hide _____ and recording _____
  - _____ fields and gaps
- Access to a _____ involves
  - _____ delay if other accesses are pending
  - \_\_\_\_\_: move the heads
  - _____ latency
  - Data _____
  - _____ overhead
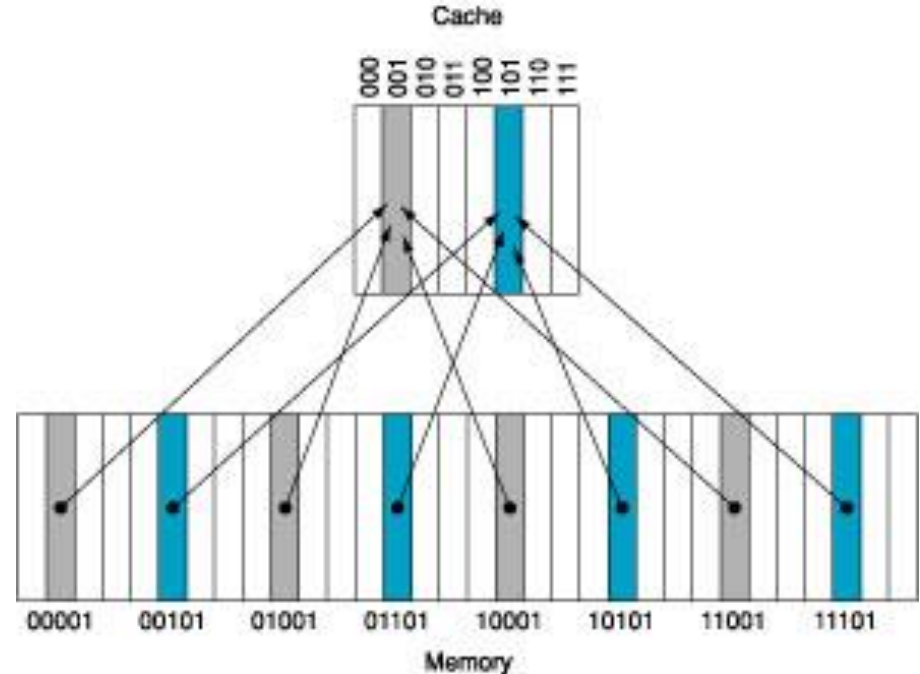
# 5.2 Disk Access Example

- Given
  - 512B sector, 15,000rpm, 4ms average seek time, 100MB/s transfer rate, 0.2ms controller overhead, idle disk

- Average read time

- If actual average seek time is 1ms

# 5.2 Disk Performance Issues

- Manufacturers quote _____ seek time
  - Based on _____ _____ seeks
  - _____ and ____ _____lead to smaller _____ average seek times
- Smart disk _____ allocate _____ sectors on disk
  - Present _____ sector interface to host
  - SCSI, ATA, SATA
- Disk drives include _____
  - _____ sectors in anticipation of access
  - Avoid _____ and _____ _____

# 5.3 Burning Question

- How do we know whether a data item is in the cache?

- If it is, how do we _____ it?

- The simplest scheme is that each item can be placed in _____ one place (_____ mapping).

- Mapping

# 5.3 Accessing a Cache

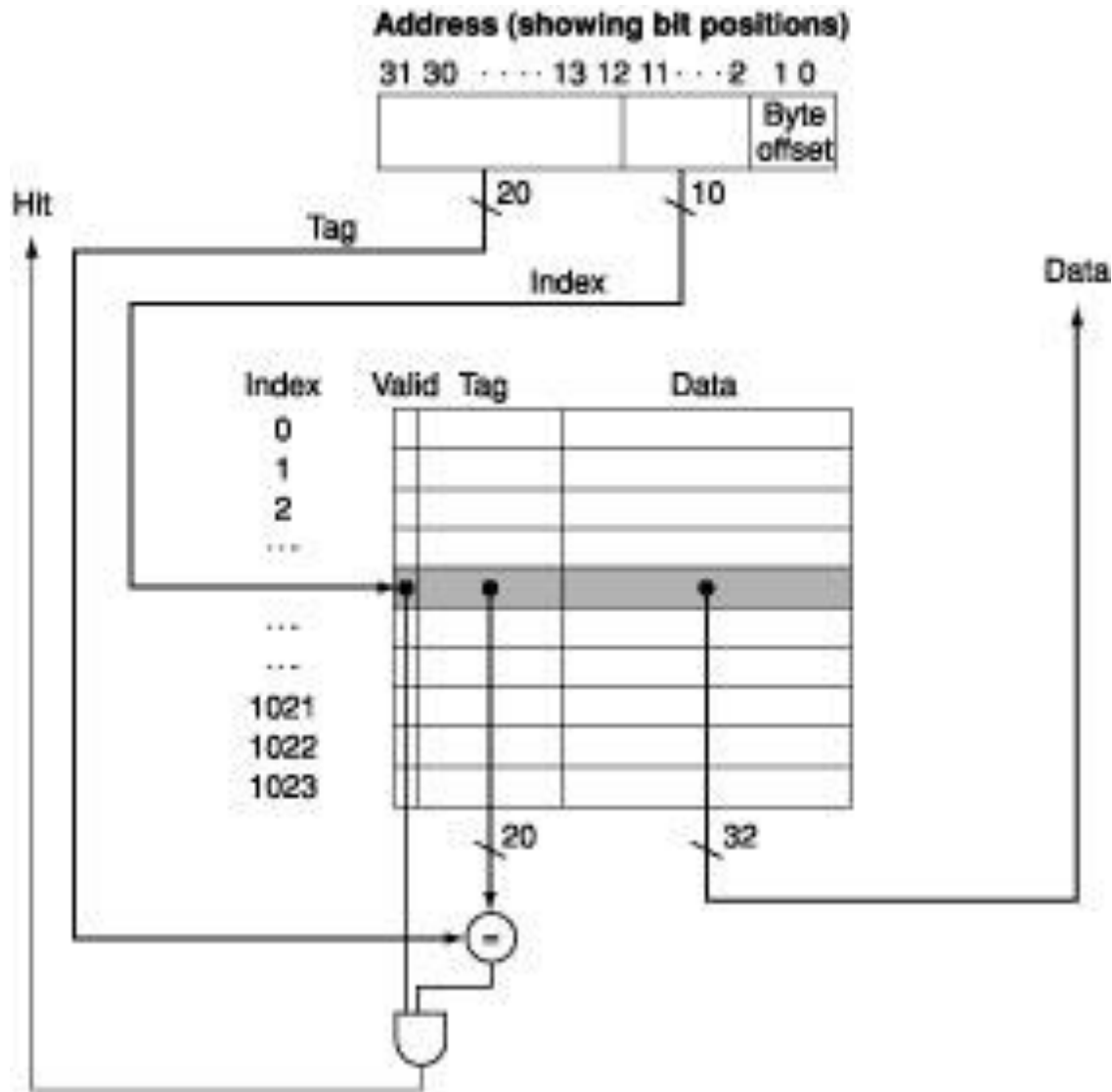| Index | V | Tag | Data |
|-------|---|-----|------|
| 000 | | | |
| 001 | | | |
| 010 | | | |
| 011 | | | |
| 100 | | | |
| 101 | | | |
| 110 | | | |
| 111 | | | |

22
26
22
26
16
3
16
18
16

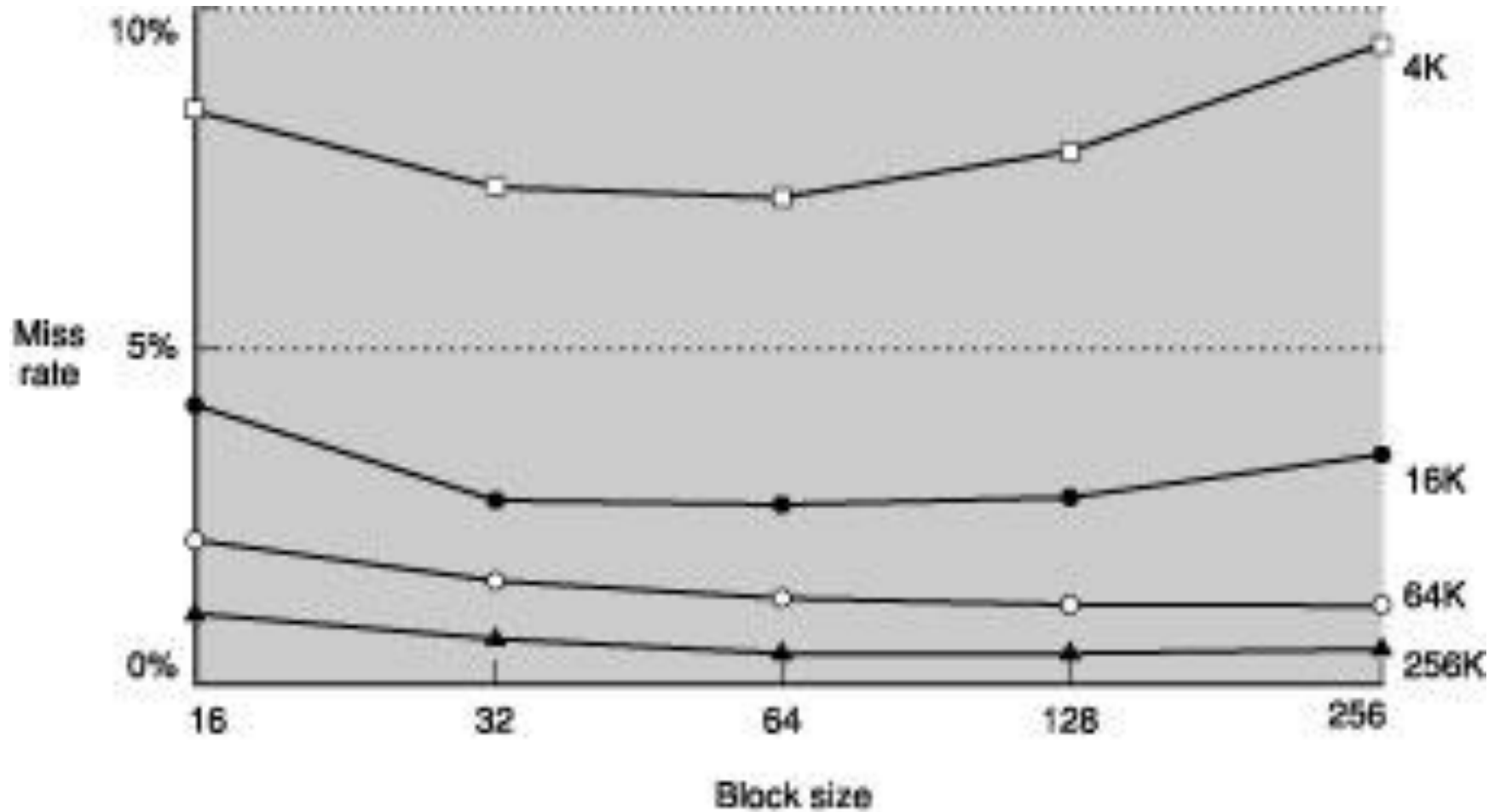# 5.3 Mapping Implemented in Hardware

# 5.3 Total Storage Required

Example: How many total bits are required for a direct-mapped cache with 16 KB of data and four-word blocks, assuming a 32-bit address?

# 5.3 Mapping an Address to a Multiword Cache Block

Consider a cache with 64 blocks and a block size of 16 bytes. What block number does byte address 1200 map to?

# 5.3 Miss Rate versus Block Size
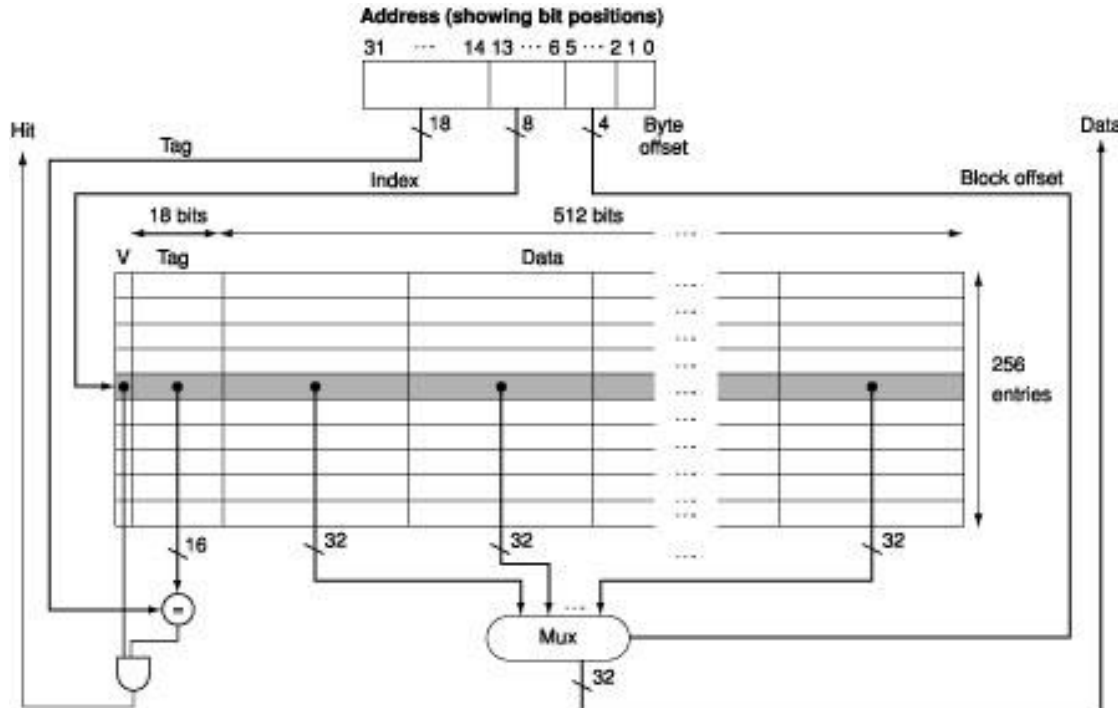
# 5.3 Handling Cache Misses

Instruction Cache Miss

1. Send the original PC value (current PC – 4) to the memory.

2. Instruct main memory to perform a read and wait for the memory to complete its access.

3. Write the cache entry, putting the data from memory in the data portion of the entry, writing the upper bits of the address into the tag field and turn the valid bit on.

4. Restart the instruction execution at the first step, which will refetch the instruction, this time finding it in the cache.

# 5.3 Handling Writes

- Suppose on a store instruction, we wrote the data into only the data cache (and not _____

  _____).

- Then the cache and main memory are said to be

  _____

- Solution A: _____ (_____

  _____)

  - Problem: _____

  - Remediation: _____

- Solution B: _____ (_____

  _____)

# 5.3 An Example Cache



This processor has a ___ stage pipeline.

When operating at peak speed, the processor can request both an _____ word and a _____ word on every clock cycle.

Separate _____ and _____ caches are used, each with ___ words and ___-word blocks.

For writes, the FastMATH offers both _____ and _____, letting the ____ decide.

The Intrinsity FastMATH Processor is a fast _____ processor that uses the MIPS architecture and a _____ cache implementation.

# 5.4 Measuring and Improving Cache Performance

- CPU time = (_____ + _____) x _____

- Memory-stall clock cycles = _____ + _____

- Read-stall cycles =

- Write-stall cycles =

- Memory-stall clock cycles =

- Calculating Cache Performance

    - $i_{miss}$ = 2 %, $d_{miss}$ = 4 %, $CPI_{perfect}$ = 2, miss penalty = 100 cycles, 36 % loads and stores

# 5.4 Impact of Increased Clock Rate

- Suppose the processor in the previous example _____ its clock rate, making the miss penalty _____.

- Total miss cycles per instruction = _____.

- Total CPI

- Performance with fast clock compared to performance with slow clock
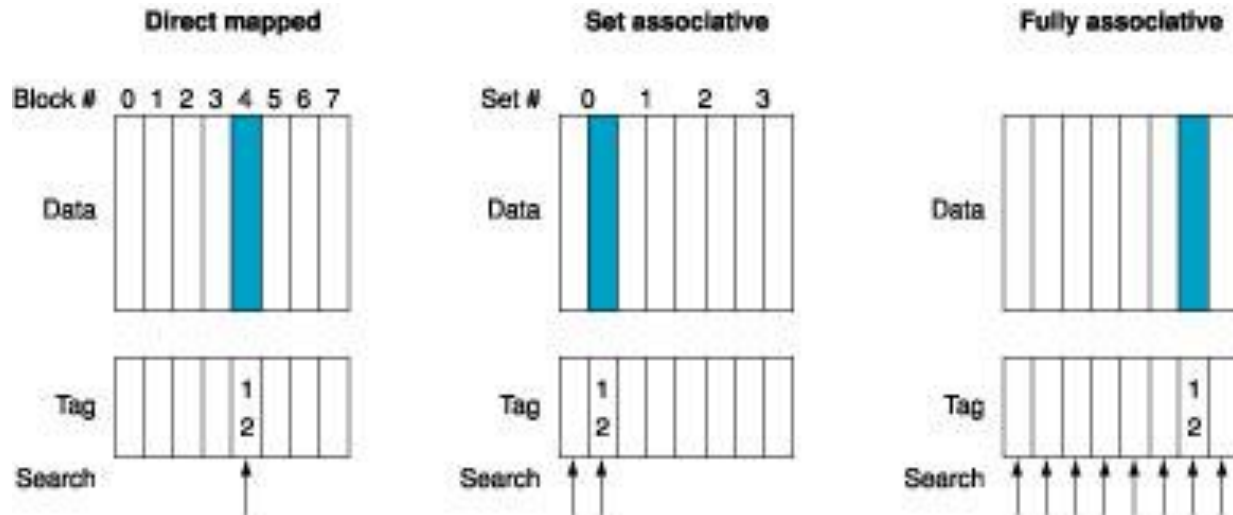
# 5.4 Average Memory Access Time

- To capture the fact that the time to _____ ____ for both ____ and _____ affects performance, designers sometimes use average memory access time (AMAT) as a way to examine _____ _____ _____.

- AMAT = Time for a hit + Miss rate x Miss penalty

- Find the AMAT for a processor with CT = 1 ns, Miss penalty = 20 cycles, Miss rate = 0.05/instruction, Cache access time = 1 cycle. Assume that read and write miss penalties are the same and ignore other write stalls.

# 5.4 Flexible Placement Reduces Cache Misses

- One Extreme - direct mapped -

- Middle Range - set associative -

- Other Extreme - fully associative -

- Set associative mapping

# 5.4 Conceptual View of Set Associativity

- One Extreme - direct mapped -

- Middle Range - set associative -

- Other Extreme - fully associative -

- Set associative mapping

# 5.4 Pseudo-Implementation View of Set Associativity

# 5.4 Misses and Associativity

Look at three small caches (four one word blocks): Address sequence: 0, 8, 0, 6, 8

  a. fully associative     b. two-way set associative     c. direct mapped

# 5.4 Locating a Block in the Cache



Which block do we replace?

# 5.4 Tag Size Considerations

- Size of Tags versus Set Associativity

    For a cache with 4K blocks, a 32-bit address with 0 bits for block and byte offsets, find the #sets, #tag bits for 1, 2, 4 and fully associative organizations

# 5.4 Performance of Multilevel Caches

- Example:
  - $CPI_{base} = 1.0$, CR = 4 GHz
  - $Mem_{access} = 100$ ns, $L1inst_{miss} = 2$ %
  - $L2_{access} = 5$ ns, $L2_{miss}$ per instruction = 0.5 %

# 5.4 Interactions with Software

- Misses depend on memory access patterns
  - Algorithm behavior
  - Compiler optimization for memory access

# 5.4 Software Optimization via Blocking

Goal: maximize accesses to data before it is replaced

Consider inner loops of DGEMM:

Blocked algorithms operate on submatrices or blocks rather than entire rows or columns of an array

```
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;             /* C[i][j] = cij */
}
```

# 5.4 Array Access Patterns

older accesses

new accesses

# 5.4 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 }
```

# 5.4 Blocked DGEMM Access Pattern

# 5.5 Dependable Memory Hierarchy

- Two states of service

    1. Service accomplishment -

    2. Service interruption –

- Transitions from 1 to 2 are _____, transitions from 2 to 1 are _____.

- Failures can be _____ or _____.

- Reliability is a measure of the continuous service accomplishment, the metric is _____.

- Availability is a measure of the service accomplishment with respect to the alternation between the two states of accomplishment and interruption.

# 5.5 MTTF vs. AFR for Disks

- Some disks today are quoted to have a 1,000,000-hour MTTF. As 1,000,000 hours is 114 years, it would seem like they practically never fail. Warehouse scale computers that run Internet services such as Search might have 50,000 servers. Assume each server has 2 disks. Use AFR to calculate how many disks we would expect to fail per year.

# 5.5 Availability Considerations

- To increase _____, you can improve the _____ of components or design systems to _____ operation in the presence of components that have _____.
- Three techniques
  - Fault avoidance
  - Fault tolerance
  - Fault forecasting
- We also need to work on decreasing _____

# 5.5 Single Error Detection - Parity

- Hamming distance
  - Number of ___ that are _____ between two bit patterns
- Minimum distance = __ provides single bit error detection
  - e.g. _____ code

- Minimum distance = __ provides _____ error correction, ___ ____ error detection

# 5.5 Encoding Single Error Correcting Hamming Code

- To calculate Hamming code:
  - Number bits from 1 on the left
  - All bit positions that are a _____ __ __ are _____ bits
  - Each _____ bit checks certain _____ bits:

| Bit position | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Encoded date bits | | p1 | p2 | d1 | p4 | d2 | d3 | d4 | p8 | d5 | d6 | d7 | d8 |
| Parity bit coverate | p1 | X | | X | | X | | X | | X | | X | |
| | p2 | | X | X | | | X | X | | | X | X | |
| | p4 | | | | X | X | X | X | | | | | X |
| | p8 | | | | | | | | X | X | X | X | X |

# 5.5 Decoding Single Error Correcting Hamming Code

- _____ of parity bits indicates which bits are __ _____
  - Use numbering from _____ procedure
  - E.g.
    - Parity bits = _____ indicates ____ _____
    - Parity bits = _____ indicates bit ____ was flipped

# 5.5 SEC/DEC Hamming Code

- Add an additional parity bit for the _____ ___ ($p_n$)

- Make Hamming distance = __

- Decoding:

  - Let H = SEC parity bits

    - H even, $p_n$ even, _____
    - H odd, $p_n$ odd, _____
    - H even, $p_n$ odd, _____
    - H odd, $p_n$ even, _____

- Note: ECC DRAM uses SEC/DEC with 8 bits protecting each 64 bits

# 5.6 Virtual Machines Redux

- First developed in the 1960s, they have remained an important part of _____ computing and have recently gained popularity due to
  - Increasing importance of _____ and _____
  - The failures in _____ and _____ of standard operating systems
  - The _____ of a single computer among many related users
  - The dramatic increase in ___ ____ of processors
- Broadest Definition
  - Includes basically all _____ methods that provide a standard software interface, like the _____
- Our Definition
  - Provide a complete _____ environment at the _____ level

# 5.6 Virtual Machine Basics

- System virtual machines present the illusion that users have an _____ _____ to themselves, including a copy of the _____ _____.

- With a VM, multiple OSes all _____ the _____ resources.

- The software that supports VMs is called a _____ _____ _____ (VMM) or _____.

- The underlying hardware is called the _____, sharing resources among the _____ VMs.

# 5.6 Virtual Machine Ancillary Benefits

- Our interest is primarily in improving _____

- Other benefits include

  - _____ _____: a typical deployment might be some OSes running legacy OSes, many running the current stable OS release, and a few testing the next OS release.

  - _____ _____: Consolidate the number of servers. Some VMMs support migration of a running VM to a different computer, either to balance load or to evacuate from failing hardware.
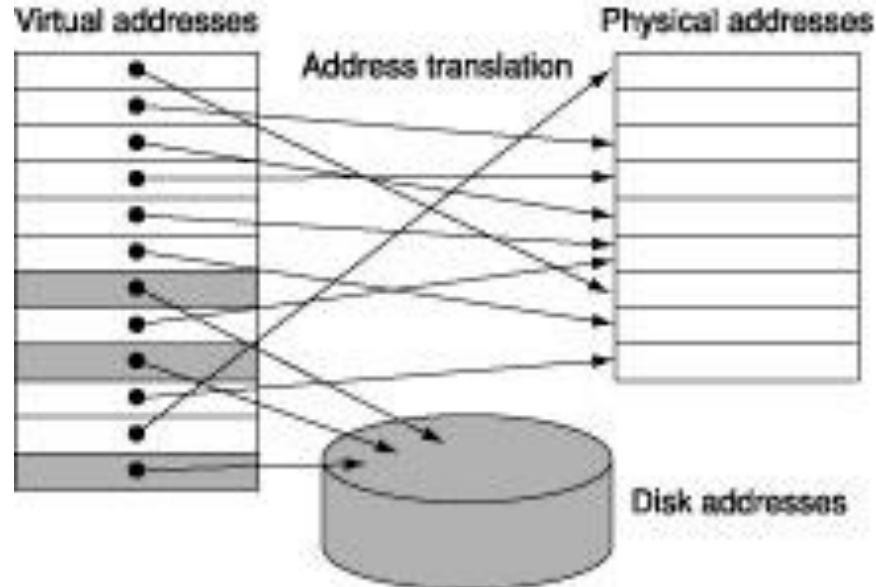
# 5.6 Requirements of a Virtual Machine Monitor

- Guest software should behave on a VM exactly as if it were running on the _____ _____, except for _____ behavior or limitations of _____ _____ shared by multiple VMs.

- Guest software should not be able to _____ _____ of real system resources directly.

- At least ___ processor modes, _____ and ____.

- A _____ subset of instructions available only in _____ mode, all system _____ must be controllable only via these instructions.

# 5.7 Virtual Memory

- The ___ _____ can act as a _____ for the _____ storage.

- Historically, two motivations for virtual memory

  – _____

  – _____

- Virtual memory implements the _____ of a program's address space to _____.

- This translation process enforces _____ of a program's address space from other programs.
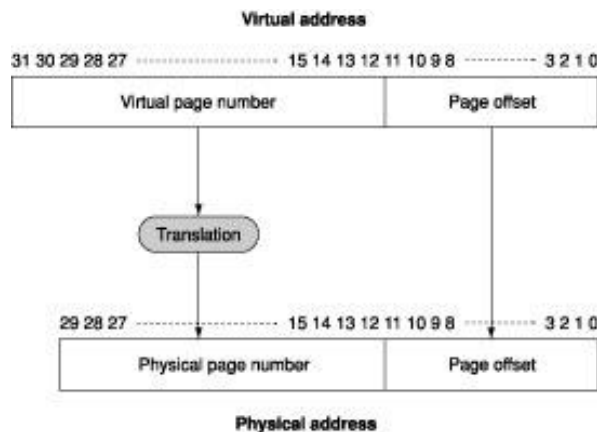
# 5.7 Virtual Memory Terminology

- A virtual memory _____ is called a _____.

- A virtual memory _____ is called a _____

  _____.

- Each _____ address is translated to a _____ address.

- This process is called _____.
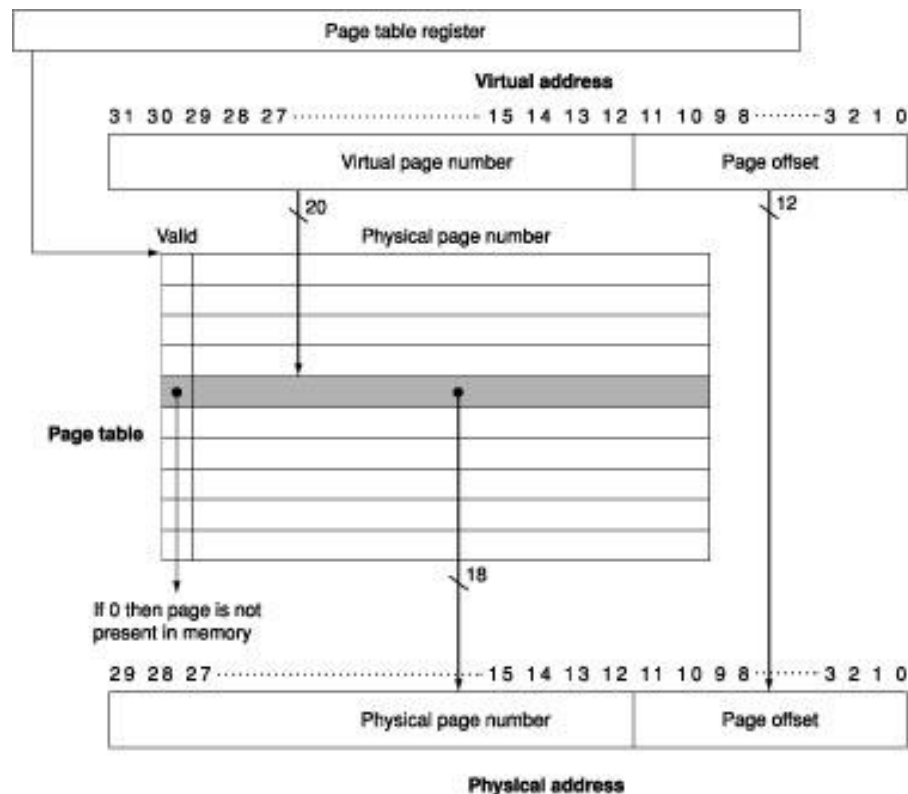
# 5.7 Virtual Memory Facts

• A virtual _____ is broken into a virtual _____ _____ and a _____ _____.



Virtual address

31 30 29 28 27 ···················· 15 14 13 12 11 10 9 8 ········· 3 2 1 0

| Virtual page number | Page offset |

Translation

29 28 27 ···················· 15 14 13 12 11 10 9 8 ········ 3 2 1 0

| Physical page number | Page offset |

Physical address

• A page fault takes _____ of cycles to process
  – Pages should be _____ enough to _____ the high access time, though _____ systems are going smaller.
  – _____ _____ placement of pages is _____.
  – Page faults can be handled in _____.
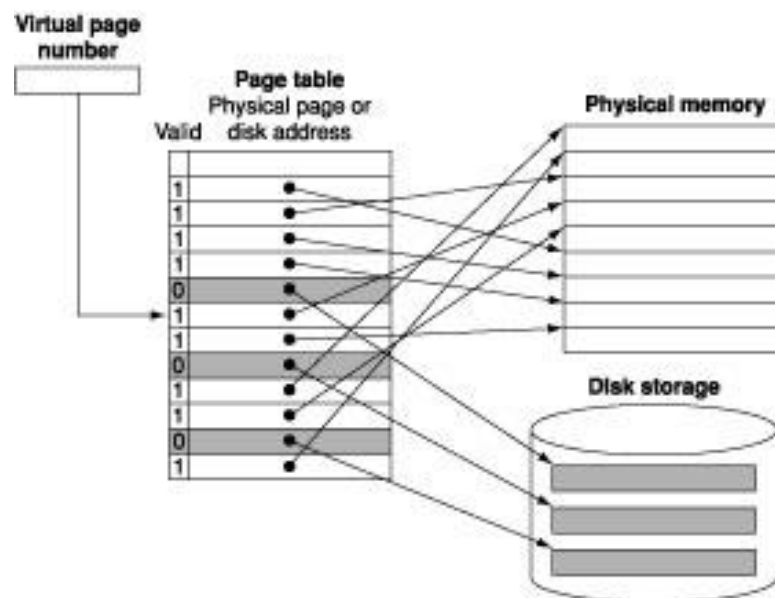  – Virtual memory uses _____.

# 5.7 Virtual Memory Mapping

- Pages are _____ by using a _____ _____ that _____ memory.
- Each _____ has its own _____ _____.
- The _____ _____ register points to the _____ of the _____ _____.
- A _____ table is too _____, _____ page tables are used.
- The _____ of a _____ consists of the _____ _____ _____, _____ _____ and _____.

# 5.7 Page Faults

- The _____ _____ manages page replacement.
- The _____ _____ usually creates the _____ __ ___ for all of the pages of a process when it creates the process, this space is called _____ _____

- A data structure records where ____ ____ is stored on disk. Another data structure tracks which _____ and which ____ _____ use each _____ ____.

- On a page fault, the _____

- ____ page is evicted.

- Consider 10, 12, 9, 7, 11, 10, then 8

# 5.7 Making Address Translation Fast: The TLB

- With virtual memory, you need ___ memory accesses, one extra for the _____.
- Add a _____ to keep track of _____ translations.
- It's called a _____ _____ _____ (TLB).
- A TLB _____ may or may not be a _____ _____

TLB characteristics
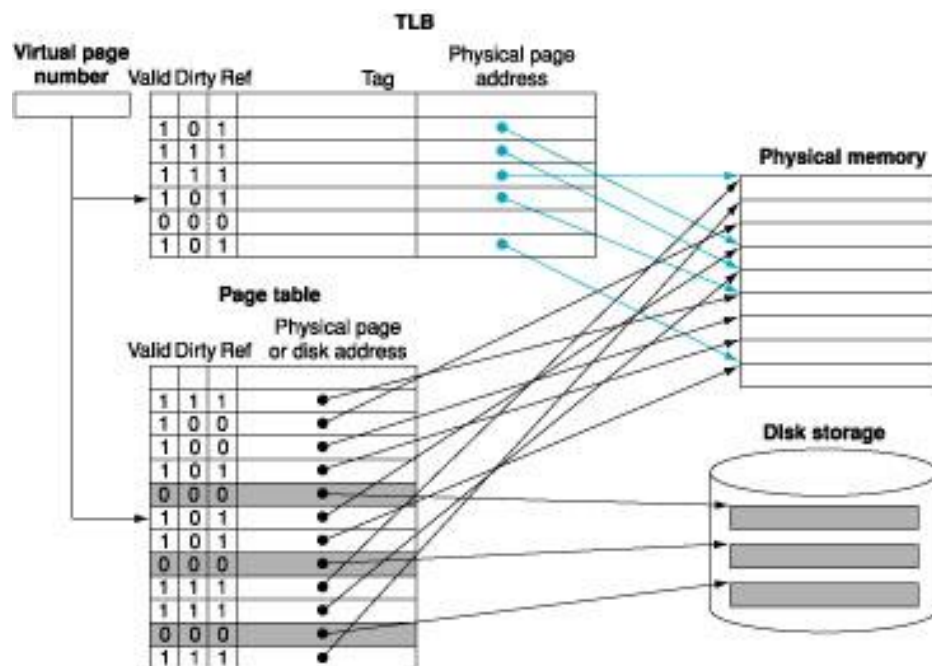    size: _____
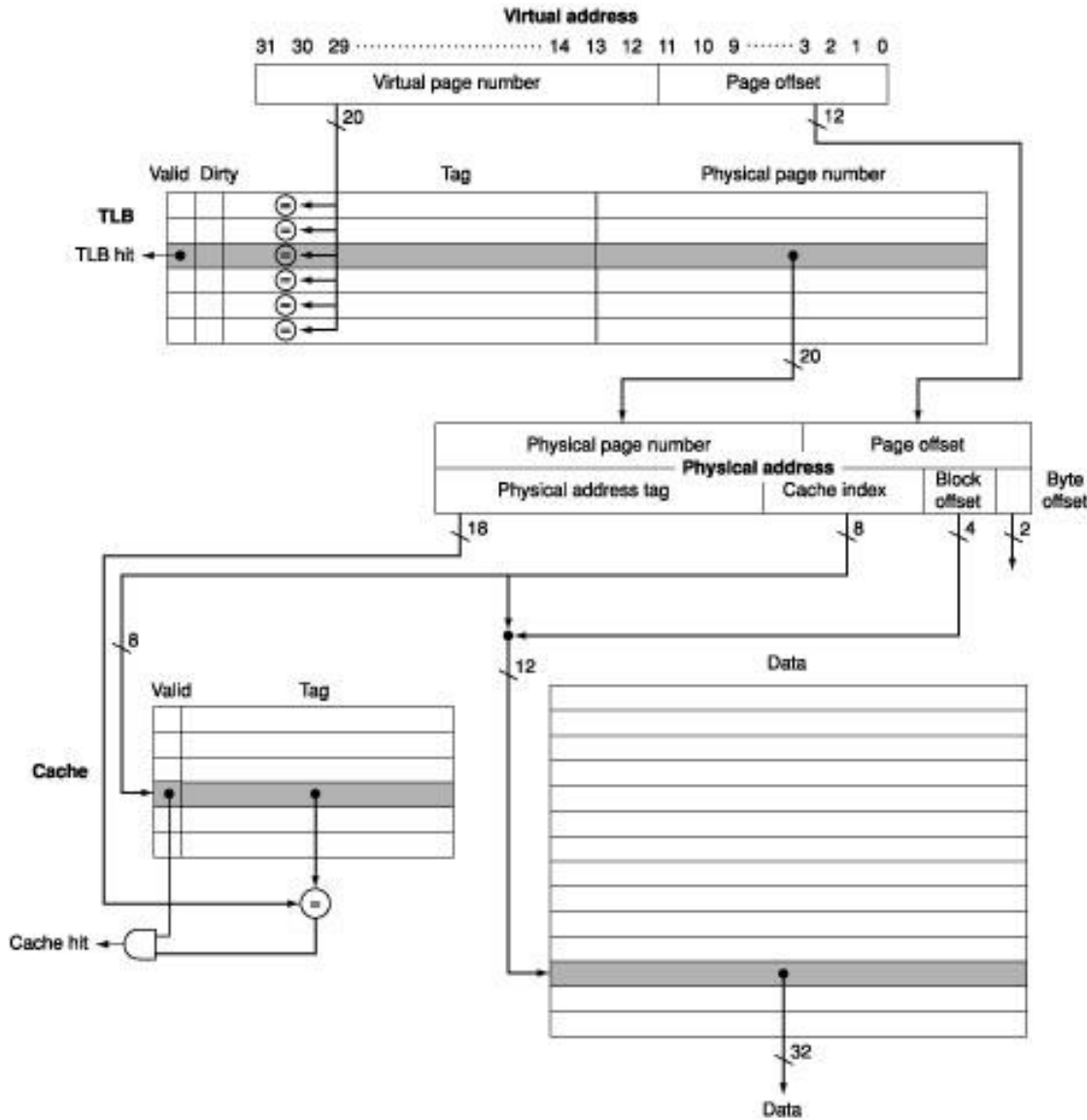    block size: _____
    hit time: _____
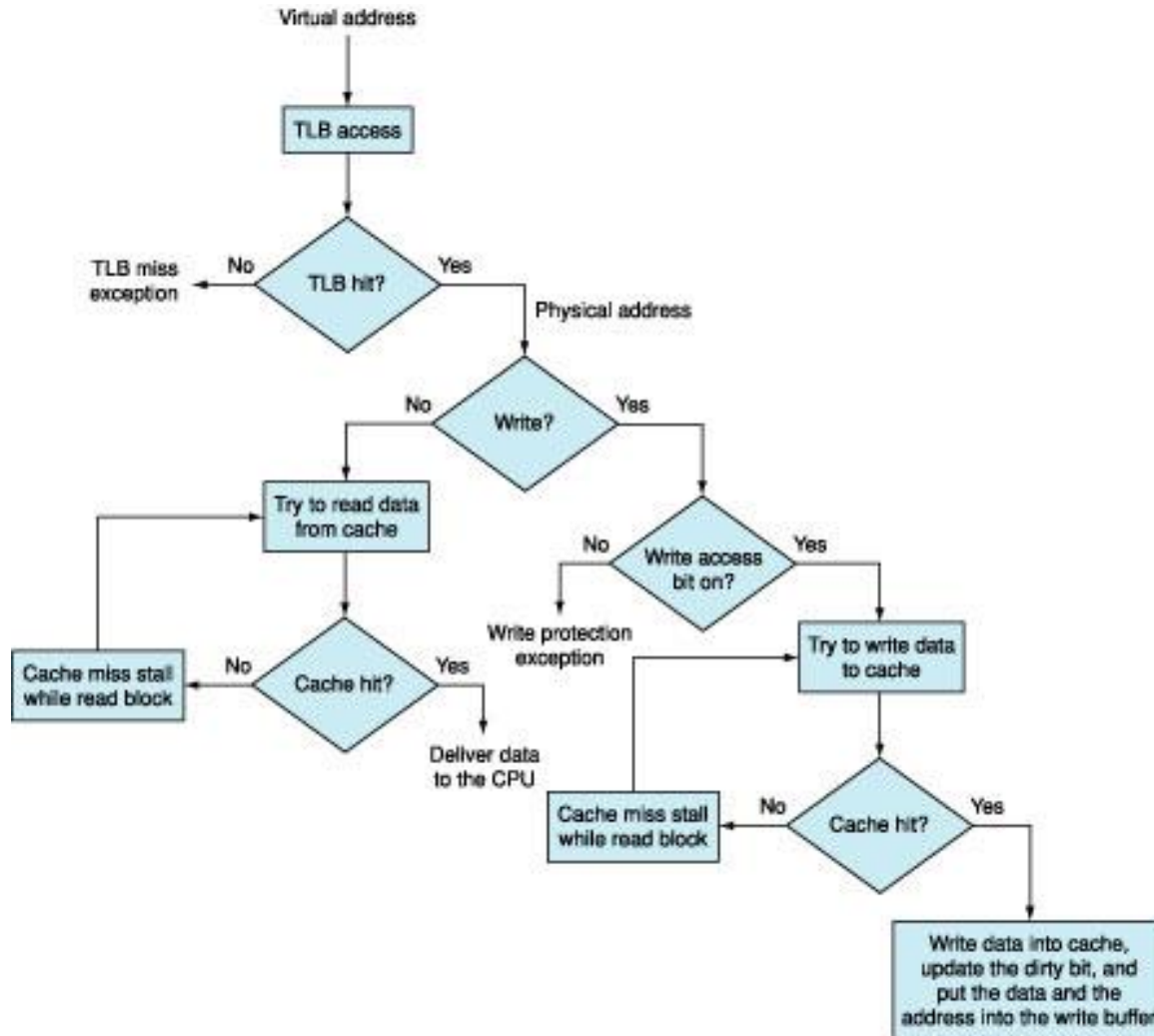    miss penalty: _____
    miss rate: _____

# 5.7 The Intrinsity FastMATH TLB

# 5.7 Processing a read or write-through

# 5.4 Overall Operation of a Memory Hierarchy

| TLB | Page Table | Cache | Possible? If so, under what circumstance? |
|---|---|---|---|
| miss | miss | miss | TLB misses and is followed by a page fault; after retry, data must miss in cache |
| miss | miss | hit | Impossible: data cannot be allowed in cache if the page is not in memory |
| miss | hit | miss | TLB misses, but entry found in page table; after retry, data misses in cache |
| miss | hit | hit | TLB misses, but entry found in page table; after retry, data is found in cache |
| hit | miss | miss | Impossible: data cannot be allowed in cache if the page is not in memory |
| hit | miss | hit | Impossible: data cannot be allowed in cache if the page is not in memory |
| hit | hit | miss | Possible, though the page table is never really checked if TLB hits |

# 5.7 Implementing Protection with Virtual Memory
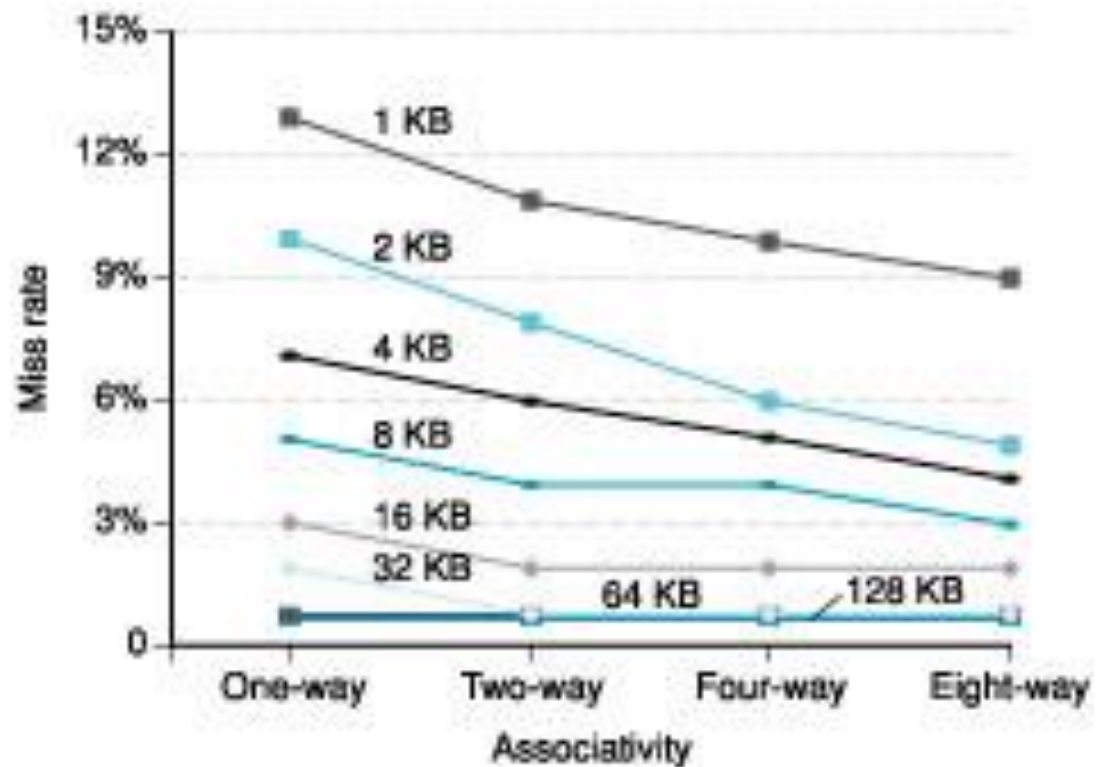
- Hardware Must Provide
  - At least two modes
    - _____
    - _____
  - Provide a portion of the _____ _____ that a user process can _____ but not _____
    - _____
  - Provide mechanisms whereby the processor can move between modes
    - _____
    - _____
- Software Can Help
  - Place the _____ _____ in the _____ address space of the _____

# 5.7 Summary

- Pages are made _____ to take advantage of _____ locality and reduce the _____ rate.

- The mapping between virtual addresses and physical addresses, which is implemented in a _____ _____, is made _____ _____so that a virtual page can be placed _____ in main memory.

- The _____ _____uses techniques, such as LRU and a reference bit, to choose which pages to _____.

# 5.8 Where can a Block be Placed?

- A _____ of Associativities is possible
- Advantage: _____ associativity _____ miss rates.
- Disadvantage: Increasing _____ increases ___ and _____ ____.

# 5.8 How is a Block Found?

- Cache
  - _Small_ degrees of associativity are used because _large_ degrees are _expensive_

- Virtual Memory
  - _Full_ _associativity_ makes sense because
    - Misses are _very_ _expensive_
    - _Software_ can implement _sophisticated_ replacement schemes
    - Full map can be easily _indexed_
    - _Large_ items means small number of _mappings_

# 5.8 Replace Which Block on a Cache Miss?

- Cache

  – _____
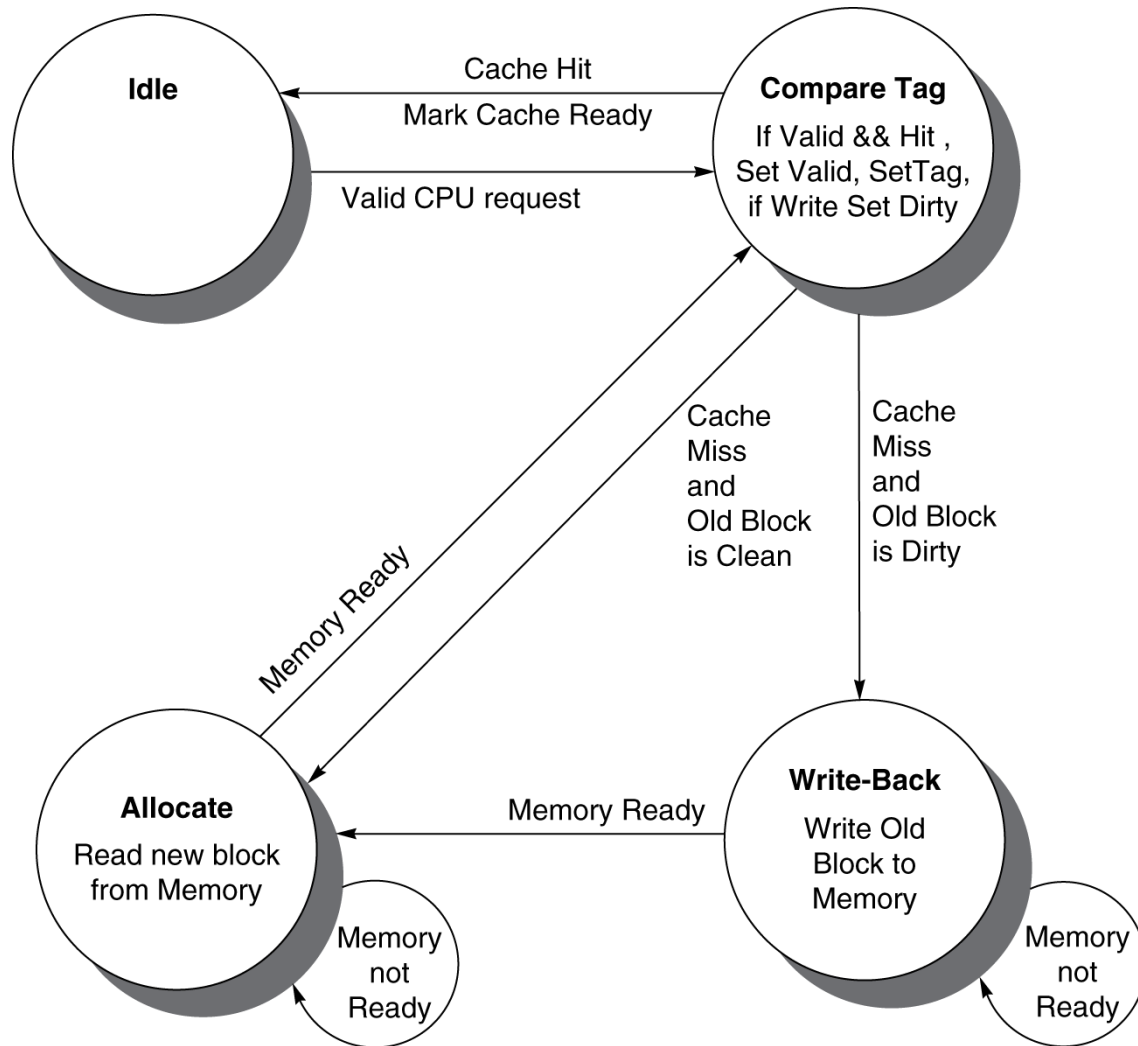
  – ___

- Virtual Memory

  – ___

# 5.8 What Happens on a Write?

- Caches - write-back is the _____ of the _____
- Virtual memory – _____ uses write-back

# 5.8 The Three Cs

- Cache misses occur in three categories:

  – Compulsory misses -

  – Capacity misses -

  – Conflict misses -

- The _____ in designing memory hierarchies is that every _____ that _____ improves the ____ rate can also _____ affect overall performance.

# 5.9 Finite State Cache Controller

# 5.10 Cache Coherence

- Multiple processors commonly _____ an _____ _____
- They may bring _____ into _____ and then write to their _____copies
- These local copies are likely
  - Different from ____ _____
  - Different from the _____ _____ value
- Example

| Time step | Event | CPU A's cache | CPU B's cache | Memory |
|---|---|---|---|---|
| 0 | | | | 0 |
| 1 | CPU A reads X | 0 | | 0 |
| 2 | CPU B reads X | 0 | 0 | 0 |
| 3 | CPU A writes 1 to X | 1 | 0 | 1 |

# 5.10 Snooping

- Every cache has the _____ _____ of the block along with the block
- The caches are all _____ via some _____ mechanism (____ or _____).
- All cache controllers _____ or _____ on the medium to see whether or not they have a _____ of the _____ requested.
- The _____ CPU broadcasts an _____ of the block.
- Example

| CPU activity | Bus activity | CPU A's cache | CPU B's cache | Memory |
|---|---|---|---|---|
| | | | | 0 |
| CPU A reads X | Cache miss for X | 0 | | 0 |
| CPU B reads X | Cache miss for X | 0 | 0 | 0 |
| CPU A writes 1 to X | Invalidate for X | 1 | | 0 |
| CPU B read X | Cache miss for X | 1 | 1 | 1 |

# 5.16 Concluding Remarks

- The difficulty of building a memory system to keep pace with faster processors is underscored by the fact that the raw material for main memory, DRAMs, is essentially the same in all computers, fast or slow.

- Multilevel caches facilitate cache optimizations.

  - _____

  - _____

    _____

- Other trends

  - _____

  - _____