

High Performance Computing

M	T	W	T	F	S	S
Page No.:						
Date:	29-03-23				YOUVA	

Unit I. Introduction to parallel computing

Insem

30

- a) What is Parallel computing & Explain Motivating parallelism.
- a) Write short note on levels of parallelism.
- a) Explain Dichotomy of Parallel computing platforms
- a) Explain N-wide superscalar Architecture
- a) Explain the communication costs in parallel machines.
- a) Differentiate between SIMD and MIMD
- a) What is implicit parallelism?
- a) Explain the Pipelining and superscalar execution & Explain with example
- a) Explain example of a two way superscalar implementation of instructions?
- a) What is VLIW Design Architecture? Explain in details with Architecture design & Diagram?
- a) Explain typical SIMD & MIMD Processors architecture & SIMT Architecture?

(a) What is parallel computing? Explain Motivating parallelism?

- A parallel computer is a set of processors that are able to work cooperatively to solve a computational problem.
- Parallel computing refers to the process of executing several processors on application or computation simultaneously.
- Generally, it is a kind of computing architecture where the large problems break into independent, smaller, usually similar parts that can be processed in one go.
- It is done by multiple CPUs communicating via shared memory, which combines results upon completion.
- Parallel computing also helps in faster application processing and task resolution by increasing the available computation power of systems.

Problem

Instructions

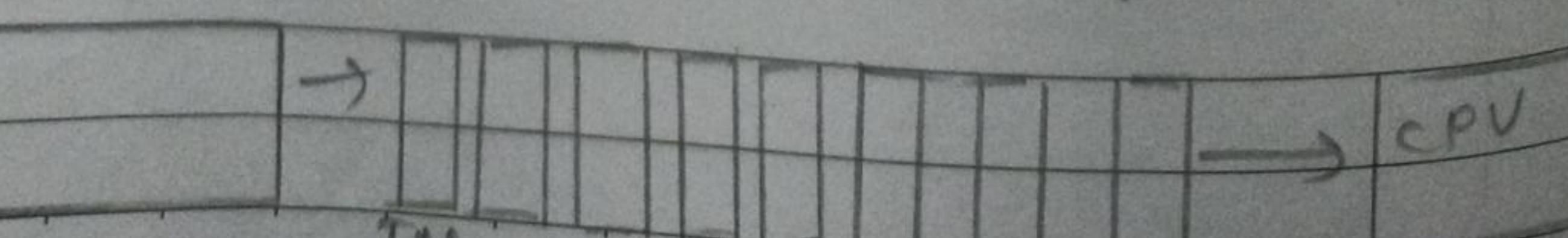
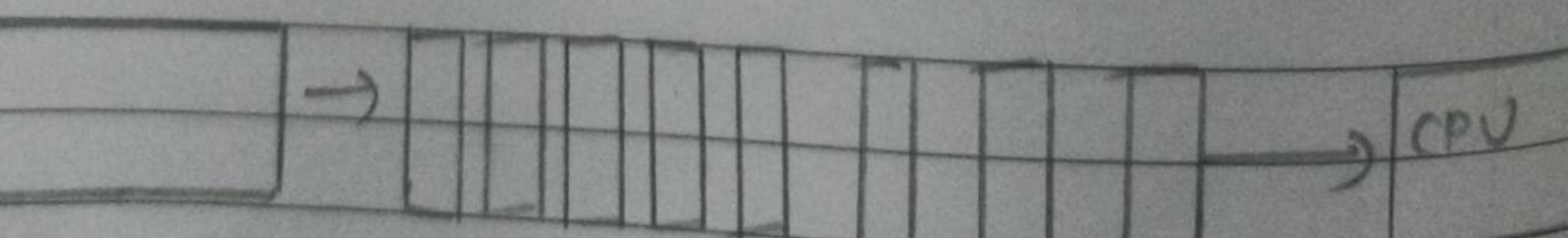
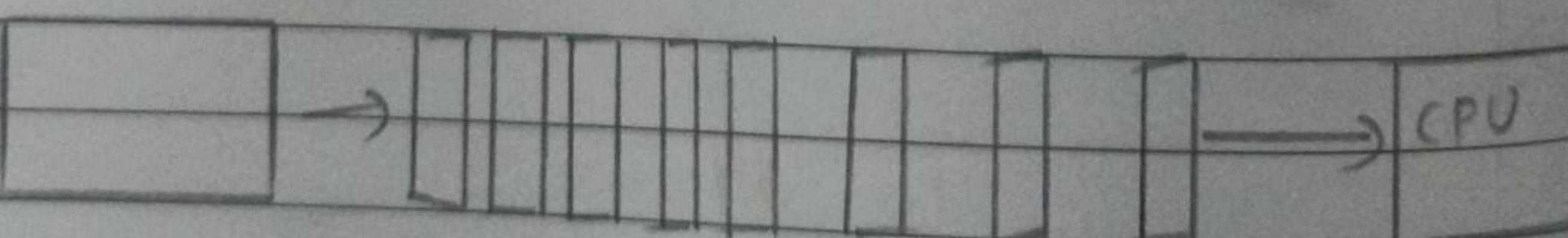
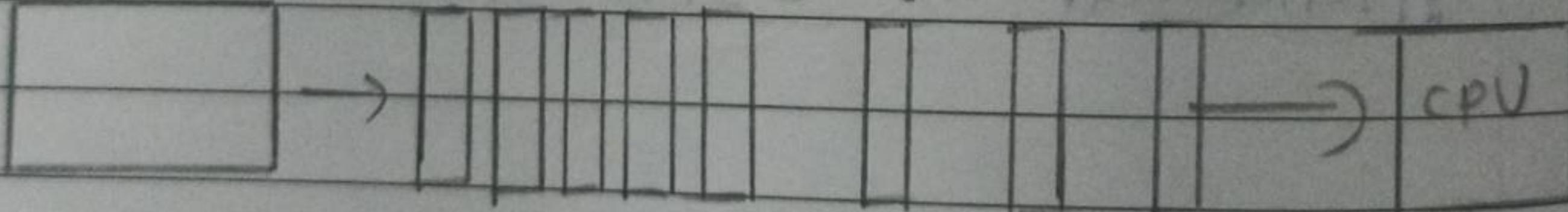


Fig: Parallel

- ① A large problem is broken into discrete parts that can be solved concurrently.
- ② Each part is further broken down to a series of instructions.
- ③ Instructions from each part execute simultaneously on different CPUs.

The Types of the levels of parallelism are.

- ① Bit-level parallelism
- ② Instructional parallelism
- ③ Data level / loop level parallelism
- ④ Task / Functional / control level parallelism.

Applications -

- ① one of the primary applications of parallel computing is Databases and Data mining.
- ② The real-time simulation of systems is another use of parallel computing.
- ③ The Technologies such as Networked videos and multimedia.
- ④ science and Engineering.
- ⑤ collaborative work environments
- ⑥ The concept of parallel computing is used by augmented reality, advanced graphics and virtual reality.

Motivating Parallelism -

- The part of parallelism in accelerate multiplying speeds has been known for numerous decades.
- Its part in provide multiplicity of data paths and improved admittance to storage elements has been important in profitable applications.
- The scalable show and lower cost of parallel stages is replicated in the wide diversity of application.
- Emerging parallel hardware and software has conventionally been time and effort demanding.
- The emergence of standardized parallel training situations, libraries and hardware have considerably summary time to (parallel) explanation,

a) → Write short note on levels of parallelism.
 The types of the levels of parallelism are
 ① Data parallelism
 ② Task parallelism
 ③ Bitlevel parallelism
 ④ Instruction-level parallelism

① Data parallelism - Data parallelism means concurrent execution of the same task on each multiple computing core.

② Task parallelism - Task parallelism means concurrent execution of the different task on multiple computing cores.

Task parallelism is the form of parallelism in which the tasks are decomposed into subtasks. Then each subtask is allocated for execution and the execution of subtasks is performed concurrently by processors.

③ Bit level parallelism - The form of parallel computing in which every task is dependent on processor word size. In terms of performing a task on large-sized data, it reduces the number of instructions the processor must execute.

④ Instruction-level parallelism

In a single CPU clock cycle, the processor decides in the instruction-level parallelism how many instructions are implemented at the same time.

The software approach in instruction-level parallelism functions on static parallelism, where the computer decides which instructions to execute simultaneously.

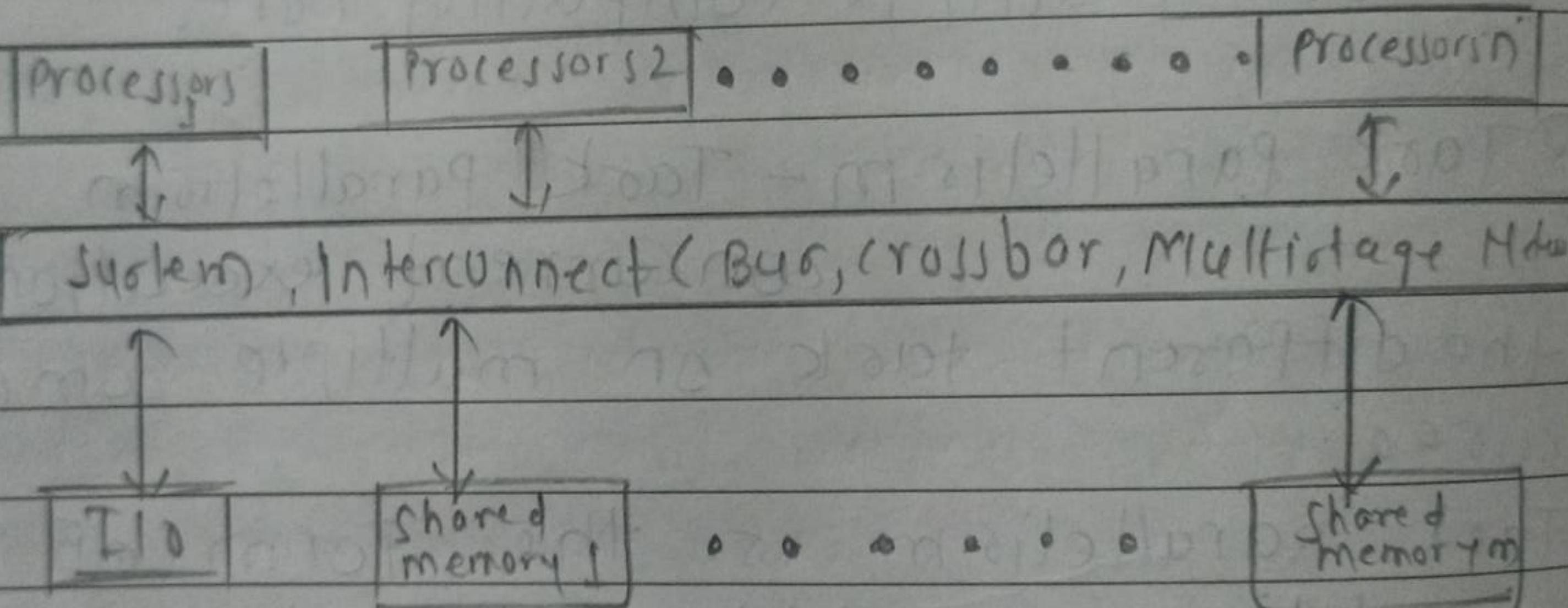


Fig. Parallel computing Architecture

Q) Explain Implicit parallelism in detail.

- - Implicit parallelism allows programmers to write their programs without any concern about the exploitation of parallelism.
- The compiler, runtime system and the underlying hardware play an important role in exploiting the parallelism implicitly.
- The parallelism is transparent to the programmer so the programmer will write the standard sequential program without adapting any special parallel constructs.
- It is the job of underlying systems to figure out the parallelism from the sequential code, with the help of different techniques.
- ← various implicit parallel mechanisms used in pipelining and superscalar execution.

(Q) Explain VLIW Architecture.

- VLIW stands for Very-long Instruction word (VLIW) Architecture.
- It is an appropriate alternative for exploiting instruction-level parallelism (ILP) in programs especially for performing more than one basic instruction at a time.
- These processors include various functional units, fetch from the instruction cache a very-long Instruction word including various primitive instructions and dispatch the whole VLIW for parallel implementation.
- The main goal of VLIW is to remove the complicated instruction scheduling and parallel dispatch that appears in most modern microprocessors.
- A VLIW processor should be quicker and less costly than a comparable RISC chip.

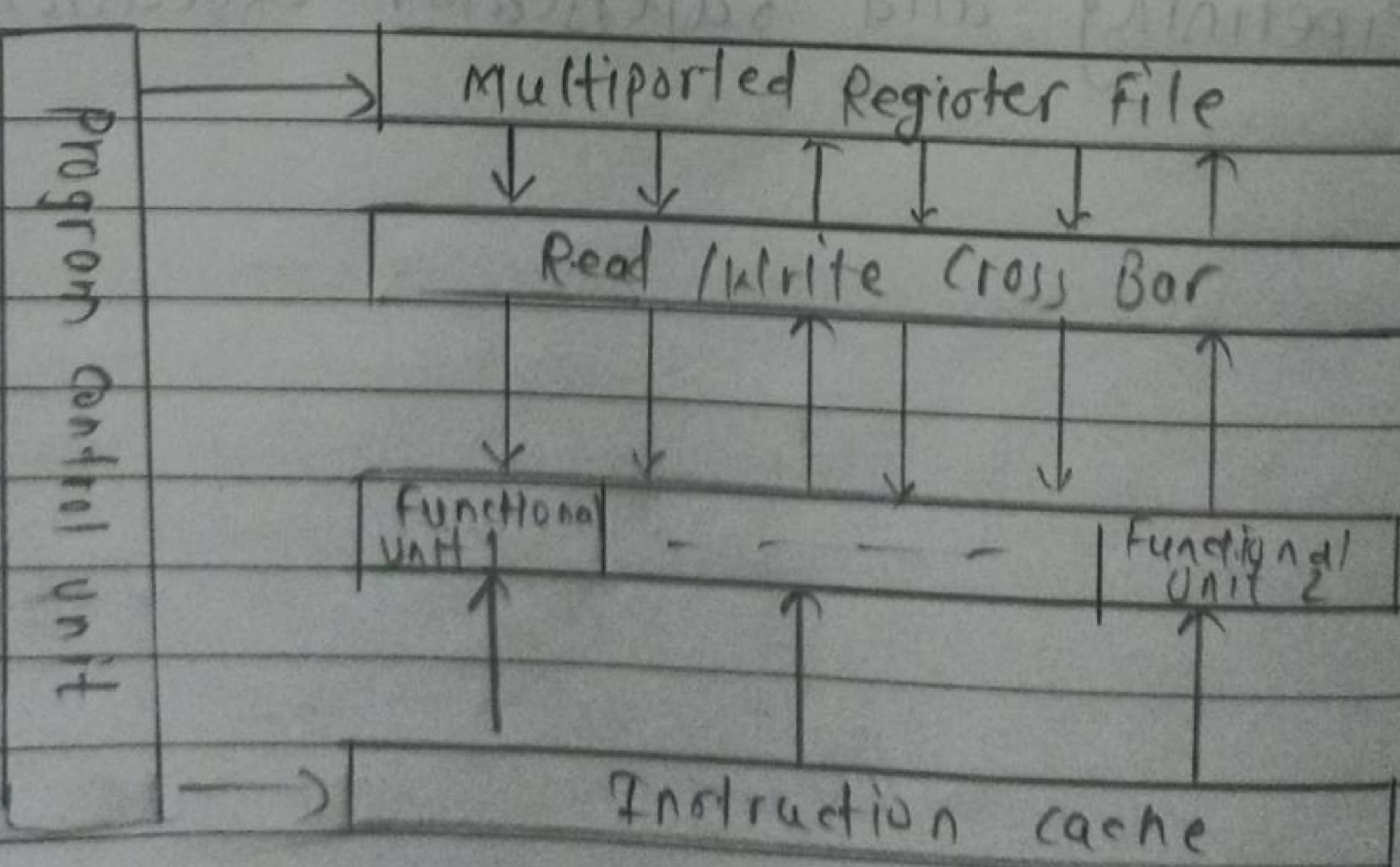


fig. VLIW Architecture

Advantages of VLIW Architecture.

- It can improve performance
- It is used to potentially scalable.
- More implementation units can be inserted and so more instructions can be overflowing into the VLIW instruction

Disadvantages

- It can be used by a new programmer required.
- The program should keep track of instruction scheduling
- It can increase memory use.
- It is used for high power consumption.

(Q) Explain dichotomy of parallel computing platforms.

→ A dichotomy is based on the logical and physical organization of parallel platforms.

- physical organisation is the actual hardware organisation of a platform whereas logical organisation refers to a programmer's view of the platform.
- From programmer's perspective the two important components of parallel computing are.

① control structure - The various way of expressing parallel tasks is known as control structure.

② Communication model - The mechanisms for specifying interaction between the parallel tasks is called as communication model.

(a) Explain typical SIMD & MIMD processors architecture.

→ SIMD and MIMD are types of computer architectures that are used to improve the performance of certain types of computational tasks.

- Michael Flynn classified the computer organization into SIMD and MIMD, where SIMD stands for single Instruction Multiple data, MIMD stands for (multiple Instruction Multiple data)

- SIMD short for single Instruction Multiple data, computer architecture can execute a single instruction on multiple data streams.

- MIMD computer architectures can execute several instructions on multiple data streams

SIMD - is a form of parallel computer architecture that is categorized under

Flynn's classification given by Michael Flynn.

- In SIMD architecture, a single instruction is applied to several data streams.

SIMD consists of a single control signal that is used to call several isolated processing units. Therefore, all the processing units accept the same instruction from the control unit and use it on separate element of data.

The SIMD organization uses shared memory unit which is divided into different modules.

SIMD As a result the memory unit can interact with all the processing units simultaneously. Since the SIMD architectures use a single copy of instruction on multiple data streams it requires less memory. SIMD architectures are particularly effective for tasks that can be easily parallelized such as Image processing, video encoding and decoding.

PE - Processing Element

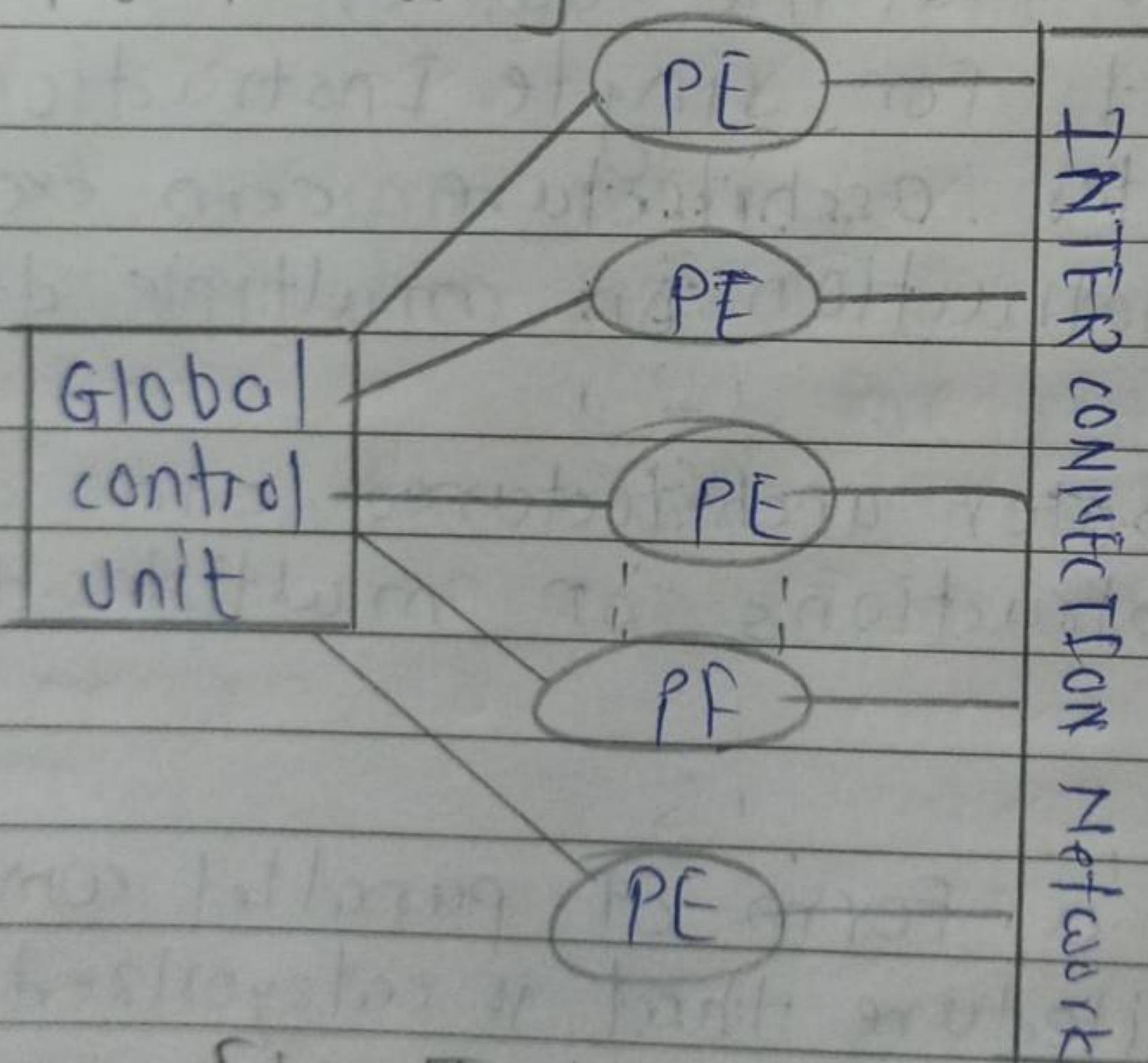


Fig. Typical SIMD Architecture.

MIMD - The MIMD architecture consists of multiple instructions and data streams. Therefore MIMD architecture requires multiple processing units. MIMD systems are considered to have most complex organisation, but they are highly efficient.

- The SIMD architecture does not need any addition control unit which reduces the cost of the system.
- It ^{also} provides efficient execution of the conditional statement such as if/else statements.
- SIMD architectures are more flexible and are better suited for tasks that require more complex and varied computation such as general-purpose computing and AI Applications.

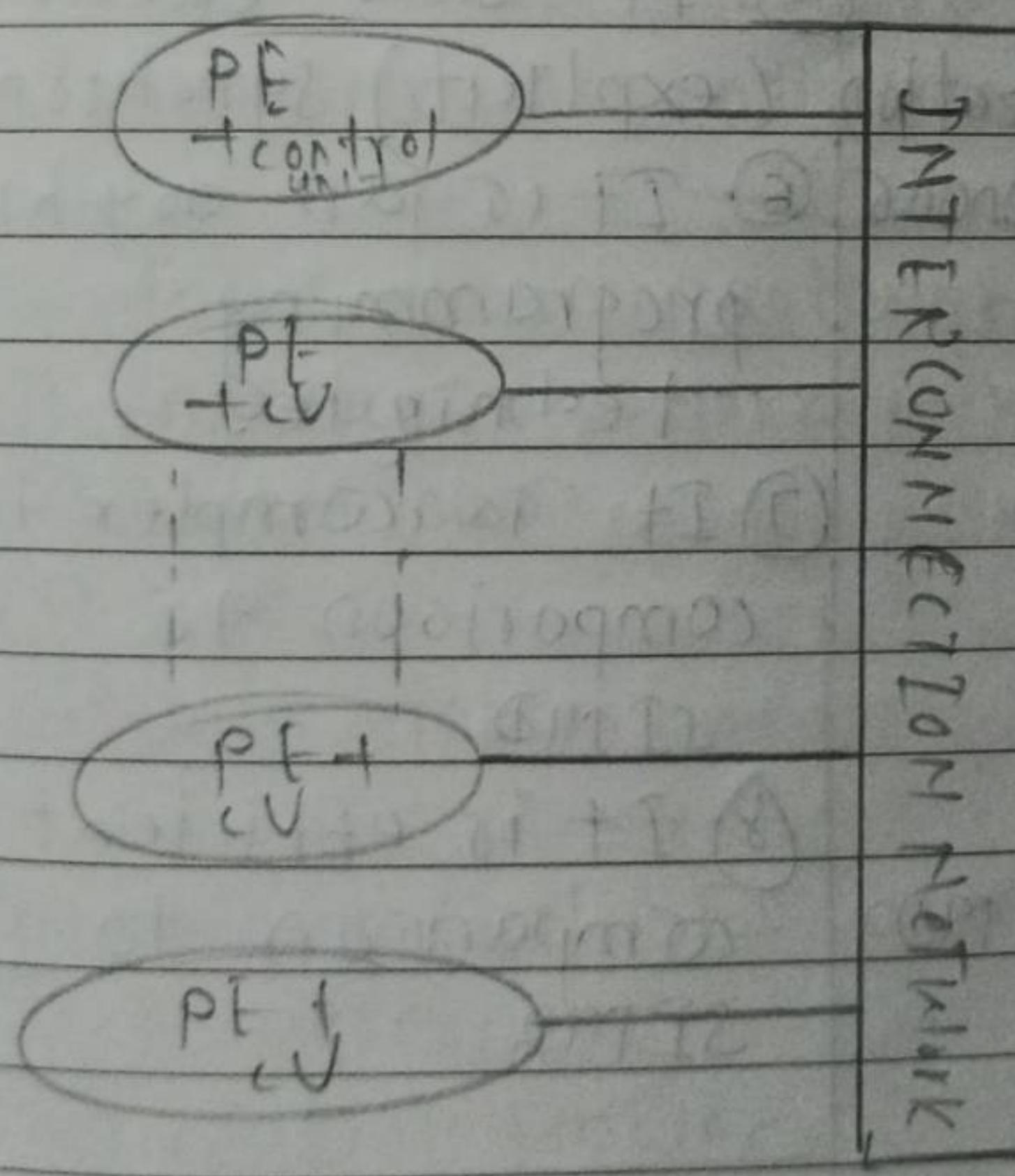


Fig. SIMD Architecture.

a) difference between SIMD & MIMD

SIMD

- ① SIMD stands for single Instruction multiple Data.
- ② It requires less memory.
- ③ It is inexpensive in comparison to MIMD.
- ④ It has a single decoder.
- ⑤ It uses latent (tacit) synchronization.
- ⑥ It is a synchronous programming technique.
- ⑦ It is simple in comparison to MIMD.
- ⑧ It is not as efficient as MIMD in terms of performance.

MIMD

- ① MIMD stands for multiple Instruction multiple Data.
- ② It requires more memory.
- ③ It is expensive in comparison to SIMD.
- ④ It contains multiple decoders.
- ⑤ It uses accurate (explicit) synchronization.
- ⑥ It is an asynchronous programming technique.
- ⑦ It is complex in comparison to SIMD.
- ⑧ It is efficient in comparison to SIMD.

(a) Explain Data flow diagrams

→ Functional modeling is represented through a hierarchy of DFDs.

- The DFD is a graphical representation of a system that shows the inputs to the system, the processing upon the inputs, the outputs of the system as well as the internal data stores.

→ DFD is A^{data} flow diagram is a graph which shows the flow of data values from their sources in objects through processes that transform them to their destinations on other objects".

The four main parts of DFD are,

① processes

② Data Flows

③ Actors

④ Data stores.

other parts of DFD are,

① constraints

② control flows.

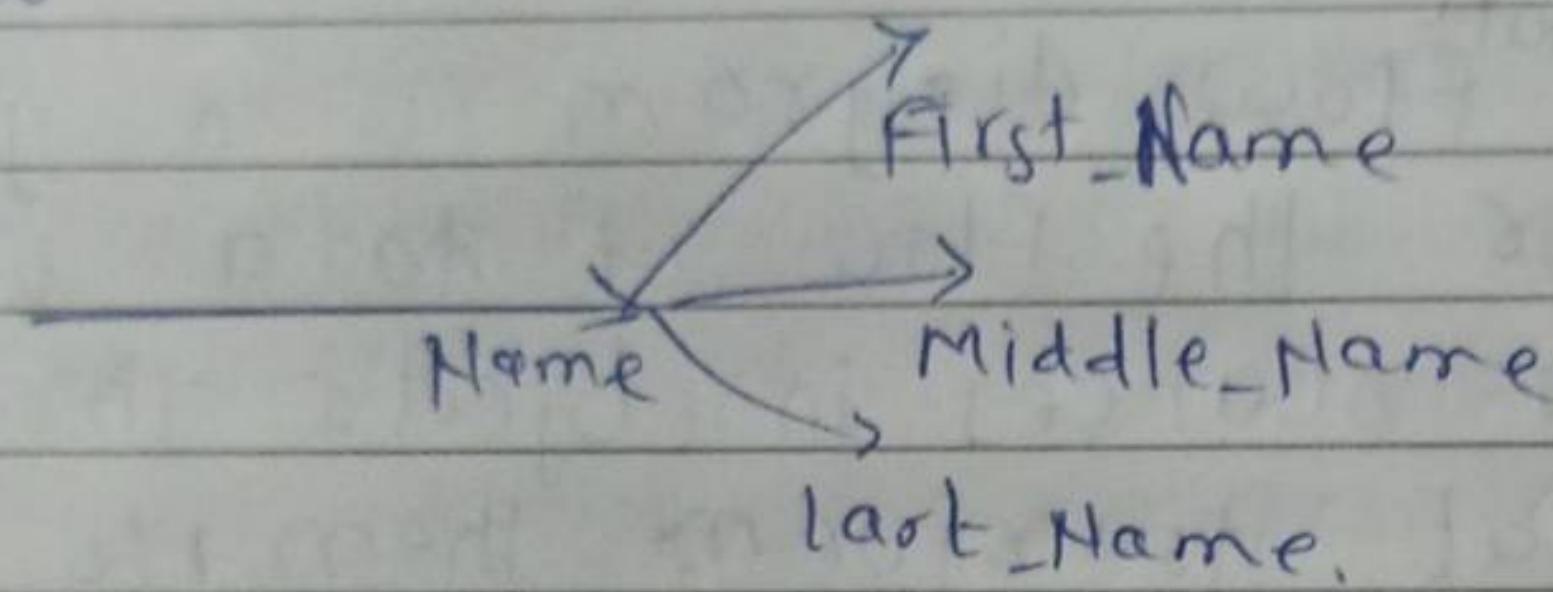
Features of a DFD

① Processes - processes are the computational activities that transform data values.

- A process may be further divided into smaller components. The lowest level process may be a simple function.

② Dataflows -

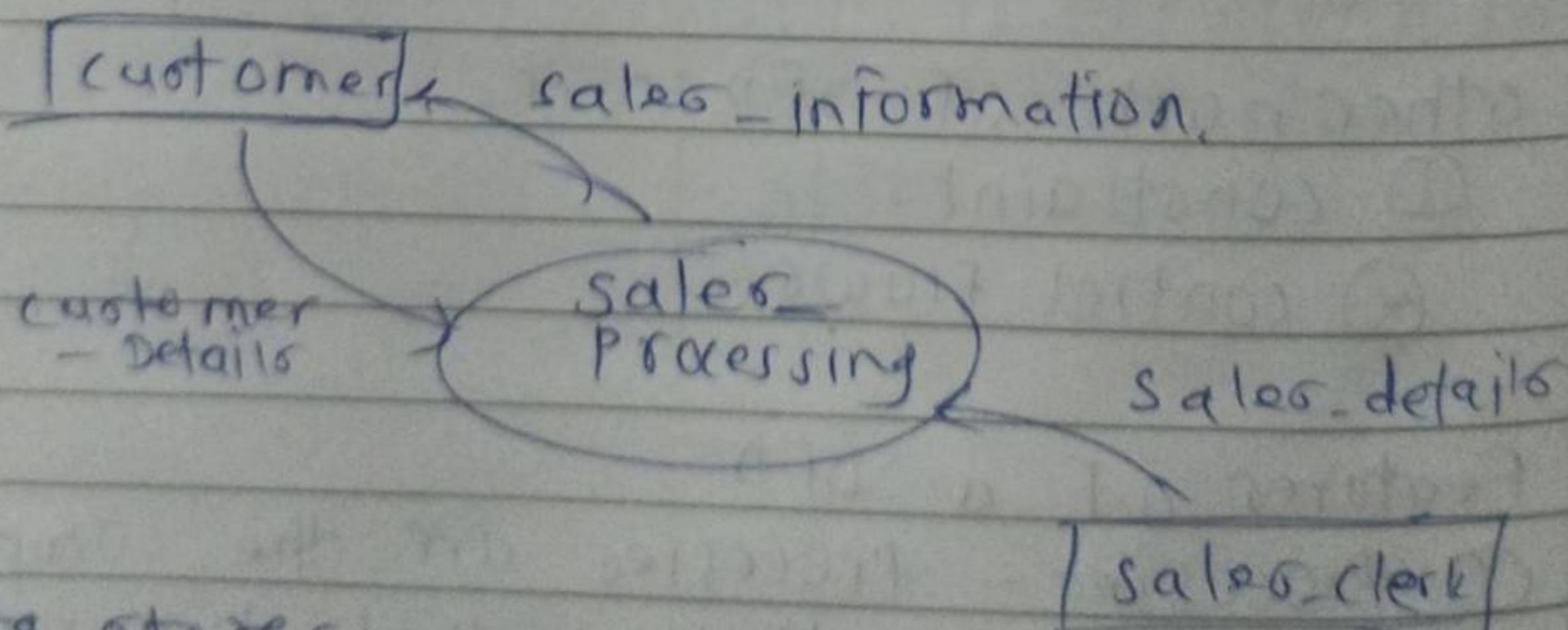
- Data flows represent the flow of data between two processes.
- It could be between an actor and a process.
- A data flow is represented by a directed arc or an arrow, labeled with name of the data item that it carries.



③ Actors

Actors are the active objects that interact with the system by either producing data and inputting them to the system.

- An Actor is represented by a rectangle.



④ Data stores

- Data stores are the passive objects that act as a repository of data.
- A data store is represented by two parallel lines containing the name of the data store.

- (a) N-wide Superscalar Architectures.
- The combination of temporal parallelism and data parallelism by issuing several instructions simultaneously in each cycle, to improve the speed of a processor, is called superscalar processing.
 - In an n-issue superscalar, n instructions are fetched, decoded, executed and committed per cycle.
 - \$

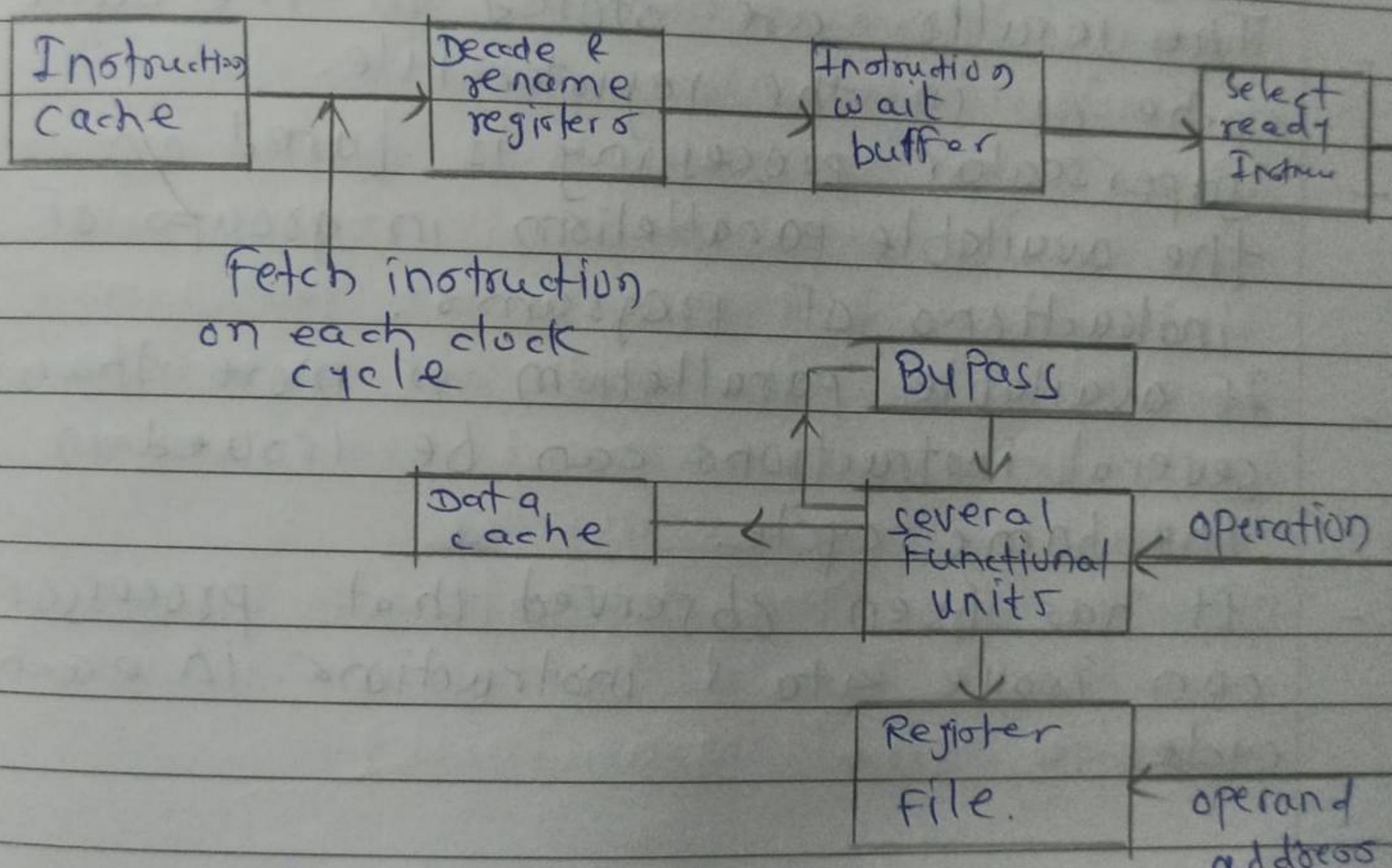


Fig. Superscalar processor pipeline stages

- The fetch unit fetches several instructions in each clock cycle from the instruction cache.
- The decode unit decodes the instructions and renamer registers so that false dependencies are eliminated and only true dependencies remain.

- The instructions are sent to an instruction buffer where they are kept temporarily till the source operands are identified and become available.
- Ready instructions are dispatched based on the availability of functional units
- The operands are fetched either from the register file or from bypass path from earlier instructions which produce these operands.
- The results are stored in the data cache or in the register file.
- superscalar processing is based on the available parallelism in groups of instructions of programs.
- If available parallelism is more than several instructions can be issued in the same cycle.
- It has been observed that processors can issue 4 to 5 instructions in each cycle.