

Computer Organization & Architecture

Unit - 4

Memory hierarchy - memory technologies - cache memory - measuring & improving cache performance

* Memory Hierarchy

- consists of multiple levels of memory with different speeds & sizes
- Faster memories are smaller and more expensive per bit

Speed	Processor	Size	Cost (/bit)	Current Technology
Fastest	<u>Registers</u>	smallest	highest	SRAM
	<u>Memory</u>			DRAM
Slowest	<u>Memory</u>	biggest	lowest	magnetic disk

- At the top of the hierarchy are registers that are matched in speed to the CPU.
- There are a small no. of registers in a processor, of the order of a few hundred or less.
- At the bottom of the hierarchy are the secondary memories such as magnetic disks / tapes where the cost per stored bit is small in terms of money and electrical power, but the access time is very long compared to registers.
- As we move up the hierarchy, greater performance is realized at a greater cost.

* Terminology

1. Block : The minimum unit of information that can be either present or not present in the two-level hierarchy is called a block or a line.
2. Hit : If the data requested by the processor appears in some block in the upper level, it is called a hit.
3. Miss : If the data is not found in the upper level, the request is called a miss. The lower level in the hierarchy is then accessed to retrieve the block containing the requested data.
4. Hitrate / hit ratio: fraction of memory accesses found in the upper level, used to measure the performance of the memory hierarchy.
5. missrate: $1 - \text{hitrate}$ = fraction of memory accesses not found in the upper level
6. Hit Time : Time to access the upper level of the memory hierarchy, which includes the time needed to determine whether the access is a hit or a miss.
7. Miss Penalty : Time to replace a block in the upper level with the corresponding block from the lower level, plus the time to deliver this to the processor.

* Principle of Locality

→ states that programs access a relatively small portion of their address space at any instant of time. There are 2 types of locality:

- (a) Temporal Locality - The principle stating that if a data location is referenced, it is likely to be referenced again soon.
- (b) Spatial Locality - The locality principle stating that if a data location is referenced, data locations with nearby addresses will likely be referenced soon.

* Memory Technology

There are 4 primary technologies used today in memory hierarchies:

- (i) DRAM (Dynamic RAM) \Rightarrow for main memory
- (ii) SRAM (Static RAM) \Rightarrow for caches
- (iii) Flash Memory \Rightarrow non-volatile, secondary memory in mobile device
- (iv) Magnetic Disk \Rightarrow largest & slowest

A. SRAM Technology

- \rightarrow are integrated circuits that are memory arrays, with a single access port that can provide either a read or write
- \rightarrow SRAMs have fixed access times to any data.
- \rightarrow do not need to refresh \Rightarrow access time is close to cycle time
- \rightarrow as long as there is power, the value can be indefinitely retained.

B. DRAM Technology

- \rightarrow The value kept in a cell is stored as a charge in a capacitor, a single transistor is used to access this stored charge (hence, ^{capacitor} cache).

- As DRAMs use the charge on a capacitor, it cannot be kept indefinitely, and must periodically be refreshed. (hence called dynamic)
- Refreshing → reading the contents & writing it back
- Some DRAMs have clocks to synchronize w/ the processor - called Synchronous DRAMs / SDRAM.
adv = eliminates time for memory & processor to synchronize
- Also a version called Double Data Rate (DDR) SDRAM ⇒ data transfers on both the rising & falling edge of the clock, meaning there is twice the bandwidth.

C. Flash Memory

- a type of electrically erasable programmable read only memory (EEPROM)
- used in mobile devices
- wear leveling : spread writes by remapping blocks that have been written many times to less trodden blocks.

D. Disk Memory

- a collection of platters that rotate on a spindle.
- have magnetic recording material on both sides.
- has a movable arm w/ a small electromagnetic coil
- tracks - disk surface divided into concentric circles
- sectors : track divided into sectors that contain the information

Seek Time: time taken to move head to desired track

Rotational Latency: time taken for the desired sector to rotate under the read / write head , at the correct track.
called rotational latency.

*Cache Memory

- represent the level of memory hierarchy between the processor and the main memory
 - The fundamental idea of cache organization is that by keeping the most frequently accessed instruction and data in the fast cache memory, the average memory access time will be reduced.
 - Typical mapping between address and cache locations for a cache is computed as:
- $$\text{block} = \text{block address} \% \text{ no. of blocks in the cache}$$
- If the no. of entries in the cache is a power of 2, then the modulo can be computed by using the low-order $\log_2(\text{cache size in blocks})$

*Accessing a cache

- Each cache location can contain the contents of a number of different memory locations
- A tag is used to identify whether the word in the cache corresponds to the requested word.
- The tag needs only to have the upper portion of the address, corresponding to the bits that are not used to index the cache
(The lower 3 bit index field of the address selects the block)

→ A valid bit is used in the tables to indicate that the associated block in the hierarchy contains valid data.

Calculating Tag Size & Cache Size

If cache block size = k

$$2^n = k$$

| n = no. of bits for the index

block size = 2^m words

| m bits for each word

if the address length is A (say 32)

the no. of bits in the tag is $A - (n+m+2)$

More Formulae

* Cache Mapping Techniques

Direct Mapping

→ A cache structure in which each memory location is mapped to exactly one location in the cache.

→ calculated as:

$$(\text{Block address}) \bmod (\text{no. of blocks in cache})$$

→ Contention may occur.

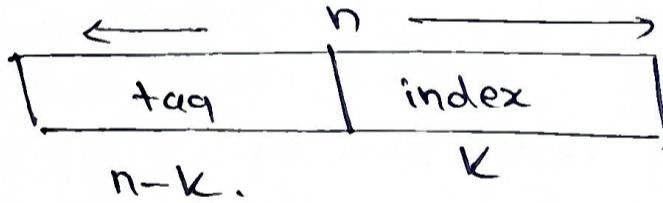
(i) when the cache is full

(ii) when more than one memory block is mapped to a given cache block position.

Solution: new blocks override the current resident

In general: 2^k words in cache memory

2^n words in main memory



∴ Each word in the cache consists of the data word & the associated tag.

Advantages → easy to implement

Disadvantages → not very flexible

→ hit ratio can drop considerably if 2 or more words whose address has the same index, but diff. tags are repeatedly accessed.

* Fully Associative

- a block can be placed in any location in the cache
 - To find a given block in a fully associative cache, all entries have to be searched because a block can be placed in any one.
- Advantage → more flexible than direct mapping
- Disadvantage → increases hardware cost
→ only suitable for caches w/ a small no. of blocks

* Set Associative

- there are a fixed no. of locations where each block can be placed.
- there are a fixed no. of locations where each block can be placed. An n-way set associative cache with n locations for a block is called n-way set associative.
- An n-way set associative cache consists of a no. of sets each of which has n-blocks.

A set containing a memory block is given by

(Block no) modulo (no. of sets in cache)

Numericals

- ① A digital computer has a memory unit of $64K \times 16$ and a cache memory of 1K words. The cache uses direct mapping with a block size of four words
- How many bits are there in the tag, index, block & word field of the address?
 - How many bits are there in each word of the cache, and how are they divided into? Include a valid bit
 - How many blocks can be accommodated?

(9)

$$\text{Main memory} = 64K \times 16$$

$= 64K$ words of 16 bits each

$$\text{Main memory size} = 64K \text{ words}$$

$$= 2^6 \times 2^{10}$$

$$= 2^{16} \text{ words}$$

Address = 16 bits

word size = 16 bits each

Data = 16 bits

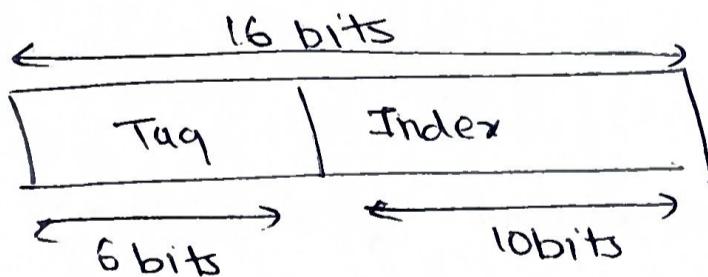
$$\text{Cache memory} = 1K$$

(length of data = 16 bits,

due to direct mapping)

$$= 2^{10} \text{ words}$$

Index = 10 bits



No. of bits in a block = ?

$$\text{Total no. of blocks} = 1K = 1024$$

size of a block = 4 words

$$\text{No. of blocks} = \frac{1024}{4} = 256 = 2^8$$

no. of bits in a block = 8

$$\text{No. of words in a block} = 4 \\ = 2^2$$

| no. of bits to access word = 2

b. Total size

$$= \text{valid bit} + \text{tag} + \text{address}$$

$$= 1 + 6 + 16$$

= 23 bits in each word of cache memory

c. cache can accommodate 256 blocks.

Q) A two-way set associative cache memory uses blocks of 4 words. The cache can accommodate a total of 2048 words from the main memory. The main memory size is 128K x 32

a. Formulate all the pertinent information required to construct the cache memory. (Tag, index, data, blocks, words)

b. What is the size of the cache memory?

Ans: Main memory:-

$$128K \times 32 = 128 \times \text{words of } 32 \text{ bits each}$$

$$\text{address} = 128K$$

$$= 2^7 \times 2^{10}$$

$$= 2^{17}$$

| address = 17 bits

| data = 32 bits

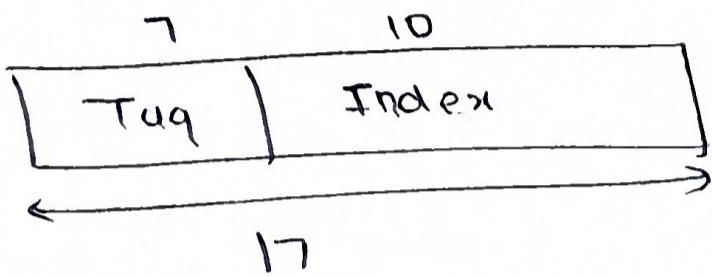
Cache memory: 2048 words

2-way set associative

= 1024 words

= 2^{10} words

$\boxed{\text{index} = 10 \text{ bits}}$



$\boxed{\text{Tag} = 7 \text{ bits}}$

bits required to represent:

$$\text{a word} = 2^2 = \boxed{2 \text{ bits}}$$

$$\text{no. of blocks} = \frac{1024 \text{ words}}{4 \text{ words}}$$

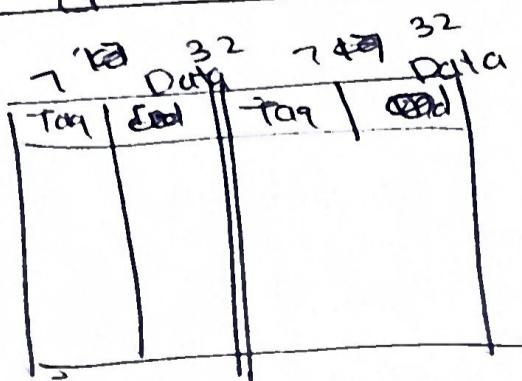
$$= 256 \text{ blocks}$$

$$\text{block} + \text{word} = \text{index}$$

$$= 2^8$$

$\boxed{\text{block} = 8 \text{ bits}}$

Sized the cache memory



to find total cache size

use tag & data

$$= 1024 \times 78$$

$$= 1K \times 78$$

$\boxed{\text{cache memory size} = 78K}$

③

How many 128×8 RAM chips are needed to provide a

memory capacity of 2048 bytes

a How many lines of address bus must be used to access 2048 bytes of memory. How many of these lines will be common to all chips?

b How many lines must be decoded for the chip select? Specify the size of the decoder?

Main memory = 128×8 RAM

$$\text{address} = 128 = 2^7$$

$$\boxed{\text{address} = 7 \text{ bits}}$$

$$\boxed{\text{data} = 8 \text{ bits}}$$

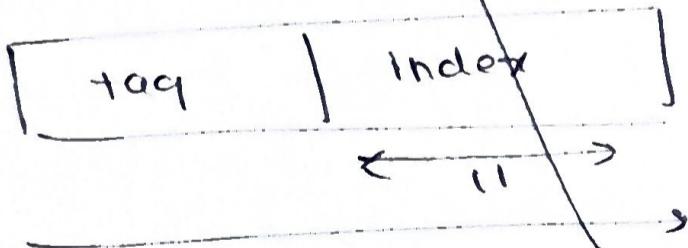
Cache memory = 2048 bytes

$$= 2^{11}$$

$$\boxed{\text{index} = 11 \text{ bits}}$$

1 byte = 8 bits

block word



(a) given 128×8 bits

required mem. capacity = 2048

$$\frac{\text{no. of chips}}{128 \times 8} = 16 \text{ chips}$$

(b) To access 2048 bytes = $2^{11} \rightarrow 11$ address lines

no. of chip bits to distinguish among 16 bits = 4

$$\text{Common lines} = 11 - 4 = 7$$

need a 4×16 decoder

(4) ~~Check~~ A computer uses RAM chips of 1024×1 capacity.

(a) How many chips are needed to provide a memory capacity of 1004 bytes

(b) How many chips are needed to provide a memory capacity of $16KB$.

assume that it is 1024×8

Ans: 1024×8

(5) A direct mapped cache has a capacity of $16KB$ and a line length of 32 bytes. How many bits are used to determine the byte that a memory of referenced within a cache line, how many bits are used to select the line in the cache that may contain the data.

$$\text{Cache size} = 16KB$$

$$= 2^4 \times 2^{10} = 2^{14}$$

$$= 14 \text{ bits}$$

$$\text{Line length} = 32 \text{ bytes}$$

$$\therefore \frac{2^{14}}{2^5} = 2^9 \Rightarrow 9 \text{ bits to select line}$$

$$\text{Word offset} = \log_2 32 = 5$$

$\Rightarrow 5$ bits to determine which byte is being referenced

(6) If a cache has 64 byte cache lines, how long does it take to fetch a catchline, if the main memory takes 2 cycles to respond to each memory request & return 8 bytes of data in response to each request?

$$\text{Ans no of memory requests needed} = \frac{64}{8} = 8$$

$$\text{no. of cycles per req} = 2$$

$$\text{Total no. of cycles} = \underline{\underline{64}}$$

(7) An address space is specified by the 24 bits & the corresponding memory space by 16 bits

(a) How many words in address & mem. space

(b) If the page has 8k words, how many pages & blocks are there in the system?

$$\text{Ans words in address space} = 2^{24} = \underline{\underline{2^{20} \times 8}} \text{ 16M words}$$

$$\text{memory space} = 2^{16} = 64K \text{ words}$$

$$\text{pages} = 16 / 2^10 = 8K \text{ pages}$$

$$\text{blocks} = 64K / 2^10 = 32 \text{ blocks}$$

(8) The logical address space in a computer system consists of 128 segments. Each segment has upto 32 pages of 4k words each. Physical memory has 4k blocks of 4words each. Find logical & physical add. formats

Logical Address

seq	page	word
-----	------	------

Physical

Block	word
-------	------

Computer Organization and Architecture

Unit 4

Additional Numericals

- ① Given ROM size of 512 kB, estimate the no. of address lines and data lines needed to access the memory

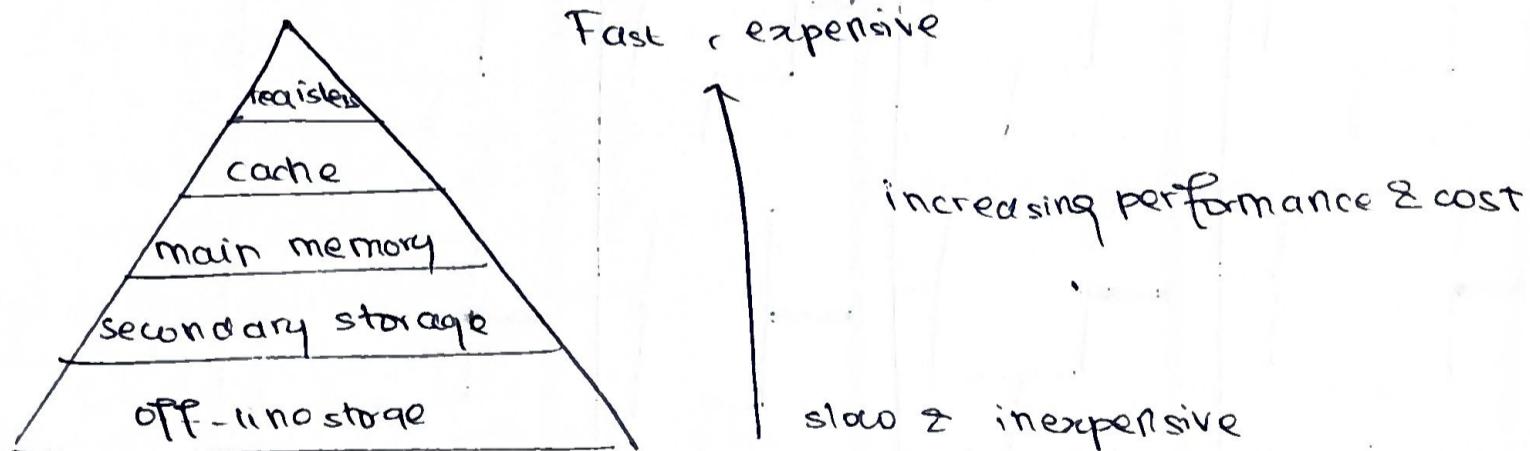
$$512 \text{ kB} = 2^9 \times 2^{10}$$

$$= 2^{19}$$

= 19 bits for Address line

8 bits for data line

- ② Show structure of memory hierarchy



- ③ Construct a static RAM memory of size $2M \times 32$ using 512 \times 8 static memory

RAM
Memory: $2M \times 32$

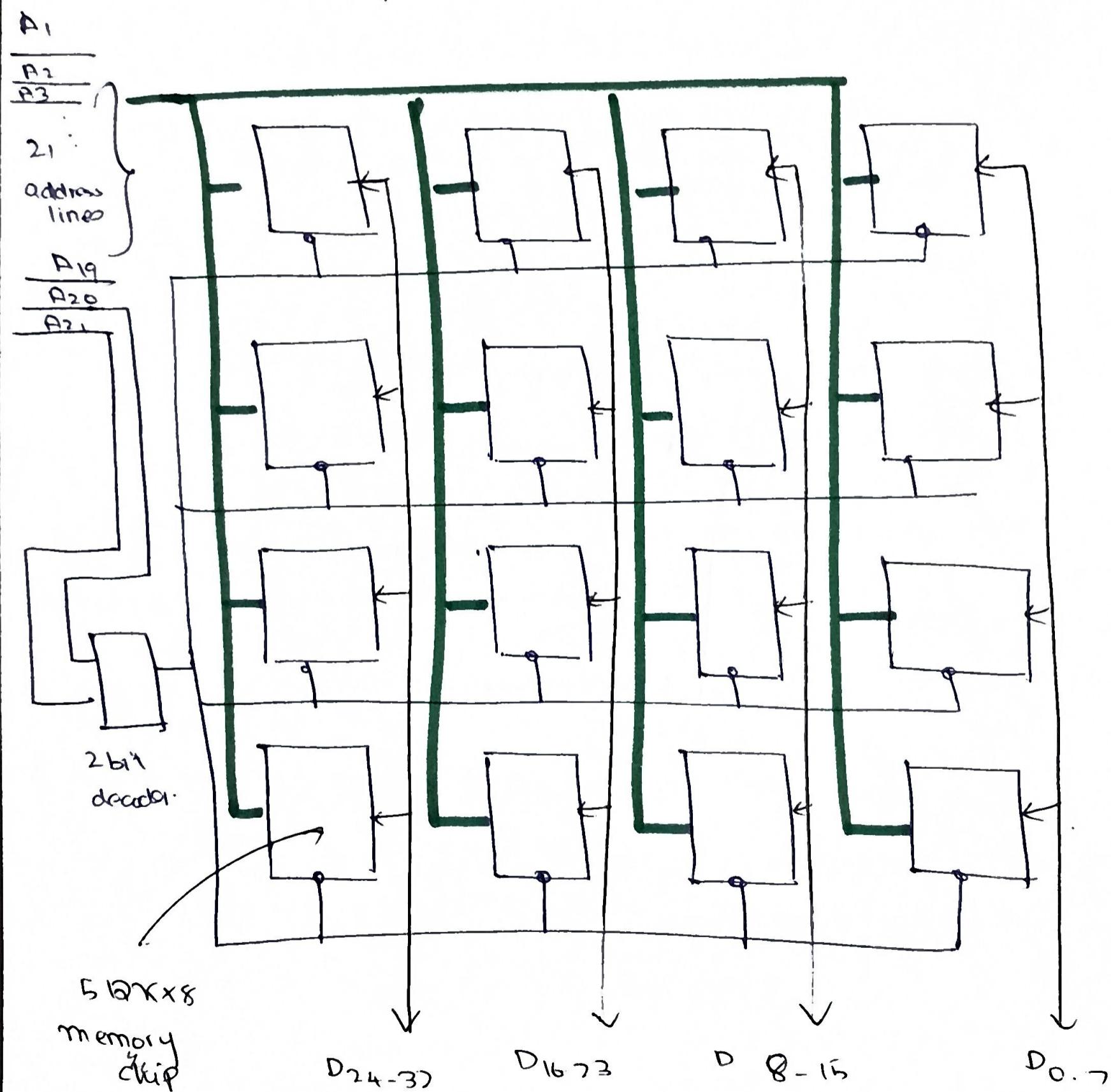
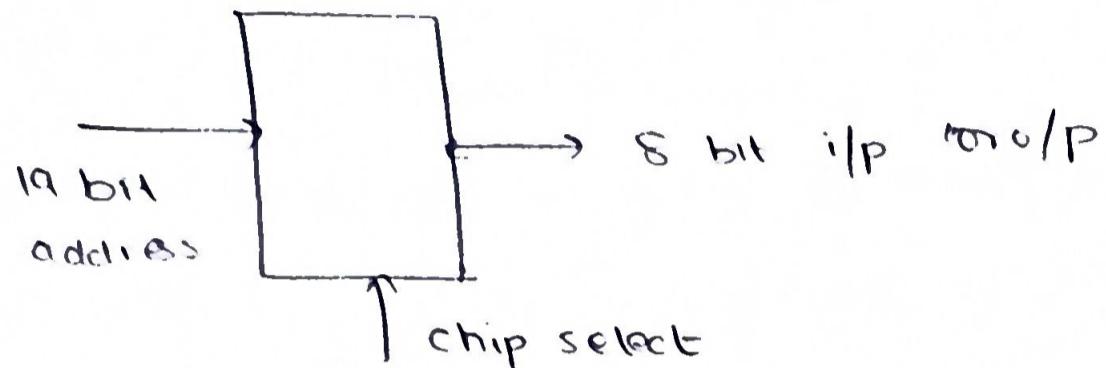
$$= 2 \times 2^{20} \times 2^5 = 2^{26}$$

static memory: 512 \times 8
 $= 2^9 \times 2^{10} \times 2^3 = 2^{22}$

$$\text{no. of blocks} = \frac{2^{26}}{2^{22}} = 2^4$$

= 16 blocks

512 K x 8 memory chip



Computer Organization and Architecture

Unit 4

* Measuring and Improving Cache Performance

Measurement: done using hit and miss rates

$$\text{CPU Time} = \left(\text{CPU execution clock cycles} + \underbrace{\text{memory stall cycles}}_{\text{from cache misses}} \right) \times \text{clock cycle time}$$

Memory stall cycles can also be defined as read-stall cycles + write-stall cycles

- Improvement : reduce the miss rate → reduce the probability that 2 different memory blocks will contend for the same cache location

reduce the miss penalty → add an additional level to the hierarchy, called multilevel caching

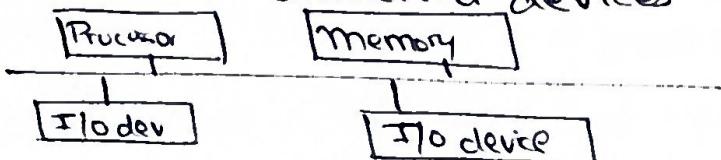
* Virtual memory + TLB - refer OS notes

* Accessing I/O Devices

→ done with a bus

↳ a shared communication link, which uses one set of wires to connect multiple subsystems.

→ enables all connected devices to share information



→ Buses consist of 3 sets of lines

(a) Address line : processor places a unique address to an I/O device on address lines

(b) Data: data placed on data lines, carry info between src & dest

(c) Control: used to signal requests and acknowledgement, and to indicate what kind of info is on the data lines

* Input - Output System

is of 2 types:

(i) memory-mapped I/O

(ii) programmed I/O

A. Memory-Mapped I/O

→ I/O devices and memory share the same address space

→ portions of address space are assigned to I/O devices

→ reads and writes to those portions are interpreted as commands to the I/O

B. Programmed I/O

→ requires all I/O operations to be executed under the direct control of the CPU.

→ data transfer happens between 2 registers - one is the CPU register and the other is attached to the I/O device. I/O devices does not have access to main memory.

→ CPU executes several instructions to transfer data from I/O to memory

1 Interrupts

- causes a CPU to temporarily transfer control from its current program to another program
- An interrupt signal is sent via a dedicated ~~process~~ line to the processor
- Processor executes the interrupt service routine (ISR)
- All registers, flags, PC values are stored onto a stack
- Time required to save status & restore contributes to Interrupt Latency

hardware interrupt - power supply failure

software interrupt - division by 0

Steps

Assert interrupt request line



Complete current instruction



Put current PC value onto stack



Load PC value w/ addr. of ISR



Execute return from interrupt



Continue where interrupted

* Handling Interrupts

1. Maskable Interrupt - can be masked under software control
2. Non maskable Interrupt - cannot be masked under software control
3. Single Line Interrupt - all I/O ports share a single interrupt request line
has a polling software that checks the load state of each device
4. Multilevel Interrupt - I/O device has individual interrupt pins
no polling required - save time
5. Vectored Interrupts - Each device assigned an interrupt vector.
Interrupts are handled by the processor after the device acknowledges w/ the interrupt vector as id
6. Interrupt Nesting - preemption of low priority interrupt by another high priority interrupt
7. Daisy Chaining - INTR common to all devices
int. req. line
INTA
ack. line propagates in a serial manner
closest device (serially) is of higher priority
8. Daisy Chaining w/ a Priority Group - combine daisy chaining and interrupt nesting to form priority group.
Each priority group has daisy chaining within.

9. Pipelined Interrupts - difficult to find interrupting instruction in pipelines , precise vs. imprecise interrupts

* Direct Memory Access

- used when large blocks of data have to be transferred at high speeds
- DMA allows transfer directly between I/O and memory, w/o minimal intervention from the processor.
- To initiate the transfer of a block of words, the processor sends the following data:
 - (i) starting address
 - (ii) word count
 - (iii) r/w mode
 - (iv) a control to start

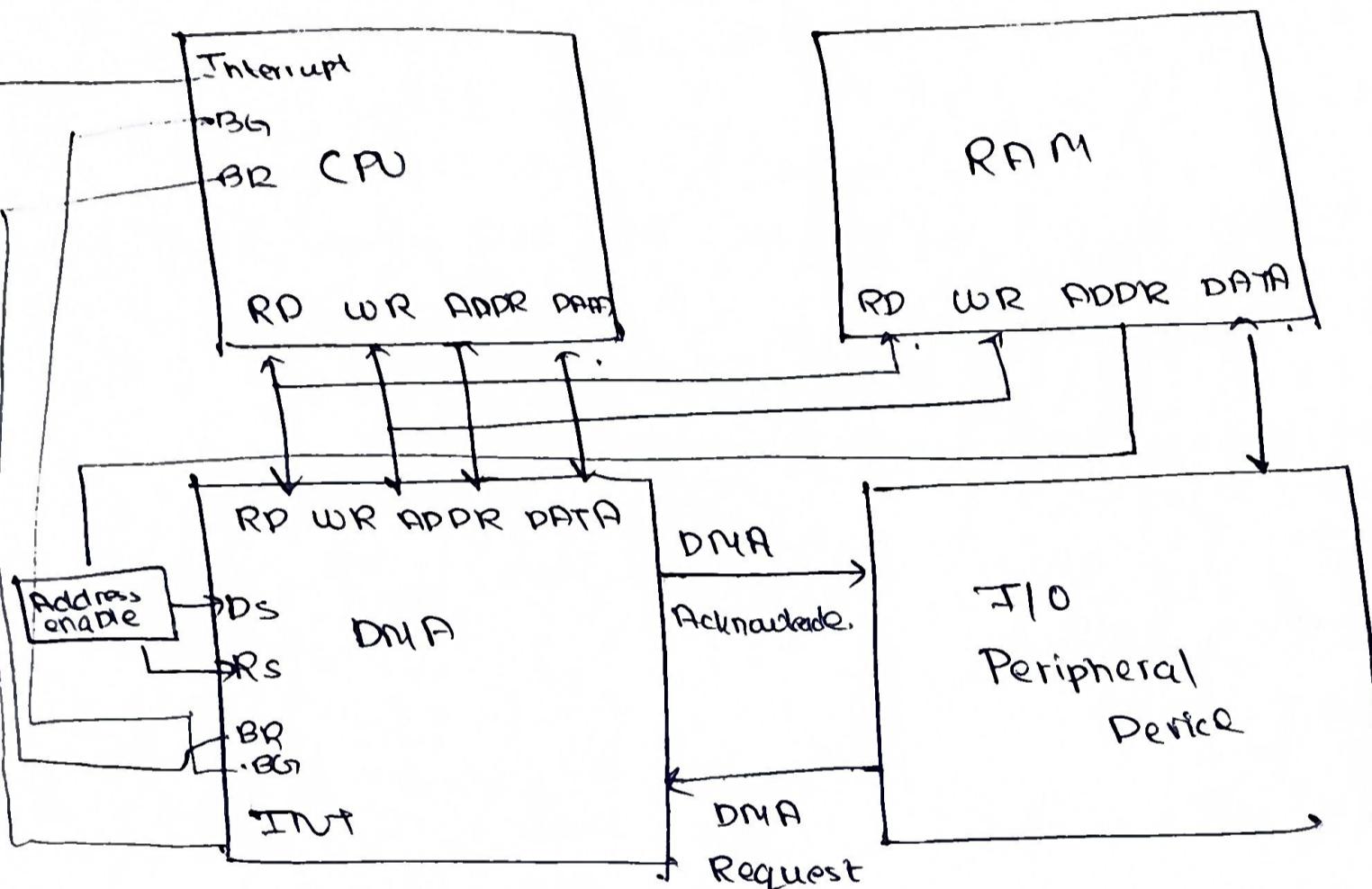
Steps

- Bus request (BR) is used by the DMAC (DMA controller) to request CPU to relinquish control of the bus.
- CPU activates the bus grant (BG)
- DMA takes control of the buses to conduct memory transfer w/o processor intervention.

DMA Transfer can be:

- (i) Burst mode - a block sequence consisting of one or more memory words transferred continuously

(ii) Cycle Stealing : allows DMA controller to transfer one data word at a time, after which control of buses is given back to the CPU. CPU delays ops for 1 clock cycle to allow the DMA transfer.



Request sent by peripheral

activate BR - relinquish control

CPU responds w/ BG line - buses are disabled.

DMA puts value of addr register in address bus

$BG_1 = 0 \Rightarrow RD, WR = \text{Input}$

$BG_1 = 1 \Rightarrow RD, WR = \text{Output}$

for every word, increments addr, decrements counter

word count = 0 \Rightarrow remove bus request

informs CPU via an interrupt

→ Bus Arbitration

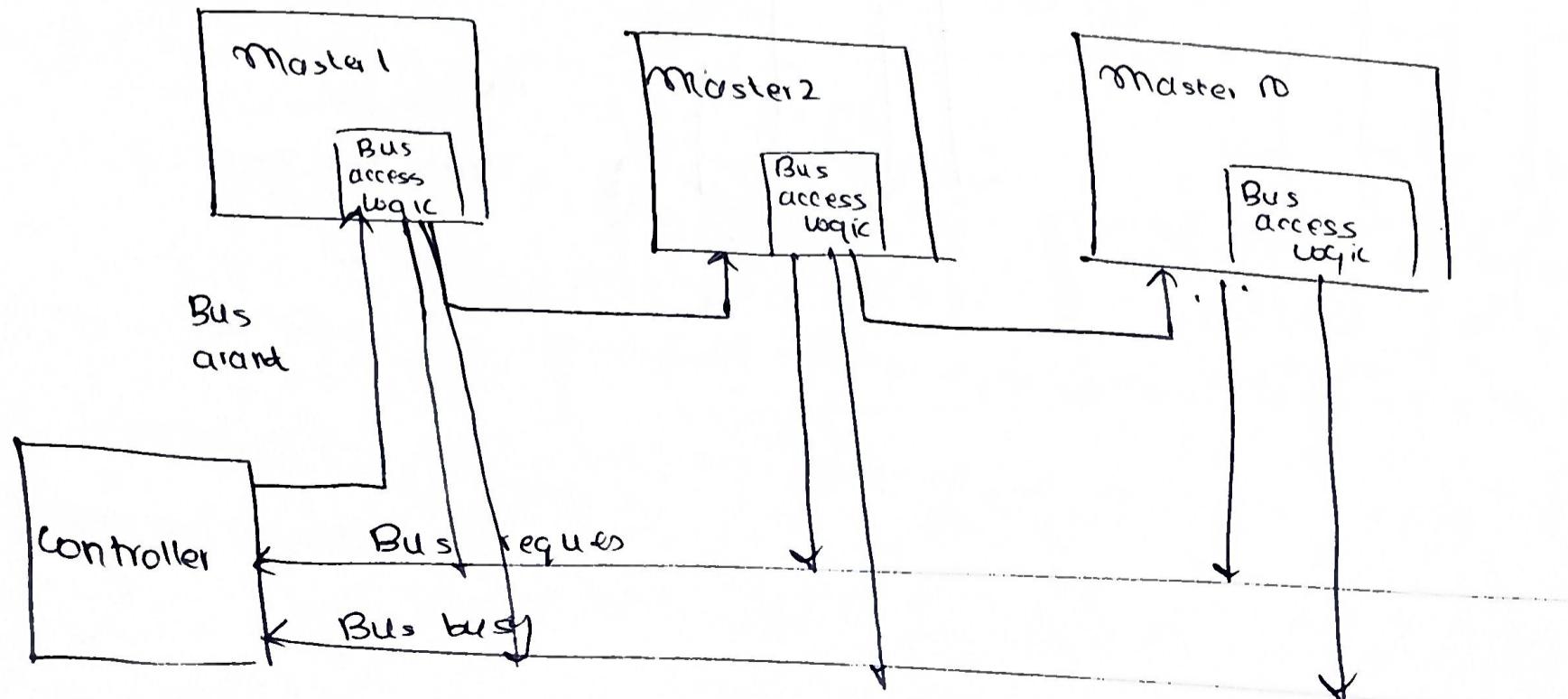
Bus Master: Device that initiates data transfer on the bus

Bus Arbitration: The process by which the next device to become master is selected

Types of Bus Arbitration: (i) Centralized
(ii) Decentralized (distributed)

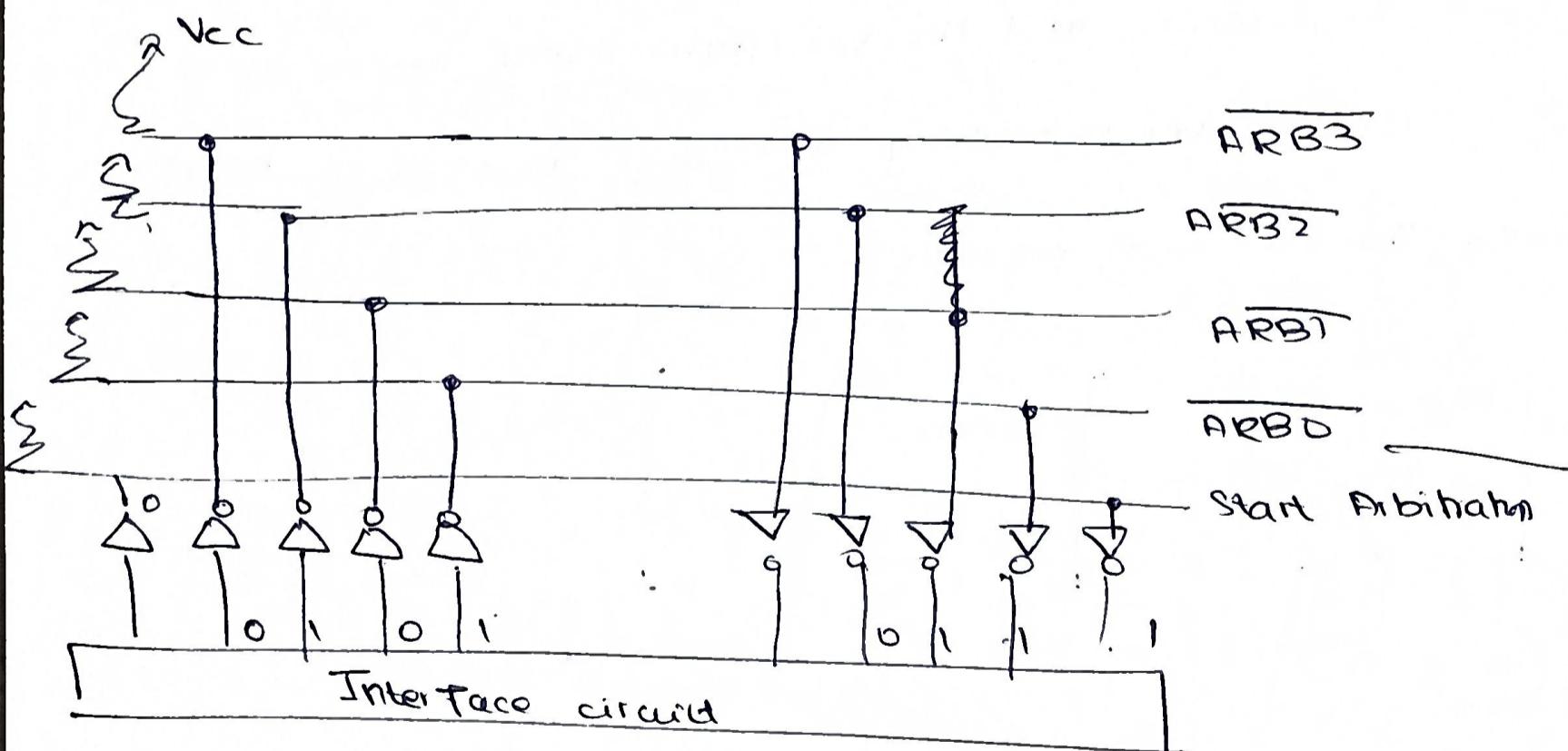
Centralized Bus Arbitration

- In this method, a single bus arbiter performs the required arbitration
- Daisy Chaining - until the bus grant signal comes across the first master who is making a request for access, it travels serially through each master.



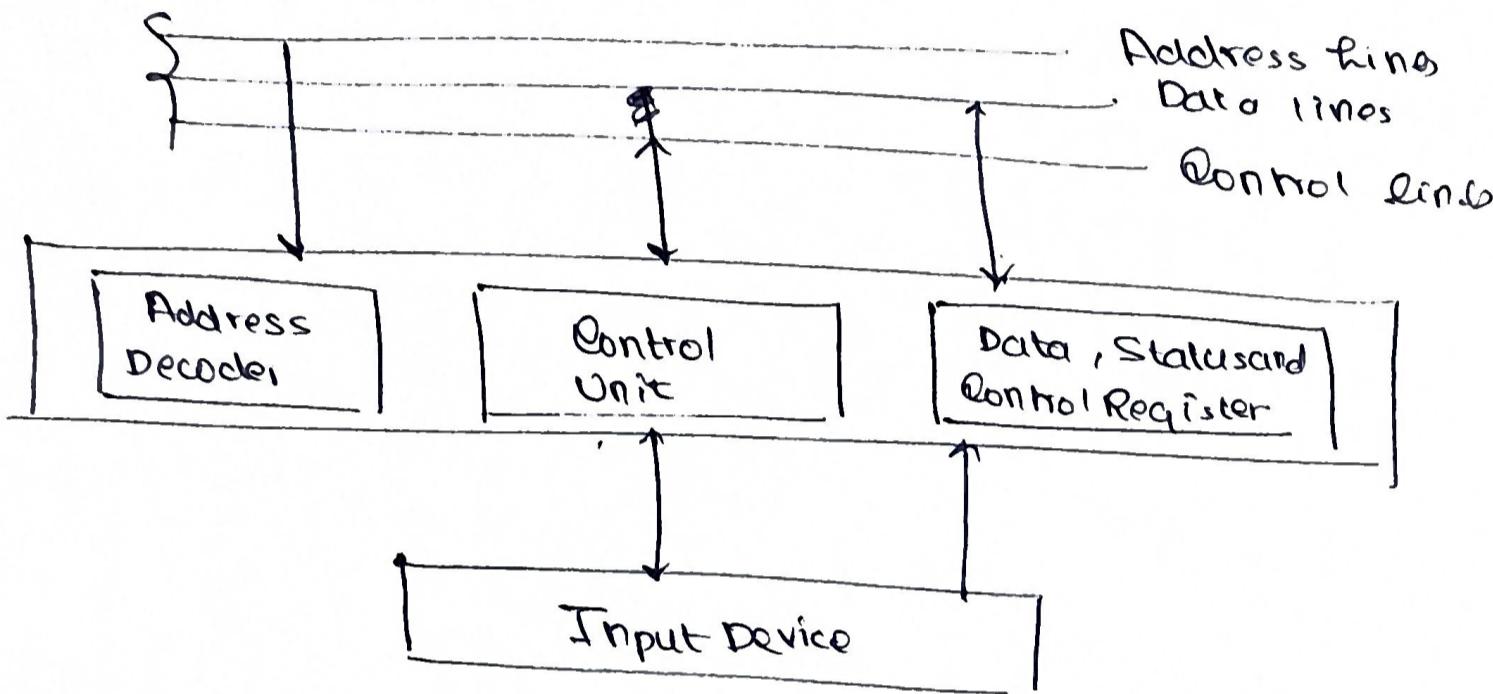
Distributed Arbitration

- No single processor or controller controls the network.
- All arbitors participate actively to decide who will be the bus master.
- Each ~~bus~~ device on the bus is given a 4 bit identification no.
- One or more devices request the bus, and place their id no. on the collector lines.
- The device w/ the highest ID no. is chosen.



* Interface Circuit

- circuitry that is designed to link I/O devices to the processor.
- acts as a mediator to make the computer communicate with the I/O modules



Features

has a data register: stores data temporarily while the data is exchanged between the I/O and processor

has a status register: bits indicate whether I/O device is free for transmission or not

has a control register: indicates type of op to be performed

User programs are prevented from using I/O operations directly because the OS does not provide direct access to the address space.