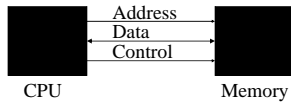


## The Main Memory Unit

- CPU and memory unit interface



- CPU issues address (and data for write)
- Memory returns data (or acknowledgment for write)

1

## Memories: Design Objectives

- Provide adequate storage capacity
- Four ways to approach this goal
  - Use of number of different memory devices with different cost/performance ratios
  - Automatic space-allocation methods in hierarchy
  - Development of virtual-memory concept
  - Design of communication links

2

## Memories: Characteristics

- Location: Inside CPU, Outside CPU, External
- Performance: Access time, Cycle time, Transfer rate
- Capacity: Word size, Number of words
- Unit of Transfer: Word, Block
- Access: Sequential, Direct, Random, associative
- Physical Type: Semiconductor, Magnetic, Optical

3

## Memories: Basic Parameters

- Cost:  $c=C/S$  (\$/bit)
- Performance:
  - Read access time ( $T_a$ ), access rate ( $1/T_a$ )
- Access Mode: random access, serial, semi-random
- Alterability:
  - R/W, Read Only, PROM, EPROM, EEPROM
- Storage:
  - Destructive read out, Dynamic, Volatility
- Hierarchy:
  - Tape, Disk, DRUM, CCD, CORE, MOS, BiPOLAR

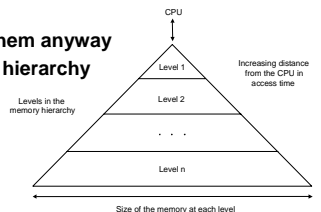
4

## Exploiting Memory Hierarchy

- Users want large and fast memories!

SRAM access times are 1 - 25ns at cost of \$100 to \$250 per Mbyte.  
 DRAM access times are 60-120ns at cost of \$5 to \$10 per Mbyte.  
 Disk access times are 10 to 20 million ns at cost of \$.10 to \$.20 per Mbyte.

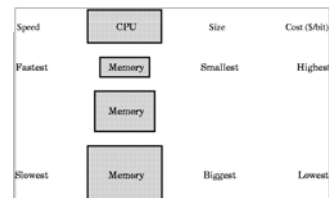
- Try and give it to them anyway
  - build a memory hierarchy



5

## Advantage of Memory Hierarchy

- Decrease cost/bit
- Increase capacity
- Improve average access time
- Decrease frequency of accesses to slow memory



6

## Memories: Review

- **SRAM:**
  - value is stored on a pair of inverting gates
  - very fast but takes up more space than DRAM
- **DRAM:**
  - value is stored as a charge on capacitor
  - very small but slower than SRAM (factor of 5/10)



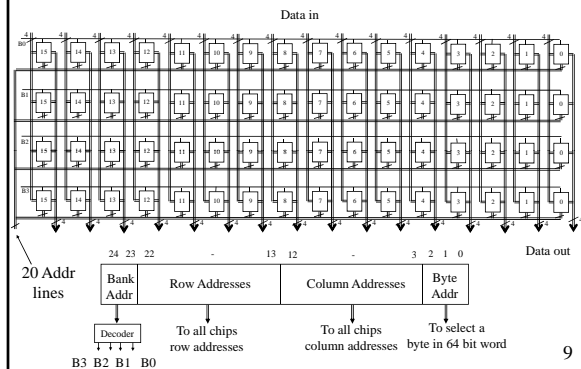
7

## Memories: Array Organization

- Storage cells are organized in a rectangular array
- The address is divided into row and column parts
- There are  $M (=2^r)$  rows of  $N$  bits each
- The row address ( $r$  bits) selects a full row of  $N$  bits
- The column address ( $c$  bits) selects  $k$  bits out of  $N$
- $M$  and  $N$  are generally powers of 2
- Total size of a memory chip =  $M \cdot N$  bits
  - It is organized as  $A=2^{r+c}$  addresses of  $k$ -bit words
- To design an  $R$  addresses  $W$ -bit words memory, we need  $\lceil R/A \rceil \cdot \lceil W/k \rceil$  chips

8

## 4Mx64-bit Memory using 1Mx4 memory chip :



9

## Locality

- A principle that makes memory hierarchy a good idea
- If an item is referenced
  - temporal locality: it will tend to be referenced again soon
  - spatial locality: nearby items will tend to be referenced soon.
- Why does code have locality?
- Our initial focus: two levels (upper, lower)
  - block: minimum unit of data
  - hit: data requested is in the upper level
  - miss: data requested is not in the upper level

10

## Memory Hierarchy and Access Time

- $t_i$  is time for access at level  $i$ 
  - on-chip cache, off-chip cache, main memory, disk, tape
- $N$  accesses
  - $n_i$  satisfied at level  $i$
  - a higher level can always satisfy any access that is satisfied by a lower level
  - $N = n_1 + n_2 + n_3 + n_4 + n_5$
- Hit Ratio
  - number of accesses satisfied/number of accesses made
  - Could be confusing
  - For example for level 3 is it  $n_3/N$  or  $(n_1+n_2+n_3)/N$  or  $n_3/(N-n_1-n_2)$
  - We will take the second definition

11

## Average Access Time

- $t_i$  is time for access at level  $i$
- $n_i$  satisfied at level  $i$
- $h_i$  is hit ratio at level  $i$ 
  - $h_i = (n_1 + n_2 + \dots + n_i) / N$
- We will also assume that data are transferred from level  $i+1$  to level  $i$  before satisfying the request
- Total time =  $n_1 \cdot t_1 + n_2 \cdot (t_1 + t_2) + n_3 \cdot (t_1 + t_2 + t_3) + n_4 \cdot (t_1 + t_2 + t_3 + t_4) + n_5 \cdot (t_1 + t_2 + t_3 + t_4 + t_5)$
- Average time = Total time/ $N$
- $t(\text{avr}) = t_1 + t_2 \cdot (1-h_1) + t_3 \cdot (1-h_2) + t_4 \cdot (1-h_3) + t_5 \cdot (1-h_4)$
- Total Cost =  $C_1 \cdot S_1 + C_2 \cdot S_2 + C_3 \cdot S_3 + C_4 \cdot S_4 + C_5 \cdot S_5$

12

## Cache

- Two issues:
  - How do we know if a data item is in the cache?
  - If it is, how do we find it?
- Our first example:
  - block size is one word of data
  - "direct mapped"

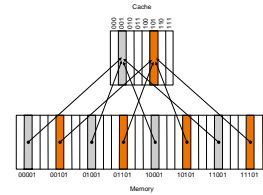
For each item of data at the lower level, there is exactly one location in the cache where it might be.

e.g., lots of items at the lower level share locations in the upper level

13

## Direct Mapped Cache

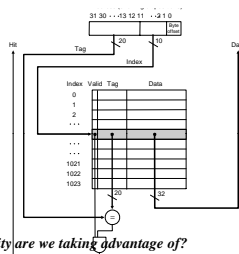
- Mapping:
  - address is modulo the number of blocks in the cache



14

## Direct Mapped Cache

- For MIPS:

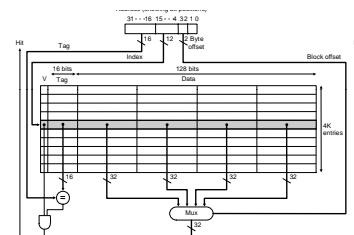


What kind of locality are we taking advantage of?

15

## Direct Mapped Cache

- Taking advantage of spatial locality:



16

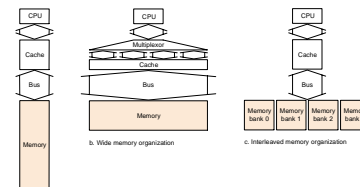
## Hits vs. Misses

- Read hits
  - this is what we want!
- Read misses
  - stall the CPU, fetch block from memory, deliver to cache, restart
- Write hits:
  - can replace data in cache and memory (write-through)
  - write the data only into the cache (write-back the cache later)
- Write misses:
  - read the entire block into the cache, then write the word

17

## Hardware Issues

- Make reading multiple words easier by using banks

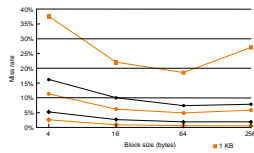


- It can get a lot more complicated...

18

## Performance

- Increase in block size tend to decrease miss rate:



- Use split caches (more spatial locality in code)

| Program | Block size in words | Instruction miss rate | Data miss rate | Effective combined miss rate |
|---------|---------------------|-----------------------|----------------|------------------------------|
| gcc     | 1                   | 5.15%                 | 2.1%           | 5.4%                         |
|         | 4                   | 2.0%                  | 1.7%           | 1.9%                         |
| spice   | 1                   | 1.2%                  | 1.3%           | 1.2%                         |
|         | 4                   | 0.3%                  | 0.6%           | 0.4%                         |

19

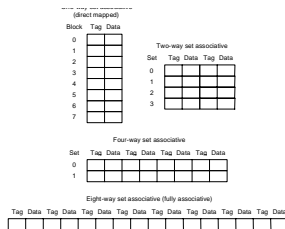
## Performance

- Simplified model:
  - execution time=(execution cycles + stall cycles)\*cct
- stall cycles= #of instructions\*miss ratio\*miss penalty
- Two ways of improving performance:
  - decreasing the miss ratio
  - decreasing the miss penalty

What happens if we increase block size?

20

## Decreasing miss ratio with associativity

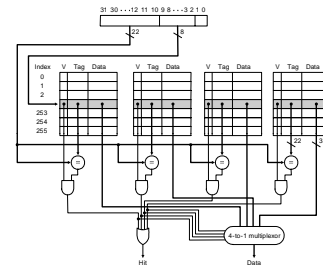


Compared to direct mapped, give a series of references that:

- results in a lower miss ratio using a 2-way set associative cache
- results in a higher miss ratio using a 2-way set associative cache assuming we use the "least recently used" replacement strategy

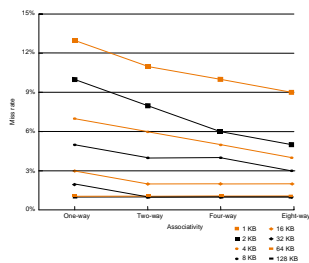
21

## An implementation



22

## Performance



23

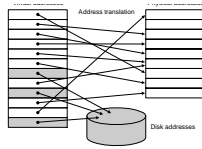
## Decreasing miss penalty with multilevel caches

- Add a second level cache:
  - often primary cache is on the same chip as the processor
  - use SRAMs to add another cache above primary memory (DRAM)
  - miss penalty goes down if data is in 2nd level cache
- Example:
  - CPI of 1.0 on a 500Mhz machine with a 5% miss rate, 200ns DRAM access
  - Adding 2nd level cache with 20ns access time decreases miss rate to 2%
- Using multilevel caches:
  - try and optimize the hit time on the 1st level cache
  - try and optimize the miss rate on the 2nd level cache

24

## Virtual Memory

- Main memory can act as a cache for the secondary storage (disk)

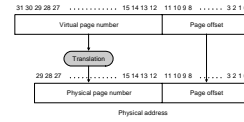


- Advantages:
  - illusion of having more physical memory
  - program relocation
  - protection

25

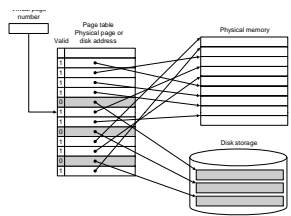
## Pages: virtual memory blocks

- Page faults: the data is not in memory, retrieve it from disk
  - huge miss penalty, thus pages should be fairly large (e.g., 4KB)
  - reducing page faults is important (LRU is worth the price)
  - can handle the faults in software instead of hardware
  - using write-through is too expensive so we use writeback



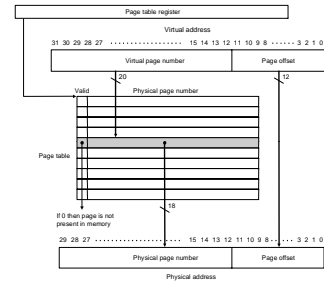
26

## Page Tables



27

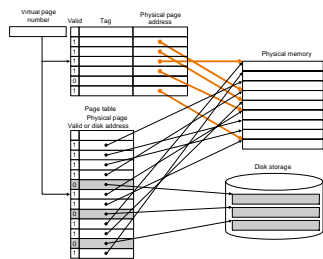
## Page Tables



28

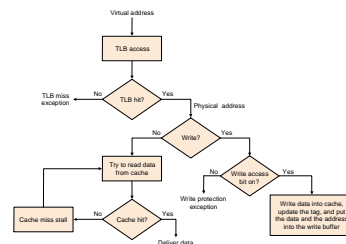
## Making Address Translation Fast

- A cache for address translations: translation lookaside buffer



29

## TLBs and Caches



30

## Replacement Policies

- Replacement Policies in Multi-way Set Associative caches
  - Random: Replace any line arbitrarily
  - Least Recently Used (LRU): Find the least recently used line to replace
  - Keep Most Recently Used (MRU): Keep the last used line in the set and replace any other randomly
- LRU performs the best
- MRU does equally well

31

## LRU Scheme

- We explain LRU with an example of a 4-way set associative cache
- Associate a 2-bit counter with each line (log k bit for k-way cache)
- Initially all lines are invalid
- For a miss bring a new line in an invalid line, make it valid, set its counter to zero, increment all other counters
  - If no invalid line, replace the line with counter value = 3, set its counter to zero, increment all other counters
- For a hit, set the accessed line's counter to zero and increment counters of those lines whose values is smaller than the accessed line
- Try this algorithm for an examples where lines read are 0, 64, 128, 64, 192, 256, 128, 0, 256, 192, 64...
  - There are 64 lines in each cache and it is 4-way set associative

32

## Reading or Writing a Memory word

- Check the address in TLB
- If not there, get the physical translation and also store the entry in TLB
  - Penalty 40-50 cycles
- If page itself is not present, page fault occurs
  - Read the page, update page table and TLB
  - Penalty 100's of thousands cycles
- Once physical address is there if there, perform read or write in cache
- If cache miss
  - Read the line in cache for read
  - May need to replace a dirty or clean line
    - Penalty 20-40 cycles
  - For Write read the line if write allocate, else write around
- If cache hit read or write in cache
  - Also write in main memory if write through

33

## A Big Example

- Instruction Frequency: LW(20%), SW(10%), R(50%), BR(15%), J(5%)
- Branch Penalty: 3 cycles on 20% mis-predictions =  $15 \cdot 0.20 \cdot 3 = 9$  cycles
- Data Cache 1: Miss rate 10% (of load/store), write back, write around, 50% dirty replacement, penalty for reading or writing a line 20 cycles
  - Load penalty =  $20 \cdot 0.10 \cdot 0.50 \cdot 20 + 20 \cdot 0.10 \cdot 0.50 \cdot (20+20) = 60$  cycles
  - Store Penalty = 0 (because of write around, otherwise will be 30)
- Data Cache 2: Miss rate 5% (of load/store), write back, write allocation, 50% dirty replacement, penalty for reading or writing a line 100 cycles
  - Load penalty =  $20 \cdot 0.05 \cdot 0.5 \cdot 100 + 20 \cdot 0.05 \cdot 0.5 \cdot (100+100) = 150$  cycles
  - Store Penalty =  $10 \cdot 0.05 \cdot 0.5 \cdot 100 + 10 \cdot 0.05 \cdot 0.5 \cdot (100+100) = 75$  cycles
- TLB: Miss Rate 2% (of load/store), Miss Penalty 100 cycles
  - Total Penalty =  $(20+10) \cdot 0.02 \cdot 100 = 60$  cycles
- Page faults: 0.01% (of load/store), Penalty 300,000 cycles
  - Total Penalty =  $(20+10) \cdot 0.0001 \cdot 300,000 = 900$  cycles
- Total Time =  $100+9+60+150+75+60+900 = 1354$  cycles, or CPI=13.54
- Notice that miss rates can be specified per instruction or per load/store

34

## Misses and Replacement Policies

- 3 C Misses
  - Compulsory: Miss will have to occur on first read (or write)
  - Capacity: A line is replaced and then brought back
  - Conflict: a miss occurs as some other line is occupying that line
- Example Suppose we read line a first time (no line is in cache), then read line b that replaces line a, and then read line a again
- The first and second misses are compulsory, second miss is also capacity and conflict, and the third miss is capacity (and also conflict)
- The terminology can be confusing here
  - The first read is always classified as compulsory
  - The replacement and read back is conflict if there was place in cache elsewhere but you had to bring it at that place due to mapping
  - If there was no place at all then it is capacity miss (like cache is full in a fully associative cache)

35

## Virtual Memory: Other Translation Schemes

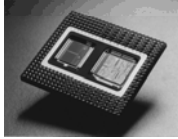
- In a single-level translation
  - 32 bit virtual address
  - 4KB Page size (12 bit address in each page)
  - Leaves 20-bit page address => 1 Million Pages => 4MB for Table
- One alternate is to only have a limited size page table with Hi and Lo Checks
  - But program use many addresses segments
- Alternate is to have a two level page table
- Divide page addresses in two parts of 10 bits each
  - There are 1K tables of 1K entries each (total is still 1M entries)
  - Most significant 10 bits points to a table (with 1K entries, each 4 bytes long, a total of 4KB that fits in a page) that contains the address of that part of table
  - Least significant 10 bits are used to access a particular entry in the selected table
- We only need to keep the first table (pointing to real tables) and some of the second level tables in memory

36

## Modern Systems

- **Very complicated memory systems:**

| Characteristic   | Intel Pentium Pro   | PowerPC 604   |
|------------------|---|---|
| Virtual address  | 32 bits   | 52 bits   |
| Physical address | 32 bits   | 32 bits   |
| Page size        | 4 KB-4 MB   | 4 KB, selectable, and 256 MB  |
| TLB organization | A TLB for instructions and a TLB for data<br>Both four-way set associative<br>Pseudo-LRU replacement<br>Instruction TLB: 32 entries<br>Data TLB: 64 entries<br>TLB misses handled in hardware | A TLB for instructions and a TLB for data<br>Both two-way set associative<br>LRU replacement<br>Instruction TLB: 128 entries<br>Data TLB: 128 entries<br>TLB misses handled in hardware |



| Characteristic      | Intel Pentium Pro                 | PowerPC 604                       |
|---------------------|-----------------------------------|-----------------------------------|
| Cache organization  | Split instruction and data caches | Split instruction and data caches |
| Cache size          | 8 KB each for instructions/data   | 16 KB each for instructions/data  |
| Cache associativity | Four-way set associative          | Four-way set associative          |
| Replacement         | Approximated LRU replacement      | LRU replacement                   |
| Block size          | 32 bytes                          | 32 bytes                          |
| Write policy        | Write-back                        | Write-back or write-through       |

37

## Some Issues

- Processor speeds continue to increase very fast
  - much faster than either DRAM or disk access times
- Design challenge: dealing with this growing disparity
- Trends:
  - synchronous SRAMs (provide a burst of data)
  - redesign DRAM chips to provide higher bandwidth or processing
  - restructure code to increase locality
  - use prefetching (make cache visible to ISA)

38