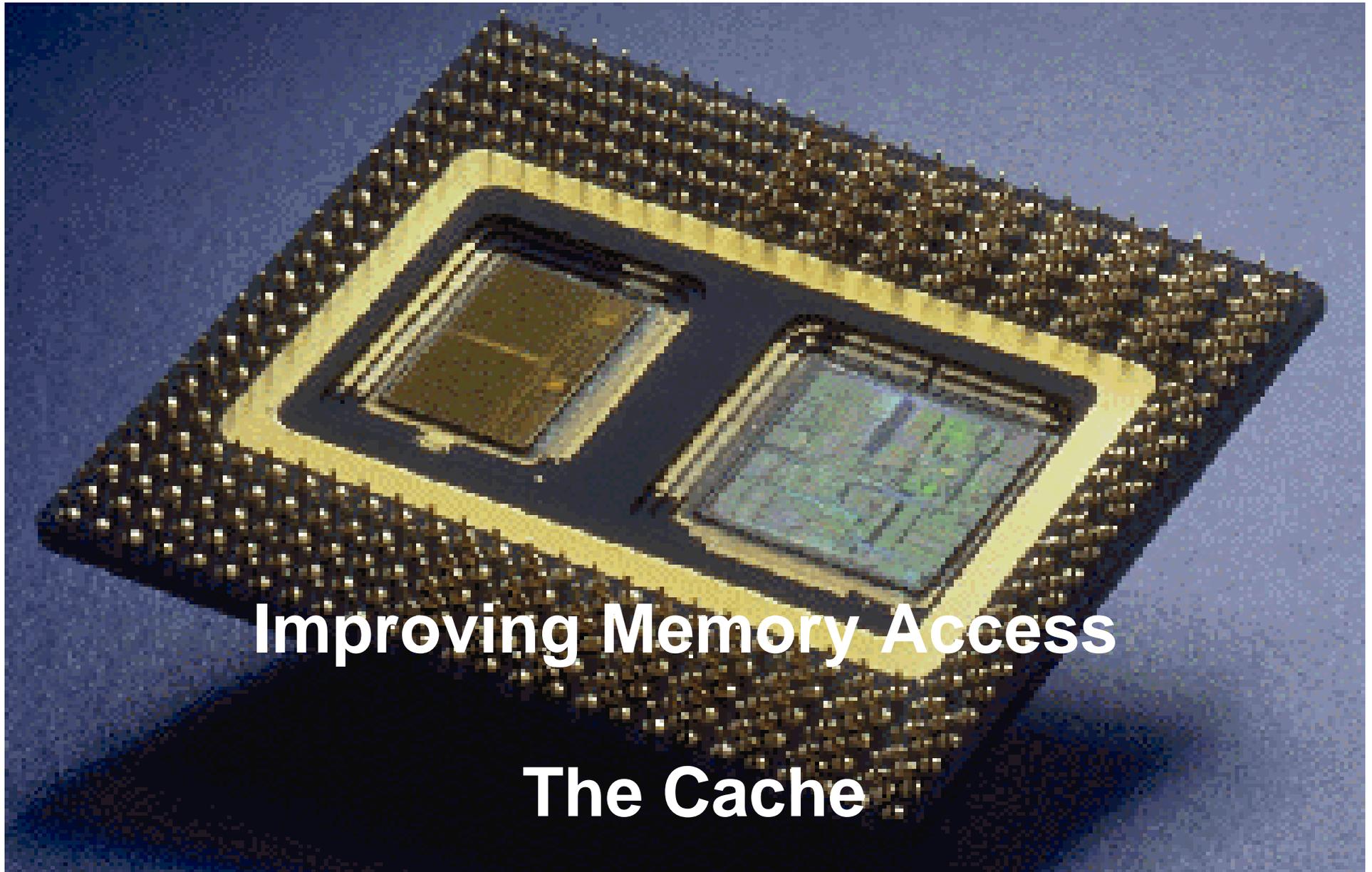


EECS 322 Computer Architecture



Improving Memory Access

The Cache

Pipelining and the cache (Designing...,M.J.Quinn, '87)



Instruction Pipelining is the use of pipelining to allow more than one instruction to be in some stage of execution at the same time.

Ferranti ATLAS (1963):

- Pipelining reduced the average time per instruction by 375%
- Memory could not keep up with the CPU, needed a cache.

Cache memory is a small, fast memory unit used as a buffer between a processor and primary memory

Principle of Locality



- **Principle of Locality**

states that programs access a relatively small portion of their address space at any instance of time

- Two types of locality

- Temporal locality (locality in time)

If an item is referenced, then

the same item will tend to be referenced soon

“the tendency to reuse recently accessed data items”

- Spatial locality (locality in space)

If an item is referenced, then

nearby items will be referenced soon

“the tendency to reference nearby data items”

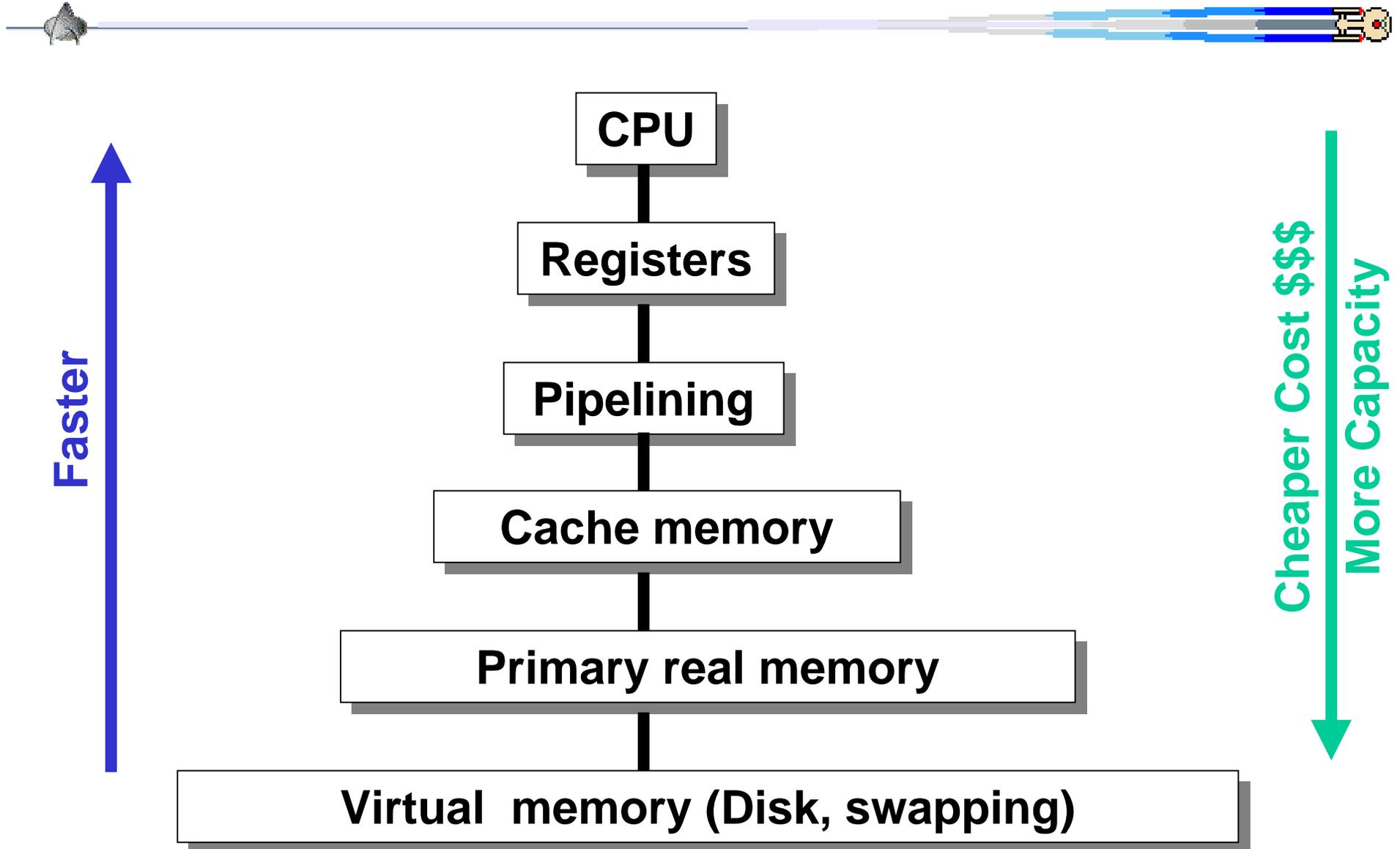
Memories Technology and Principle of Locality



- **Faster Memories are more expensive per bit**
- **Slower Memories are usually smaller in area size per bit**

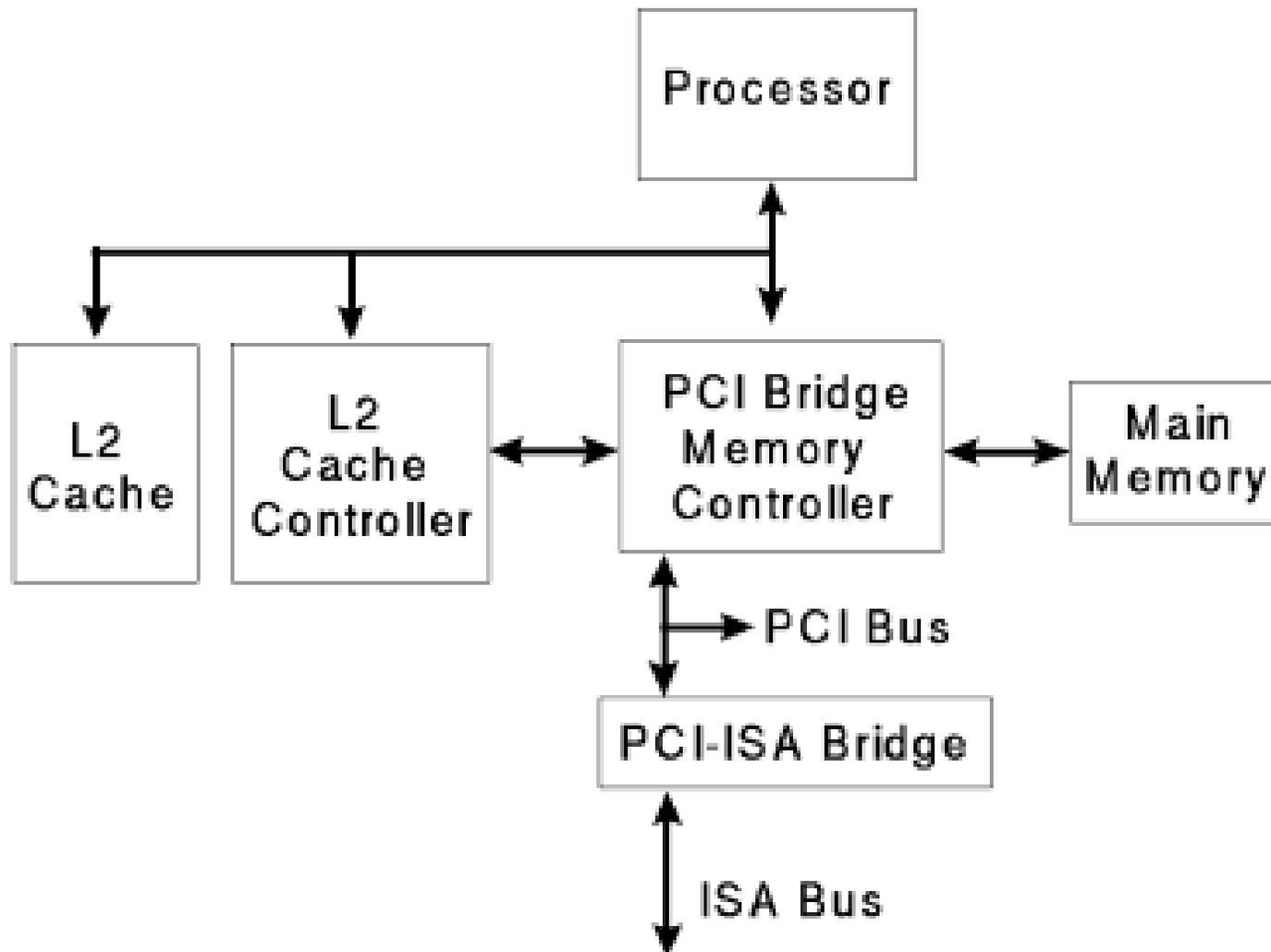
Memory Technology	Typical access time	\$ per Mbyte in 1997
SRAM	5-25 ns	\$100-\$250
DRAM	60-120 ns	\$5-\$10
Magnetic Disk	10-20 million ns	\$0.10-\$0.20

Memory Hierarchy



Basic Cache System

Figure 1. Basic Cache System



Cache Terminology



A hit if the data requested by the CPU is in the upper level

Hit rate or **Hit ratio**

is the fraction of accesses found in the upper level

Hit time

is the time required to access data in the upper level
= <detection time for hit or miss> + <hit access time>

A miss if the data is not found in the upper level

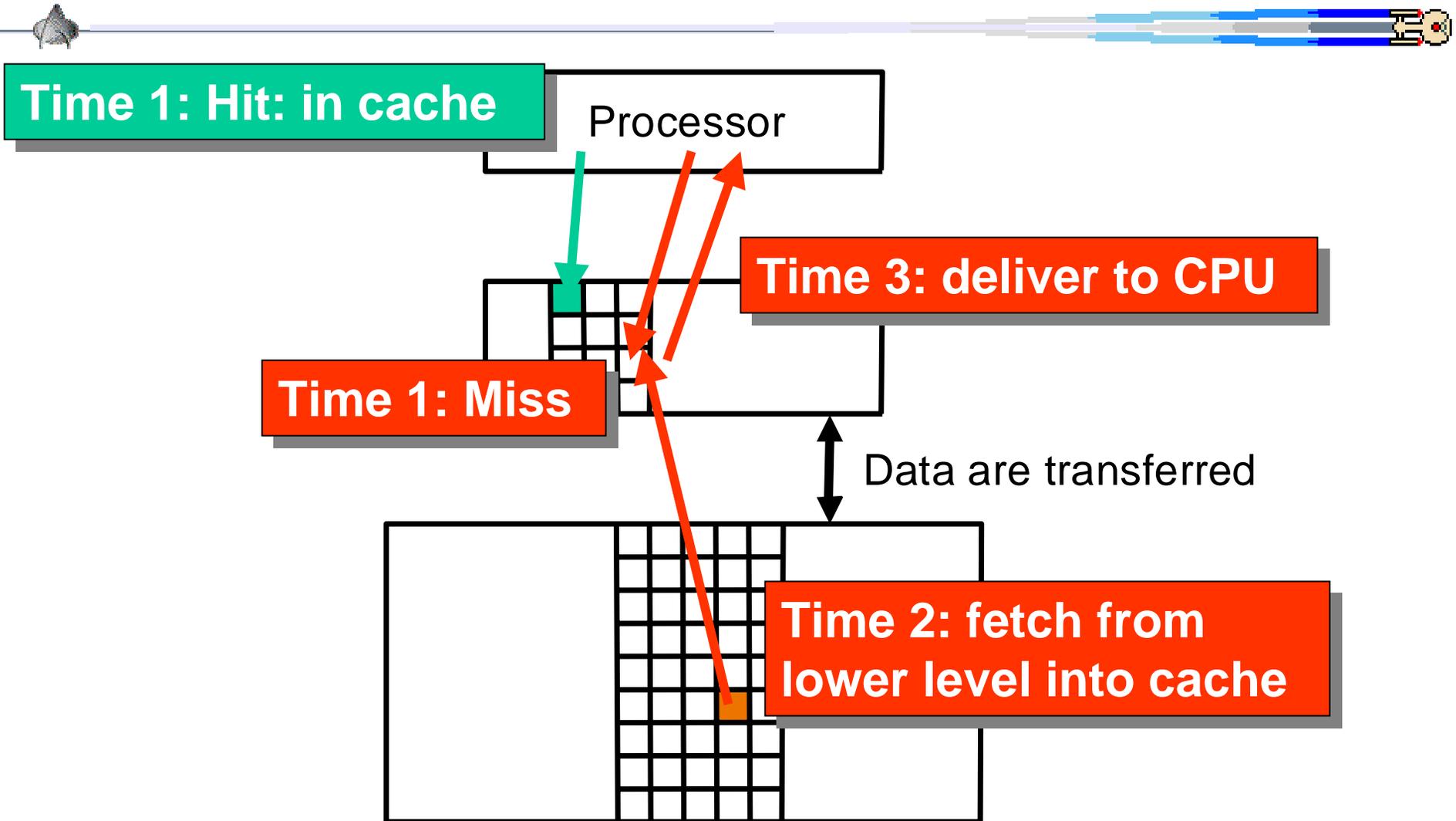
Miss rate or **(1 – hit rate)**

is the fraction of accesses not found in the upper level

Miss penalty

is the time required to access data in the lower level
= <lower access time> + <reload processor time>

Cache Example



Hit time = Time 1

Miss penalty = Time 2 + Time 3

Cache Memory Technology: SRAM

see page B-27



- **Why use SRAM (Static Random Access Memory)?**

- **Speed.**

The primary advantage of an SRAM over DRAM is speed.

The fastest DRAMs on the market still require **5 to 10 processor clock cycles** to access the first bit of data.

SRAMs can operate at processor speeds of 250 MHz and beyond, with access and **cycle times equal to the clock cycle** used by the microprocessor

- **Density.**

when 64 Mb DRAMs are rolling off the production lines, the largest SRAMs are expected to be only 16 Mb.

see reference: <http://www.chips.ibm.com/products/memory/sramoperations/sramop.html>

Cache Memory Technology: SRAM

(con't)



- **Volatility.**

Unlike DRAMs, SRAM cells do not need to be refreshed.
SRAMs are available 100% of the time for reading & writing.

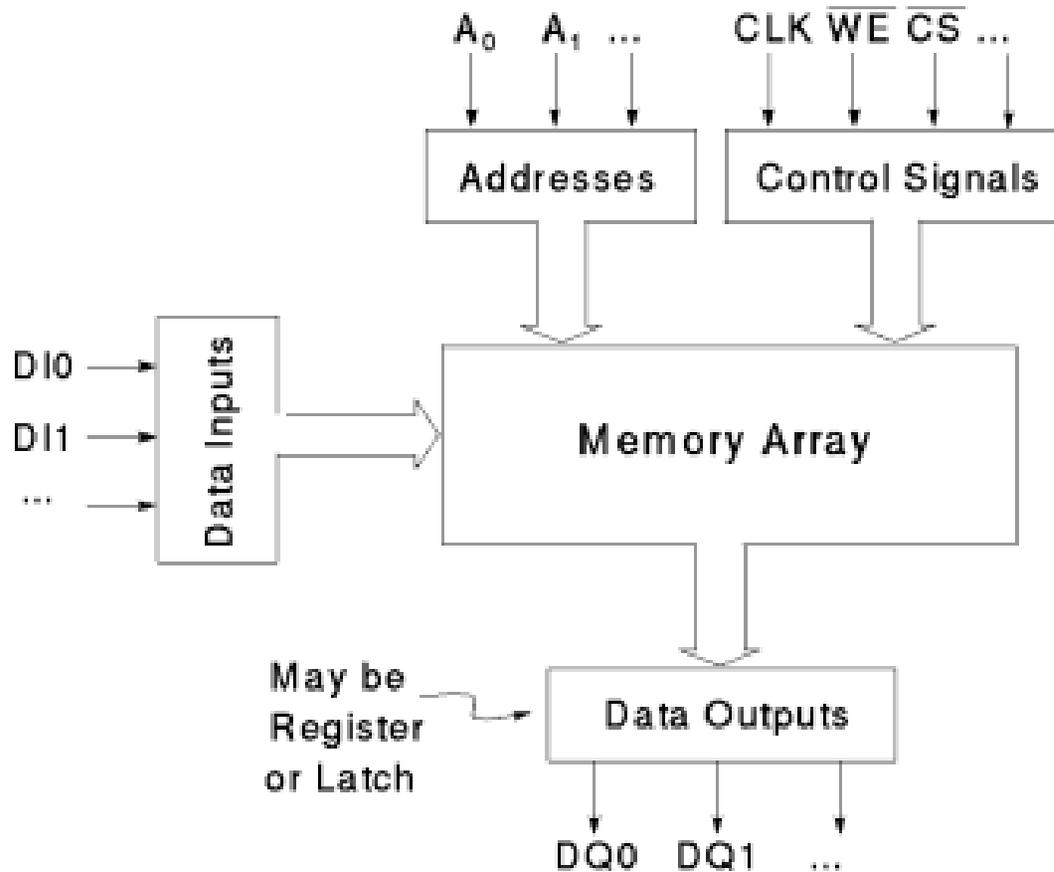
- **Cost.**

If cost is the primary factor in a memory design,
then DRAMs win hands down.

If, on the other hand, performance is a critical factor,
then a well-designed SRAM is an effective cost
performance solution.

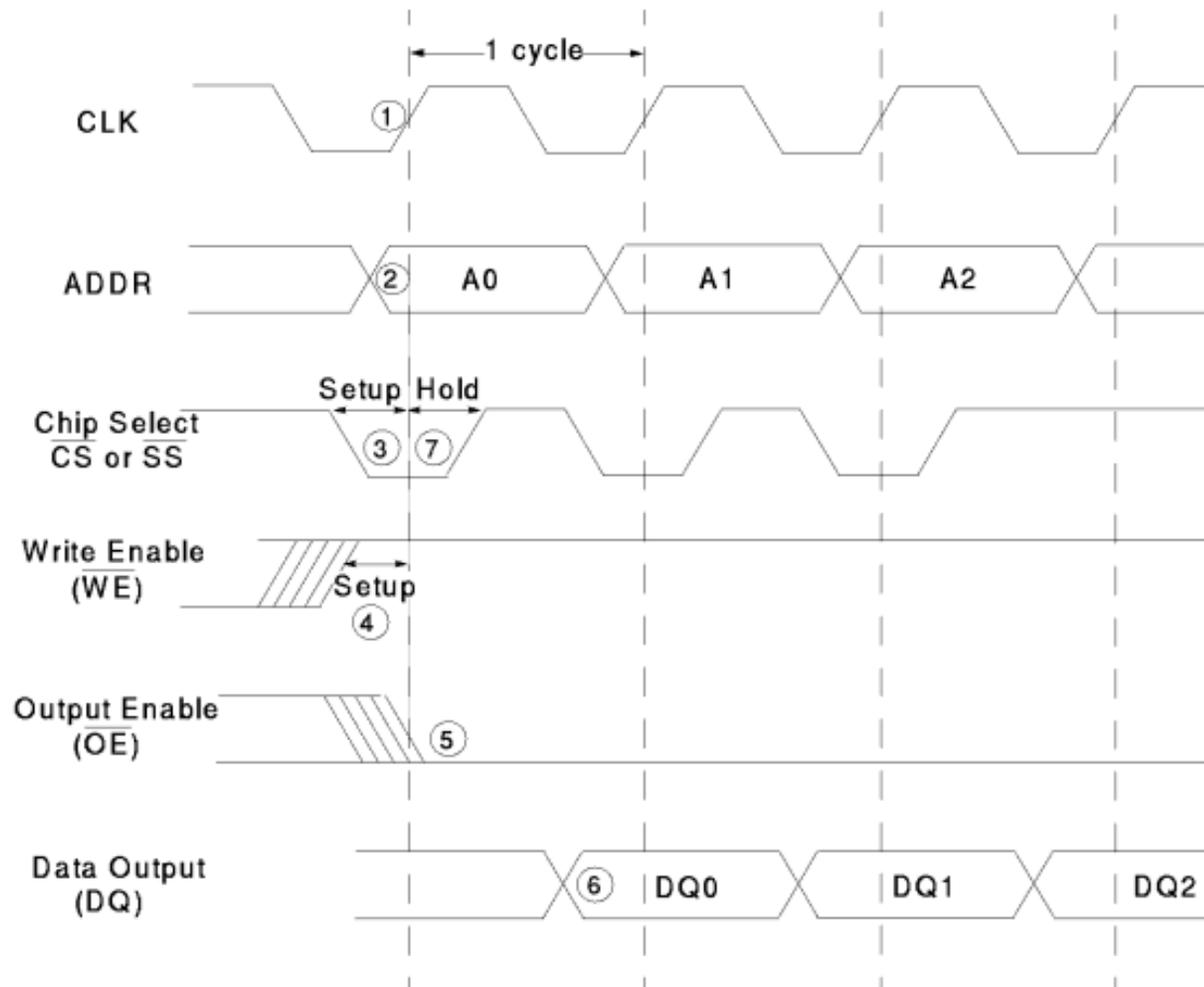
Cache Memory Technology: SRAM Block diagram

Figure 2. Simplified Block Diagram of a Synchronous SRAM



Cache Memory Technology: SRAM timing diagram

Figure 4. Reading from Memory (Flow Thru mode)

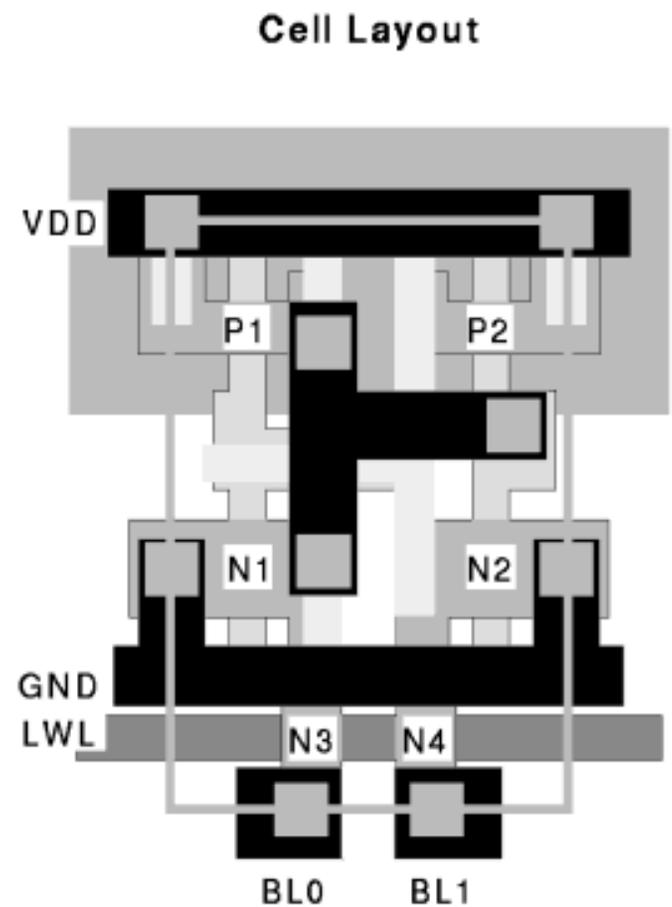
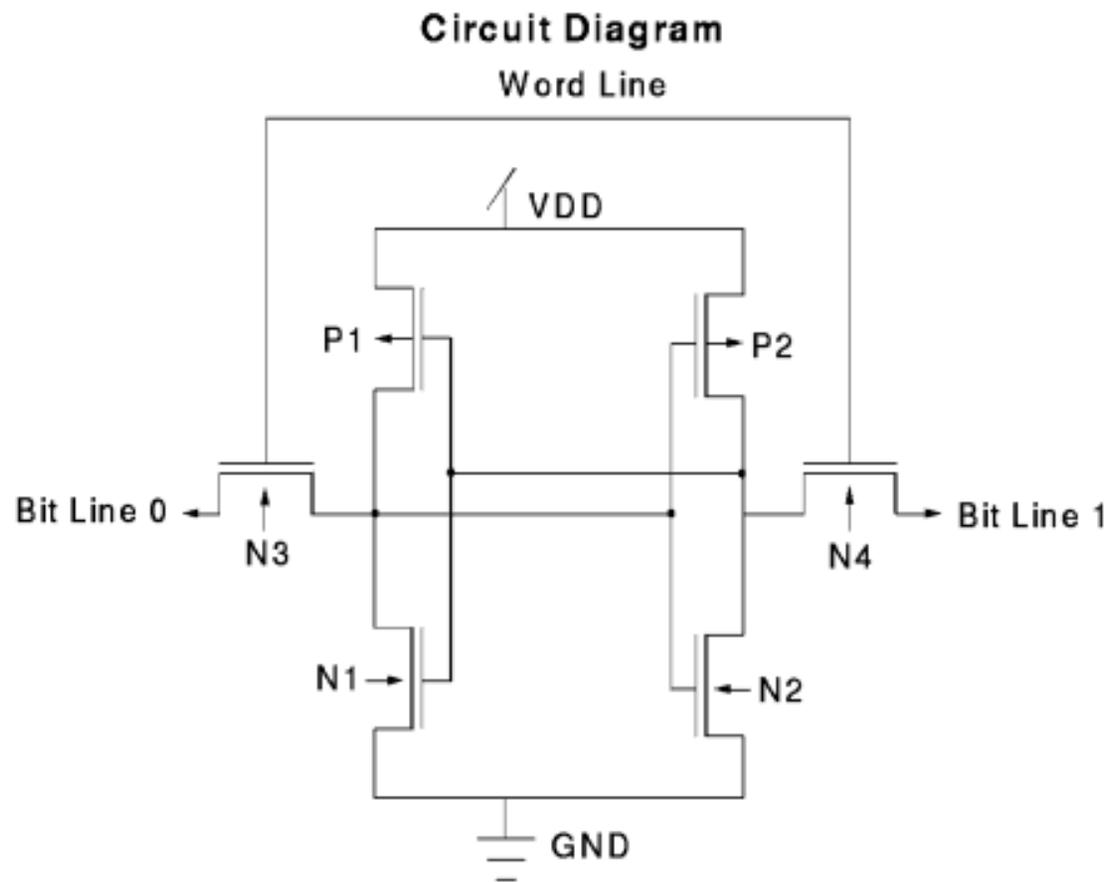


Note: DQ0 is the data associated with Address 0 (A0). DQ1 is the data associated with Address 1 (A1).

Cache Memory Technology: SRAM 1 bit cell layout

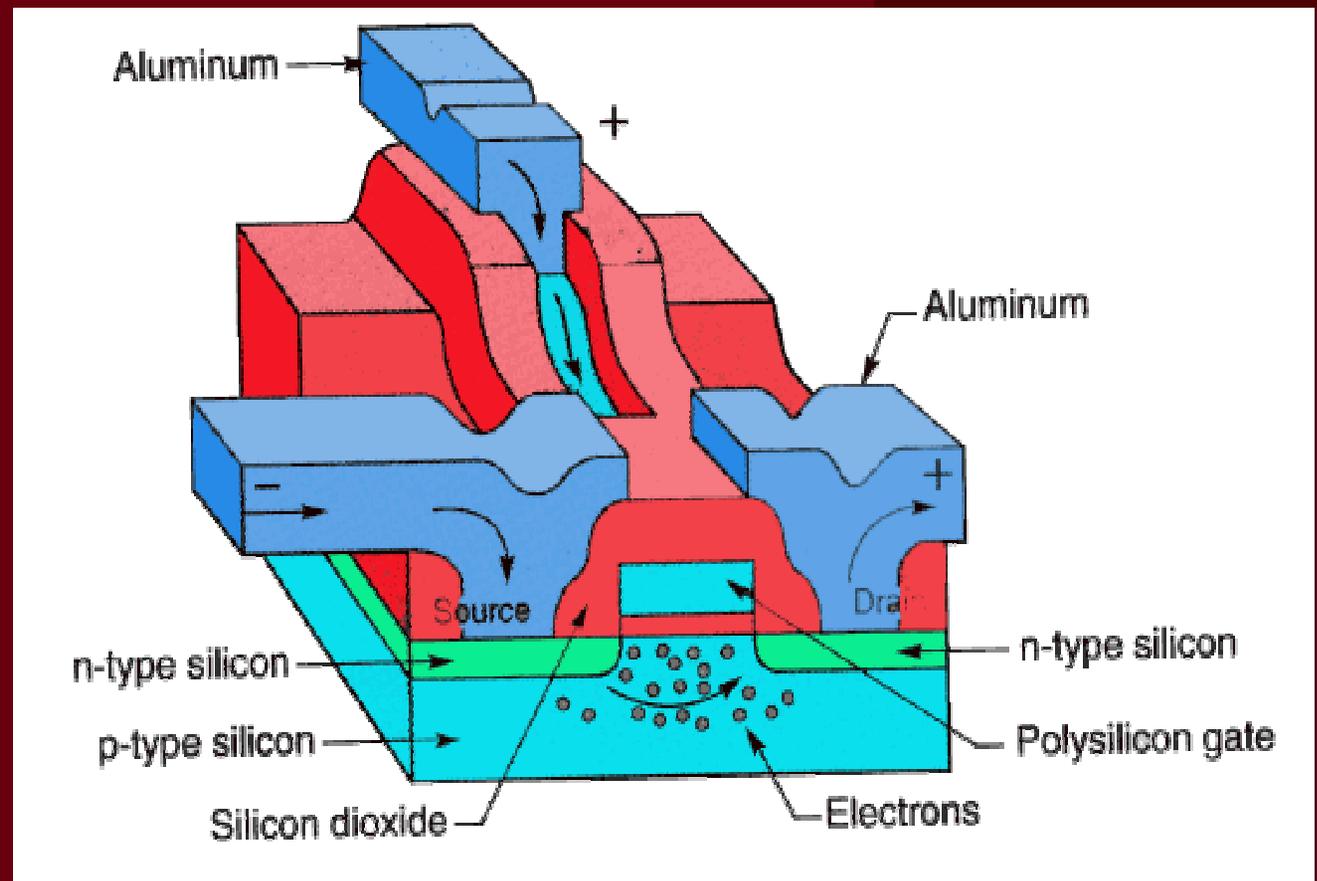


Figure 3. IBM's 6-Transistor Memory Cell



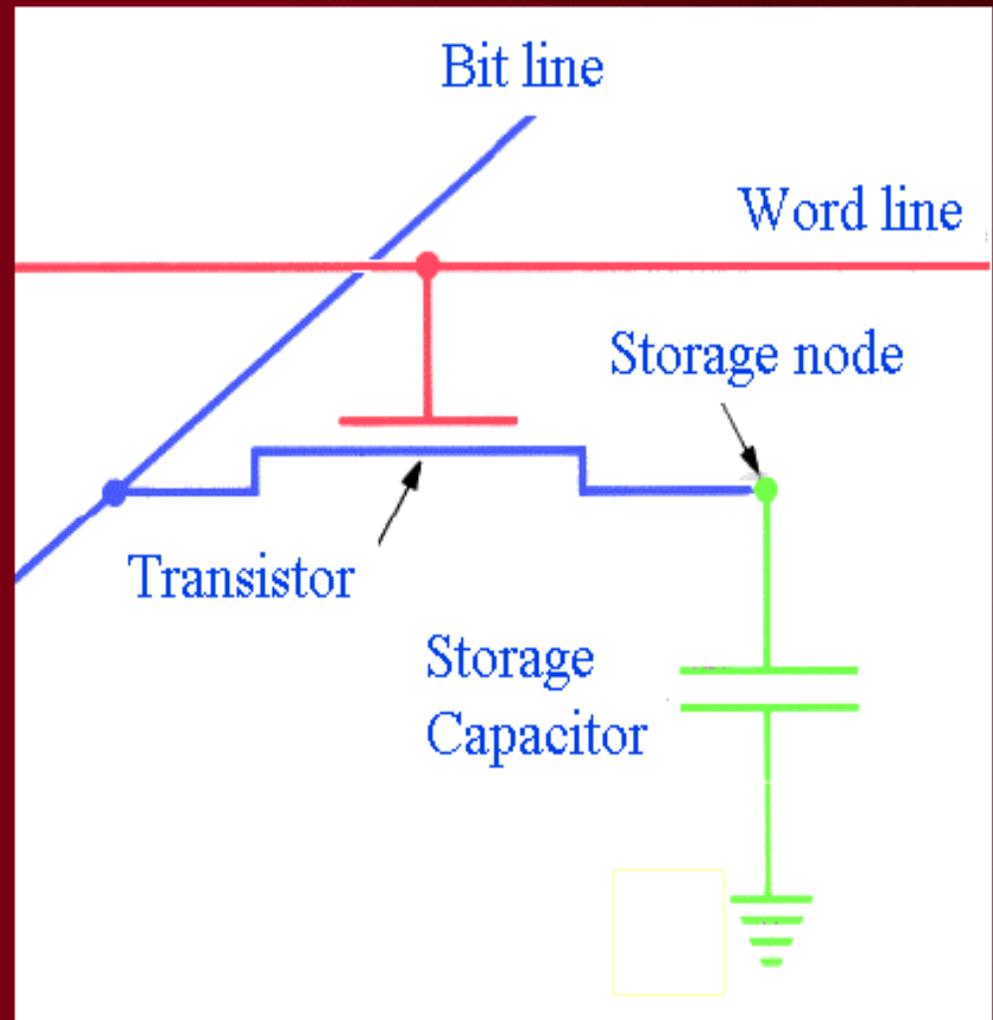
Real transistor

- 3-D structure
- Real materials



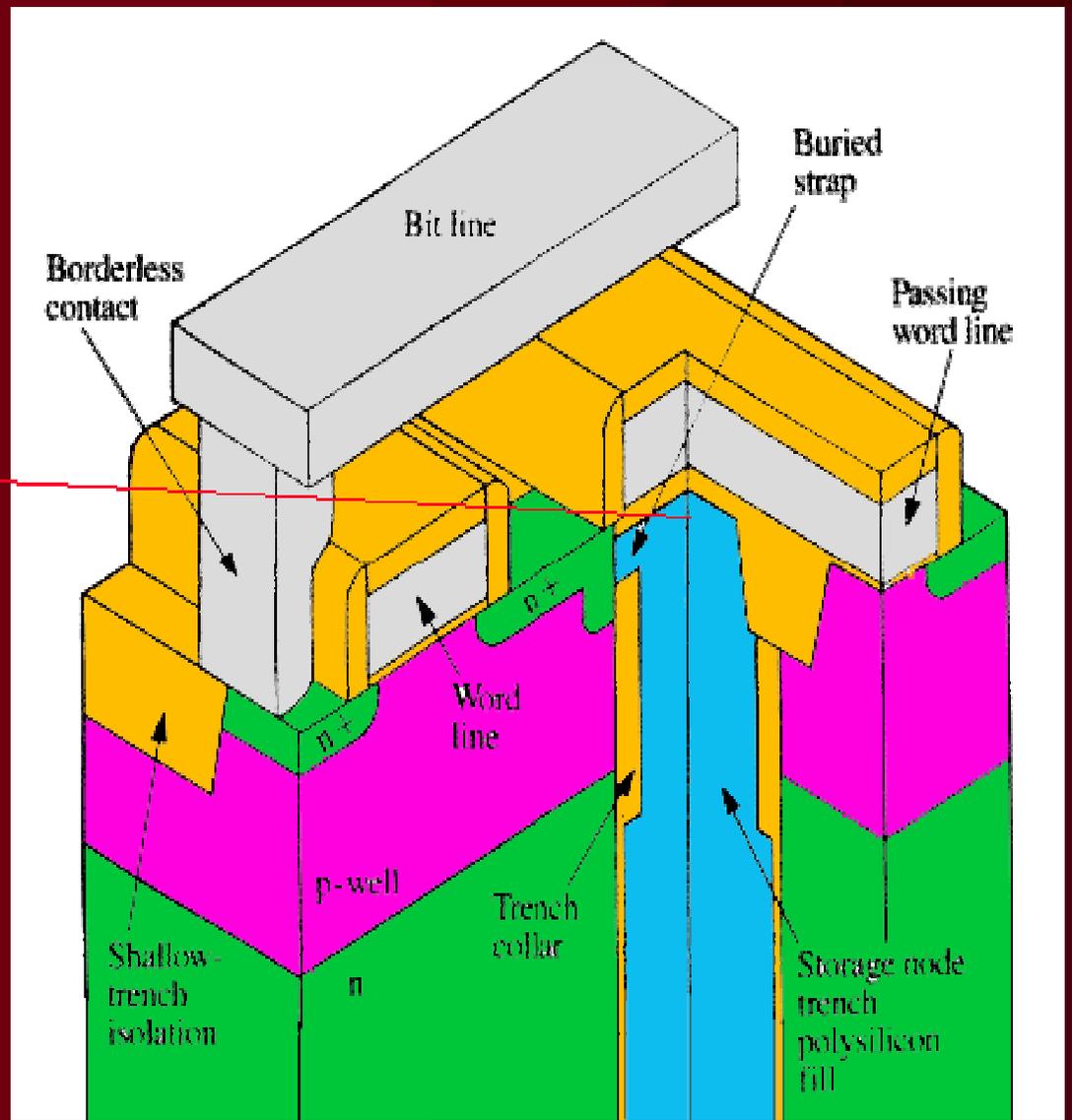
Basic DRAM design

- DRAM replaces all but one transistors of flip-flop with a capacitor
- => smaller!
- Capacitor stores information
- Charge leakage requires periodic refreshment (sense & rewrite)



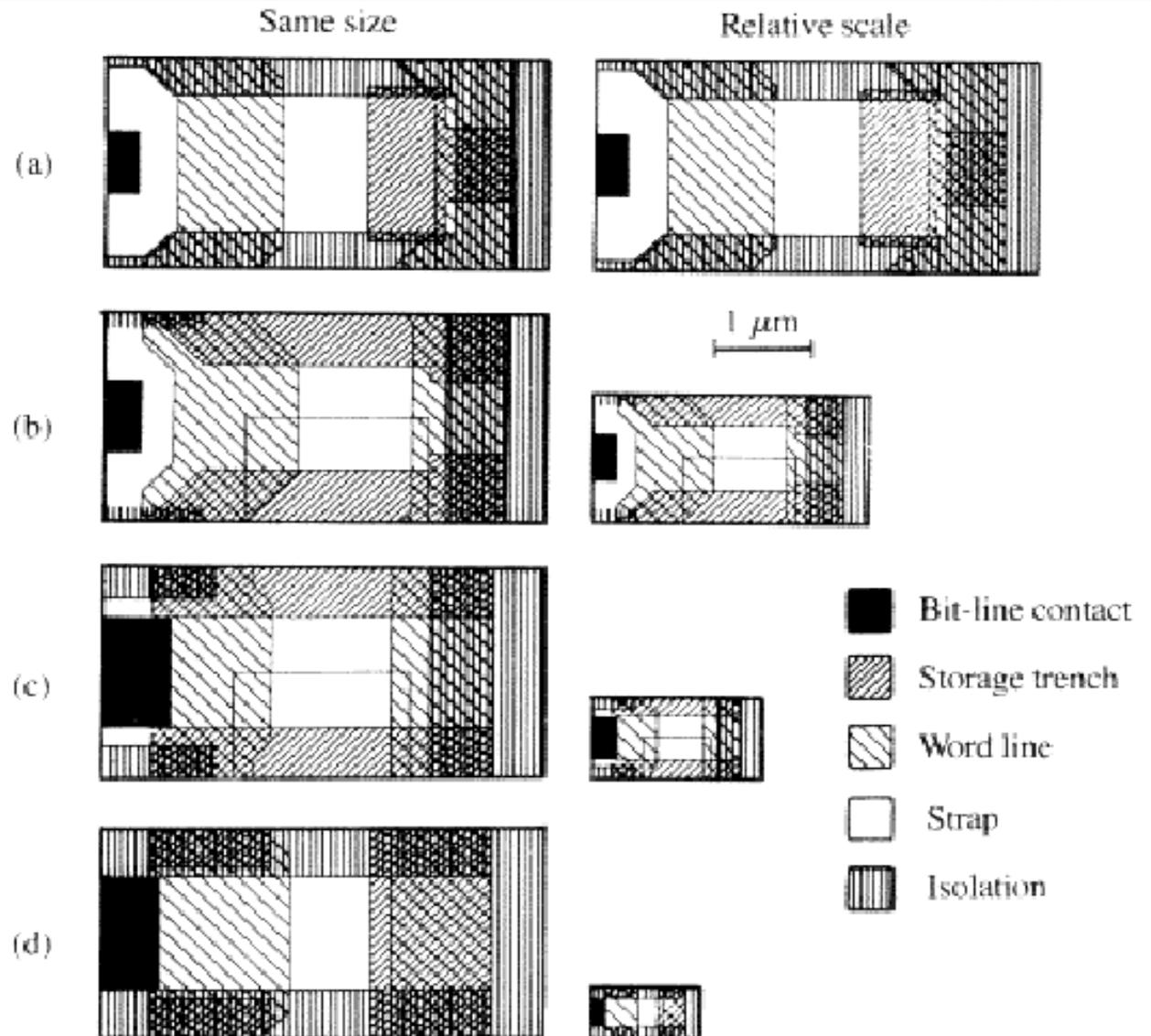
256Mb DRAM

- Increased vertical integration
- Word line passes over capacitor and contact
- Cell area $\sim 0.5\mu\text{m}^2$
- Capacitor area smaller - dielectric must be thinner
- => higher quality dielectric required

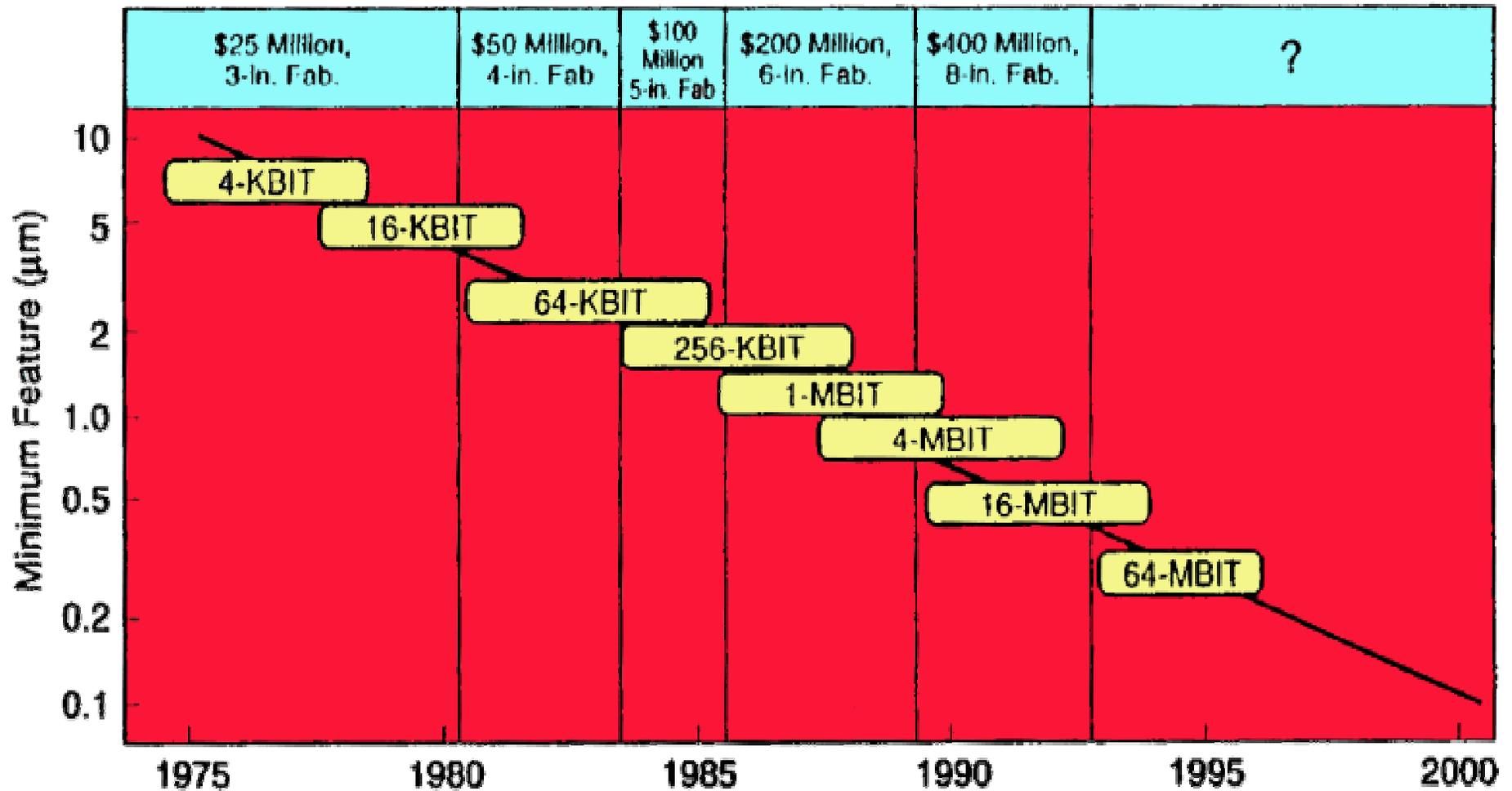


Memory Technology: DRAM Evolution

DRAM evolution (II)



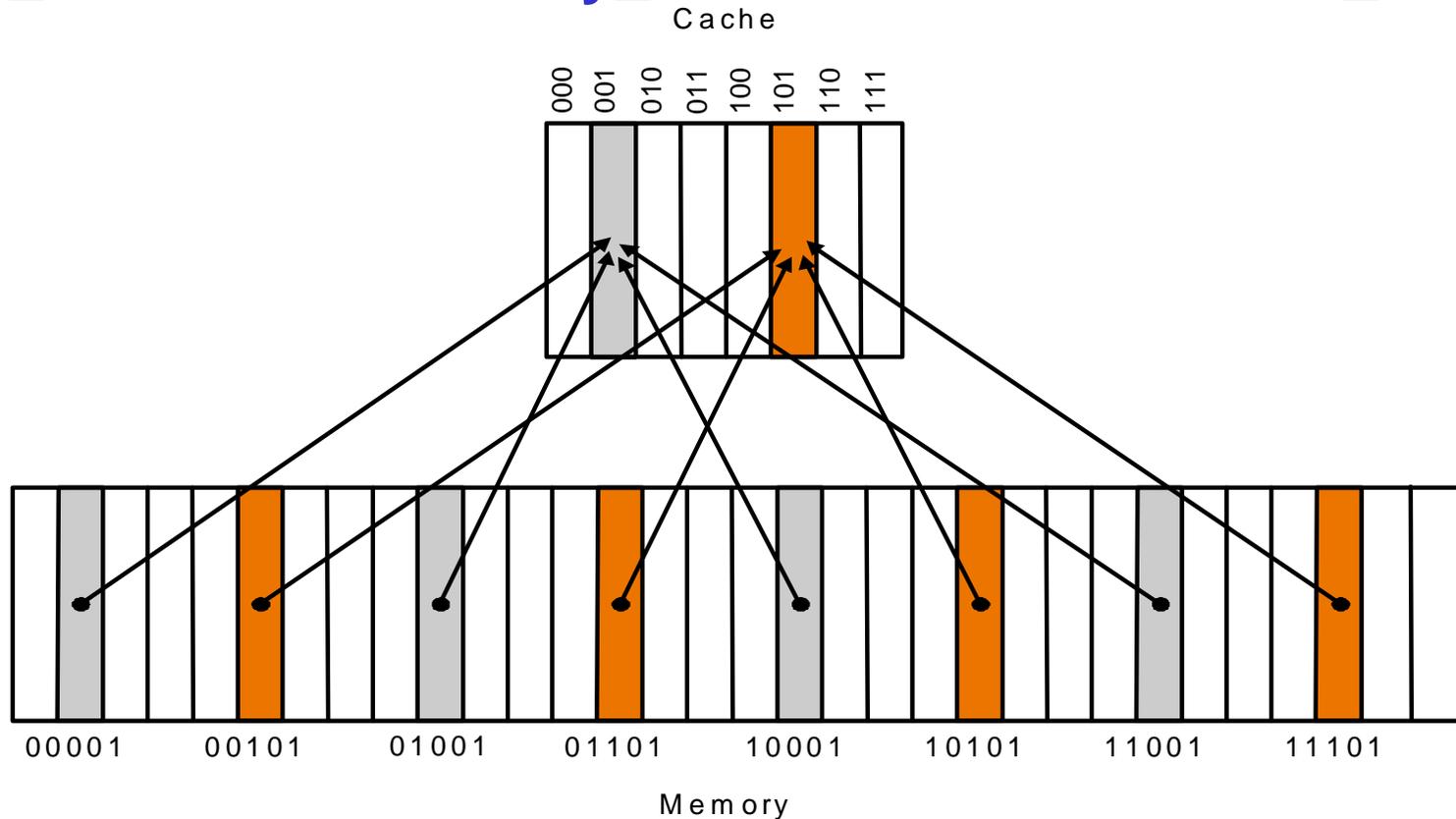
DRAM development



Direct Mapped Cache



- **Direct Mapped:** assign the cache location based on the address of the word in memory
- **cache_address = memory_address modulo cache_size;**



Observe there is a Many-to-1 memory to cache relationship

Direct Mapped Cache: Data Structure



There is a **Many-to-1 relationship** between memory and cache

How do we know whether the data in the cache corresponds to the requested word?

tags

- contain the address information required **to identify** whether a word in the cache corresponds to the requested word.
- tags need only to contain the **upper portion** of the memory address (often referred to as a **page address**)

valid bit

- indicates whether an entry contains a valid address

Figure 7.6

Direct Mapped Cache: Temporal Example



<p>lw \$1,10 110 (\$0)</p> <p>lw \$2,11 010 (\$0)</p> <p>lw \$3,10 110 (\$0)</p>	<div style="background-color: red; color: white; padding: 5px; margin-bottom: 5px;">Miss: valid</div> <div style="background-color: red; color: white; padding: 5px; margin-bottom: 5px;">Miss: valid</div> <div style="background-color: green; color: white; padding: 5px;">Hit!</div>	<p>lw \$1,22(\$0)</p> <p>lw \$2,26(\$0)</p> <p>lw \$3,22(\$0)</p>
--	--	---

Index	Valid	Tag	Data
000	N		
001	N		
010	Y	11	Memory[11010]
011	N		
100	N		
101	N		
110	Y	10	Memory[10110]
111	N		



Figure 7.6

Direct Mapped Cache: Worst case, always miss!

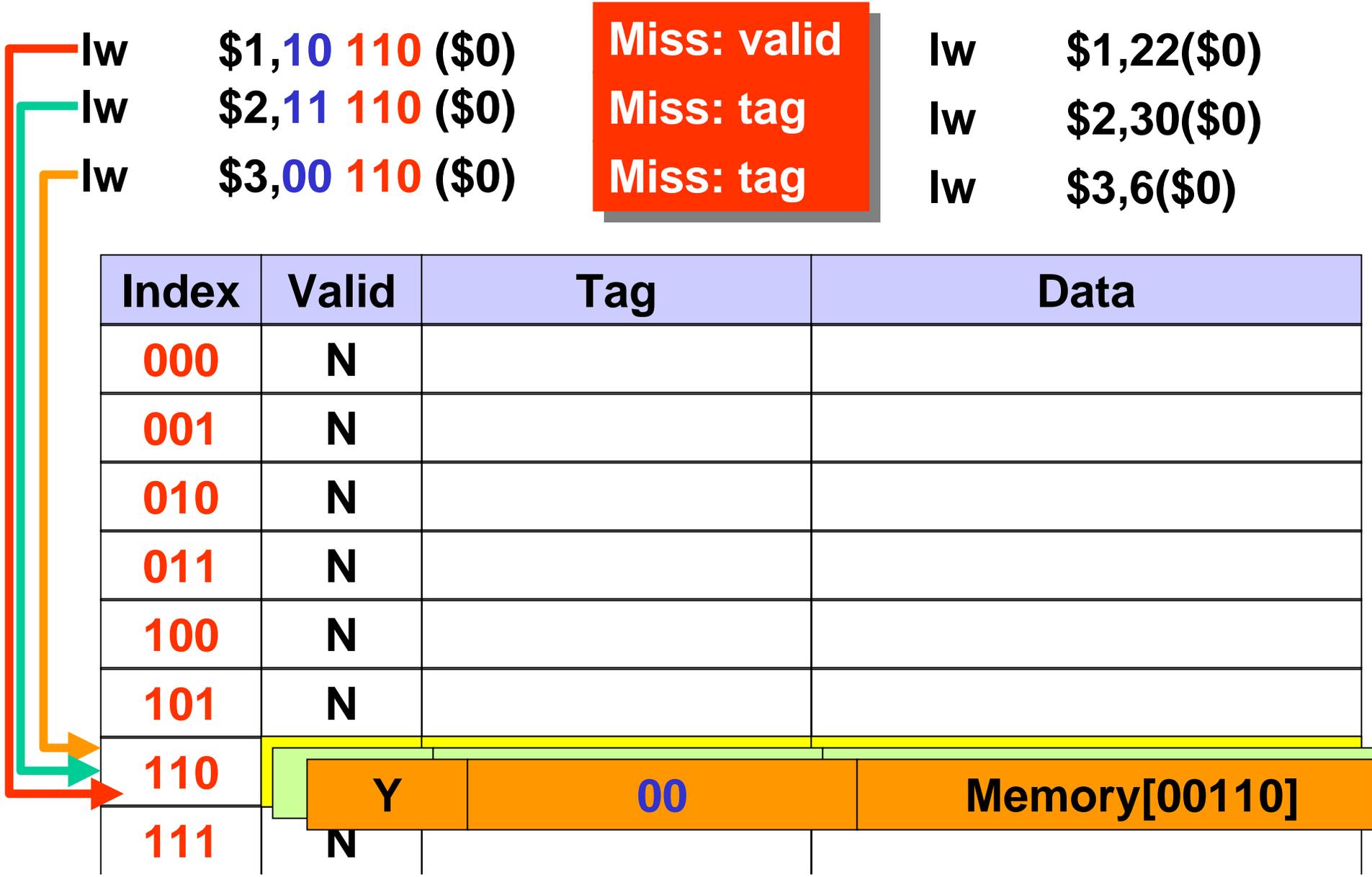
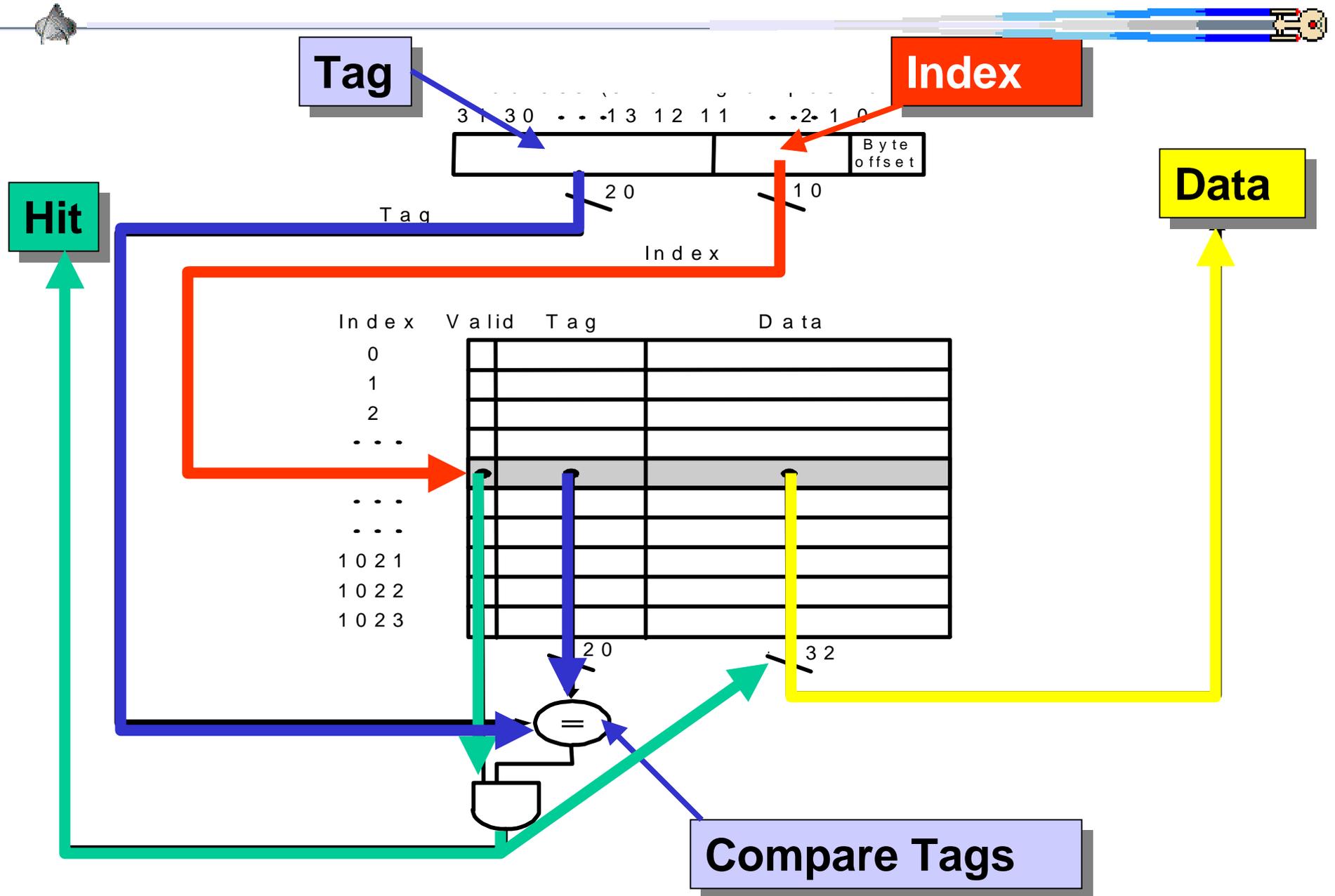
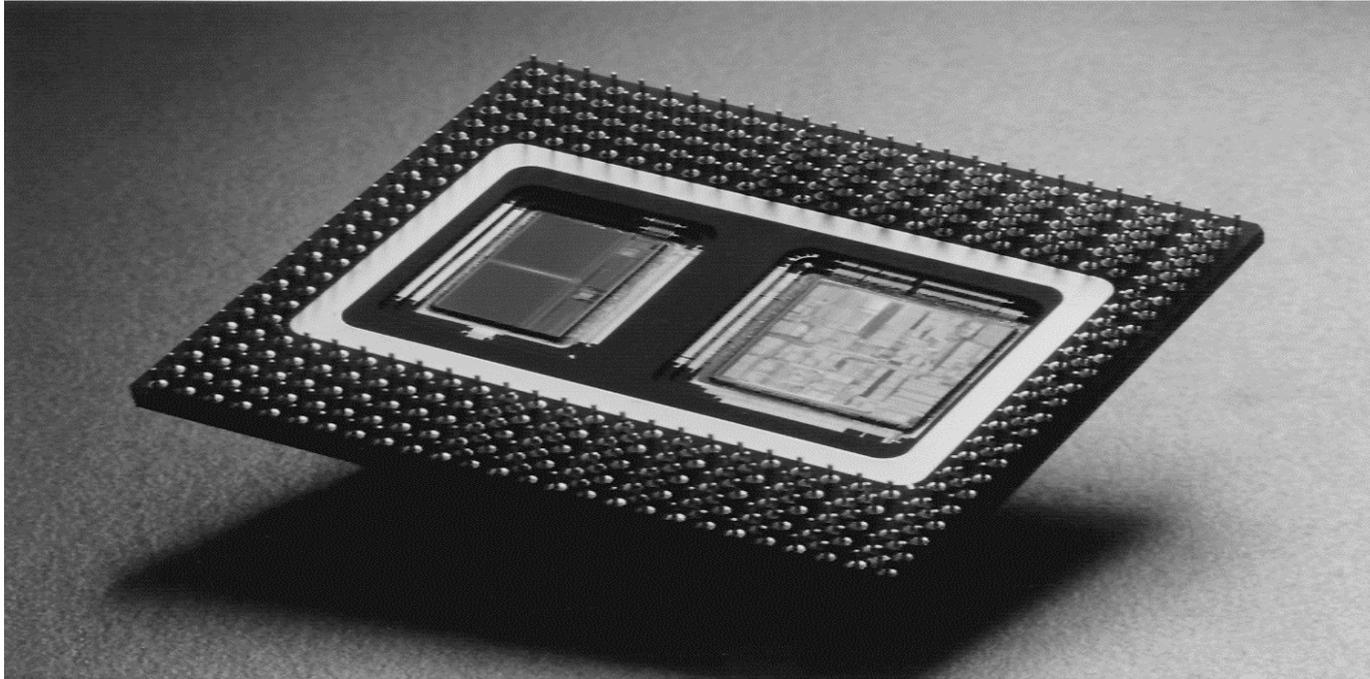


Figure 7.7

Direct Mapped Cache: Mips Architecture



Modern Systems: Pentium Pro and PowerPC



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