

DIRECT MEMORY ACCESS

that has
an instruction to transfer input or output data is executed
only after the processor determines that I/O device is ready.
To do this considerable overhead is incurred if processor
has to keep intervening.

To transfer large blocks of data at high speed, a
special control unit maybe provided to allow transfer of
the entire data block b/w external device and main memory
without the continuous intervention by processor \rightarrow DMA

DMA transfers are performed by a control circuit that is
part of I/O interface called DMA Controller

DMA controller performs operations that would normally
require processor to carry out by Processor

To initiate transfer of block of words, processor sends
Starting address, no. of words in block & direction of
transfer.

DMA controller proceeds to do all of it. While processor can
do something else.

Once job is done, it ~~sends~~ an interrupt to inform processor

31 30		1 0	
IRQ IE	R/W	DONE	
Starting Address			Registers in DMA

2 registers are used for storing starting address & word
count. The third has status & control flags.

R/W : determines direction of transfer bit = 1 Read opn

Done : set to 1 if task is completed and ready to receive
another command.

IE (Interrupt Enable) : set to 1 after task is completed to let processor know
IRQ : set to 1 when it has requested an interrupt.

DMA controller modes

1) Cycle stealing: Since processor originates most memory access
cycles, DMA controller is said to be 'stealing'
memory cycles from processor

2) Block/Purge mode: DMA controller may be given exclusive access
to transfer a block of data without interruption

BUS ARBITRATION

Need: A conflict may arise if both the processor and DMA
controller or 2 DMA controllers try to use the bus at
same time to access main memory. \therefore To resolve this,
we need an arbitration procedure

Bus Master: Device that's allowed to initiate data transfers on
the bus at any time.

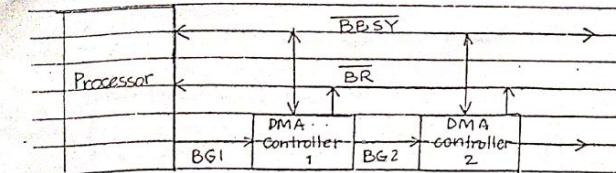
" Bus Arbitration is the process by which next device to
become bus master is selected & bus mastership transferred to it!"

2 approaches

i) Centralized : single bus arbiter performs reqd arbitration

ii) Distributed : All devices participate in selection of next
bus master

CENTRALIZED ARBITRATION



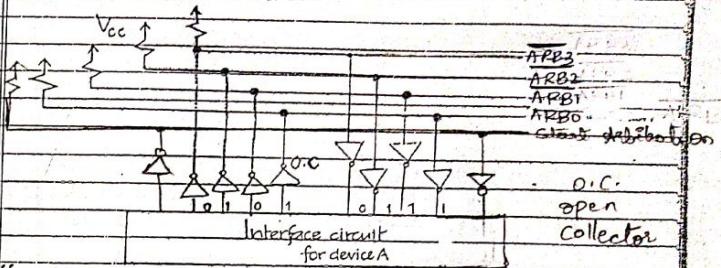
BR : Bus Request

BBSY : Busy Busy

BG : Bus Grant BG1, BG2

- 1) Processor is normally the bus master unless it grants membership to DMA controllers
- 2) DMA controller sends BR signal if it needs to be the bus master.
- 3) The BR line is open line. When it's activated, the processor activates Bus Grant signal BG1 indicating to DMAC that they may use the bus when it becomes free.
- 4) This signal is in Daisy chain so if the first one isn't (DMAC) requesting, it's passed on.
- 5) Current bus master indicates that it's using the bus by activating 'busy busy' BBSY line. To prevent other devices from using the bus at the same time.

DISTRIBUTED ARBITRATION

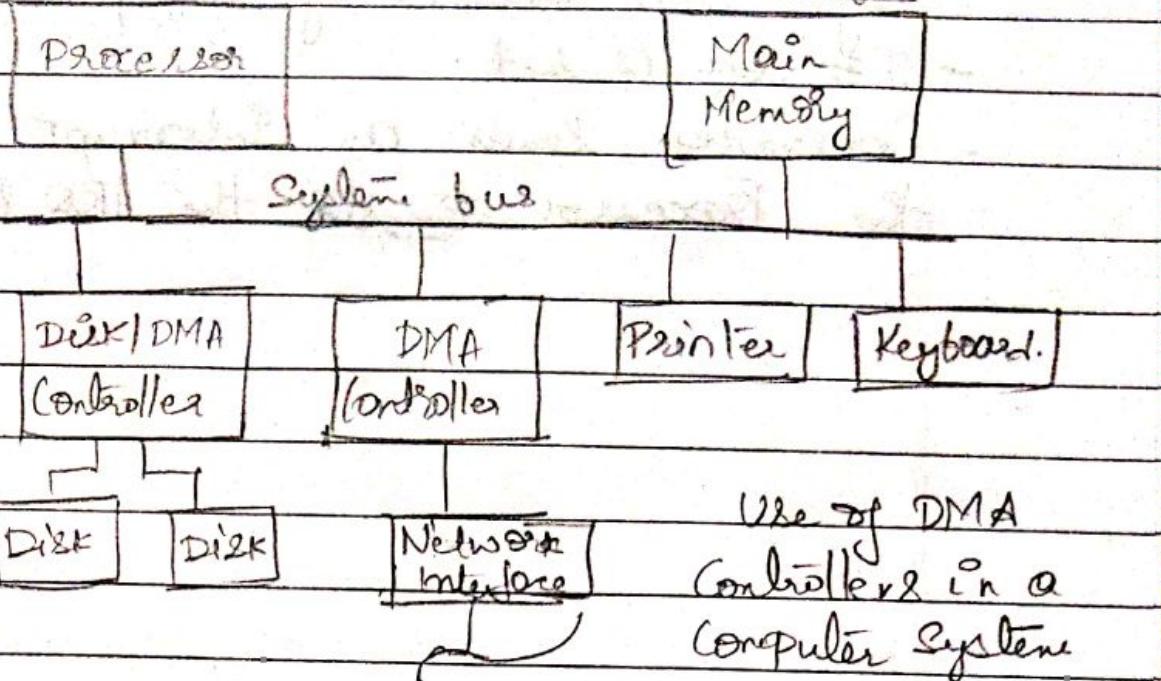


- 1) Each device on bus is assigned a 4-bit ID
- 2) When a device requests the bus, they assert Start-Arbitration signal and place their 4-bit ID mes on 4 open collector lines.
- 3) Winner is selected as a result of interaction among the signals transmitted over these lines by all contenders.
- 4) Net outcome is that the code on lines represents request that has highest ID no. (ID number)

Advantages

- 1) Higher reliability as operation of bus isn't dependent on a single device.

Use of DMA Controllers in a Computer System



- A example of a computer system how DMA (DMAC) Controllers may be used .
- A DMAC connects a high-speed link to the computer bus.
- The disk controller also has DMA Capability & provides two DMA channels.
- To start a DMA transfer a block of data from the main memory to one of the disks, a program writes → Address
→ Word Count Information to registers of the corresponding channel of the disk Controller .
- When the DMA transfer is completed, it is recorded in the status and control register of DMA channel by setting the DDONE Bit.
IE bit is set.
- Controller sends an interrupt signal (request) to the Processor . Sets IRQ bit

BUSES

The processor, main memory & I/O devices can be interconnected by means of a common bus whose primary function is to provide a communications path for data transfer.

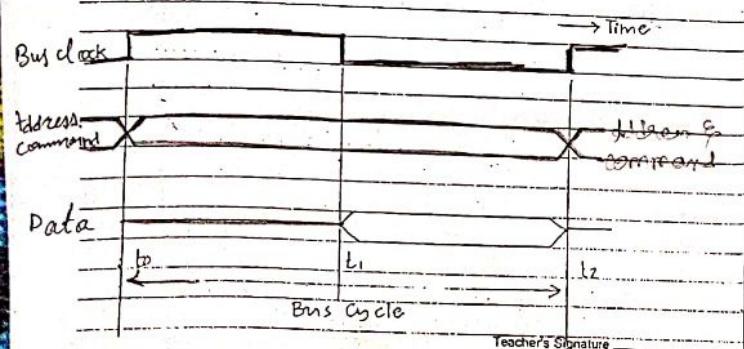
Bus protocol is the set of rules that govern the behaviour of various devices connected to bus.

Bus lines can be grouped into 3 types: data, address & control lines.

Control Signals specify whether R/W. They also carry the timing information. Read is high to 1 & write is low to 0
 Master \rightarrow initiates data transfers (Processor or DMA)
 Slave \rightarrow Device addressed by master

Synchronous Bus

In a synchronous bus, all devices derive timing info from a common clock line. Each equally spaced pulses define intervals which constitute a bus-cycle, during which one data transfer can take place.



Sequence of events during i/p operation

- 1) Master places device address on address lines & sends appropriate command on control lines at time t_0 .
- 2) Info travels over bus determined by its physical & electrical characteristics.
- 3) Clock pulse width must be longer to allow all devices to decode address, control signals, so that the addressed device can respond at t_1 . (t_1-t_0)
- 4) Addressed slave places requested i/p data on data lines at time t_1 .
- 5) At t_2 , master captures data & stores it in register buffer.
 Similarly, master places o/p data on data lines when it transmits address & command info. At t_2 , addressed slave captures data & loads it into buffer.

Asynchronous Bus

Based on use of a 'handshake' between the master and the slave.

Principle: The master places the address & command info on the bus. It then indicates to all devices by activating Master-ready line. This causes all devices on bus to decode the address. The selected slave performs read operation & informs processor with Slave-ready line. Master then removes its signal from bus.

Teacher's Signature

Master Ready

Slave Ready

Data

$t_0 \dots t_1$ $t_2 \dots t_3$ $t_3 \dots t_4$ $t_4 \dots t_5$

← Bus cycle →

to Master places address & command info on bus & All devices on bus begin to decode it.

t_1 - Master sets the Master-ready line to 1 (informs all I/O devices)
Delay $t_1 - t_0$ is to allow for any 'skew' that may occur

t_2 - The selected slave, having decoded the address & info performs the read input operation & sets slave-ready to 1.
($t_2 - t_1$ depends on distance b/w master & slave.)

t_3 - Slave-ready signal arrives at the master. After a delay eqvt to max. bus skew, it captures data & drops the Master-ready signal.

t_4 - Master removes the address & command info from the bus.

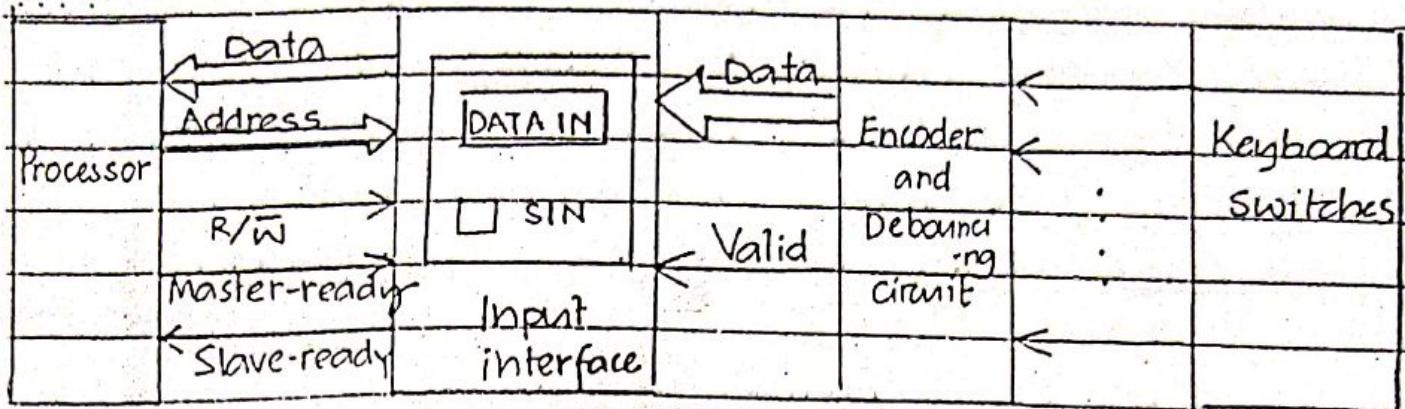
t_5 - When device interface sees MR signal as 0, removes data & SR signal. This completes i/p transfer.

Teacher's Signature

MR - Master Ready
SR - Slave Ready

PARALLEL PORT

Keyboard - Processor Connection



When a key is pressed, switch closes establishing a path for electrical signal

Signal is detected by Encoder that generates ASCII code for corresponding character.

Debouncer ensures I/O routine waits long enough to take one input as one.

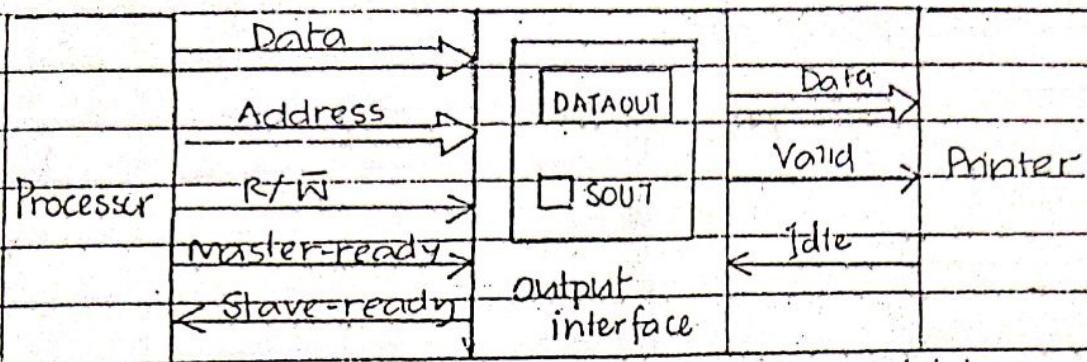
O/P of encoder is in bits that represent encoded character & control signal called Valid \rightarrow indicates a key being pressed.

ASCII code is loaded into DATAIN & SIN is set to 1

O/P of DATAIN connected to datalines of bus that are turned on when processor issues a read instruction that selects register. Inside the contents of the DATAIN register.

Interface circuit is connected to asynchronous bus. Master ready and Slave ready signals.
R/W = Read & Write transfers.

b) Printer-Processor Connection



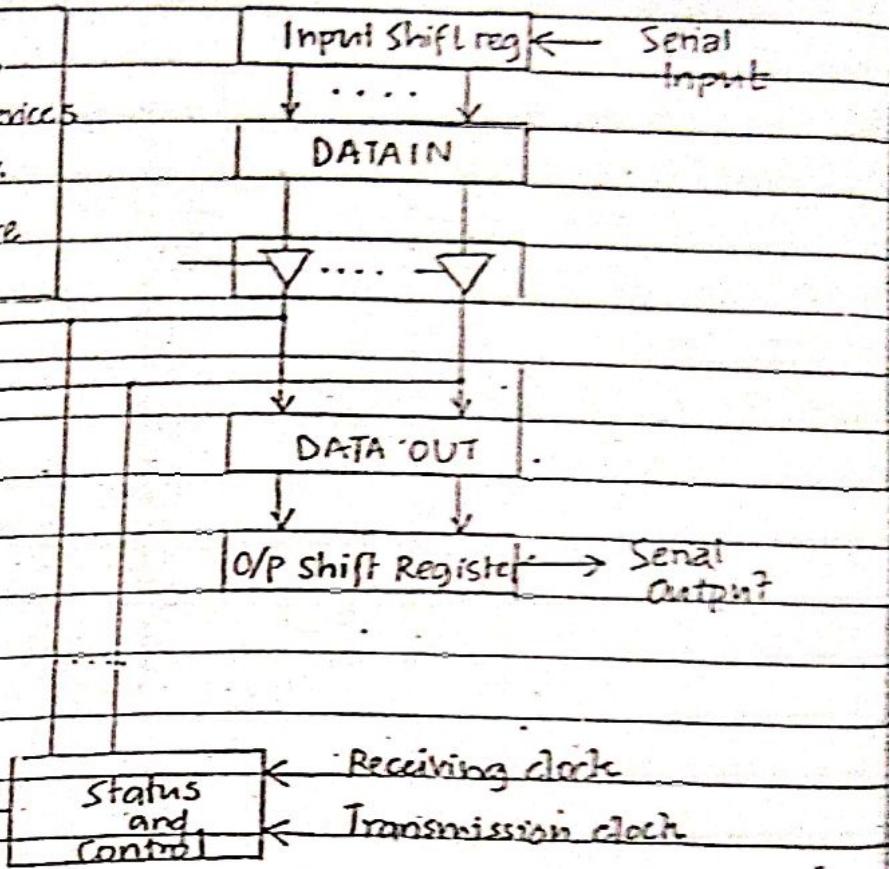
- 1) Printer operates under control of handshake signals Valid, Idle
- 2) When it's ready to accept a character, it asserts Idle signal.
- 3) The interface ckt can then place a new character on the data lines & activate Valid signal
- 4) In response, printer starts printing new character & negates Idle
- 5) SOUT flag is set to 1 when printer is ready to accept another character and cleared out when a new character is loaded into DATAOUT by the Processor.

SERIAL PORT

A serial port is used to connect processor to I/O devices that require transmission of data 1-bit at a time.
 Key Feature : Its interface circuit is capable of communicating bit-serially on device side & bit-parallel on bus-side.

It requires fewer wires,
 transmission is suited for devices
 which are away from each other.
 Data rate : depends on nature
 of devices connected

D₇ —
 ... :
 D₀ —
 My address →
 RST →
 RxD → Chip and
 RTW → Register
 Ready → Select
 accept ←



Includes DATAIN & DATAOUT. The i/p shift register accepts serial input. When all 8 bits are received, contents are loaded into DATAIN parallel.

As soon as data are transferred to DATAIN, shift register starts accepting next 8-bit character

Double buffer : Transfer of second character can begin as soon as first character is loaded into DATAIN.

Similarly, o/p data in DATAOUT register are loaded into o/p shift register, from which bits are shifted out to I/O device.

The status flags SIN & SOUT have similar functions

Teacher's Signature

* Explain DATAOUT, DATAIN, SIN & SOUT.

~~SPEED, SIZE, AND COST~~

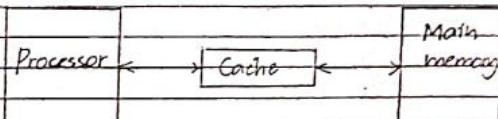
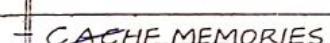
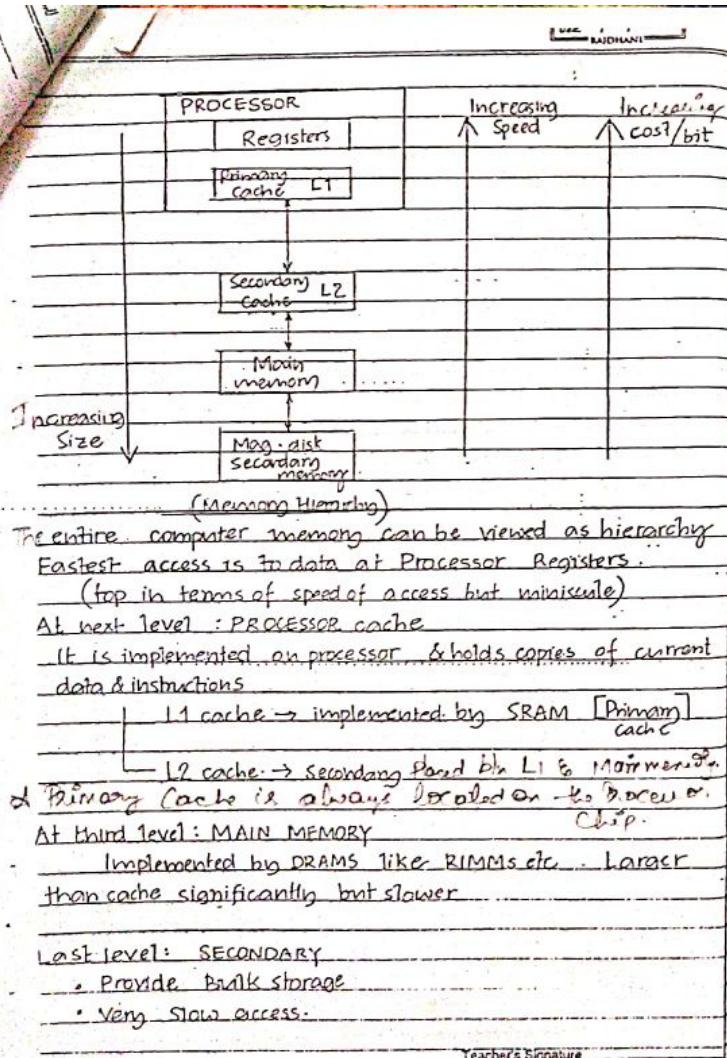
1. A very fast memory can be implemented if SRAM chips are used. But these chips are very expensive. Thus for cost reasons it's impractical to build large memory (Alternatively DRAM is used)
2. Dynamic RAM can be implemented at reasonable cost but affordable size is still small compared to demand.
3. Secondary storage (mag. disks) are used to implement large memory. However they're much slower.

CONCLUSION

1. Huge amt. of cheap storage → mag. disks
2. Large & affordable main memory → DRAM
3. Cache memories → SRAM.

(P.T.O)

mag. disk (Magnetic disk)



The speed of main memory is limited by electronic & packaging constraints. ∴ An efficient solution is to use 'cache memory' that makes main memory appear faster to processor than it really is.

Operations

Locality of Reference

It manifests in 2 ways

i) Temporal: means that a recently executed instruction is likely to be executed again

ii) Spatial : means that institutions in close proximity of recently executed one are likely to be executed soon.

So if active segments like these can be placed in a fast cache memory, then total execution time can be reduced significantly.

Note :-

Load-through or early restart :- Word sent to the processor as soon as it is read from main memory.

Write with different word not in Cache _____
during write up.

- When a Read Request is received from the Processor,
 - the contents of a block of memory words containing the location specified are transferred into the Cache one word at a time.
 - When the program references any of the locations in this block, the desired contents are directly read from Cache.

Read or write bit: Processor issues RD or RW request. The Cache Control Circuitry determines whether the requested word currently exists in the Cache. If it does, the Read or Write op_p is performed on the appropriate Cache location.

Write through Protocol:- The Cache location & Main memory location are updated simultaneously.

Write back or copy back:- Update only the Cache location & to mark it as updated with associated flag bit called dirty or modified bit. The main memory location updated, the block containing marked is to be removed from Cache.

Read Miss:- When the addressed word in the Read op_p is not in the Cache.

MAPPING FUNCTIONS

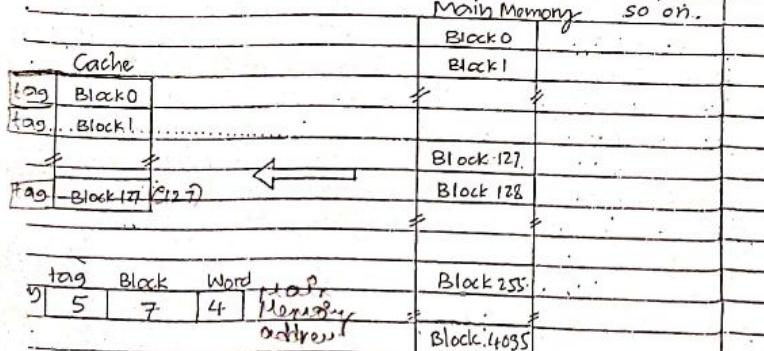
"The correspondence between the main memory blocks and those in cache is specified by mapping function. It decides which block should contain what and which block needs replacing."

Direct Mapping:

Simplest way to determine cache locations in which to store memory blocks. In this technique,

block j of main memory maps onto block j modulo 128 of cache.

Blocks 0, 128, 256 etc of main memory \rightarrow Block 0 of cache
1, 129, 257 etc \rightarrow Block 1 and so on.



Contention of a block may arise for long programs which is solved by overwriting old block.

Placement determined by memory address.

Contention may arise over jobs to cache not full.
e.g. Program starts with block 1 & continue in block 129 possibly after a branch.

Block 129 must be transferred to Block 1 position in Cache.

Memory address can be divided into 3 fields -

- Lower 4 bits \rightarrow Selection of 16 words in a block
 - 7-bit field \rightarrow determines cache position
 - 5-bits \rightarrow tag bits
- CONTENTION PROBLEM**
- EASY TO IMPLEMENT BUT NOT VERY FLEXIBLE

Advantage

Disadvantage

b) Associative Mapping

Much more 'flexible'. Main memory blocks can be placed into any cache position.

- In this case, 12 tag bits are reqd to identify a resident block.
- It gives complete freedom in choosing the cache location in which to place the memory block.
- Space in cache can be used more EFFICIENTLY.
- COSTLIER. Needs associative search & Replacement algo.

Main memory address

12 bits \rightarrow tag

11 bits \rightarrow word

Main memory

Block 0

Block 1

Cache

Tag

Block 0

Tag

Block 1

Tag

Block 127

Block i

Block i

Tag Word

12

4

Main mem.

addr

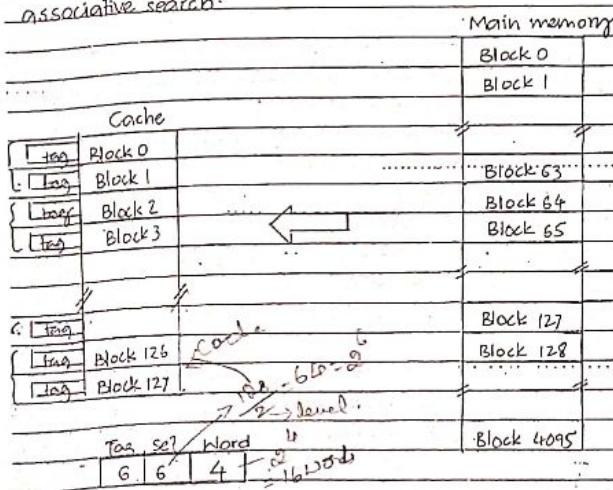
Block 4095

Teacher's Signature

Associative Mapping

It's a combination of direct- and associative-mapping
blocks of cache are grouped into sets & the mapping
allows a block to reside in only block of specific set

- contention problem eliminated
- hardware cost reduced by decreasing size of associative search:



Having 64 sets \Rightarrow 6 bits field for tag.

The control bit \rightarrow VALID BIT provided to indicate if block contains valid data.

Cache-coherence problems.

REPLACEMENT ALGORITHMS

In cache memories where cache controller has to decide which instructions to keep that'll boost processor speed the best, several replacement algorithms come into picture.

i) When a block is being overwritten, it's sensible to overwrite the one that has gone the longest time unused. It's called LRU (Least Recently Used) block & the technique is LRU replacement algorithm.

- Cache controller tracks references to all blocks
- Extensively used.
- Can even lead to poor performance often
- ∴ performance is improved by a small amount of randomness in decision.

ii) Another reasonable rule is to remove "oldest" block from a full set.

Since it does not take into account the recent pattern of access, it's generally not very effective.

iii) Simplest: Randomly choose a block to overwrite.
QUITE EFFECTIVE too.

PERFORMANCE CONSIDERATIONS

Key factors in commercial success of a computer are efficiency and cost. A common measure of success is performance ratio. $\frac{Performance}{Efficiency}$
 Hit Rate & Miss Penalty (Processor or Bus on chip).

HIT RATE and MISS PENALTY

An excellent indicator of effectiveness of particular implementation of memory hierarchy is the success rate of accessing info at various levels.

HIT RATE: No. of hits stated as a fraction of all attempted accesses.

MISS RATE: No. of misses stated as a fraction of all attempted accesses.

MISS PENALTY: It is the extra time needed to bring desired information into cache after a failed attempt (miss).

* It can be reduced by efficiently implementing memory hierarchy.

* Also if load-through approach is used.

How to improve hit rate?

Increase size of cache \rightarrow but costly

Increase obj. block size \rightarrow miss penalty increases

in Block sizes that're neither very small nor very large (16-128 bytes) give best results

Average access time experienced by processor:

$$T_{avg} = hC + (1-h)M$$

h - hit ratio

M - Miss penalty

C - time to access info in cache.

Teacher's Signature

Teacher's Signature

Caches on Processor chip

If both L1 & L2 caches are used;

- L1 should be designed to allow fast access as it'll have large effect on clock rate
- L2 cache can be slower but larger to ensure high hit rate.
 \rightarrow SRAM chip.

AVERAGE ACCESS TIME

$$T_{avg} = h_1 C_1 + (1-h_1) h_2 C_2 + (1-h_1)(1-h_2) M$$

where

$h_1 \rightarrow$ hit rate in L1 cache

$h_2 \rightarrow$ hit rate in L2 cache

$C_1 \rightarrow$ time to access info in L1

$C_2 \rightarrow$ time to access info in L2

$M \rightarrow$ time to access info from main memory.

Q1: Two Operate parallel for P-230 to data
 62060, P-3.8 P-10 (Part 1)

Cache for data to Path - PART 1

\rightarrow A practical way to speed up access of cache
 ex. to access in the bus one word circulates normally,

Interleaving

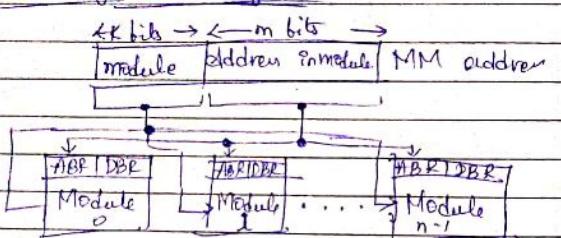
- If main memory of a computer is structured as cols of physically separate modules, each module has its own

- 1) ABR (Address Block Register)
- 2) DBR (Data Buffer Registers)

So, memory access opns may proceed in more than one module at the same time.

→ Thus, the aggregate rate of transmission of words to / from the main memory can be increased.

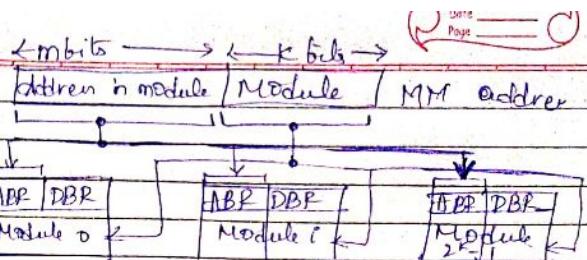
Methods of address layout:



Consecutive words in a module.

- High order k bits of memory address determine the module.
- Low order m bits - word within module.
- When a block of words is transferred from MM to Cache, only one module is busy at a time.

(2)



- Consecutive words in consecutive modules.
- Consecutive addresses can be located in consecutive modules.
- While transferring a block of data several memory modules can be kept busy at the same time.
- Faster access.
- Higher avg utilisation of the memory system as a whole.