

Fundamentals of Computer Architecture

Madiba Hudson-Quansah

CONTENTS

CHAPTER 1	TERMINOLOGY	PAGE 3
1.1	Computer Hardware	3
1.2	Key Characteristics of Hardware	3
1.3	Operating System Fundamentals Execution of Programs — 4	3
1.4	System Architecture Von Neumann Architecture — 4 • Harvard Architecture — 4	4
1.5	Components of a Computer System	4
1.6	Assembly Language	5
1.7	General Terminology	5
CHAPTER 2	BASIC BUILDING BLOCKS	PAGE 7
CHAPTER 3	TRANSFORMATION HIERARCHY	PAGE 8
CHAPTER 4	CLASSES OF COMPUTERS	PAGE 9
4.1	Internet of Things / Embedded Computers	9
4.2	Personal Mobile Devices (PMD)	10
4.3	Desktop Computing	10
4.4	Servers Availability — 11 • Scalability — 11 • Throughput — 11	10
4.5	Clusters / Warehouse Scale Computers (WSC) WSCs vs Supercomputers — 11	11
CHAPTER 5	CLASSES OF PARALLELISM AND PARALLEL ARCHITECTURES	PAGE 12
CHAPTER 6	DEFINING COMPUTER ARCHITECTURE	PAGE 14
6.1	Instruction Set Architecture	14

CHAPTER 7	ENERGY EFFICIENCY AND SYSTEMS ARCHITECTURE	PAGE 15
7.1	Energy Efficiency	15
	Metrics — 15	
	7.1.1.0.1 Power Usage Effectiveness (PUE)	15
	7.1.1.0.2 Green Energy Usage	16
	7.1.1.0.3 Performance Per Watt (PPW)	16
	7.1.1.0.4 Energy Star Rating	16
	7.1.1.0.5 Carbon Footprint	16
	7.1.1.0.6 Dynamic Voltage and Frequency Scaling (DVFS)	16
	7.1.1.0.7 Idle Power Consumption	16
	7.1.1.0.8 Energy Consumption Index (ECI)	16
	7.1.1.0.9 Energy Efficiency Ratio (EER)	17
	7.1.1.0.10 Power Supply Efficiency (PSE)	17
	7.1.1.0.11 Renewable Energy Factor (REF)	17
7.2	Importance of Efficient Designs in Modern Systems	17
CHAPTER 8	KEY DESIGN PRINCIPLES	PAGE 18
8.1	Modularity	18
8.2	Scalability	18
8.3	Reliability	18
8.4	Maintainability	19
8.5	Balancing Trade-offs	19

Chapter 1

Terminology

1.1 Computer Hardware

Definition 1.1.1: Hardware

The physical components of a computer system that can be seen and touched.

- Central Processing Unit (CPU) - Instruction sets and execution.
- Memory - RAM / ROM, caching mechanisms and memory hierarchy.
- Input / Output (IO)- Peripherals, buses and device controllers.
- Motherboards and Chipsets - Physical layout and connectivity of components.

1.2 Key Characteristics of Hardware

- Physical Components - Physical parts of a computer system.
- Electronic Circuits - Electrical circuits that perform functions such as processing data, storing information and facilitating communication between other components.
- Peripheral Devices - Devices that are connected to the computer system to provide additional functionality.
- Physical Infrastructure - Servers, routers, switches.
- Assembly of Components - The physical assembly of components to form a functional computer system.
- Firmware - Software that is embedded in hardware components for low level control of the specific component.
- Mechanical Parts - Physical parts of a computer system that are not electronic.

1.3 Operating System Fundamentals

Definition 1.3.1: Process

An instance of a program running on a computer system. A process differs from a program by having its own memory space and CPU time, i.e. A program becomes a process when it is loaded into memory and executed by the CPU.

- Process Management - Processing handling - Scheduling, Multitasking, Threads.
- Memory Management - Virtual Memory, Paging, Segmentation.
- File Systems (FS) - How data is stored, accessed, and organized on ROM.

- I/O Management - Managing data transfer between peripherals and the CPU.
- Security - Authentication, Authorization, Encryption, and Firewalls.

1.3.1 Execution of Programs

Program execution from source files goes through the following steps:

1. Compilation - Source code is compiled into machine code or an intermediate language.
2. Linking - Libraries and dependencies are linked to the executable file.
3. Loading - The executable file is loaded into memory by a loader software and executed by the CPU.
4. CPU Time - The CPU scheduler allocates CPU time to the process for execution.
5. Process Running - The process is executed by the CPU and performs the required operations.
6. Process Termination - The process is terminated and the resources are released.

1.4 System Architecture

- Von Neumann Architecture - CPU, Memory, IO, and Bus.
- Harvard Architecture - Separate memory for data and instructions.
- System Buses - Communication between different components of the system, PCI, USB.
- Interrupts and Handling - Managing interrupts from hardware devices for processing.

1.4.1 Von Neumann Architecture

Definition 1.4.1: Von Neumann Architecture

A computer architecture that is based on the concept of a stored-program computer. The Von Neumann architecture consists of a CPU, Memory, IO, and a Bus. The CPU fetches instructions from memory, decodes them, and executes them and then fetches the data the instructions operate on from the same memory usually at another memory address/location. The Von Neumann architecture is used in most modern computers.

1.4.2 Harvard Architecture

Definition 1.4.2: Harvard Architecture

A computer architecture that has separate memory for data and instructions. The Harvard architecture allows the CPU to fetch data and instructions simultaneously, which can improve performance. The Harvard architecture is used in some embedded systems and microcontrollers.

1.5 Components of a Computer System

Definition 1.5.1: System

A collection of components that work together to perform complex computational tasks, manage resources, or provide specific services. These components are interconnected and interdependent on each other.

- CPU - Data Path and Control Unit

Data Path - Arithmetic and Logic Unit (ALU), Registers, and Cache.

Control Unit - Instruction Fetch, Decode, and Execute. Controls the flow of electrons / data to perform operations in the CPU.

- Memory - Speed Size trade-off. The faster the memory the smaller the size. Memory hierarchy.
 1. Registers - Fastest memory, used to store data that is currently being processed.
 2. Cache - Faster than RAM, used to store frequently accessed data. L1, L2, L3.
 3. RAM - Random Access Memory, used to store data that is currently being processed.
 4. ROM - Read Only Memory, used to store firmware and boot instructions.
 - SSD - Solid State Drive, faster than HDD, used to store data.
 - HDD - Hard Disk Drive, slower than SSD, used to store data.
 - Optical Drives - CD, DVD, Blu-ray, used to store data.
 - Magnetic Tapes - Used for long term storage of data / Archival.
- Input / Output

1.6 Assembly Language

Definition 1.6.1: Assembly

A low-level programming language that is used to write programs that are executed by the CPU. Assembly language is specific to the CPU architecture and provides a way to directly control the hardware components of the system. Maps mnemonics to machine code instructions. Example **ADD**, **MOV**, **JUMP**

Assembly language is less productive than high-level languages but provides more control over the hardware components of the system. Assembly language programs are translated into machine code by an assembler and executed by the CPU.

1.7 General Terminology

Definition 1.7.1: Transistor

The fundamental building block of modern electronic devices, acting as a switch for electrical signals.

Definition 1.7.2: Latency

The speed at which memory can be accessed. The time taken for a CPU to access memory.

Definition 1.7.3: Throughput

The amount of data that can be processed in a given amount of time. The number of instructions that can be executed per second.

Definition 1.7.4: Cache Hit Ratio

The percentage of memory accesses that are found in the cache. A high cache hit ratio indicates that the cache is effectively storing frequently accessed data.

Definition 1.7.5: RISC

Reduced Instruction Set Computer. A computer architecture that uses a small set of simple instructions that can be executed quickly. RISC architectures are designed to be efficient and fast by simplifying the instruction set and reducing the complexity of the CPU.

Definition 1.7.6: ARM

Advanced RISC Machine. A family of RISC architectures that are widely used in embedded systems, smartphones, tablets, and other devices. ARM processors are known for their low power consumption and high performance.

Theorem 1.7.1 Moore's Law

The number of transistors on a microchip doubles approximately every two years, resulting in an exponential increase in computing power. Moore's Law has been a driving force behind the rapid advancement of computer technology.

Theorem 1.7.2 Amdahl's Law

Amdahl's Law is a formula that describes the theoretical speedup of a program when running on multiple processors. The formula states that the speedup of a program is limited by the fraction of the program that can be parallelized. Amdahl's Law is used to analyse the performance of parallel computing systems.

$$S(n) = \frac{1}{(1 - P) + \frac{P}{n}}$$

Where $S(n)$ is the speedup of the program running on n processors, P is the fraction of the program that can be parallelized.

Definition 1.7.7: Dennard Scaling

Dennard scaling is a principle that states that as transistors get smaller, their power density remains constant. This principle has allowed for the continued increase

Improvements in processor performance has slowed down due to the following factors

- Transistors no longer getting much better because of the slowing of Moore's Law 1.7 and the end of Dennard Scaling 1.7.
- The unchanging power budgets for microprocessors.
- The replacement of the single processor with several energy-efficient processors.
- The limits to multiprocessing to achieve Amdahl's Law 1.7.

Chapter 2

Basic Building Blocks

- Electrons - Negatively charged particles that flow through electrical circuits.
- Voltage - The difference in electric potential between two points in an electrical circuit.
- Current - The flow of electrons through an electrical circuit.

Chapter 3

Transformation Hierarchy

- Problem
- Algorithm
- Program /Language
- System Software
- Software / Hardware Interface
 - ISA
 - * RISC - Modular relatively simple instructions, energy efficient.
 - * CISC - Complex instructions, more powerful, less energy efficient.
- Micro-architecture
- Logic
- Devices
- Electrons

Chapter 4

Classes of Computers

Rapid advancements in computer technology have led to the development of different classes of computers that are optimized for specific tasks and applications. These classes of computers can be broadly categorized into the following categories:

- Internet of Things / Embedded Computers
- Personal Mobile Devices
- Desktop Computing
- Servers
- Clusters/Warehouse Scale Computers

Feature	Personal Mobile Device	Desktop	Server	Clusters	IOT
Price of System	\$100 - \$1000	\$300 - \$2500	\$5000 - \$10,000,000	\$100,000 - \$200,000,000	\$10 - \$100,000
Price of Microprocessor	\$10 - \$100	\$50 - \$500	\$200 - \$2000	\$50 - \$250	\$0.01 - \$100
Critical System Design Issues	Cost, energy, media performance, responsiveness	price-performance, energy, graphics performance	Throughput, availability, scalability, energy	Price-performance, Throughput, energy proportionality	Price, energy, application-specific performance

Table 4.1: Summary of Classes of Computers

4.1 Internet of Things / Embedded Computers

Definition 4.1.1: Embedded Computer

A computer system that is designed to perform a specific task or function. Embedded computers are used in a wide range of applications, including consumer electronics, industrial automation, and automotive systems. Embedded computers are typically small, low-power devices that are optimized for a specific task.

Definition 4.1.2: Internet of Things (IOT)

Embedded computers that are connected to the internet, typically wirelessly, collecting useful data about their environs and interacting with the physical world. Some examples of IOT devices include smart thermostats, smart watches, smart cars, and smart homes.

Embedded computers have the widest spread of processing power and cost. They include 8-bit to 32-bit processors that may cost a penny, and high end 64-bit processors for cars and network switches that cost \$100. Price is a key factor in the design of computers for embedded systems.

4.2 Personal Mobile Devices (PMD)

Definition 4.2.1: Personal Mobile Device

A small, portable computing device that is designed to be used on the go. Personal mobile devices include smart-phones, tablets, and wearable devices. These devices are optimized for mobility, battery life, and connectivity.

Cost is also a key factor in the design of PMDs, with energy efficiency and media performance also being critical. Applications of PMDs are often web-based and media oriented. Processors in PMDs are often also considered to embedded computers because of their low power consumption and small size but are separated due to their ability to run externally developed software and share many features with desktop computers.

Responsiveness and predictability are key characteristics for media-applications, often requiring real-time performance.

Definition 4.2.2: Real Time Performance

A segment of an application that has an absolute maximum execution time. For example, in playing a video on a PMD, the time to process each frame is limited since the processor must accept and process the next frame shortly.

In other applications another requirement arises where the average time for a particular task is constrained as well as the number of instances when some maximum time is exceeded. This is known as **Soft Real Time Performance**.

4.3 Desktop Computing

Definition 4.3.1: Desktop Computer

A personal computer that is designed to be used on a desk or table. Desktop computers are typically larger and more powerful than personal mobile devices, with more storage, memory, and processing power. Desktop computers are used for a wide range of applications, including gaming, multimedia production, and office work.

The desktop market seeks to optimize price-performance chiefly, with energy and graphics performance also being critical.

Definition 4.3.2: Price-Performance

The combination of performance, measures in terms of compute performance and graphics performance, and the price of the computer system.

4.4 Servers

Definition 4.4.1: Server

A computer system that is designed to provide services or resources to other computers on a network. Servers are used for a wide range of applications, including web hosting, email, file storage, and database management. Servers are typically more powerful than desktop computers and are optimized for throughput, availability, and scalability.

The server market is chiefly characterized by the maximization of availability, scalability and throughput.

4.4.1 Availability

Definition 4.4.2: Availability

The percentage of time that a server is operational and available to provide services. Availability is a critical factor for servers that are used in mission-critical applications, such as e-commerce websites and financial systems.

4.4.2 Scalability

Definition 4.4.3: Scalability

The ability of a server system to grow in response to increasing demand for the services they support for an expansion in functional requirements.

The ability of a server to scale up computing capacity, memory, storage, and the I/O bandwidth is critical for servers that are used in applications that require high performance and reliability.

4.4.3 Throughput

Definition 4.4.4: Throughput

The overall performance of a server system, measured in terms of the number of requests it can handle per second or the amount of data it can process in a given time period.

4.5 Clusters / Warehouse Scale Computers (WSC)

Definition 4.5.1: Cluster

A collection of desktop computers or servers connected by local area networks to act as one unified computing resource. With each node, a separate computer system, running its own operating system and communicating via a network protocol.

Price-Performance is also critical to WSCs as they are so large with the majority of the cost of a warehouse associated with power and cooling the computers inside the warehouse. Availability is also crucial for WSCs as the cost for downtime is very high.

4.5.1 WSCs vs Supercomputers

Supercomputers are designed to optimize floating-point performance and are used for scientific and engineering, running large communication-intensive batch programs that can run for weeks at a time. WSCs are designed to emphasize interactive applications, large-scale storage, dependability, and high internet bandwidth.

Chapter 5

Classes of Parallelism and Parallel Architectures

Definition 5.0.1: Parallelism

The ability to perform multiple tasks simultaneously. Parallelism can be achieved at different levels of a computer system, including instruction level, task level, and data level.

There are mainly two kinds of parallelism in applications:

Data-Level Parallelism (DLP) - The ability to perform the same operation on multiple data elements simultaneously.

Task-Level Parallelism (TLP) - The ability to perform tasks simultaneously and independently.

Computer hardware can be designed to exploit these two kinds of parallelism in four major ways:

Instruction Level Parallelism (ILP) - Exploits DLP with the use of compilers to optimize code to perform tasks like pipelining, and speculative execution.

Vector architectures, graphics processing units (GPUs) and multimedia instruction sets - Exploits DLP by applying a single instruction to a collection of data in parallel.

Thread-level parallelism - Exploits DLP or TLP in a hardware model that allows for interaction between parallel threads.

Request-Level parallelism - Exploits TLP among largely decoupled tasks specified by the programmer of operating system.

With Flynn's Taxonomy, all computers can be classified into categories based on the way they handle parallelism in the instruction and data streams. The four categories are:

Single Instruction, Single Data (SISD) - One instruction stream and one data stream are processed at a time. The uni-processor, or a single-core computer. Can perform ILP.

Single Instruction, Multiple Data (SIMD) - The same instruction is executed by multiple processors using different data streams. SIMD computers exploit DLP by applying the same operations to multiple items of data in parallel. Each processor has its own data memory, but there is a single instruction memory and control processor that fetches and dispatches instructions.

Multiple Instruction, Single Data (MISD) - Multiple processors execute different instructions on the same data stream. MISD computers are not common and are not used in practice.

Multiple Instruction, Multiple Data (MIMD) - Each processor fetches its own instructions and operates on its own data and targets TLP. MIMD computers are more flexible than SIMD computers and can be used for a wider range of applications, but is inherently more complex and expensive than SIMD.

This taxonomy is not mutually exclusive, and many computers can be classified into more than one category. For example, a GPU can be classified as both SIMD and MIMD.

Chapter 6

Defining Computer Architecture

6.1 Instruction Set Architecture

Definition 6.1.1: Instruction Set Architecture (ISA)

The boundary between software and hardware. The actual instructions the computer/processor receives to process data

Class of ISA - Most modern ISAs are classified as general-purpose register architectures, where operands are either registers or memory locations.

Memory Addressing - All modern computers use byte addressing to access memory operands, with some architectures like ARMv8, requiring objects to be aligned. An access to an object of size s bytes at byte address A is aligned if $A \bmod s = 0$, alignment allows for generally faster accessing.

Addressing Modes - In addition to specifying registers and constant operands addressing modes specify the address of a memory object. For example RISC-V have addressing modes including Register, Immediate, Displacement.

Types and size of operands - Most ISAs support operand sizes of 8-bit (ASCII characters), 16-bit (Unicode character / Half word), 32-bit (integer / word), 64-bit (long integer/double word), IEEE 754 floating point 32-bit (single precision), and 64-bit (double precision).

Operations -

Chapter 7

Energy Efficiency and Systems Architecture

Definition 7.0.1: System Architecture

The way in which the components of a computer system are organized and connected. System architecture includes the design of the CPU, memory, IO, and bus, as well as the way in which these components interact with each other.

7.1 Energy Efficiency

Definition 7.1.1: Energy Efficiency

The ability of a computer system to perform tasks using the least amount of energy possible. Energy efficiency is important for reducing the environmental impact of computing and lowering the cost of operating computer systems.

(DVFS)

Definition 7.1.2: Performance

The speed at which a computer system can perform tasks. Performance is typically measured in terms of the number of instructions executed per second or the amount of data processed per second.

$$\text{Performance} = \frac{1}{\text{Execution Time}}$$

7.1.1 Metrics

7.1.1.0.1 Power Usage Effectiveness (PUE)

Definition 7.1.3: Power Usage Effectiveness (PUE)

Used to evaluate the energy efficiency of a data centre. The ratio of the total facility power to the power used by the equipment in a data centre. PUE is used to measure the energy efficiency of a data centre, with lower values indicating higher efficiency. Always greater than 1.

$$\text{PUE} = \frac{\text{Total Facility Power}}{\text{IT Equipment Power}}$$

7.1.1.0.2 Green Energy Usage

Definition 7.1.4: Green Energy Usage

The percentage of energy consumed by a computer system that comes from renewable sources, such as solar, wind, or hydroelectric power. Reflects the system's environmental sustainability.

7.1.1.0.3 Performance Per Watt (PPW)

Definition 7.1.5: Performance Per Watt (PPW)

The amount of performance that can be achieved per watt of power consumed. PPU is used to measure the energy efficiency of a computer system, with higher values indicating higher efficiency.

$$\text{PPW} = \frac{\text{Computational performance}}{\text{Power Consumed}}$$

7.1.1.0.4 Energy Star Rating

Definition 7.1.6: Energy Star Rating

A certification program indicating that a computer or other electronic devices meet specific energy efficiency guidelines set by the U.S. Environmental Protection Agency (EPA). Devices with a higher star rating are more energy efficient.

7.1.1.0.5 Carbon Footprint

Definition 7.1.7: Carbon Footprint

The total amount of greenhouse gas emissions, measured in Carbon Dioxide equivalents (CO_2), produced directly or indirectly by a computer system.

7.1.1.0.6 Dynamic Voltage and Frequency Scaling (DVFS)

Definition 7.1.8: Dynamic Voltage and Frequency Scaling

Allows for the adjustment of a system's voltage and frequency based on workload requirements. Measures how effectively a system can scale its power consumption in response to varying workloads.

7.1.1.0.7 Idle Power Consumption

Definition 7.1.9: Idle Power Consumption

The power consumption of a computer system when it is in an idle or low-activity state. Low idle power consumption is vital especially in situations where the system spends significant time in idle states.

7.1.1.0.8 Energy Consumption Index (ECI)

Definition 7.1.10: Energy Consumption Index

The total energy consumed by a computer system over a specific period of time, often normalized to a standard unit of time. Provides a holistic view of the energy consumption of a system and useful in long term sustainability assessments.

7.1.1.0.9 Energy Efficiency Ratio (EER)

Definition 7.1.11: Energy Efficiency Ratio

Measures the energy efficiency of a computer system by comparing the useful work output to the energy input. Commonly used in terms of cooling units within data centres.

$$\text{EER} = \frac{\text{Useful Work Output}}{\text{Energy Input}}$$

7.1.1.0.10 Power Supply Efficiency (PSE)

Definition 7.1.12: Power Supply Efficiency (PSE)

Assesses the efficiency of a computer system's power supply unit. The higher the better

$$\text{PSE} = \frac{\text{Output Power}}{\text{Input Power}} \times 100$$

7.1.1.0.11 Renewable Energy Factor (REF)

Definition 7.1.13: Renewable Energy Factor (REF)

Quantifies the proportion of a data centre's energy that comes from renewable sources.

7.2 Importance of Efficient Designs in Modern Systems

Resource Optimization - Efficient design minimizes resource usage (CPU, Memory, energy), while maximizing performance.

Cost Saving - Optimization resource utilization leads to less overall resources being using saving cost and increasing the overall return on investment (ROI) for organizations.

Environmental Impact - Promotes sustainability by reducing energy consumption and carbon footprint.

User Experience - Enhances system responsiveness and reliability leading to a better user experience.

Chapter 8

Key Design Principles

8.1 Modularity

Definition 8.1.1: Modularity

Involves breaking a system down into independent, interchangeable and cohesive parts, with each part performing a specific function. This allows for independent development, testing and maintenance of different parts of system reducing complexity, i.e. Abstraction of various of the systems various operations.

Modularity is important as it allows for enhanced flexibility, scalability, re usability and maintainability of a system.

8.2 Scalability

Definition 8.2.1: Scalability

A system's ability to handle an increasing workload by adding resources without significantly affecting performance. A well modularized system can be easily scaled by replicating or adding more modules targeting the specific performance bottleneck.

Scalability ensures that a system can grow and adapt to changing requirements and effectively utilize available resources as workload increases.

8.3 Reliability

Definition 8.3.1: Reliability

A measure of a system's ability to perform its intended function consistently and accurately with minimal downtime over time. Reliability is achieved through redundancy / backup systems, fault tolerance, and error detection and recovery..

Reliability is critical for mission-critical systems, such as servers, where downtime can result in significant consequences.

8.4 Maintainability

Definition 8.4.1: Maintainability

The ease at which a system can be repaired, updated, or modified. A well-designed system is easy to maintain and requires minimal effort to fix issues or add new features. Maintainability is achieved through good documentation, debugging tools, and modular design.

Maintainability reduces downtime, the cost of ownership, and extends the life of a system.

8.5 Balancing Trade-offs