

MODULE-1

- **The State of Computing**
- **Multiprocessors and Multicomputer**
- **Multivector and SIMD Computers, PRAM and VLSI Models**
- **Program and Network Properties**
- **Conditions of Parallelism, Program Partitioning and Scheduling**
- **Program Flow Mechanisms, System Interconnect Architectures**
- **Principles of Scalable Performance, Performance Metrics and Measures**
- **Parallel Processing Applications**
- **Speedup Performance Laws**
- **Scalability Analysis and Approaches.**

1. The State of Computing

- Computers have gone through two major stages of development: mechanical and electronic. Prior to 1945, computers were made with mechanical or electromechanical parts. The earliest mechanical computer can be traced back to 500 BC in the form of the abacus used in China.

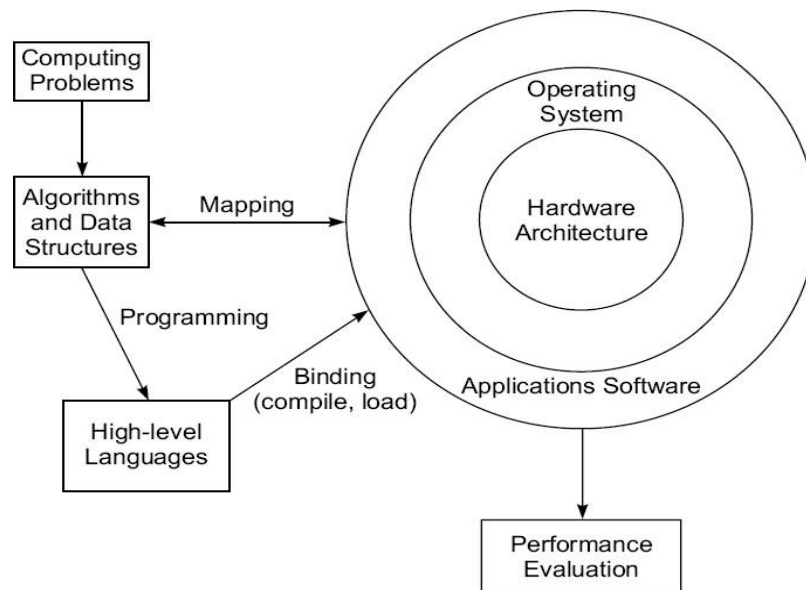
Computer Generations:

- Over the past five decades, electronic computers have gone through five generations of development. Each of the first three generations lasted about 10 years. The fourth generation covered a time span of 15 years. We have just entered the fifth generation with the use of processors and memory devices with more than 1 million transistors on a single silicon chip.

Table 1-1 Five **Generations** of Electronic Computers

Generation	Technology and Architecture	Software and Applications	Representative Systems
First (1945-54)	Vacuum tubes and relax memories, CPU driven by PC and accumulator, fixed-point arithmetic.	Machine/assembly languages, single user, no sub -routine linkage, programmed I/O using CPU.	ENIAC , Princeton IAS, IBM 701.
Second (1955-64)	Discrete transistors and core memories, floating-point arithmetic, I/O processors, multiplexed memory access.	HLL used with compilers, subroutine libraries, batch processing monitor .	IBM 7090, CDC 1604, Univac LARC .
Third (1965-74)	Integrated circuits (SSI/-MSI), microprogramming, pipelining, cache, and lookahead processors.	Multiprogramming and time-sharing OS, multiuser applications.	IBM 360/370, CDC 6600, TI-ASC, PDP-8.
Fourth (1975-90)	LSI/VLSI and semiconductor memory, multiprocessors, vector supercomputers, multicomputers.	Multiprocessor OS, languages, compilers, and environments for parallel processing.	VAX 9000, Cray X-MP, IBM 3090, BBN TC2000.
Fifth (1991-present)	ULSI/VHSIC processors, memory, and switches, high-density packaging, scalable architectures.	Massively parallel processing, grand challenge applica tions, heterogeneous processing.	Fujitsu VPP500, Cray/MPP , TMC/CM-5, Intel Paragon.

1.1.1 Elements of Modern Computers

**Fig. 1.1** Elements of a modern computer system

- Computing Problem
- Algorithms and Data Structure
- High Level languages
- Hardware resources
 - ✓ System Software
 - ✓ Application Software
 - ✓ Operating System
- Compiler Support
 - **Computing Problems**
 - The use of a computer is driven by real-life problems demanding fast and accurate solutions. Depending on the nature of the problems, the solutions may require different computing resources.
 - For numerical problems in science and technology, the solutions demand complex mathematical formulations and tedious integer or floating-point computations.
 - For alpha numerical problems in business and government, the solutions demand accurate transactions, large database management, and information retrieval operations.
 - For artificial intelligence (AI) problems, the solutions demand logic inferences and symbolic manipulations.
 - These computing problems have been labeled *numerical computing*, *transaction processing*, and *logical reasoning*.
 - Some complex problems may demand a combination of these processing modes.
 - **Hardware Resources**
 - A modern computer system demonstrates its power through coordinated efforts by hardware

resources, an operating system, and **application** software.

- Processors, memory, and peripheral devices form the hardware core of a computer system.
- Special hardware interfaces are often built into I/O devices, such as terminals, workstations, optical page scanners, magnetic ink character recognizers, modems, file servers, voice data entry, printers, and plotters.
 - These peripherals are connected to mainframe computers directly or through local or wide-area networks.
 - Operating System
- An effective operating system manages the allocation and deallocation of resources during the execution of user programs.
- Beyond the OS, application software must be developed to benefit the users.
- Standard benchmark programs are needed for performance evaluation.
- Mapping is a bidirectional process matching algorithmic structure with hardware architecture, and vice versa.
- Efficient mapping will benefit the programmer and produce better source codes.
- The mapping of algorithmic and data structures onto the machine architecture includes processor scheduling, memory maps, interprocessor communications, etc.
- These activities are usually architecture-dependent.
 - **System Software Support**
- Software support is needed for the development of efficient programs in high-level languages. The source code written in a HLL must be first translated into object code by an optimizing compiler.
- The *compiler* assigns variables to registers or to memory words and reserves functional units for operators.
- An *assembler* is used to translate the compiled object code into machine code which can be recognized by the machine hardware. A *loader* is used to initiate the program execution through the OS kernel.

- **Compiler Support**

There are three compiler upgrade approaches:

- **Preprocessor:** A preprocessor uses a sequential compiler and a low-level library of the target computer to implement high-level parallel constructs.
- **Precompiler:** The precompiler approach requires some program flow analysis, dependence

checking, and limited optimizations toward parallelism detection.

- **Parallelizing Compiler:** This approach demands a fully developed parallelizing or vectorizing compiler which can automatically detect parallelism in source code and transform sequential codes into parallel constructs.

1.1.2 Evolution of Computer Architecture

- The study of computer architecture involves both hardware organization and programming/software requirements. As seen by an assembly language programmer, computer architecture is abstracted by its instruction set, which includes opcode (operation codes), addressing modes, registers, virtual memory, etc. From the hardware implementation point of view, the abstract machine is organized with CPUs, caches, buses, microcode, pipelines, physical memory, etc. Therefore, the study of architecture covers both instruction-set architectures and machine implementation organizations.

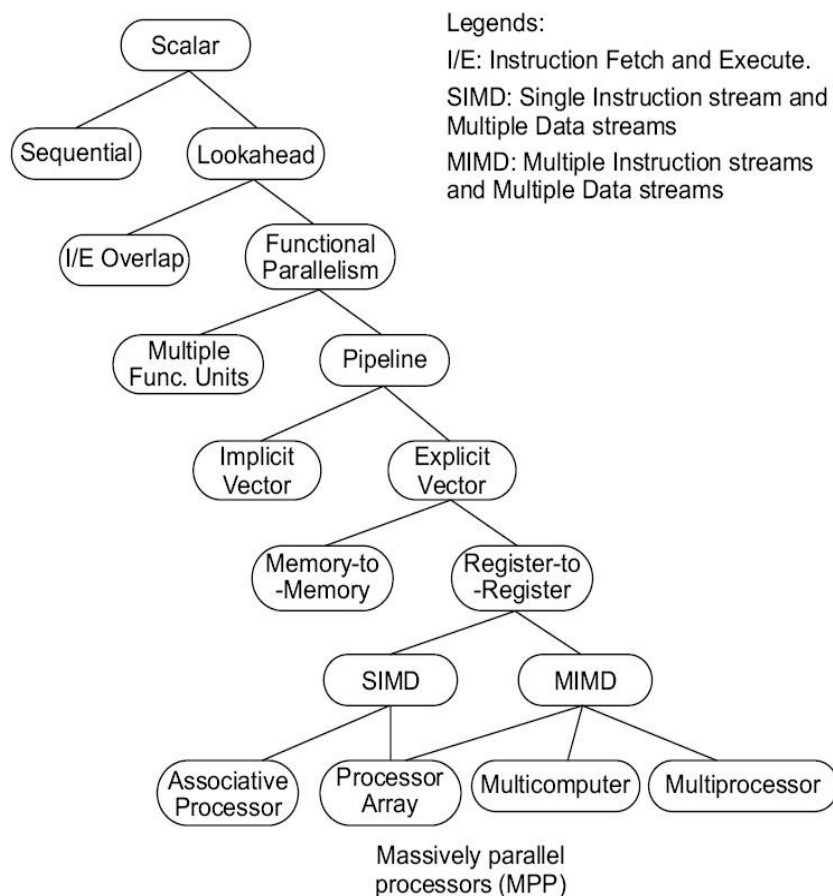


Fig. 1.2 Tree showing architectural evolution from sequential scalar computers to vector processors and parallel computers

Lookahead, Parallelism, and Pipelining

Lookahead techniques were introduced to prefetch instructions in order to overlap I/E (instruction fetch/decode and execution) operations and to enable functional parallelism. Functional parallelism was supported by two approaches:

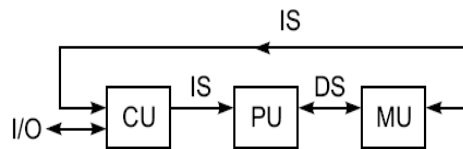
1. using multiple functional units simultaneously,
2. to practice pipelining at various processing levels.

Flynn's Classification

Michael Flynn (1972) introduced a classification of various computer architectures based on notions of instruction and data streams.

1. **SISD** (Single Instruction stream over a Single Data stream) computers
2. **SIMD** (Single Instruction stream over Multiple Data streams) machines
3. **MIMD** (Multiple Instruction streams over Multiple Data streams) machines.
4. **MISD** (Multiple Instruction streams and a Single Data stream) machines

1. SISD (Single Instruction stream over a Single Data stream) computers

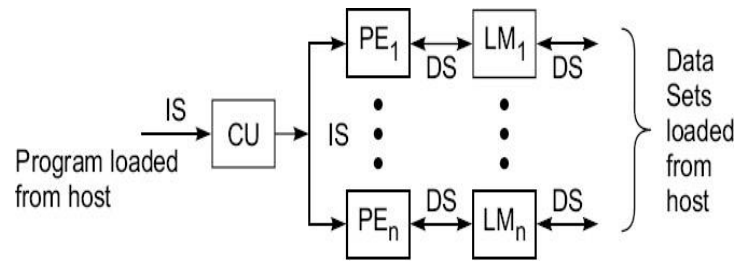


(a) SISD uniprocessor architecture

- Conventional sequential machines are called SISD computers.
- They are also called scalar processor i.e., one instruction at a time and each instruction have only one set of operands.
- Single instruction: only one instruction stream is being acted on by the CPU during any one clock cycle
- Single data: only one data stream is being used as input during any one clock cycle
- Deterministic execution
- Instructions are executed sequentially.
- This is the oldest and until recently, the most prevalent form of computer

□ Examples: most PCs, single CPU workstations and mainframes.

2. SIMD (Single Instruction stream over Multiple Data streams) machines



(b) SIMD architecture (with distributed memory)

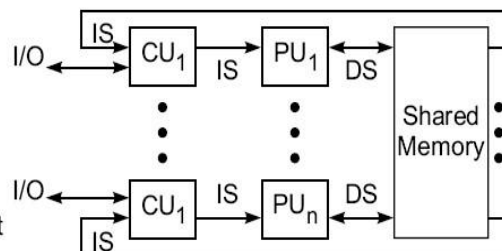
A type of parallel computer

- **Single instruction:** All processing units execute the same instruction issued by the control unit at any given clock cycle.
- **Multiple data:** Each processing unit can operate on a different data element. The processors are connected to shared memory or interconnection network providing multiple data to processing unit.
- This type of machine typically has an instruction dispatcher, a very high-bandwidth internal network, and a very large array of very small-capacity instruction units.
- Thus single instruction is executed by different processing unit on different set of data.
- Best suited for specialized problems characterized by a high degree of regularity, such as image processing and vector computation.
- Synchronous (lockstep) and deterministic execution.
- Two varieties: Processor Arrays e.g., Connection Machine CM-2, Maspar MP-1, MP-2 and Vector Pipelines processor e.g., IBM 9000, Cray C90, Fujitsu VP, NEC SX-2, Hitachi S820

3. MIMD (Multiple Instruction streams over Multiple Data streams) machines.

Captions:

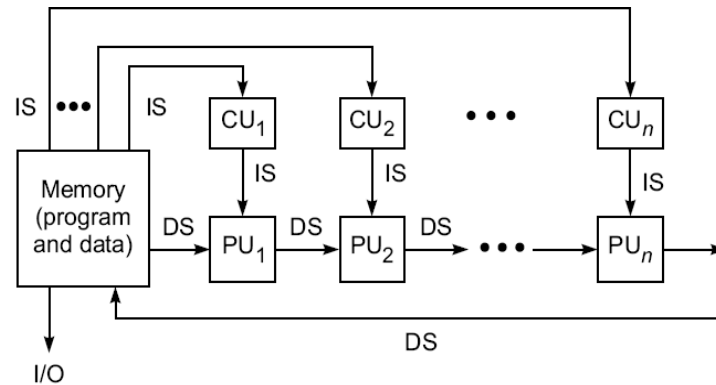
CU = Control Unit
 PU = Processing Unit
 MU = Memory Unit
 IS = Instruction Stream
 DS = Data Stream
 PE = Processing Element
 LM = Local Memory



(c) MIMD architecture (with shared memory)

- A single data stream is fed into multiple processing units.
- Each processing unit operates on the data independently via independent instruction streams.
- A single data stream is forwarded to different processing unit which are connected to different control unit and execute instruction given to it by control unit to which it is attached.

- Thus in these computers same data flow through a linear array of processors executing different instruction streams.
 - This architecture is also known as **Systolic Arrays** for pipelined execution of specific instructions. Some conceivable uses might be:
 1. multiple frequency filters operating on a single signal stream
 2. multiple cryptography algorithms attempting to crack a single coded message.
4. MISD (Multiple Instruction streams and a Single Data stream) machines



(d) MISD architecture (the systolic array)

Fig. 1.3 Flynn's classification of computer architectures (Derived from Michael Flynn,

- **Multiple Instructions:** Every processor may be executing a different instruction stream
- **Multiple Data:** Every processor may be working with a different data stream, multiple data stream is provided by shared memory.
- Can be categorized as loosely coupled or tightly coupled depending on sharing of data and control.
- Execution can be synchronous or asynchronous, deterministic or non-deterministic
- There are multiple processors each processing different tasks.
- Examples: most current supercomputers, networked parallel computer "grids" and multi-processor SMP computers - including some types of PCs.

Development Layers

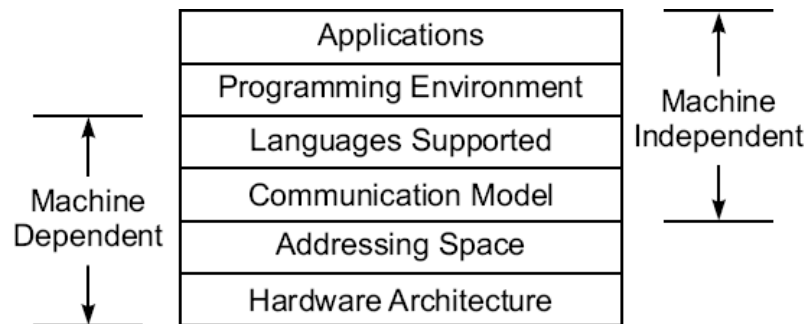


Fig. 1.4 Six layers for computer system development (Courtesy of Lionel Ni, 1990)

Development Layers A layered development of parallel computers is illustrated in Fig. 1.4.

- **Hardware configurations** differ from machine to machine, even those of the same model.
- **Address Space** of a processor in a computer system varies among different architectures. It depends on the memory organization, which is machine-dependent. These features are up to the designer and should match the target application domains.
- We want to develop **Application Programs** and **Programming Environments** which are machine-independent. Independent of machine architecture, the user programs can be ported to many computers with minimum conversion costs.
- **High-level languages** and **Communication Models** depend on the architectural choices made in a computer system. From a programmer's viewpoint, these two layers should be architecture-transparent.
- At present, Fortran, C, Pascal, Ada, and Lisp are supported by most computers.
- However, the **Communication Models**, shared variables versus message passing, are mostly machine-dependent. The Linda approach using *tuple spaces* offers architecture transparent Communication model for parallel computers.
- Application programmers prefer more architectural transparency. However, kernel programmers have to explore the opportunities supported by hardware.
- As a good computer architect, one has to approach the problem from both ends.
- The compilers and OS support should be designed to remove as many architectural constraints as possible from the programmer.

1.1.3 System Attributes affecting Performance

Clock Rate and CPI

- The CPU (or simply the *processor*) of today's digital computer is driven by a clock with a constant

cycle time (τ in nanoseconds).

- The inverse of the cycle time is the *clock rate* ($f = 1/\tau$ in megahertz). The size of a program is determined by its *instruction count* (I_c), in terms of the number of machine instructions to be executed in the program.
- Different machine instructions may require different numbers of clock cycles to execute. Therefore, the *cycles per instruction* (CPI) becomes an important parameter for measuring the time needed to execute each instruction.
- For a given instruction set, we can calculate an *average* CPI over all instruction types, provided we know their frequencies of appearance in the program.
- An accurate estimate of the average CPI requires a large amount of program code to be traced over a long period of time.
- Unless specifically focusing on a single instruction type, we simply use the term CPI to mean the average value with respect to a given instruction set and a given program mix.

Performance Factors

- Let I_c be the number of instructions in a given program, or the instruction count.
- The CPU time (T in seconds/program) needed to execute the program is estimated by finding the product of three contributing factors:

$$T = I_c \times \text{CPI} \times \tau \quad (1.1)$$

- The execution of an instruction requires going through a cycle of events involving the instruction fetch, decode, operand(s) fetch, execution, and store results.
- In this cycle, only the instruction decode and execution phases are carried out in the CPU.
- The remaining three operations may be required to access the memory. We define a *memory cycle* as the time needed to complete one memory reference.
- Usually, a memory cycle is k times the processor cycle τ .
- The value of k depends on the speed of the memory technology and processor-memory interconnection scheme used.
- The CPI of an instruction type can be divided into two component terms corresponding to the total processor cycles and memory cycles needed to complete the execution of the instruction.
- Depending on the instruction type, the complete instruction cycle may involve one to four memory references (one for instruction fetch, two for operand fetch, and one for store results). Therefore we can rewrite Eq. 1.1 as follows;

$$T = I_c \times (p + m \times k) \times \tau \quad (1.2)$$

where p is the number of processor cycles needed for the instruction decode and execution,
 m is the number of memory references needed,
 k is the ratio between memory cycle and processor cycle,
 I_c is the instruction count,
 r is the processor cycle time.

Equation 1.2 can be further refined once the CPI components (p, m, k) are weighted over the entire instruction set.

System Attributes

- The above five performance factors (I_c, p, m, k, τ) are influenced by four system attributes: instruction-set architecture, compiler technology, CPU implementation and control, and cache and memory hierarchy, as specified in Table 1.2.
- The instruction-set architecture affects the program length (I_c) and processor cycle needed (p). The compiler technology affects the values of I_c, p and the memory reference count (m).
- The CPU implementation and control determine the total processor time ($p \cdot \tau$) needed.
- Finally, the memory technology and hierarchy design affect the memory access latency ($k \cdot \tau$). The above CPU time can be used as a basis in estimating the execution rate of a processor.

Table 1.2 Performance Factors Versus System Attributes

System Attributes	Performance Factors				
	Instr. Count,	Average Cycles per Instruction, CPI			Processor Cycle Time, T
		Processor Cycles per Instruction, p	Memory References per Instruction, m	Memory-Access Latency, k	
Instruction-set Architecture	X	X			
Compiler Technology	X	X	X		
Processor Implementation and Control		X			X
Cache and Memory Hierarchy				X	X

MIPS Rate

- Let C be the total number of clock cycles needed to execute a given program.
- Then the CPU time in Eq. 1.2 can be estimated as $T = C \times \tau = C/f$.
- Furthermore, $CPI = C/I_c$ and $T = I_c \times CPI \times \tau = I_c \times CPI/f$. The processor speed is often measured in

terms of *million instructions per second* (MIPS).

- We simply call it the MIPS rate of a given processor. It should be emphasized that the MIPS rate varies with respect to a number of factors, including the clock rate (f), the instruction count (I_c), and the CPI of a given machine, as defined below:

$$\text{MIPS rate} = \frac{I_c}{T \times 10^6} = \frac{f}{\text{CPI} \times 10^6} = \frac{f \times I_c}{C \times 10^6} \quad (1.3)$$

- Based on Eq. 1.3, the CPU time in Eq. 1.2 can also be written as $T = I_c \times 10^{-6} / \text{MIPS}$.
- Based on the system attributes identified in Table 1.2 and the above derived expressions, we conclude by indicating the fact that the MIPS rate of a given computer is directly proportional to the clock rate and inversely proportional to the CPI.
- All four system attributes, instruction set, compiler, processor, and memory technologies, affect the MIPS rate, which varies also from program to program.

Throughput Rate

- Number of programs a system can execute per unit time, called the *system throughput* W_s (in programs/second).
- In a multiprogrammed system, the system throughput is often lower than the *CPU throughput* W_p defined by:

$$W_p = \frac{f}{I_c \times \text{CPI}} \quad (1.4)$$

- Note that $W_p = (\text{MIPS}) \times 10^6 / I_c$ from Eq. 1.3- The unit for W_p is programs/second.
- The CPU throughput is a measure of how many programs can be executed per second, based on the MIPS rate and average program length (I_c).
- The reason why $W_s < W_p$ is due to the additional system overheads caused by the I/O, compiler, and OS when multiple programs are interleaved for CPU execution by multiprogramming or timesharing operations.
- If the CPU is kept busy in a perfect program-interleaving fashion, then $W_s = W_p$. This will probably never happen, since the system overhead often causes an extra delay and the CPU may be left idle for some cycles.

Programming Environments

- Programmability depends on the programming environment provided to the users.

- Conventional computers are used in a sequential programming environment with tools developed for a uniprocessor computer.
- Parallel computers need parallel tools that allow specification or easy detection of parallelism and operating systems that can perform parallel scheduling of concurrent events, shared memory allocation, and shared peripheral and communication links.

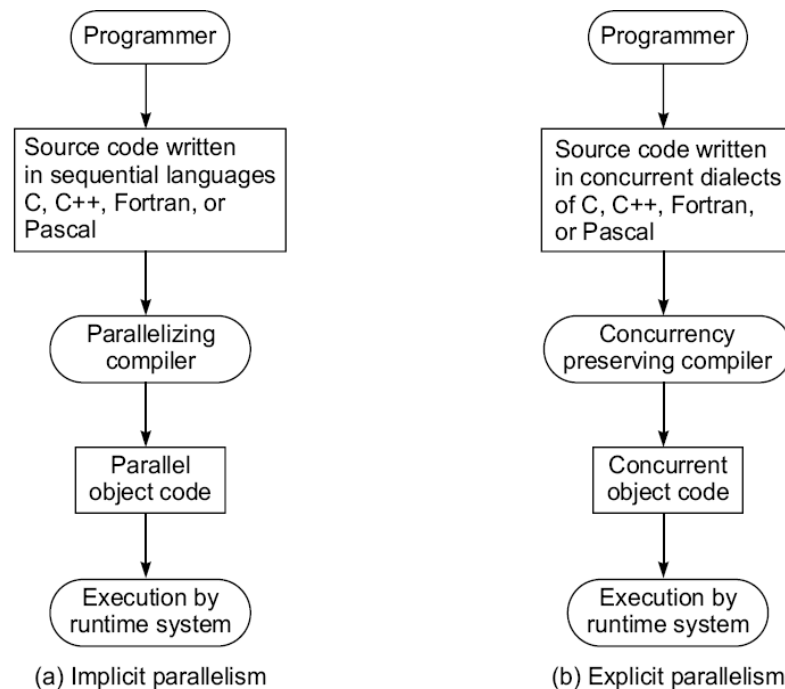


Fig. 1.5 Two approaches to parallel programming (Courtesy of Charles Seitz; adapted with permission from

Implicit Parallelism

- An implicit approach uses a conventional language, such as C, Fortran, Lisp, or Pascal, to write the source program.
- The sequentially coded source program is translated into parallel object code by a parallelizing compiler.
- The compiler must be able to detect parallelism and assign target machine resources. This compiler approach has been applied in programming shared-memory multiprocessors.
- With parallelism being implicit, success relies heavily on the "intelligence" of a parallelizing compiler.
- This approach requires less effort on the part of the programmer.

Explicit Parallelism

- The second approach (Fig. 1.5b) requires more effort by the programmer to develop a source program **using parallel** dialects of C, Fortran, Lisp, or Pascal.

- Parallelism is explicitly specified in the user programs.
- This will significantly reduce the burden on the compiler to detect parallelism.
- Instead, the compiler needs to preserve parallelism and, where possible, assigns target machine resources.

1.2 Multiprocessors and Multicomputer

Two categories of parallel computers are architecturally modeled below. These physical models are distinguished by having a shared common memory or unshared distributed memories.

1. Shared-Memory Multiprocessors

There are 3 shared-memory multiprocessor models:

- i. Uniform Memory-access (UMA) model,
- ii. Non uniform-Memory-access (NUMA) model
- iii. Cache-Only Memory Architecture (COMA) model.

These models differ in how the memory and peripheral resources are **shared** or distributed.

i. Uniform Memory-Access (UMA) model

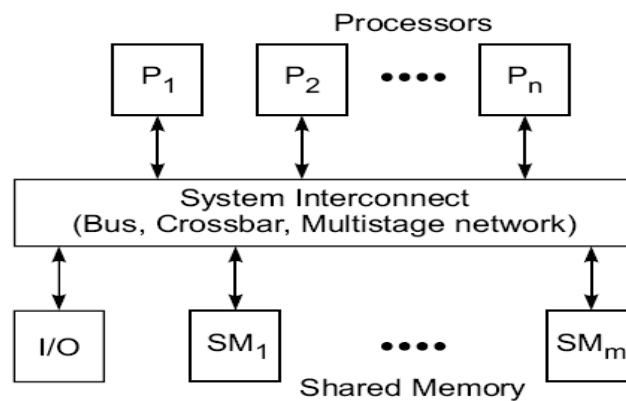


Fig. 1.6 The UMA multiprocessor model

- In a UMA multiprocessor model (Fig. 1.6), the physical memory is uniformly shared by all the processors.
- All processors have equal access time to all memory words, which is why it is called uniform memory access.
- Each processor may use a private cache. Peripherals are also shared in some fashion.
- Multiprocessors are called *tightly coupled systems* due to the high degree of resource sharing. The system interconnect takes the form of a common bus, a crossbar switch, or a multistage network.

- Most computer manufacturers have *multiprocessor* (MP) extensions of their *uniprocessor* (UP) product line.
- The UMA model is suitable for general-purpose and timesharing applications by **multiple** users. It can be used to speed up the execution of a single large program in time-critical applications. To coordinate parallel events, synchronization and communication among processors are done through using shared variables in the common memory.
- When all processors have equal access to all peripheral devices, the system is called a **symmetric multiprocessor**. In this case, all the processors are equally capable of running the executive programs, such as the OS kernel and I/O service routines.

ii. Non uniform-Memory-Access (NUMA) model

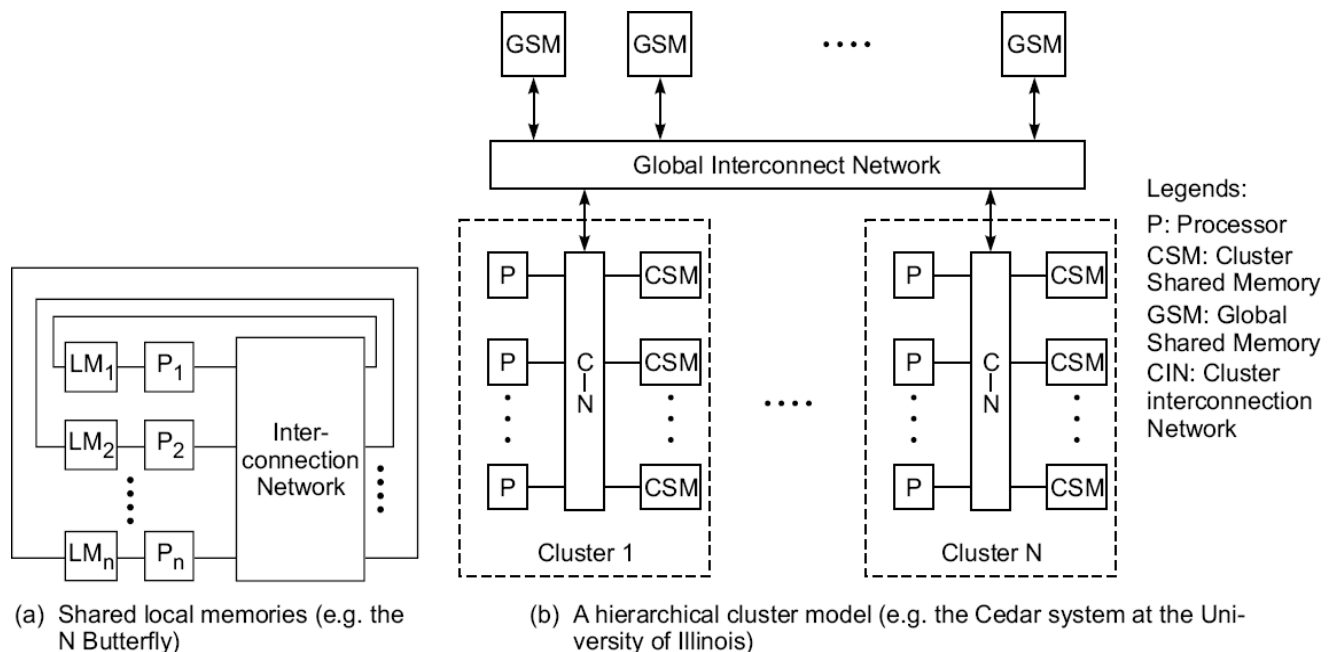


Fig. 1.7 Two NUMA models for multiprocessor systems

- A NUMA multiprocessor is a shared-memory system in which the access time varies with the location of the memory word.
- Two NUMA machine models are depicted in Fig. 1.7.
- The shared memory is physically distributed to all processors, called *local memories*.
- The collection of all local memories forms a global address space accessible by all processors.
- It is faster to access a local memory with a local processor. The access of remote memory attached to other processors takes longer due to the added delay through the interconnection network.

- The BBN TC-2000 Butterfly multiprocessor assumes the configuration shown in Fig. 1.7a.

iii. Cache-Only Memory Architecture (COMA) model

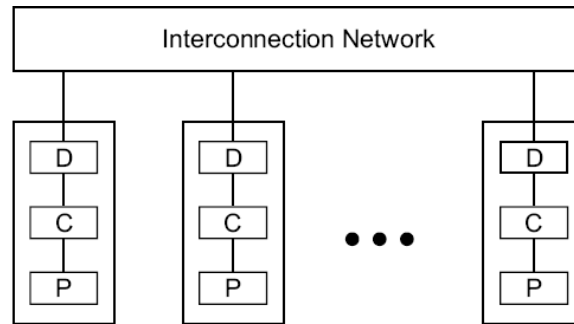


Fig. 1.8 The COMA model of a multiprocessor (P: Processor, C: Cache, D: Directory; e.g. the KSR-1)

- A multiprocessor using cache-only memory assumes the COMA model.
- Examples of COMA machines include the Swedish Institute of Computer Science's Data Diffusion Machine and Kendall Square Research's KSR-1 machine.
- The COMA model is a special case of a NUMA machine, in which the distributed main memories are converted to caches.
- There is no memory hierarchy at each processor node. All the caches form a global address space.
- Remote cache access is assisted by the distributed cache directories (D in Fig. 1.8).
- Depending on the interconnection network used, sometimes hierarchical directories may be used to help locate copies of cache blocks.
- Initial data placement is not critical because data will eventually migrate to where it will be used.

2. Distributed-Memory Multicomputers

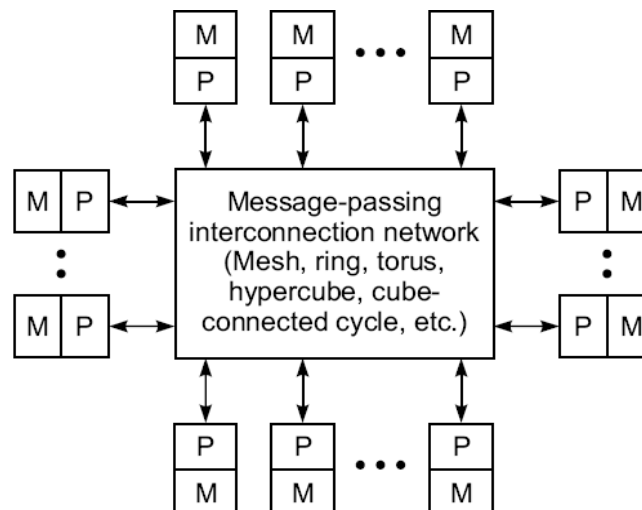


Fig. 1.9 Generic model of a message-passing multicomputer

- A distributed-memory multicomputer system is modeled in the above figure consists of multiple computers, often called *nodes*, interconnected by a message-passing network.
- Each node is an autonomous computer consisting of a processor, local memory, and sometimes attached disks or I/O peripherals.
- The message-passing network provides point-to-point static connections among the nodes.
- All local memories are private and are accessible only by local processors.
- For this reason, traditional multicomputers have been called *no-remote-memory-access* (NORMA) machines.
- However, this restriction will gradually be removed in future multi computers with distributed shared memories. Internode communication is carried out by passing messages through the static connection network.

1.3 Multivector and SIMD Computers

We can classify super computers as:

- Pipelined vector machines** using a few powerful processors equipped with vector hardware
- SIMD computers** emphasizing massive data parallelism

1.3.1 Vector Supercomputers

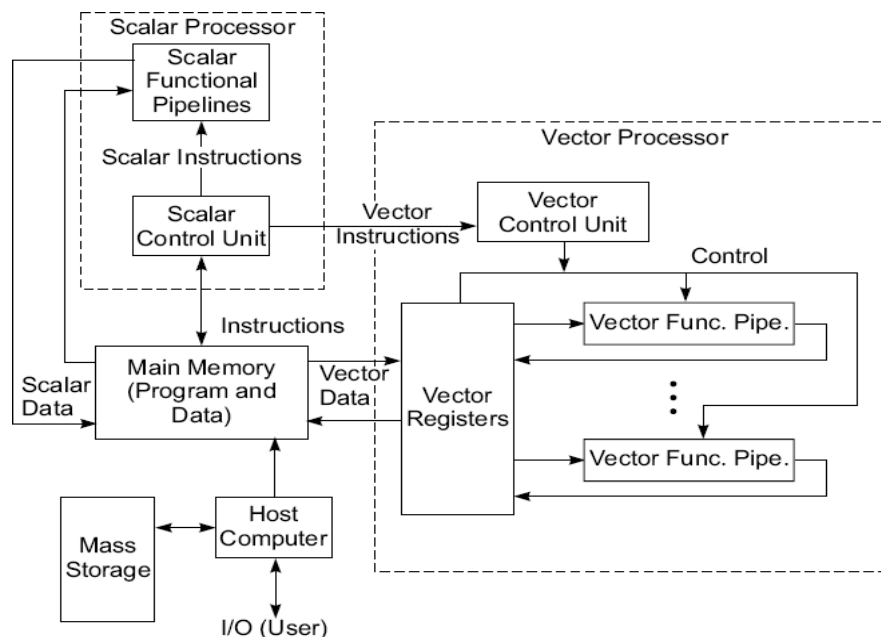


Fig. 1.11 The architecture of a vector supercomputer

- A vector computer is often built on top of a scalar processor.
- As shown in Fig. 1.11, the vector processor is attached to the scalar processor as an optional feature.
- Program and data are first loaded into the main memory through a host computer.
- All instructions are first decoded by the scalar control unit.
- If the decoded instruction is a scalar operation or a program control operation, it will be directly executed by the scalar processor using the scalar functional pipelines.
- If the instruction is decoded as a vector operation, it will be sent to the vector control unit.
- This control unit will supervise the flow of vector data between the main memory and vector functional pipelines.
- The vector data flow is coordinated by the control unit. A number of vector functional pipelines may be built into a vector processor.

Vector Processor Models

- Figure 1.11 shows a **register-to-register** architecture.
- Vector registers are used to hold the vector operands, intermediate and final vector results.
- The vector functional pipelines retrieve operands from and put results into the vector registers.
- All vector registers are programmable in user instructions.
- Each vector register is equipped with a component counter which keeps track of the component registers used in successive pipeline cycles.
- The length of each vector register is usually fixed, say, sixty-four 64-bit component registers in a vector register in a Cray Series supercomputer.
- Other machines, like the Fujitsu VP2000 Series, use reconfigurable vector registers to dynamically match the register length with that of the vector operands.

1.3.2 SIMD Supercomputers

SIMD computers have a single instruction stream over multiple data streams.

An operational model of an SIMD computer is specified by a 5-tuple:

$$M = (N, C, I, M, R)$$

where

1. N is the number of *processing elements* (PEs) in the machine. For example, the Illiac IV had 64 PEs and the Connection Machine CM-2 had 65,536 PEs.

2. **C** is the set of instructions directly executed by the *control unit* (CU), including scalar and program flow control instructions.
3. **I** is the set of instructions broadcast by the CU to all PEs for parallel execution. These include arithmetic, logic, data routing, masking, and other local operations executed by each active PE over data within that PE.
4. **M** is the set of masking schemes, where each mask partitions the set of PEs into enabled and disabled subsets.

R is the set of data-routing functions, specifying various patterns to be set up in the interconnection network for inter-PE communications.

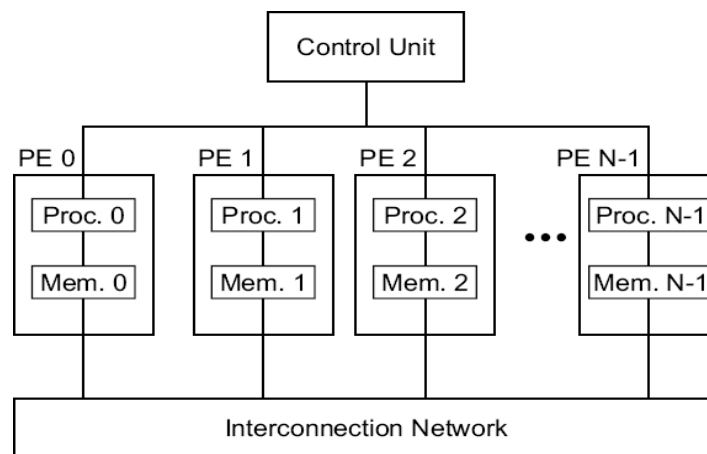


Fig. 1.12 Operational model of SIMD computers

1.4 PRAM AND VLSI MODELS

1.4.1 PRAM model (Parallel Random Access Machine)

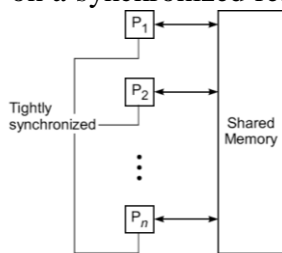
- PRAM is a theoretical model of parallel computation in which an arbitrary but finite number of processors can access any value in an arbitrarily large shared memory in a single time step.
- ✓ The Parallel Random Access Machine (PRAM) is an abstract model for parallel computation which assumes that all the processors operate synchronously under a single clock and are able to randomly access a large shared memory.
- Processors may execute different instruction streams, but work synchronously. This model assumes a shared memory, multiprocessor machine as shown:
- The machine size n can be arbitrarily large.
- The machine is synchronous at the instruction level. That is, each processor is executing its own

series of instructions, and the entire machine operates at a basic time step (cycle). Within each cycle, each processor executes exactly one operation or does nothing, i.e. it is idle.

- An instruction can be any random access machine instruction, such as: fetch some operands from memory, perform an ALU operation on the data, and store the result back in memory.
- All processors implicitly synchronize on each cycle and the synchronization overhead is assumed to be zero.
- Communication is done through reading and writing of shared variables.
- Memory access can be specified to be UMA, NUMA, EREW, CREW, or CRCW with a defined conflict policy.
- The PRAM model can apply to SIMD class machines if all processors execute identical instructions on the same cycle or to MIMD class machines if the processors are executing different instructions.
- Load imbalance is the only form of overhead in the PRAM model.

An n-processor PRAM (Fig. 1.14) has a globally addressable memory.

The shared memory can be distributed among the processors or centralized in one place. The n processors operate on a synchronized read-memory, compute, and write-memory cycle. With shared memory, the



model must specify how concurrent read and concurrent write of memory are handled.

Four memory-update options are possible:

- **Exclusive Read (ER)** — This allows at most one processor to read from any memory location in each cycle, a rather restrictive policy.
- **Exclusive Write (EW)** — This allows at most one processor to write into a memory location at a time.
- **Concurrent Read (CR)** — This allows multiple processors to read the same information from the same memory cell in the same cycle.

- **Concurrent Write (CW)** — This allows simultaneous writes to the same memory location.

In order to avoid confusion, some policy must be set up to resolve the write conflicts. Various combinations of the above options lead to several variants of the PRAM model as specified below.

PRAM Variants

There are 4 variants of the PRAM model, depending on how the memory reads and writes are handled.

- **EREW – PRAM Model (Exclusive Read, Exclusive Write):** This model forbids more than one processor from reading or writing the same memory cell simultaneously. This is the most restrictive PRAM model proposed.
- **CREW – PRAM Model (Concurrent Read, Exclusive Write);** The write conflicts are avoided by mutual exclusion. Concurrent reads to the same memory location are allowed.
- **ERCW – PRAM Model** – This allows exclusive read or concurrent writes to the same memory location.
- **CRCW – PRAM Model (Concurrent Read, Concurrent Write);** This model allows either concurrent reads or concurrent writes to the same memory location.

1.4.2 VLSI Model

Parallel computers rely on the use of VLSI chips to fabricate the major components such as processor arrays memory arrays and large-scale switching networks. The rapid advent of very large scale integrated (VLSI) technology now computer architects are trying to implement parallel algorithms directly in hardware. An AT^2 model is an example for two-dimensional VLSI chips

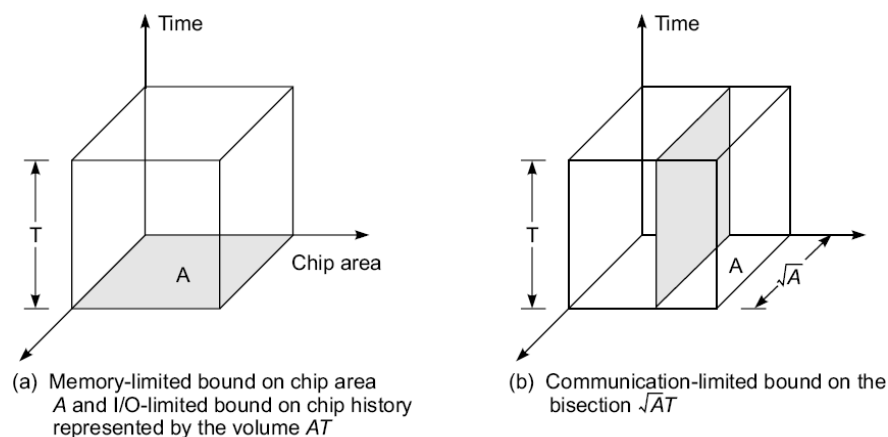


Fig. 1.15 The AT^2 complexity model of two-dimensional VLSI chips

Program & Network properties

Condition of parallelism

2.2.1 Data and Resource Dependence

- The ability to execute several program segments in parallel requires each segment to be independent of the other segments. We use a dependence graph to describe the relations.
- The nodes of a dependence graph correspond to the program statement (instructions), and directed edges with different labels are used to represent the ordered relations among the statements.
- The analysis of dependence graphs shows where opportunity exists for parallelization and vectorization.

Data dependence:

The ordering relationship between statements is indicated by the data dependence. Five type of data dependence are defined below:

1. **Flow dependence:** A statement S2 is flow dependent on S1 if an execution path exists from s1 to S2 and if at least one output (variables assigned) of S1 feeds in as input (operands to be used) to S2 also called RAW $s_1 \rightarrow s_2$ hazard and denoted as
2. **Anti dependence:** Statement S2 is antidependent on the statement S1 if S2 follows S1 in the program order and if the output of S2 overlaps the input to S1 also called RAW hazard and denoted as $s_1 \rightarrow s_2$
3. **Output dependence:** Two statements are output dependent if they produce (write) the same output variable. Also called WAW hazard and denoted as $s_1 \rightarrow s_2$
4. **I/O dependence:** Read and write are I/O statements. I/O dependence occurs not because the same variable is involved but because the same file referenced by both I/O statement.
5. **Unknown dependence:** The dependence relation between two statements cannot be determined in the following situations:
 - The subscript of a variable is itself subscribed (indirect addressing)
 - The subscript does not contain the loop index variable.
 - A variable appears more than once with subscripts having different coefficients of the loop variable.
 - The subscript is non linear in the loop index variable.

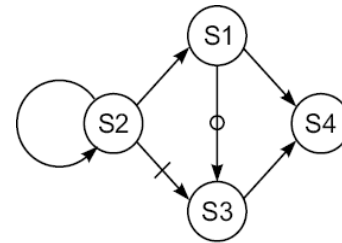
Parallel execution of program segments which do not have total data independence can produce non-deterministic results.

Example: Consider the following fragment of a program:

```

S1:  Load  R1, A      /R1 ← Memory(A) /
S2:  Add   R2, R1     /R2 ← (R1) + (R2)/
S3:  Move  R1, R3     /R1 ← (R3)/
S4:  Store  B, R1     /Memory(B) ← (R1)/

```



(a) Dependence graph

- Here the Flow dependency from **S1 to S2**, **S3 to S4**, **S2 to S2**
- Anti-dependency from **S2 to S3**
- Output dependency **S1 to S3**

```

S1: Read (4), A(I)  /Read array A from file 4/
S2: Rewind (4)     /Process data/
S3: Write (4), B(I) /Write array B into file 4/
S4: Rewind (4)     /Close file 4/

```



(b) I/O dependence caused by accessing the same file by the read and write statements

The read/write statements S1 and S2 are I/O dependent on each other because they both access the same file.

Control Dependence:

- This refers to the situation where the order of the execution of statements cannot be determined before run time.
- For example all condition statement, where the flow of statement depends on the output.
- Different paths taken after a conditional branch may depend on the data hence we need to eliminate this data dependence among the instructions.
- This dependence also exists between operations performed in successive iterations of looping procedure. Control dependence often prohibits parallelism from being exploited.

Control-independent example:

```

for (i=0; i<n; i++)
{
    a[i] = c[i];
    if (a[i] < 0) a[i] = 1;
}

```

Control-dependent example:

```

for (i=1; i<n; i++)
{
    if (a[i-1] < 0) a[i] = 1;
}

```

Control dependence also avoids parallelism to being exploited. Compilers are used to eliminate this control dependence and exploit the parallelism.

Resource dependence:

- Data and control dependencies are based on the independence of the work to be done.
- Resource independence is concerned with conflicts in using shared resources, such as registers, integer and floating point ALUs, etc. ALU conflicts are called ALU dependence.
- Memory (storage) conflicts are called storage dependence.

Bernstein's Conditions

Bernstein's conditions are a set of conditions which must exist if two processes can execute in parallel.

Notation

- I_i is the set of all input variables for a process P_i . I_i is also called the read set or domain of P_i . O_i is the set of all output variables for a process P_i . O_i is also called write set.
- If P_1 and P_2 can execute in parallel (which is written as $P_1 \parallel P_2$), then:

$$I_1 \cap O_2 = \emptyset$$

$$I_2 \cap O_1 = \emptyset$$

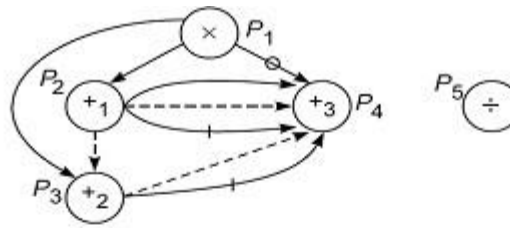
$$O_1 \cap O_2 = \emptyset$$

- In terms of data dependencies, Bernstein's conditions imply that two processes can execute in parallel if they are flow-independent, anti-independent, and output-independent.
- The parallelism relation \parallel is commutative ($P_i \parallel P_j$ implies $P_j \parallel P_i$), but not transitive ($P_i \parallel P_j$ and $P_j \parallel P_k$ does not imply $P_i \parallel P_k$).
- Therefore, \parallel is not an equivalence relation. Intersection of the input sets is allowed.

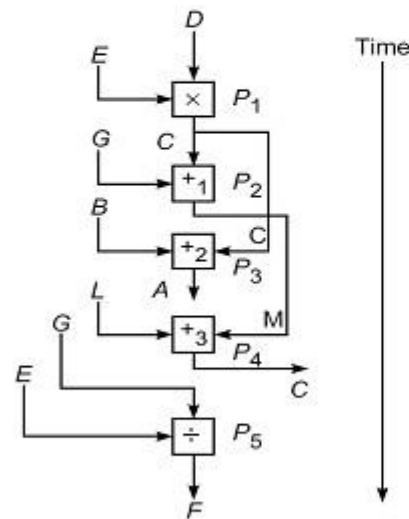
Example: Detection of parallelism in a program using Bernstein's conditions

Consider the simple case in which each process is a single HLL statement. We want to detect the parallelism embedded in the following 5 statements labeled P1, P2, P3, P4, P5 in program order.

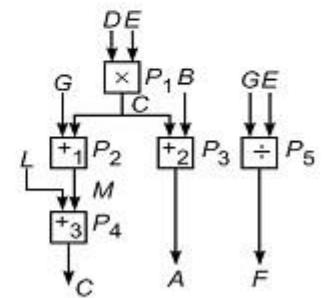
$P_1:$ $C = D \times E$
 $P_2:$ $M = G + C$
 $P_3:$ $A = B + C$
 $P_4:$ $C = L + M$
 $P_5:$ $F = G / E$



(a) A dependence graph showing both data dependence (solid arrows) and resource dependence (dashed arrows)



(b) Sequential execution in five steps, assuming one step per statement (no pipelining)



(c) Parallel execution in three steps, assuming two adders are available per step

Fig. 2.2 Detection of parallelism in the program of Example 2.2

- Assume that each statement requires one step to execute. No pipelining is considered here. The dependence graph shown in 2.2a demonstrates flow dependence as well as resource dependence. In sequential execution, five steps are needed (Fig. 2.2b).
- If two adders are available simultaneously, the parallel execution requires only 3 steps as shown in Fig 2.2c.
- Pairwise, there are 10 pairs of statements to check against Bernstein's conditions. Only 5 pairs, $P1 \parallel P5$, $P2 \parallel P3$, $P2 \parallel P5$, $P5 \parallel P3$ and $P4 \parallel P5$ can execute in parallel as revealed in Fig 2.2a if there are no resource conflicts.
- Collectively, only $P2 \parallel P3 \parallel P5$ is possible (Fig. 2.2c) because $P2 \parallel P3$, $P3 \parallel P5$ and $P5 \parallel P2$ are all possible.

Hardware and software parallelism

Hardware parallelism

- Hardware parallelism is defined by machine architecture and hardware multiplicity i.e., functional parallelism times the processor parallelism.
- It can be characterized by the number of instructions that can be issued per machine cycle.
- If a processor issues k instructions per machine cycle, it is called a k -issue processor.
- Conventional processors are *one-issue* machines.
- This provide the user the information about **peak attainable performance**.

Examples: Intel i960CA is a three-issue processor (arithmetic, memory access, branch).

IBM RS -6000 is a four-issue processor (arithmetic, floating-point, memory access, branch).

A machine with n k -issue processors should be able to handle a maximum of nk threads simultaneously.

Software Parallelism

Software parallelism is defined by the control and data dependence of programs, and is revealed in the program's flow graph i.e., it is defined by dependencies with in the code and is a function of algorithm, programming style, and compiler optimization.

Example: Mismatch between Software parallelism and Hardware parallelism

- Consider the example program graph in Fig. 2.3a. There are eight instructions (four *loads* and four *arithmetic* operations) to be executed in three consecutive machine cycles.
- Four *load* operations are performed in the first cycle, followed by two *multiply* operations in the second cycle and two *add/subtract* operations in the third cycle.
- Therefore, the parallelism varies from 4 to 2 in three cycles. The average software parallelism is equal to $8/3 = 2.67$ instructions per cycle in this example program.
- Now consider execution of the same program by a two-issue processor which can execute one memory access (*load* or *write*) and one arithmetic (*add*, *subtract*, *multiply*, etc.) operation simultaneously.
- With this hardware restriction, the program must execute in seven machine cycles as shown in Fig. 2.3b. Therefore, the *hardware parallelism* displays an average value of $8/7 = 1.14$ instructions executed per cycle.
- This demonstrates a mismatch between the software parallelism and the hardware parallelism.

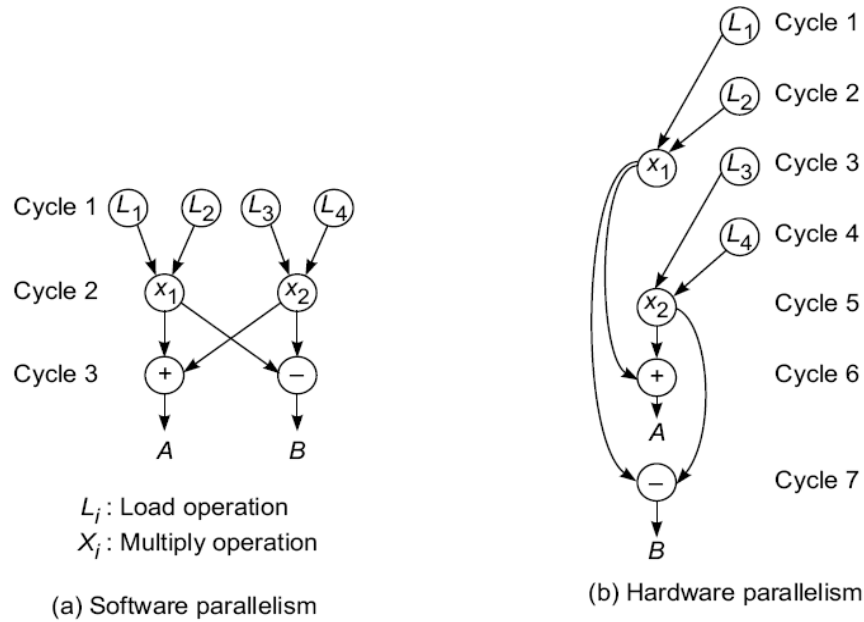


Fig. 2.3 Executing an example program by a two-issue superscalar processor

- Let us try to match the software parallelism shown in Fig. 2.3a in a hardware platform of a dual-processor system, where single-issue processors are used.
- The achievable hardware parallelism is shown in Fig 2.4. Six processor cycles are needed to execute 12 instructions by two processors.
- S1 and S2 are two inserted store operations, l5 and l6 are two inserted load operations for interprocessor communication through the shared memory.

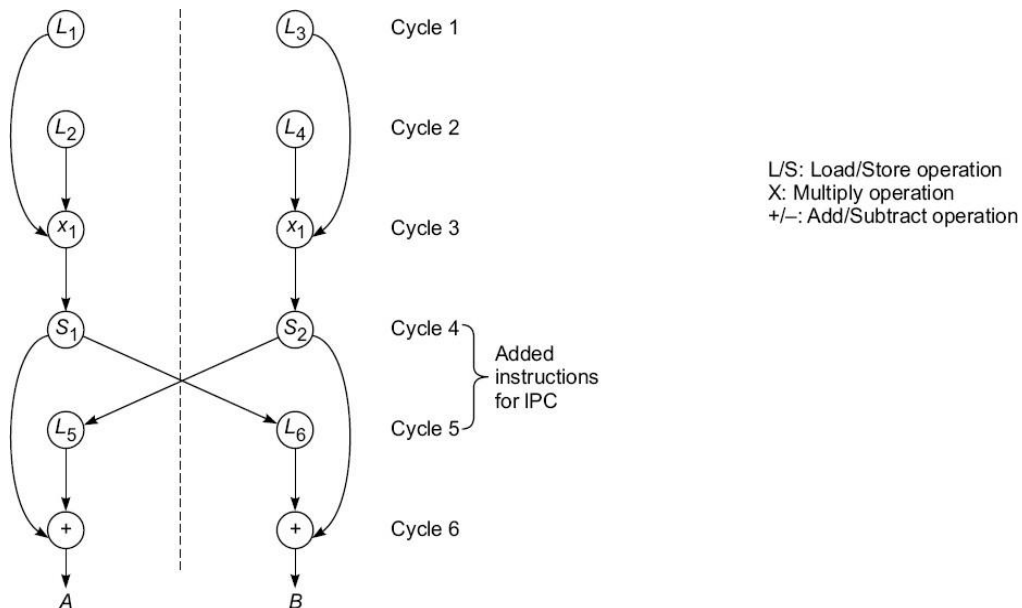


Fig. 2.4 Dual-processor execution of the program in Fig. 2.3a

The Role of Compilers

- Compilers used to exploit hardware features to improve performance. Interaction between compiler and architecture design is a necessity in modern computer development.
- It is not necessarily the case that more software parallelism will improve performance in conventional scalar processors.
- The hardware and compiler should be designed at the same time.

Program Partitioning & Scheduling

Grain size and latency

- The size of the parts or pieces of a program that can be considered for parallel execution can vary.
- The sizes are roughly classified using the term –granule size,|| or simply –granularity.||
- The simplest measure, for example, is the number of instructions in a program part.
- Grain sizes are usually described as fine, medium or coarse, depending on the level of parallelism involved.

Latency

Latency is the time required for communication between different subsystems in a computer. Memory latency, for example, is the time required by a processor to access memory. Synchronization latency is the time required for two processes to synchronize their execution. Computational granularity and communication latency are closely related.

Latency and grain size are interrelated and some general observation are

- As grain size decreases, potential parallelism increases, and overhead also increases.
- Overhead is the cost of parallelizing a task. The principle overhead is communication latency.
- As grain size is reduced, there are fewer operations between communication, and hence the impact of latency increases.
- Surface to volume: inter to intra-node comm.

Levels of Parallelism

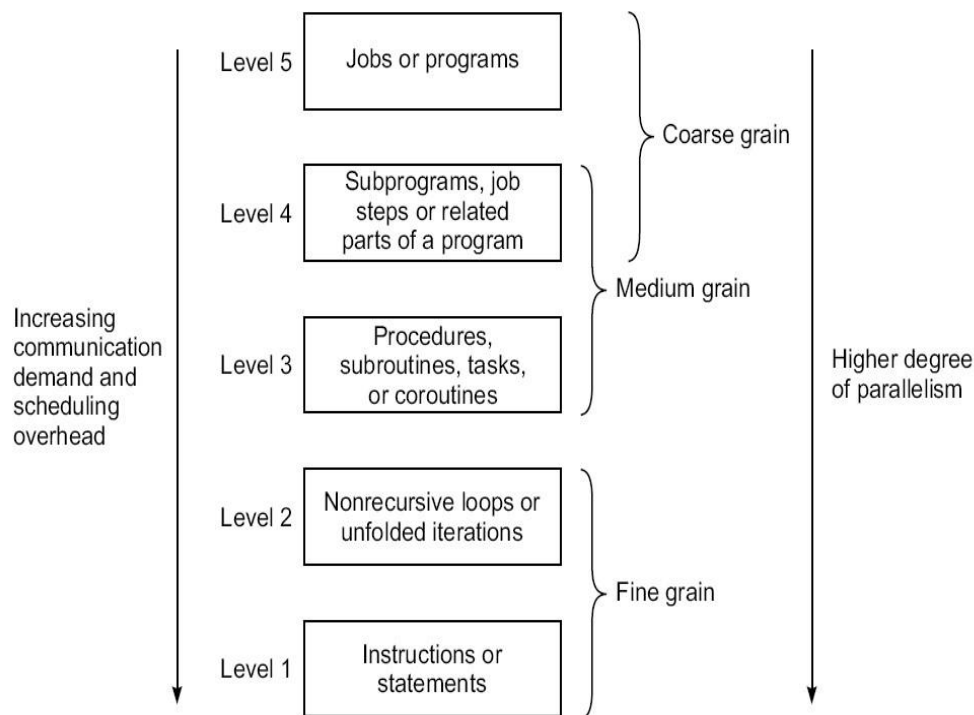


Fig. 2.5 Levels of parallelism in program execution on modern computers (Reprinted from Hwang, *Proc. IEEE*, October 1987)

Instruction Level Parallelism

- This fine-grained, or smallest granularity level typically involves less than 20 instructions per grain.
- The number of candidates for parallel execution varies from 2 to thousands, with about five instructions or statements (on the average) being the average level of parallelism.

Advantages:

There are usually many candidates for parallel execution. Compilers can usually do a reasonable job of finding this parallelism

Loop-level Parallelism

- Typical loop has less than 500 instructions. If a loop operation is independent between iterations, it can be handled by a pipeline, or by a SIMD machine.
- Most optimized program construct to execute on a parallel or vector machine.
- Some loops (e.g. recursive) are difficult to handle. Loop-level parallelism is still considered fine grain computation.

Procedure-level Parallelism

- Medium-sized grain; usually less than 2000 instructions.
- Detection of parallelism is more difficult than with smaller grains; interprocedural dependence analysis is difficult and history-sensitive.
- Communication requirement less than instruction level SPMD (single procedure multiple data) is a special case Multitasking belongs to this level.

Subprogram-level Parallelism

- Job step level; grain typically has thousands of instructions; medium- or coarse-grain level.
- Job steps can overlap across different jobs. Multiprogramming conducted at this level No compilers available to exploit medium- or coarse-grain parallelism at present.

Job or Program-Level Parallelism

- Corresponds to execution of essentially independent jobs or programs on a parallel computer.
- This is practical for a machine with a small number of powerful processors, but impractical for a machine with a large number of simple processors (since each processor would take too long to process a single job).

Communication Latency

Balancing granularity and latency can yield better performance. Various latencies attributed to machine architecture, technology, and communication patterns used.

Latency imposes a limiting factor on machine scalability.

Ex: Memory latency increases as memory capacity increases, limiting the amount of memory that can be used with a given tolerance for communication latency.

Interprocessor Communication Latency

- Needs to be minimized by system designer
- Affected by signal delays and communication patterns **Ex:** n communicating tasks may require $n(n - 1)/2$ communication links, and the complexity grows quadratically, effectively limiting the number of processors in the system.

Communication Patterns

- Determined by algorithms used and architectural support provided
- Patterns include permutations broadcast multicast conference
- Tradeoffs often exist between granularity of parallelism and communication demand.

Grain Packing and Scheduling

Two questions:

- How can I partition a program into parallel —pieces to yield the shortest execution time?
- What is the optimal size of parallel grains?

There is an obvious tradeoff between the time spent scheduling and synchronizing parallel grains and the speedup obtained by parallel execution.

One approach to the problem is called –grain packing.¶

Program Graphs and Packing (Basic concept of Program Partitioning)

- A program graph shows the structure of the program, similar to dependence graph. Each node in the program graph corresponds to a computational unit in the program.
- Grain size is measured by the number of basic machine cycles needed to execute all the operations within the node.
- Each node is denoted by, **Nodes** = { (n,s) }, where n = node name (id), **s = grain size** (larger s = larger grain size), Fine-grain nodes have a smaller grain size, and coarse-grain nodes have a larger grain size.
- **Edges** = { (v,d) }, where v = variable being –communicated,¶ and d = communication delay.
- Packing two (or more) nodes produces a node with a larger grain size and possibly more edges to other nodes.
- Packing is done to eliminate unnecessary communication delays or reduce overall scheduling overhead.

Example: Basic concept of Program Partitioning

Fig. 2.6, shows an example program graph in two different grain sizes.

Var $a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q$

Begin

- | | |
|----------------------|-----------------------|
| 1. $a := 1$ | 10. $j := e \times f$ |
| 2. $b := 2$ | 11. $k := d \times f$ |
| 3. $c := 3$ | 12. $l := j \times k$ |
| 4. $d := 4$ | 13. $m := 4 \times l$ |
| 5. $e := 5$ | 14. $n := 3 \times m$ |
| 6. $f := 6$ | 15. $o := n \times i$ |
| 7. $g := a \times b$ | 16. $p := o \times h$ |
| 8. $h := c \times d$ | 17. $q := p \times q$ |
| 9. $i := d \times e$ | |

End

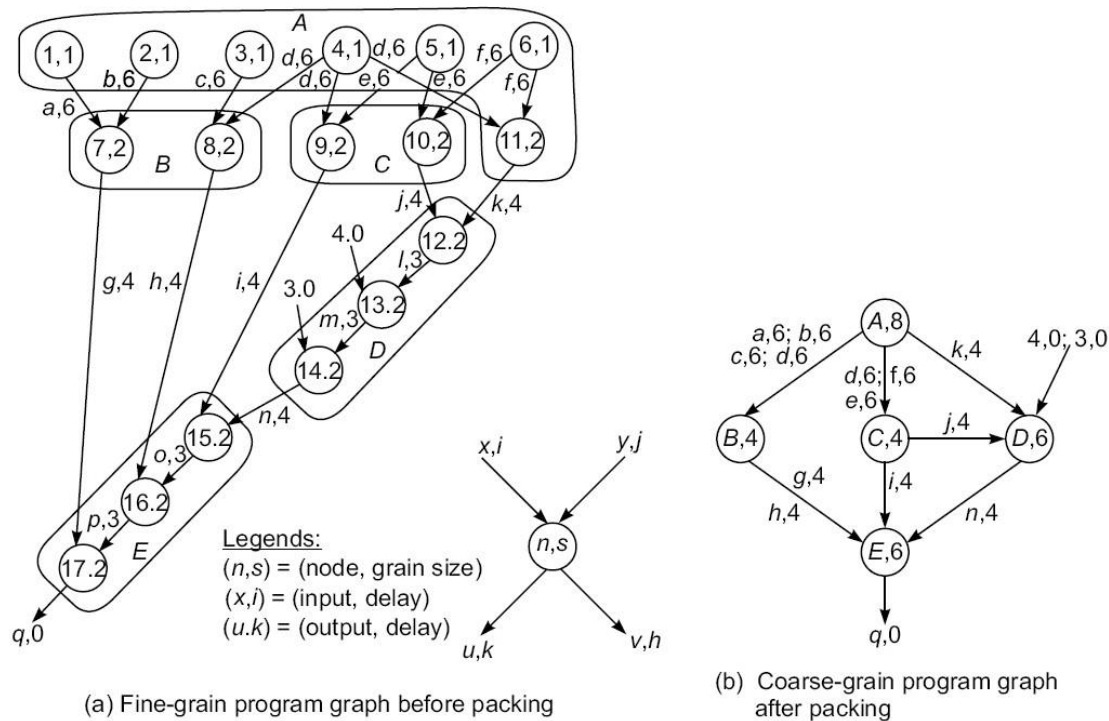


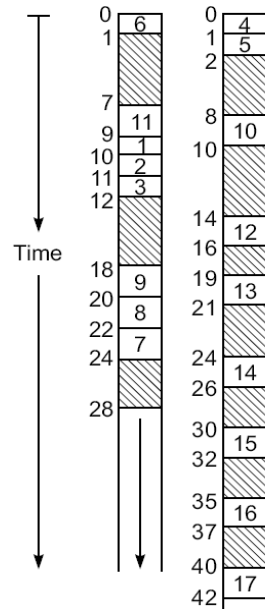
Fig. 2.6 A program graph before and after grain packing in Example 2.4 (Modified from Kruatrachue and Lewis, *IEEE Software*, Jan. 1988)

Scheduling

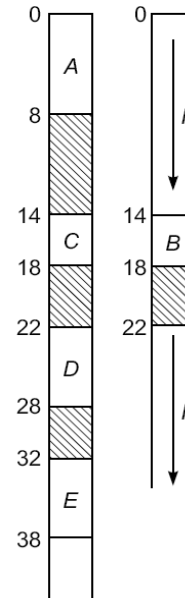
A schedule is a mapping of nodes to processors and start times such that communication delay requirements are observed, and no two nodes are executing on the same processor at the same time.

Some general scheduling goals are:

- Schedule all fine-grain activities in a node to the same processor to minimize communication delays.
- Select grain sizes for packing to achieve better schedules for a particular parallel machine.



(a) Fine grain (Fig. 2.6a)



(b) Coarse grain (Fig. 2.6b)

Fig. 2.7 Scheduling of the fine-grain and coarse-grain programs (arrows: idle time; shaded areas: communication delays)

- With respect to the fine-grain versus coarse-grain program graphs in Fig. 2.6, two multiprocessor schedules are shown in Fig. 2.7. The fine-grain schedule is longer (42 time units) because more communication delays were included as shown by the shaded area.
- The coarse-grain schedule is shorter (38 time units) because communication delays among nodes 12, 13 and 14 within the same node D (and also the delays among 15, 16 and 17 within the nodeE) are eliminated after grain packing.

Node Duplication

- Grain packing may potentially eliminate interprocessor communication, but it may not always produce a shorter schedule.
- By duplicating nodes (that is, executing some instructions on multiple processors), we may eliminate some interprocessor communication, and thus produce a shorter schedule.

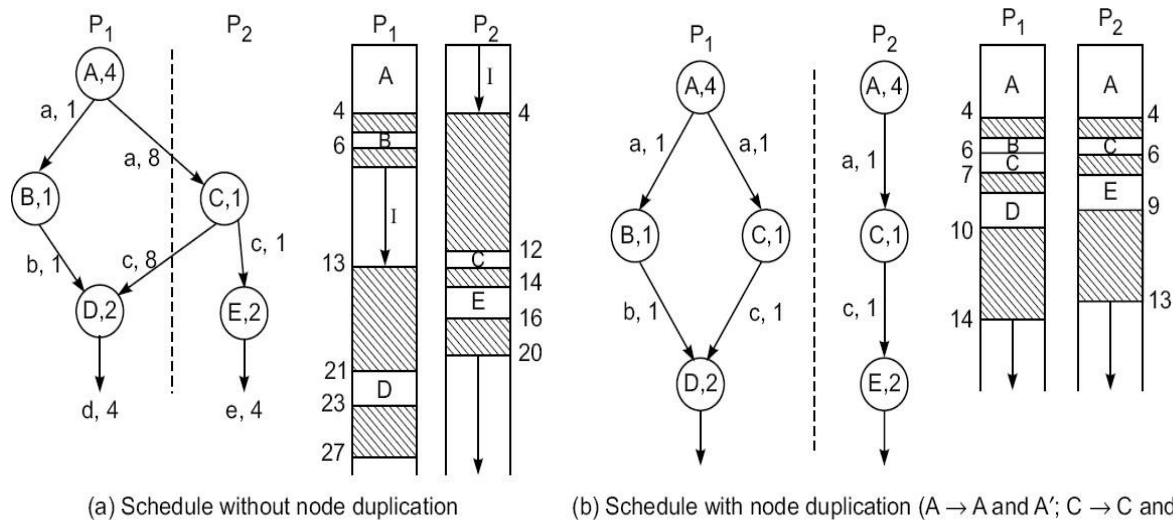


Fig. 2.8 Node-duplication scheduling to eliminate communication delays between processors (I: idle time; shaded areas: communication delays)

- Figure 2.8a shows a schedule without duplicating any of the 5 nodes. This schedule contains idle time as well as long interprocessor delays (8 units) between P1 and P2.
- In Fig 2.8b, node A is duplicated into A' and assigned to P2 besides retaining the original copy A in P1.
- Similarly, a duplicated node C' is copied into P1 besides the original node C in P2.
- The new schedule is shown in Fig. 2.8b is almost 50% shorter than that in Fig. 2.8a. The reduction in schedule time is caused by elimination of the (a, 8) and (c, 8) delays between the two processors.

Grain packing and node duplication are often used jointly to determine the best grain size and corresponding schedule.

Four major steps are involved in the grain determination and the process of scheduling optimization:\

Step 1: Construct a fine-grain program graph

Step 2: Schedule the fine-grain computation

Step 3: Perform grain packing to produce the coarse grains.

Step 4: Generate a parallel schedule based on the packed graph.

Program Flow Mechanisms

Control Flow vs. Data Flow

- In Control flow computers the next instruction is executed when the last instruction as stored in the program has been executed where as in Data flow computers an instruction executed when the data (operands) required for executing that instruction is available.

- Control flow machines used shared memory for instructions and data.
- Since variables are updated by many instructions, there may be side effects on other instructions. These side effects frequently prevent parallel processing.
- Single processor systems are inherently sequential.
- Instructions in dataflow machines are unordered and can be executed as soon as their operands are available; data is held in the instructions themselves. *Data tokens* are passed from an instruction to its dependents to trigger execution.

Program Flow Mechanisms

- **Control flow mechanism**: Conventional machines used control flow mechanism in which order of program execution explicitly stated in user programs.
- **Dataflow machines** which instructions can be executed by determining operand availability.
- **Reduction machines** trigger an instruction's execution based on the demand for its results.

Control flow machines used shared memory for instructions and data. Since variables are updated by many instructions, there may be side effects on other instructions. These side effects frequently prevent parallel processing. Single processor systems are inherently sequential.

Instructions in dataflow machines are unordered and can be executed as soon as their operands are available; data is held in the instructions themselves. *Data tokens* are passed from an instruction to its dependents to trigger execution.

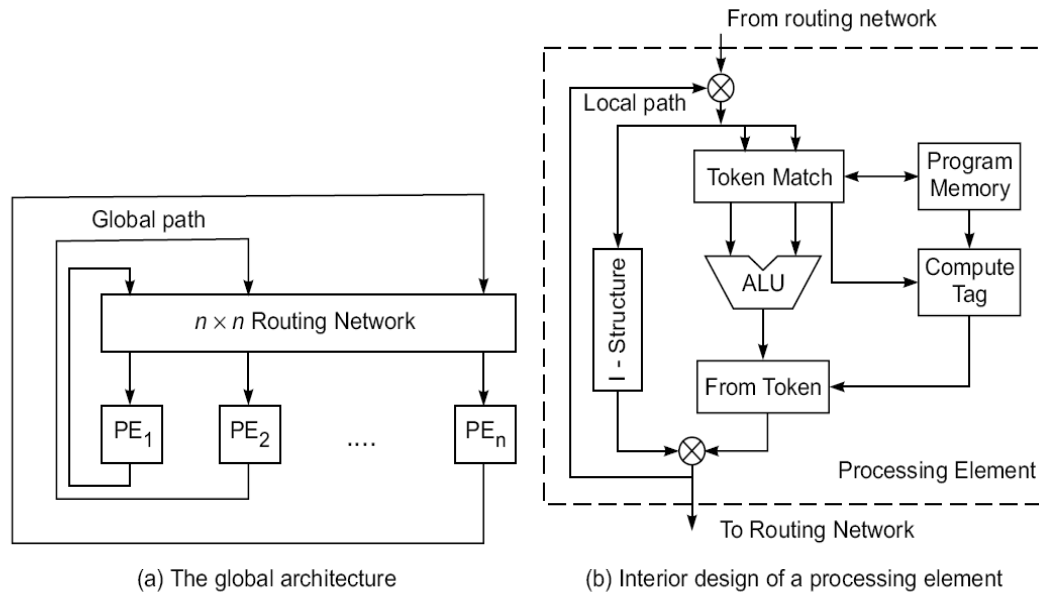
Data Flow Features

No need for

- shared memory
- program counter
- control sequencer

Special mechanisms are required to

- detect data availability
- match data tokens with instructions needing them
- enable chain reaction of asynchronous instruction execution



A Dataflow Architecture

- The Arvind machine (MIT) has N PEs and an N -by- N interconnection network.
- Each PE has a token-matching mechanism that dispatches only instructions with data tokens available.
- Each datum is tagged with
 - address of instruction to which it belongs
 - context in which the instruction is being executed
- Tagged tokens enter PE through local path (pipelined), and can also be communicated to other PEs through the routing network.
- Instruction address(es) effectively replace the program counter in a control flow machine.
- Context identifier effectively replaces the frame base register in a control flow machine.
- Since the dataflow machine matches the data tags from one instruction with successors, synchronized instruction execution is implicit.
- An I-structure in each PE is provided to eliminate excessive copying of data structures.
- Each word of the I-structure has a **two-bit tag** indicating whether the value is **empty**, **full** or **has pending read requests**.
- This is a retreat from the pure dataflow approach.
- Special compiler technology needed for dataflow machines.

Demand-Driven Mechanisms

- Demand-driven machines take a top-down approach, attempting to execute the instruction (a demand) that yields the final result.
- This triggers the execution of instructions that yield its operands, and so forth.
- The demand-driven approach matches naturally with functional programming languages (e.g. LISP and SCHEME).

Reduction Machine Models

- **String-reduction model:**
 - each demander gets a separate copy of the expression string to evaluate
 - each reduction step has an operator and embedded reference to demand the corresponding operands
 - each operator is suspended while arguments are evaluated
- **Graph-reduction model:**
 - expression graph reduced by evaluation of branches or subgraphs, possibly in parallel, with demanders given pointers to results of reductions.
 - based on sharing of pointers to arguments; traversal and reversal of pointers continues until constant arguments are encountered.

2.1 System interconnect architecture

Various types of interconnection networks have been suggested for SIMD computers. These are basically classified have been classified on network topologies into two categories namely

1. Static Networks

2. Dynamic Networks

- Direct networks for static connections
- Indirect networks for dynamic connections
- Networks are used for
 - internal connections in a centralized system among
 - processors
 - memory modules
 - I/O disk arrays
 - distributed networking of multicomputer nodes
- The goals of an interconnection network are to provide
 - low-latency

- high data transfer rate
 - wide communication bandwidth
- Analysis includes
 - latency
 - bisection bandwidth
 - data-routing functions
 - scalability of parallel architecture
- The topology of an interconnection network can be either static or dynamic. Static networks are formed of point-to-point direct connections which will not change during program execution.
- Dynamic networks are implemented with switched channels, which are dynamically configured to match the communication demand in user programs.
- Packet switching and routing is playing an important role in modern multiprocessor architecture.

Node Degree and Network Diameter:

- The number of edges (links or channels) incident on a node is called the node degree d .
- In the case of unidirectional channels, the number of channels into a node is the in degree, and that out of a node is the out degree.
- Then the node degree is the sum of the two. The node degree reflects the number of IO ports required per node, and thus the cost of a node.
- Therefore, the node degree should be kept a (small) constant, in order to reduce cost.
- The Diameter D of a network is the maximum shortest path between any two nodes.
- The path length is measured by the number of links traversed.
- The network diameter indicates the maximum number of distinct hops between any two nodes, thus providing a figure of communication merit for the network.
- Therefore, the network diameter should be as small as possible from a communication point of view.

Bisection Width:

When a given network is cut into two equal halves, the minimum number of edges (channels) along the cut is called the bisection width b . In the case of a communication network, each edge may correspond to a channel with w bit wires.

To summarize the above discussions, the performance of an interconnection network is affected by the following factors:

Functionality: refers to how the network supports data routing, interrupt handling, synchronization, request-message combining, and coherence.

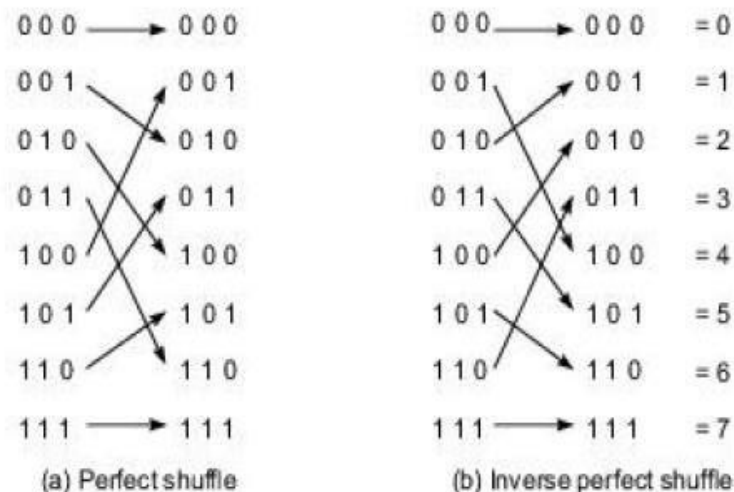
Network Latency:- This refers to the worst-case time delay for a unit message to be transferred through the network.

Bandwidth: This refers to the maximum data transfer rate, in terms of Mbps or Gbps transmitted through the network.

Hardware Complexity'—This refers to implementation costs such as those for wires, switches, connectors, arbitration, and interface logic.

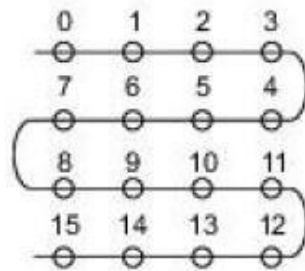
Scalability—This refers to the ability of a network to be modularly expandable with a scalable performance with increasing machine resources.

Perfect Shuffle and Exchange Perfect shuffle is a special permutation function suggested by Harold Stone (1971) for parallel processing applications. The mapping corresponding to a perfect shuffle is shown in Fig. 2.14a. Its inverse is shown on the right-hand side (Fig. 2.14b).

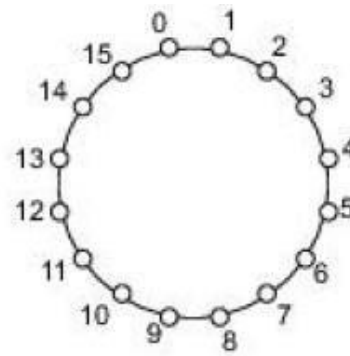


2.4.2 Static Connection Networks:

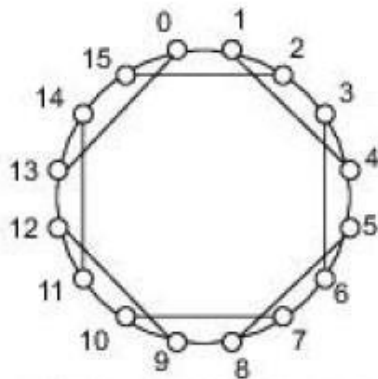
Static networks use direct links which are fixed once built. This type of network is more suitable for building computers where the communication patterns are predictable or implementable with static connections. We describe their topologies below in terms of network parameters and comment on their relative merits in relation to communication and scalability.



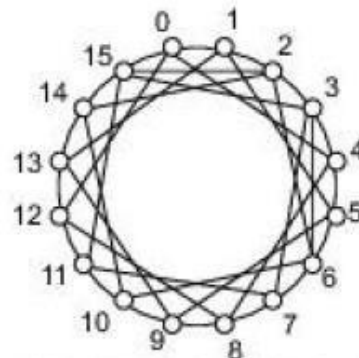
(a) Linear array



(b) Ring



(c) Chordal ring of degree 3

(d) Chordal ring of degree 4
(same as Illiac mesh)

Dynamic Connection Networks

- Dynamic connection networks can implement all communication patterns based on program demands.
- In increasing order of cost and performance, these include
 - bus systems
 - multistage interconnection networks
 - crossbar switch networks
- Price can be attributed to the cost of wires, switches, arbiters, and connectors.
- Performance is indicated by network bandwidth, data transfer rate, network latency, and communication patterns supported.

Digital Buses

- A bus system (contention bus, time-sharing bus) has
 - a collection of wires and connectors
 - multiple modules (processors, memories, peripherals, etc.) which connect to the wires
 - data transactions between pairs of modules

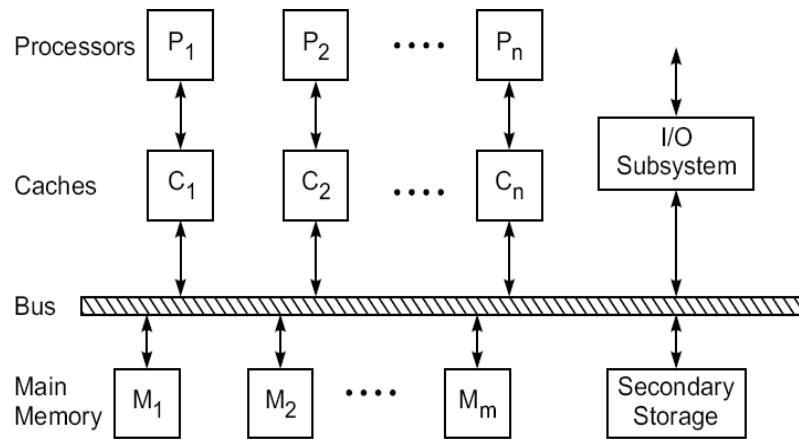


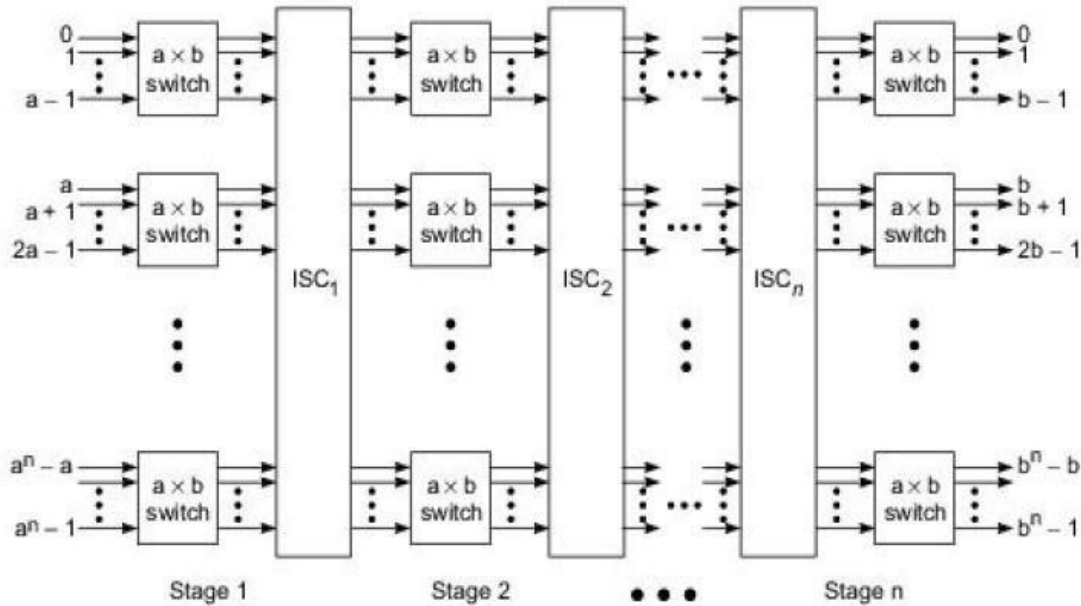
Fig. 2.22 A bus-connected multiprocessor system, such as the Sequent Symmetry S1

- Bus supports only one transaction at a time.
- Bus arbitration logic must deal with conflicting requests.
- Lowest cost and bandwidth of all dynamic schemes.
- Many bus standards are available.

Switch Modules An $a \times b$ switch module has a inputs and b outputs. A binary switch corresponds to a 2×2 switch module in which $a = b = 2$. In theory, a and b do not have to be equal. However, in practice, a and b are often chosen as integer powers of 2; that is, $a = b = 2^k$ for some $k \geq 1$.

Multistage Interconnection Networks MINs have been used in both MIMD and SIMD computers. A generalized multistage network is illustrated in Fig. 2.23. A number of $a \times b$ switches are used in each stage. Fixed interstage connections are used between the switches in adjacent stages. The switches can be dynamically set to establish the desired connections between the inputs and outputs.

Different classes of MINs differ in the switch modules used and in the kind of *interstage connection* (ISC) patterns used. The simplest switch module would be the 2×2 switches ($a = b = 2$ in Fig. 2.23). The ISC



Principles of Scalable Performance

- Performance measures
- Speedup laws
- Scalability principles

Scaling up vs. scaling down

Performance metrics and measures:

The degree of parallelism reflects the extent to which software parallelism matches hardware parallelism.

Degree of Parallelism The execution of a program on a parallel computer may use different numbers of processors at different time periods during the execution cycle. For each time period, the number of processors used to execute a program is defined as the

The plot of the D-OP as a function of time is called the parallelism pm-file of a given program.

Fluctuation of the profile during an observation period depends on the algorithmic structure, program optimization, resource utilization, and run-time conditions of a computer system. The DOP was defined under the assumption of having an unbounded number of available processors and other necessary resources.

The DOP may not always be achievable on a real computer with limited resources.

Average Parallelism

- ✓ Parallelism profiles
 - ✓ Asymptotic speedup factor
 - ✓ System efficiency,
 - ✓ utilization and quality
 - ✓ Standard performance measures
 - ✓ Degree of parallelism
 - ✓ Reflects the matching of software and hardware parallelism
- Discrete time function – measure, for each time period, the # of processors used
- Parallelism profile is a plot of the **DOP** as a function of time
- Ideally have unlimited resources

Factors affecting parallelism profiles

- Algorithm structure
- Program optimization
- Resource utilization
- Run-time conditions
- Realistically limited by # of available processors,
- memory, and other non processor resources
- Average parallelism variables
- n – homogeneous processors
- m – maximum parallelism in a profile
- D – computing capacity of a single processor (execution rate only, no overhead)
- $DOP = i$ – # processors busy during
- an observation period

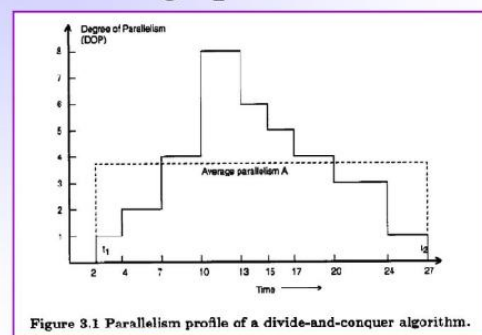
Average parallelism

- Total amount of work performed is proportional to the area under the profile curve

$$W = \Delta \int_{t_1}^{t_2} DOP(t) dt$$

$$W = \Delta \sum_{i=1}^m i \cdot t_i$$

Example: parallelism profile and average parallelism



Performance measures and Metrics

Consider n processors executing m programs in various modes want to define the mean performance of these multimode computers:

- ✓ Arithmetic mean performance
- ✓ Geometric mean performance

Harmonic mean performance

Arithmetic mean performance

$$R_a = \sum_{i=1}^m R_i / m$$

Arithmetic mean execution rate
(assumes equal weighting)

Geometric mean performance

$$R_g = \prod_{i=1}^m R_i^{1/m}$$

Geometric mean execution rate

Harmonic mean performance

$$T_i = 1 / R_i$$

Mean execution time per instruction
For program i

$$T_a = \frac{1}{m} \sum_{i=1}^m T_i = \frac{1}{m} \sum_{i=1}^m \frac{1}{R_i}$$

Arithmetic mean execution time
per instruction

Scalability Metrics

Machine size (n) : # of processors

Clock rate (f) : determines basic machine cycle

Problem size (s) : amount of computational workload. Directly proportional to $T(s,1)$.

CPU time ($T(s,n)$) : actual CPU time for execution

I/O demand (d) : demand in moving the program, data, and results for a given run

Scalability Metrics

Memory capacity (m) : max # of memory words Demanded

Communication overhead ($h(s,n)$) : amount of time for interprocessor communication, synchronization,

Computer cost (c) : total cost of h/w and s/w resources required

Programming overhead (p) : development overhead associated with an application program.

Efficiency ,Utilization and Quality

Speedup and Efficiency

- The problem size is the independent parameter

$$S(s,n) = \frac{T(s,1)}{T(s,n) + h(s,n)}$$

$$E(s,n) = \frac{S(s,n)}{n}$$

Scalable Systems

- Ideally, if $E(s,n)=1$ for all algorithms and any s and n , system is scalable
- Practically, consider the scalability of a m/c

$$\Phi(s,n) = \frac{S(s,n)}{S_1(s,n)} = \frac{T_1(s,n)}{T(s,n)}$$

Standard Performance Measures

Most of the standard measures adopted by the industry to compare the performance of various parallel computers are based on the concepts of:

Peak Performance

[Theoretical maximum based on best possible utilization of all resources]

Sustained Performance

[based on running application-oriented benchmarks]

Generally measured in units of:

MIPS [to reflect instruction execution rate]

MFLOPS [to reflect the floating-point capability]

Benchmarks

Benchmarks are a set of programs of program fragments used to compare the performance of various machines.

Machines are exposed to these benchmark tests and tested for performance.

When it is not possible to test the applications of different machines, then the results of benchmark programs that most resemble the applications run on those machines are used to evaluate the performance of machine.

Sources of Parallel Overhead

Parallel computers in practice do not achieve linear speedup or an efficiency of 1 because of parallel overhead.

The major sources of which could be:

- Inter-processor Communication*
 - Load Imbalance*
- Inter-Task Dependency*
 - Extra Computation*
- Parallel Balance Point*
- Speedup Performance Laws**

- Amdahl's Law
- [based on fixed problem size or fixed work load]
- Gustafson's Law
- [for scaled problems, where problem size increases with machine size
- i.e. the number of processors]

Amdahl's Law (1967)

- For a given problem size, the speedup does not increase linearly as the number of processors increases. In fact, the speedup tends to become saturated.
- This is a consequence of Amdahl's Law.
- According to Amdahl's Law, a program contains two types of operations:
- **Completely sequential Completely parallel**
- Let, the time T_s taken to perform sequential operations be a fraction α ($0 < \alpha \leq 1$) of the total execution time $T(1)$ of the program, then the time T_p to perform parallel operations shall be $(1-\alpha)$ of $T(1)$

Amdahl's Law:

Thus, $T_s = \alpha.T(1)$ and $T_p = (1-\alpha).T(1)$

Assuming that the parallel operations achieve linear speedup

(i.e. these operations use $1/n$ of the time taken to perform on each processor), then

$$T(n) = T_s + T_p/n =$$

Thus, the speedup with n processors will be:

$$\alpha.T(1) + \frac{(1-\alpha).T(1)}{n}$$

$$S(n) = \frac{T(1)}{T(n)} = \frac{T(1)}{\alpha.T(1) + \frac{(1-\alpha).T(1)}{n}} = \frac{1}{\alpha + \frac{(1-\alpha)}{n}}$$

Amdahl's Law

Sequential operations will tend to dominate the speedup as n becomes very large.

As $n \rightarrow \infty$, $S(n) \rightarrow 1/\alpha$

This means, no matter how many processors are employed, the speedup in this problem is limited to $1/\alpha$.

This is known as **sequential bottleneck** of the problem.

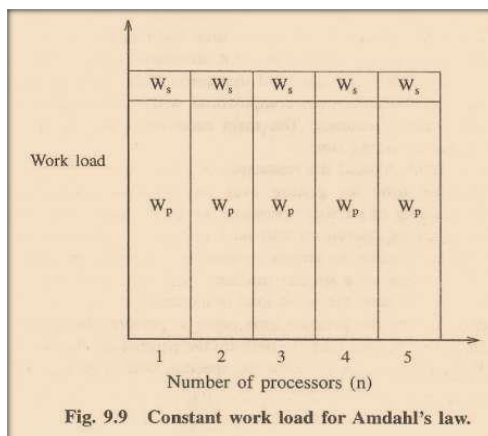
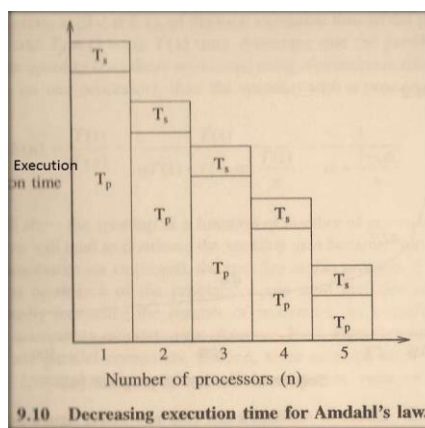


Fig. 9.9 Constant work load for Amdahl's law.



9.10 Decreasing execution time for Amdahl's law.

Gustafson's Law (1988)

It relaxed the restriction of fixed size of the problem and used the notion of fixed execution time for getting over the sequential bottleneck.

According to Gustafson's Law,

If the number of parallel operations in the problem is increased (or scaled up) sufficiently,

Then sequential operations will no longer be a bottleneck.

In accuracy-critical applications, it is desirable to solve the largest problem size on a larger machine rather than solving a smaller problem on a smaller machine, with almost the same execution time.

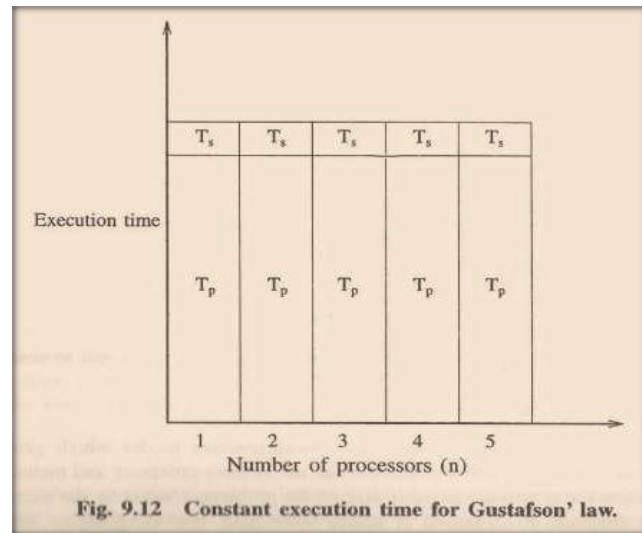
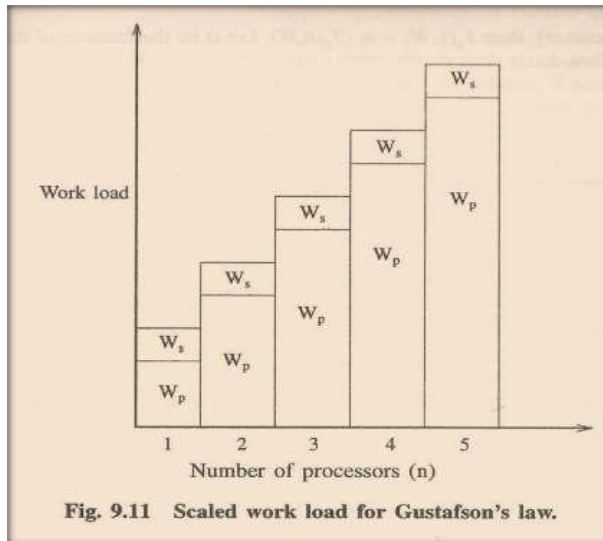
As the machine size increases, the work load (or problem size) is also increased so as to keep the fixed execution time for the problem.

Let, T_s be the constant time tank to perform sequential operations; and

$T_p(n, W)$ be the time taken to perform parallel operation of problem size or workload W using n processors;

Then the speedup with n processors is:

$$S'(n) = \frac{T_s + T_p(1, W)}{T_s + T_p(n, W)}$$



Assuming that parallel operations achieve a linear speedup

(i.e. these operations take $1/n$ of the time to perform on one processor) Then, $T_p(1, W) = n \cdot T_p(n, W)$

Let α be the fraction of sequential work load in problem, i.e.

$$\alpha = \frac{T_s}{T_s + T_p(n, W)}$$

$$S'(n) = \frac{T_s + n \cdot T_p(n, W)}{T_s + T_p(n, W)} = \alpha + n(1 - \alpha) = n - \alpha(n - 1)$$