

1) unifunction pipeline -

→ A pipeline unit with a fixed and dedicated function such as the floating point adder is called unifunctional.

e.g. → 2) The cray-1 has 12 unifunctional pipeline units for various Scaler, Vector, fixed point and floating point operations.

3) A multifunction pipe may perform different functions, either at different times or at the same time.

4) It is interconnecting different subsets of stages in the pipeline.

2) Data dependency Hazards -

→ Hazards are caused by resource usage conflicts among various instruction in the pipeline.

2) Such hazards are triggered by interinstruction dependencies.

3) Method to cope with such hazards are needed in any type of lookahead processors for either synchronous - pipeline or asynchronous pipeline.

4) There are three classes of data-dependent hazards, according to various data update patterns :
1) write after read (WAR) hazards

2) read after write (RAW) hazards

3) write after write (WAW)

→ most
with

exec

→ T

sho

per

→ sca

→ D

pro

an

→ I

sca

→ I

in

4)

Se

in

5)

e

ts

ba

bu

ts

ba

bu

ts

ba

bu

ts

- 3) Pipeline Throughput :-
- The no. of results (clocks) that can be completed by a pipeline per unit time is called its throughput.
 - This rate reflects the computing power of a pipeline. In terms of efficiency η and clock period τ of a linear pipeline we define the throughput follows:-

$$K = \frac{n}{K\tau + (n-1)\tau} = \frac{n}{\tau}$$

- When n equals the total no. of tasks being processed during an observation period $K\tau + (n-1)\tau$
- In the ideal case $\omega = 1/\tau = f$ when $n \rightarrow 1$
This means that the maximum throughput of a linear pipeline is equal to its frequency, which corresponds to one output result per clock period.

- 4) Static pipeline:-
- A static pipeline may assume only one functional configuration at a time.
 - Static pipeline can be either unifunctional or multifunctional.
 - Pipelining is made possible in static pipe only if instructions of same type are to be

- most existing computers are equipped with static pipes, either unifunctional or multifunctional.
- executed continuously.
- the function perm performed by a static pipeline should not change frequently. Otherwise, its performance may be very low.
- s) scalar pipeline-
- 1) Depending on the instruction or data types, pipeline processor can be also classified as a scalar pipeline or and vector pipeline.
 - 2) A scalar pipeline processes a sequence of scalar operands under the control of a Do loop.
 - 3) Instruction in a small Do loop are often prefetched into the instruction buffer.
 - 4) The required scalar operands for repeated scalar instructions are moved into a data cache in order to continuously supply the pipeline operands.
 - 5) The IBM System 1360, Model 91 is a typical example of a machine equipped with scalar pipelining. However, the model 91 does not have a cache.

- Q-2) Draw internal arch. of CPU pipeline structure of typical CPU.
 Q-3) Draw architecture of Cray-1

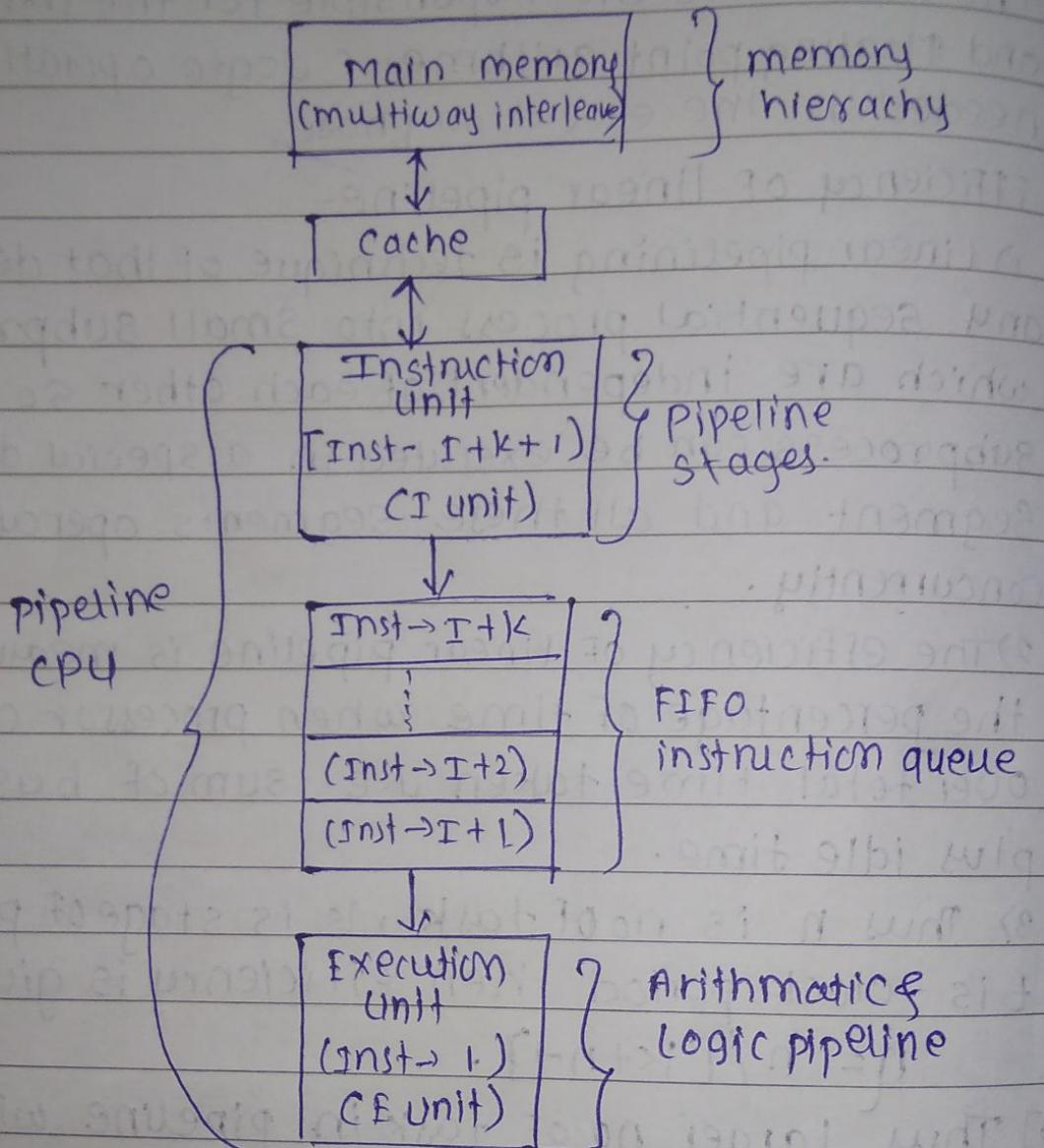


fig. The pipelined structure of typical CPU.

- Q) The central processing unit (CPU) of a modern digital computer can generally be partitioned into 3 sections → A] The instruction unit B] The instruction queue C] The execution unit.
- 2) Programs and data reside in the main memory which usually consists of interleaved memory modules.
- 3) The cache is faster storage of copies of program and data which are ready for execution.
It is used to close up the speed gap between main memory and the CPU.
- 4) The instruction unit consists of pipeline stages for instruction fetch, decode, operand addr calculation and operand fetch.
- 5) The instruction queue is a FIFO storage area for decoded instruction and fetched operands.
- 6) The execution unit may contain multiple functional pipelines for arithmetic & logical fun.

EFFICIENCY OF LINEAR PIPELINE -

- 1) Linear pipelining is technique of that decomposes any sequential process into small subprocesses which are independent of each other so that each subprocess can be executed in a special dedicated segment and all these segments operates concurrently.
 - 2) The efficiency of linear pipeline is measured by the percentage of time when processor are busy over total time taken i.e. sum of busy time plus idle time.
 - 3) Thus n is no of task, k is stage of pipeline & t is clock period then efficiency is given by $\eta = n / [k + n - 1]$.
 - 4) Thus larger no of task in pipeline will be pipeline busy hence better will be efficiency.
- The pipelined structure of a typical central processing unit →

Q-3) Draw architecture of Cray-1

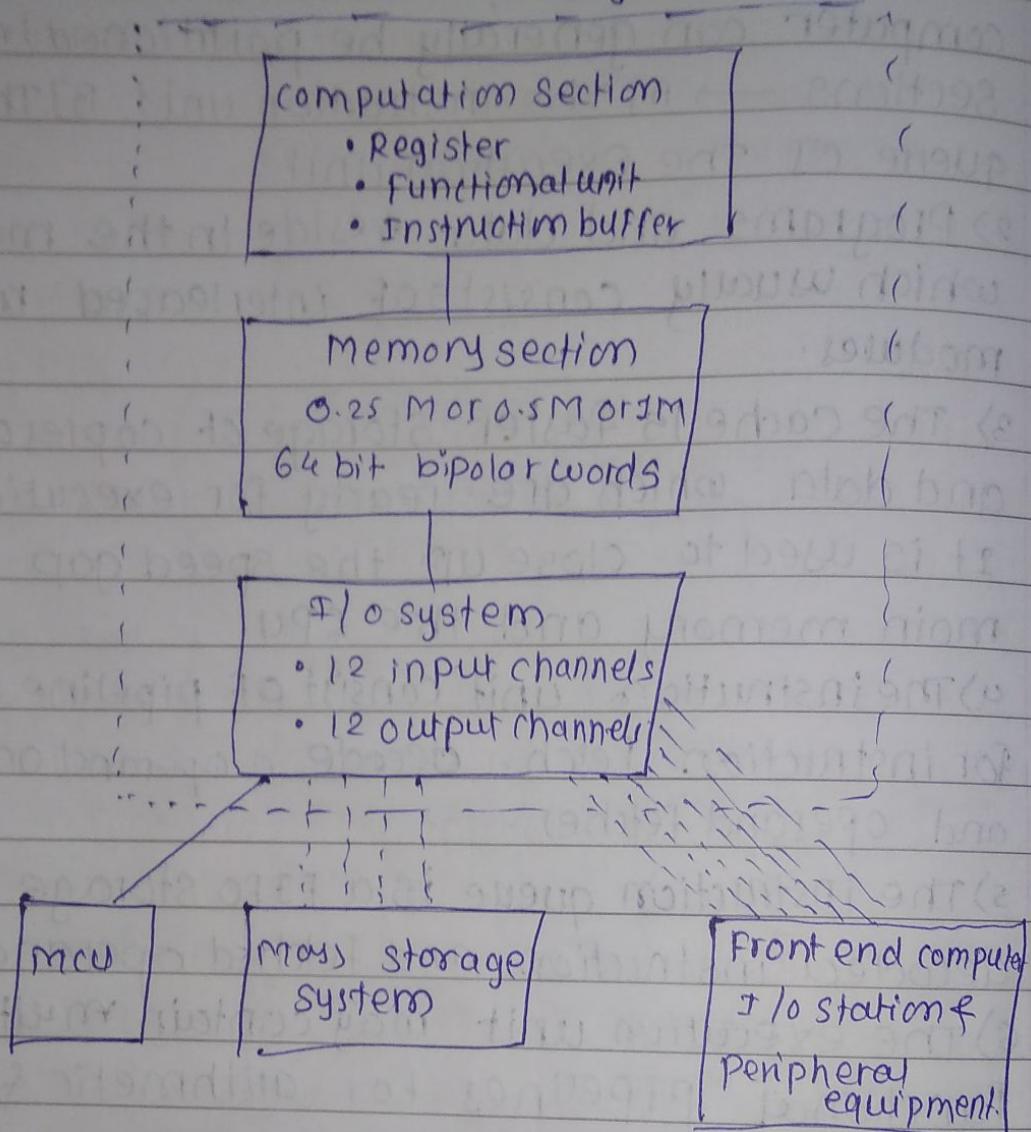


Fig. Front end system & Cray-1 connection.

Q-4

> The Cray-1 has been available as the first modern vector processing since 1976. The architecture of Cray-1 consist of no. of working registers, large instruction buffer, data buffer and 12 functional pipeline units.

- 2) with the chaining of pipeline units, interim results are used immediately once they become available.
- 3) The clock rate in the Cray-1 is 12.5 ns
- 4) The Cray is not a "stand alone" computer.
- 5) The CPU contains a computation section, a memory section and I/O section.
- 6) 24 I/O channels are connected to the front-end computer, the I/O stations, peripheral equipment, the mass storage subsystem, and maintenance control unit (MCU).
- 7) The front-end system will collect data, present it to the Cray-1 for processing, and receive output from the Cray-1 for distribution to slower devices.

Q-4 What are characteristics of vector processing?
How vector instructions are classified?

- 1) A vector operand contains an ordered set of n elements, where n is called the length of vector.
- 2) Each element in vector is scalar quantity which may be a floating point no., an integer, a logical value, or a character.
- 3) Vector instruction can be classified into four primitive types:-

$$f_1: V \rightarrow V$$

$$f_2: V \rightarrow S$$

$$f_3: V \times V \rightarrow V$$

$$f_4: V \times S \rightarrow V$$

where V and S denote a vector operand & scalar operand, resp.

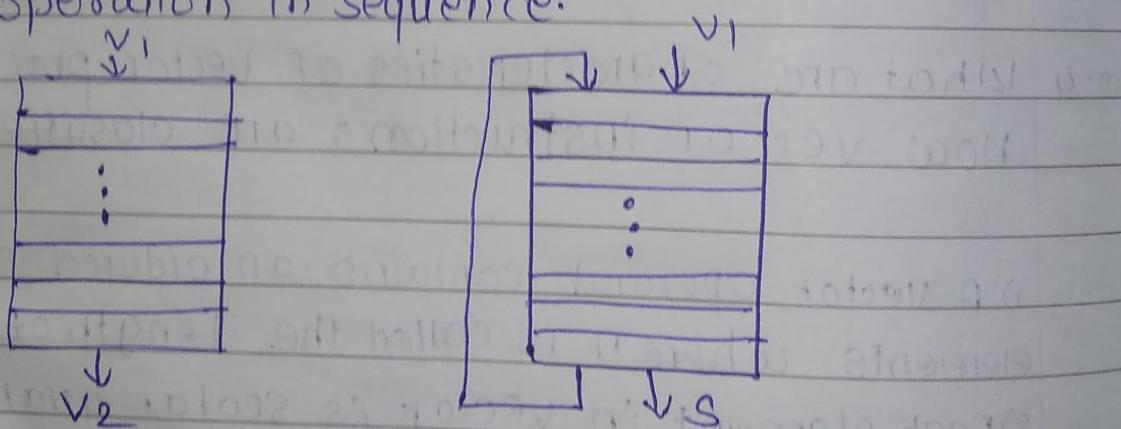
(a) The mappings f_1 and f_2 are unary operations and f_3 and f_4 are binary operations.

(b) The VSQR (vector square root) is an f_1 operation.
VSUM (vector summation) is an f_2 operation.

SNP (scalar vector product) is f_3 operation &

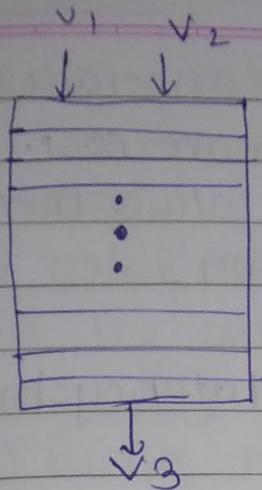
VADD (vector add) is f_4 operation.

(c) The dot product of two vectors $V_1 \cdot V_2 = \sum_i V_{1,i} V_{2,i}$ is generated by applying f_3 and then f_2 operation in sequence.

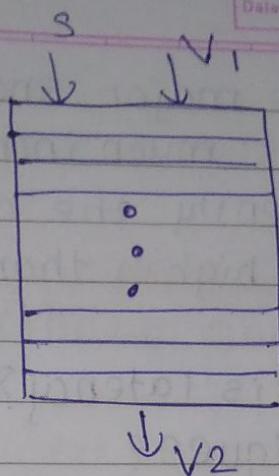


a) $f_1: V_1 \rightarrow V_2$

b) $f_2: V_1 \rightarrow S$



c) F3: $V_1 \times V_2 \rightarrow V_3$



d) F4: $S \times V_1 \rightarrow V_2$

ASS-2

Q.7 Explain Associative memory. How it's diff' from Random Access memory.

- 1) Data stored in an associative memory are addressed by their contents.
- 2) In this sense, associative memories have been known as content-addressable memory, parallel search memory and multiaccess memory.
- 3) The major advantage of associative memory over RAM is capability of performing parallel search and parallel comparison operations.
- 4) These are frequently needed in many important applications, such as the storage and retrieval of rapidly changing database, radar signal tracking, image processing.

- 5) The major shortcoming of associative memory is its much increased hardware cost.
- 6) Presently the cost of associative memory is much higher than that of RAM's.

Q-2 What is latency? What are latency hiding techniques?

1) The value of latency includes: network delays, cache miss penalty, delays caused by contentions in split transactions.

2) Interval between switches refers to the cycles between switches triggered by remote reference.

3) Latency hiding can be achieved using four methods-

1) Prefetching technique

2) Using coherent caches

3) Using relaxed memory consistency models

4) Using multiple context support.

1) Using prefetching technique -

which bring instruction or data close to the processor before they are actually needed.

2) Using coherent caches -

Supported by hardware to reduce cache misses.

3) Using relaxed memory consistency models -

By allowing buffering and pipeling of memory references.

4) Using multiple context support -

To allow a processor to switch from one context to another.

Q-3 Draw and explain CM* architecture in loosely coupled system.

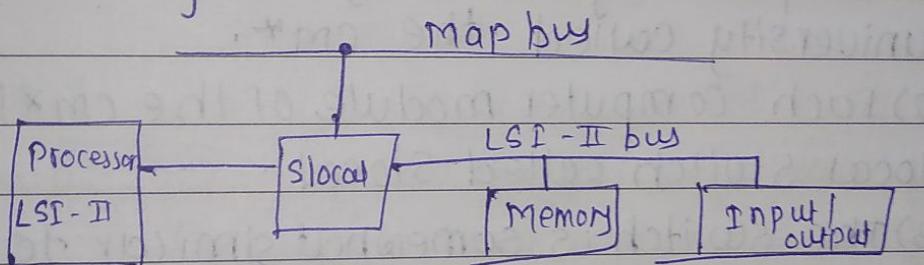


Fig. A] - A computer module.

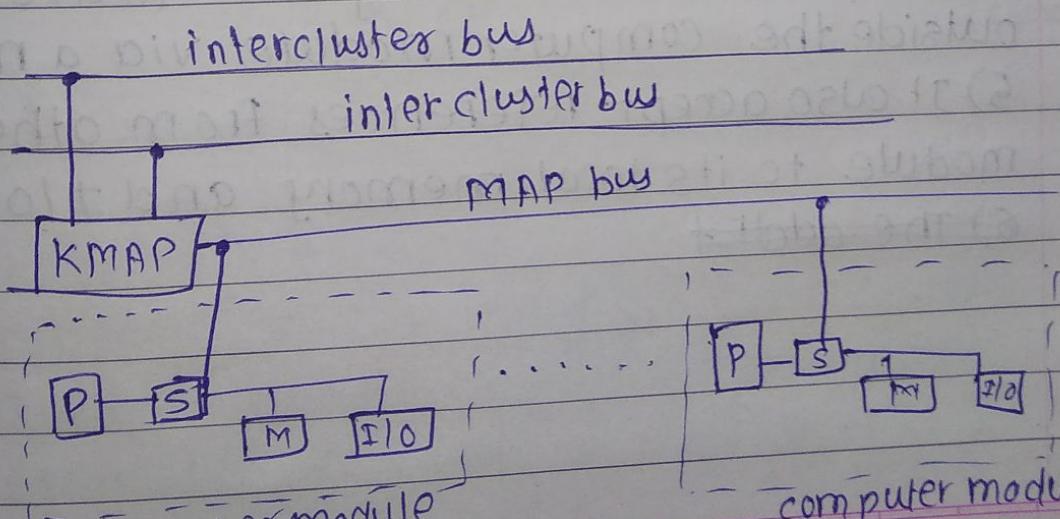
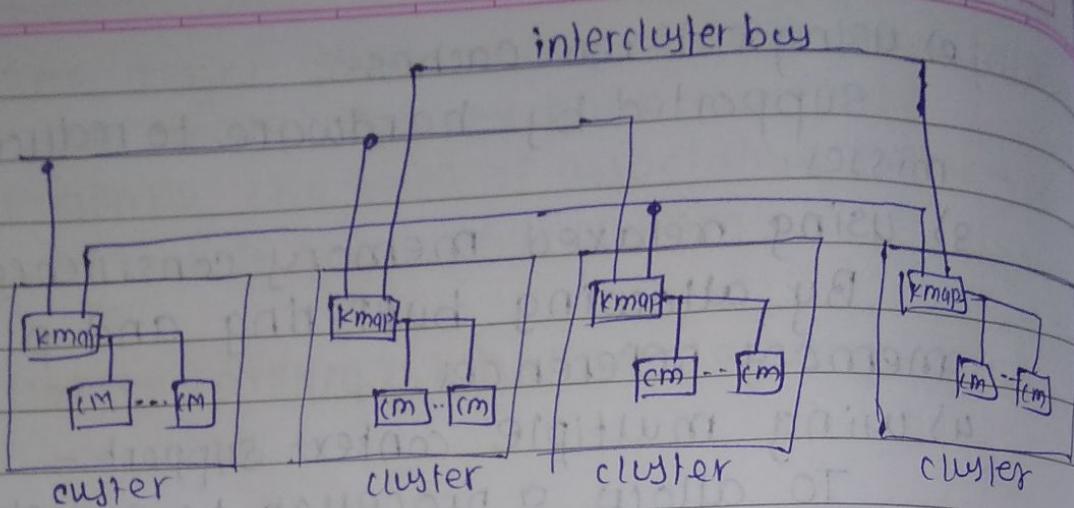


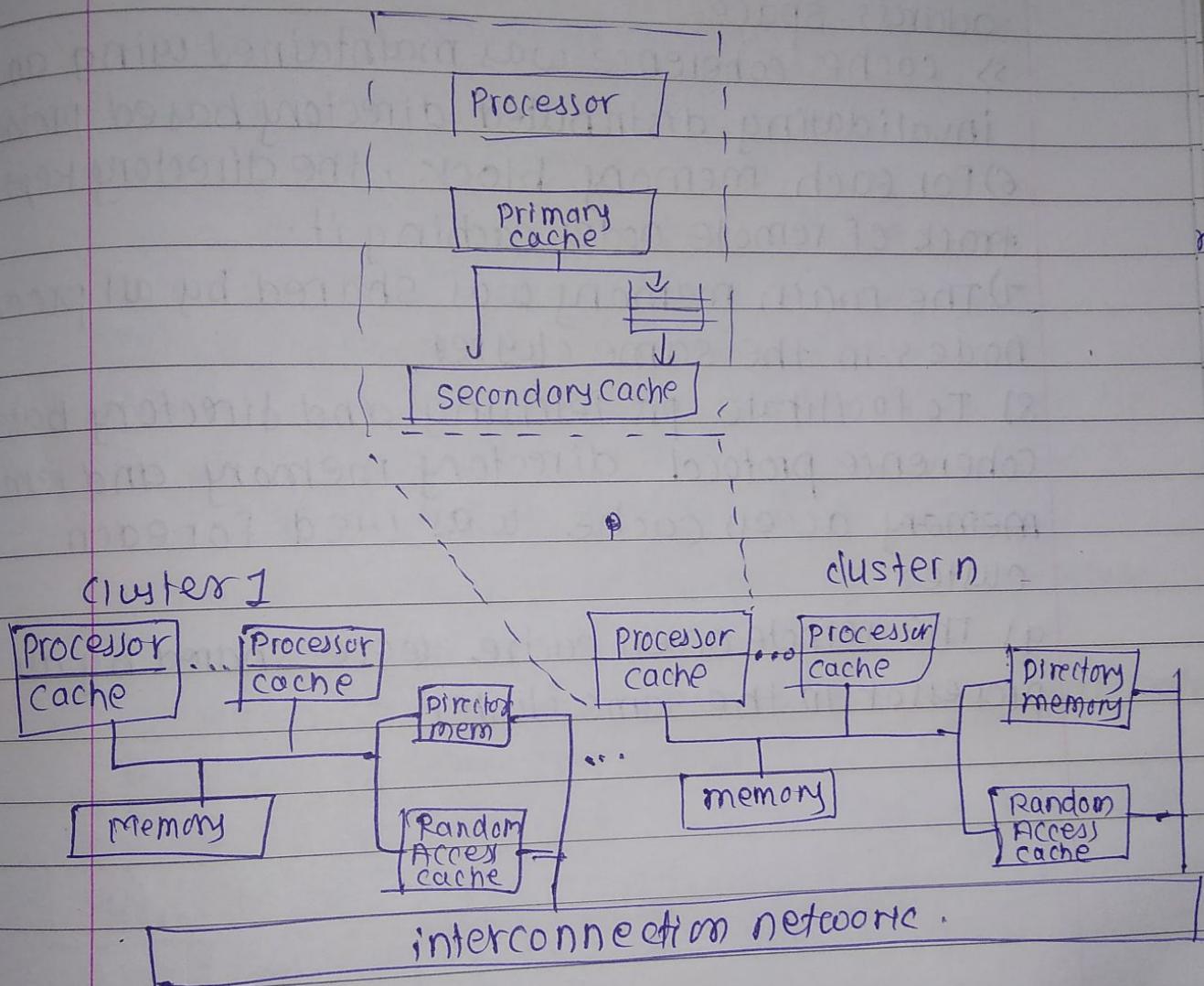
Fig. B] A cluster of computer module



c) - A network of clusters.

- 1) For a hierarchical LCS, we take example of a computer system project at Carnegie Mellon university called the CM*.
- 2) Each computer module of the CM* includes a local switch called slocal.
- 3) The switch is somewhat similar to the CAS.
- 4) The slocal intercepts and routes the processor's requests to the memory and I/O devices outside the computer module via a map bus.
- 5) It also accept references from other computer module to its local memory and I/O devices.
- 6) The address

Q.4 Draw and explain scalable coherent cache multiprocessor with distributed shared memory model.



- 1) The DASH architecture was a large-scale, cache-coherent, NUMA multi-processor system.
- 2) It consists of multiple multi-processor clusters connected through a scalable, low-latency interconnection network.

- 3) Physical memory was distributed among the processing nodes in various cluster.
- 4) The distributed memory formed a global address space.
- 5) Cache coherence was maintained using an invalidating, distributed directory based protocol.
- 6) For each memory block, the directory kept track of remote nodes caching it.
- 7) The main memory was shared by all processing nodes in the same cluster.
- 8) To facilitate pre fetching and directory based coherence protocol, directory memory and remote memory access cache was used for each cluster.
- g) The remote access cache were shared by all processor in the same cluster.

ASS - 3

Page No.:
Date: *youva*

Q-1 What are differ component of kmap in CM* architecture? State funn of each component
intercluster bus 1

intercluster bus 0

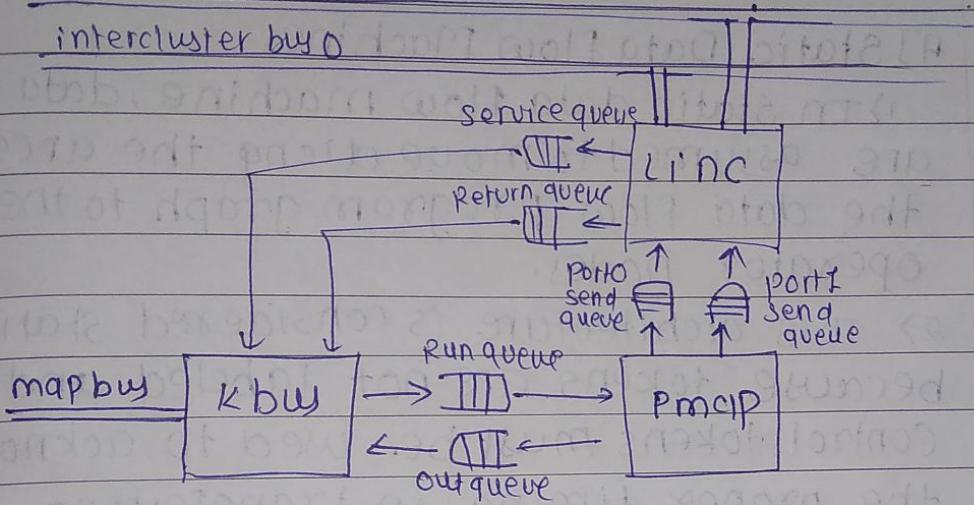


Fig. The component of Kmap in CM*

- 1) Three processors in the Kmap are the Kbus, Linc and the Pmap.
- 2) Kbus is the bus controller which arbitrates requests to the map bus.
- 3) The Linc message communication betn the Kmap and other Kmap.
- 4) The Pmap is the mapping processor which responds to requests from the Kbus and Linc; It also performs most of the request processing.
- 5) The pmap communicates with the comput modules in the cluster via map bus which is a packet switched bus.

Q-2 With block diagram explain dynamic and static data flow architecture.

A) Static Data flow Machines -

1) In static data flow machine, data tokens are assumed to move along the arcs of the data flow program graph to the operator nodes.

2) This architecture is considered static because tokens are not labeled and control tokens must be used to acknowledge the proper timing in transferring data tokens from node to node.

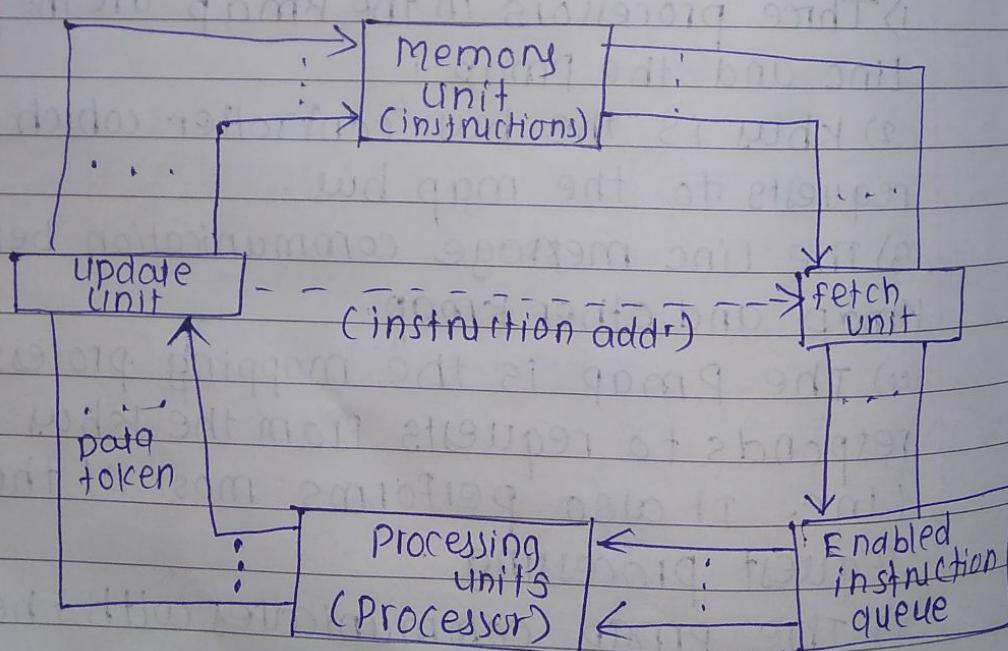


Fig. A static dataflow architecture

B) Dynamic
1) A di
tokens
exist
2) The
label a
the c
3) The
of par
comes
terms
data
mat

B) Dynamic Dataflow machine-

- 1) A dynamic data flow machine uses tagged tokens, so that more than one token can exist in arc.
 - 2) The tagging is achieved by attaching a label with each token which uniquely identifies the context of these particular tokens.
 - 3) The dynamic data flow allows greater exploitation of parallelism; however, this advantage comes at the expense of the overhead in terms of the generation of tags, larger data tokens, and complexity of the matching tokens.

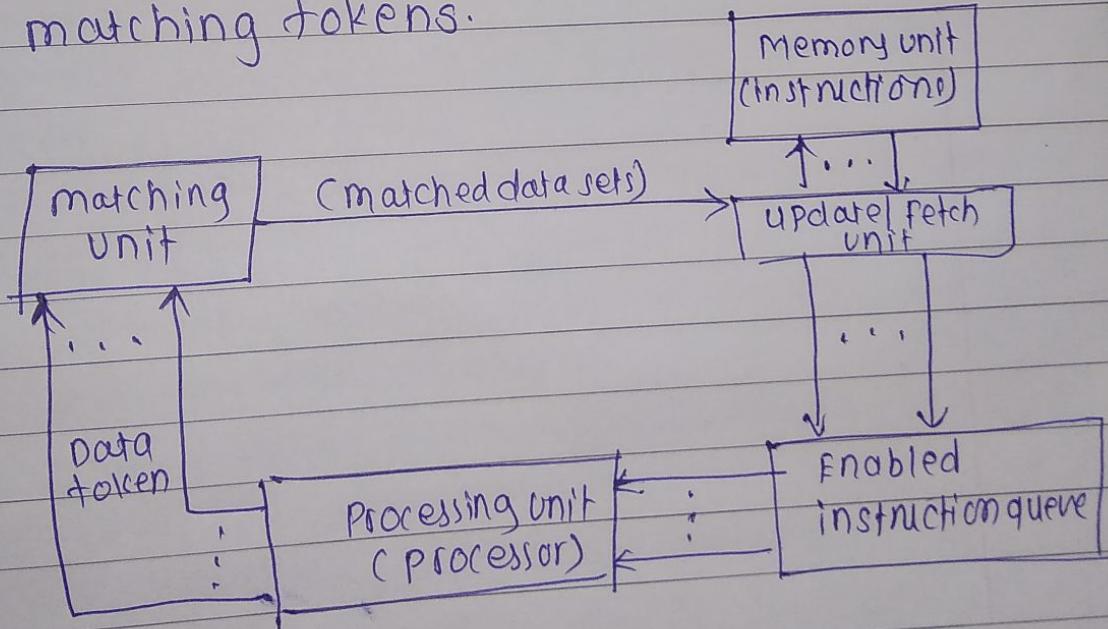


fig : Dynamic dataflow architecture.