

# RESUMO HARDWARE – Arquitetura de computadores

Prof.: Antônio Eugênio

## SUMÁRIO

1.	Introdução à Arquitetura de Computadores .....	3
1.1	Computador .....	4
1.1.1	Sistema de Computador Típico.....	5
1.1.2	Tipos de Computadores.....	7
1.2	Unidades Métricas nos Sistemas Computacionais .....	9
1.3	O Modelo de Von Neumann .....	10
1.4	Arquitetura de Harvard .....	11
1.5	Unidade Central de Processamento (CPU) .....	12
1.5.1	A Unidade de Controle .....	14
1.5.2	O decodificador de instruções.....	16
1.5.3	Unidade Lógica e Aritmética (ULA).....	17
1.5.4	Registradores.....	18
1.5.5	Memória Cache .....	18
1.5.6	Pipeline .....	18
1.6	Evolução dos Processadores .....	18
1.7	Arquiteturas Modernas e Tendências .....	19
1.8	Avaliação de Desempenho.....	22
1.8.1	Fatores de Desempenho.....	23
1.9	Unidade Lógica Aritmética (ULA) .....	23
1.10	O Modelo de Barramento do Sistema .....	23
1.10.1	Tipos de Barramentos .....	24
1.11	Clock e Clock Cycle.....	26
1.12	Conjuntos de Instruções (Instruction Set Architecture - ISA) .....	27
2.	Memória.....	30
2.1	Tipos de Memória .....	30
2.2	Características das Memórias .....	30
2.2.1	RAM (Random Access Memory).....	30
2.2.2	Cache .....	31
2.2.3	ROM (Read-Only Memory).....	32
2.2.4	EPROM (Erasable Programmable Read-Only Memory).....	32
2.2.5	EEPROM (Electrically Erasable Programmable Read-Only Memory) .....	33

2.2.6	Flash Memory .....	33
2.2.7	Registradores de Memória .....	35
2.2.8	Virtual Memory .....	36
2.3	Hierarquia de Memória.....	36
2.4	Encapsulamentos de Memória .....	37
3.	Memória Secundária .....	38
3.1	Discos Rígidos – (Hard Disk Drive – HDD) .....	38
3.2	Discos Flexíveis (Disquetes) .....	38
3.3	Discos Óticos .....	38
3.4	Drives de Estado Sólido (SSD) .....	39
3.5	Pen Drives .....	39
3.6	Cartões de Memória .....	39
4.	Placa-mãe .....	40
4.1	Componentes.....	40
4.2	Detalhes da Placa-mãe.....	41
4.2.1	Conexões e Slots .....	41
4.2.2	Chipset.....	41
4.2.3	CPU Socket.....	42
4.2.4	Slots de Memória RAM.....	42
4.2.5	Slots de Expansão .....	42
4.2.6	Conectores de Armazenamento .....	43
4.2.7	Conectores de E/S .....	44
4.2.8	BIOS e UEFI .....	44
4.2.9	Alimentação (Power Connectors).....	45
4.2.10	Área de Montagem e Layout.....	46
4.3	Padrões de Formatos de Placa-mãe.....	46
4.4	Sistema de Entrada e Saída (E/S) .....	46
5.	ADENDO 1 - “threads” de processamento .....	47
6.	ADENDO 2 - Comparativo de Processadores e Threads de Processamento .....	49

## 1. Introdução à Arquitetura de Computadores

Arquitetura de computadores é o ramo da ciência da computação e da engenharia que estuda o projeto, a organização e a funcionalidade dos sistemas computacionais, focando na forma como os diversos componentes de hardware interagem entre si para executar programas. Ela descreve como um computador é estruturado internamente, ou seja, como processador, memória, barramentos, dispositivos de entrada e saída e outros elementos estão interconectados, bem como o conjunto de instruções que o processador entende e executa.

A **Arquitetura de Computadores** refere-se aos atributos de um sistema que são visíveis ao programador e que impactam a execução de programas – por exemplo, o conjunto de instruções, a quantidade de bits usados para representar dados, mecanismos de endereçamento etc. Já a **Organização de Computadores** diz respeito à forma como os componentes operacionais são interligados e implementados internamente (unidade lógica e aritmética, unidades de controle, registradores, memórias, barramentos etc.), aspectos esses geralmente transparentes ao programador. Em outras palavras, *arquitetura define o que o computador faz*, enquanto *organização trata de como isso é realizado em nível de hardware*. Por exemplo, dois processadores podem compartilhar a mesma arquitetura (mesmo conjunto de instruções e formatos de dados), mas terem organizações internas distintas – com diferentes números de unidades de execução, tamanhos de cache ou tecnologias de implementação – sem alterar em nada o software que neles roda. Por outro lado, mudanças na arquitetura (como ocorreu na migração dos computadores Apple de CPUs PowerPC para **x86**) geralmente exigem adaptações no software, pois o conjunto de instruções e convenções básicas se altera.

Outra distinção fundamental é entre **hardware**, **software** e **firmware**. O **hardware** consiste nos componentes físicos do computador (circuitos integrados, chips de memória, dispositivos etc.), enquanto o **software** é o conjunto de programas e instruções executados pelo hardware. O **firmware** refere-se a programas específicos gravados em memória não volátil do hardware, geralmente responsáveis pelo controle básico de dispositivos (por exemplo, a BIOS/UEFI que inicializa o sistema). Esses conceitos se interrelacionam: a arquitetura define a interface de hardware e software (ex.: linguagem de máquina), e o firmware muitas vezes implementa funções de baixo nível, ficando numa camada intermediária.

A arquitetura de computadores define como componentes como CPU, memória e barramentos se organizam para executar programas. O modelo clássico de von Neumann (1945) unificou dados e instruções na mesma memória, permitindo alterar programas apenas modificando seu conteúdo na memória. Esse conceito de programa armazenado transformou os computadores, mas criou um *gargalo* quando a CPU precisa alternar entre buscar instruções e ler/escrever dados. Já a arquitetura de Harvard é um modelo de organização de computadores no qual a memória de instruções (programa) e a memória de dados são fisicamente separadas, ou seja, utilizam espaços distintos de armazenamento e barramentos independentes. Essa separação permite que o processador acesse simultaneamente uma instrução e um dado, o que resulta em melhor desempenho em comparação com a arquitetura de von Neumann, onde ambos compartilham o mesmo espaço de memória e barramentos.

Um computador moderno pode ser entendido através de níveis de abstração. Em um nível mais baixo estão os **circuitos lógicos digitais** (portas lógicas, flip-flops, etc.) que implementam funções binárias simples. Acima deles, vem o nível da **microarquitetura** (organização interna do processador), seguido pelo nível da **arquitetura do conjunto de instruções** (ISA), que é a visão do computador em termos das instruções que ele executa. Em níveis superiores situam-se o **sistema operacional** e as **linguagens de alto nível**, mas estes já extrapolam nosso foco nesta apostila. O importante é saber que a arquitetura de computadores atua como um “projeto conceitual” definindo repertório de instruções, formatos de dados, capacidades e funcionalidades vistas pelo programador, enquanto a organização e implementação concretas podem variar conforme as escolhas de projeto e tecnologia. Assim, a arquitetura de computadores é normalmente dividida em três níveis de abstração:

Nível	Descrição
<b>Arquitetura do Conjunto de Instruções (ISA)</b>	Define o <i>conjunto de instruções</i> , tipos de dados, registradores, modos de endereçamento e a interface entre hardware e software. É o “contrato” entre hardware e o compilador.
<b>Organização do Hardware (Microarquitetura)</b>	Trata de como a arquitetura é implementada fisicamente: pipelines, unidades de execução, caches, barramentos, controladores.
<b>Implementação Física</b>	Trata da realização concreta em circuitos integrados (silício), incluindo temporização, consumo de energia e desempenho térmico.

Portanto, a arquitetura de computadores é o projeto e a estruturação dos componentes físicos e lógicos de um sistema computacional. Ela engloba desde o nível mais baixo, como os transistores e circuitos eletrônicos, até a organização e funcionamento do sistema como um todo. Vamos explorar, resumidamente, os principais componentes e conceitos envolvidos na arquitetura de computadores.

## 1.1 Computador

Um computador é uma máquina complexa e altamente sofisticada que tem a capacidade de processar e manipular dados em uma escala exponencialmente rápida. Ele é a culminação de décadas de avanços em engenharia e ciência da computação, combinando hardware e software para realizar tarefas desde cálculos complexos até a manipulação de multimídia em tempo real.

Computador, de uma forma geral, é um equipamento composto de duas partes (hardware e software), que recebe dados através de um meio de entrada, processa-os sob o controle de um programa e produz resultados através de um meio de saída.

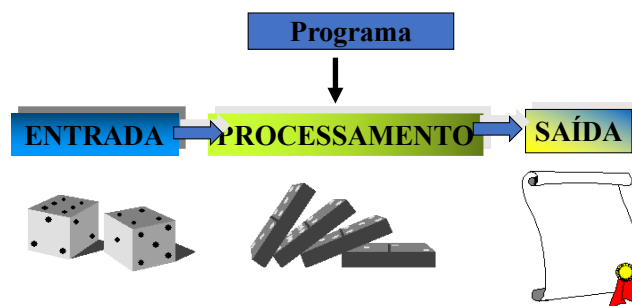


Figura 1 – Modelo de computador

O hardware compreende os componentes físicos do sistema. Executa, sob controle do software, as tarefas necessárias ao funcionamento e fluxo de dados entre os componentes do computador: UCP, Monitor de Vídeo, Teclado, Impressora, Drives, Mouse etc.

O software é o conjunto de programas (incluindo o Sistema Operacional), necessário ao funcionamento do hardware.

Em um computador moderno, existem vários níveis de máquina que desempenham funções distintas para garantir o funcionamento adequado do sistema. Estes são:

- **Nível de Transistores e Portas Lógicas:** Neste nível, o funcionamento do hardware é descrito em termos de transistores individuais e das operações lógicas que eles realizam (como AND, OR, NOT). Este é o nível mais baixo do computador formado por componentes eletrônicos e fios. As portas lógicas são implementadas usando transistores e fios de conexão.
- **Nível de Unidades Funcionais:** Aqui, os transistores são organizados em blocos que realizam funções específicas, como a ALU (Unidade Lógica e Aritmética) e os registradores. Nesse nível os registradores internos da CPU, a unidade lógica e aritmética e, a memória do computador é organizada sob a forma de unidades funcionais, de acordo com a função que desempenham para realizar as transferências de dados entre estas unidades funcionais.

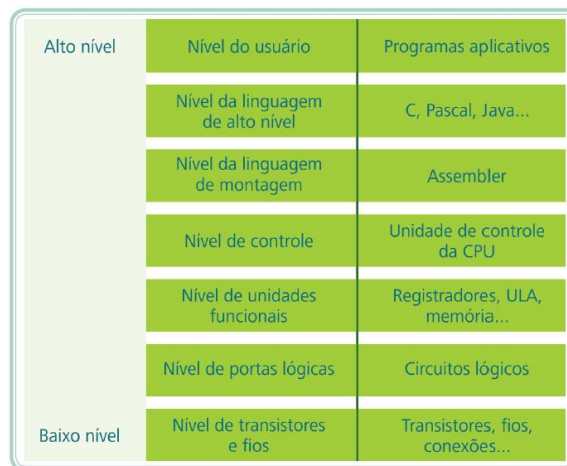


Figura 2 - Níveis de máquina num computador moderno

- **Nível de Controle:** Este nível coordena as operações das unidades funcionais e controla o fluxo de dados entre elas. Aqui a unidade de controle, que está dentro do processador, efetua as devidas transferências de dados entre os registradores, memória e dispositivos de entrada e saída. Essa transferência é feita através de sinais de controle por um circuito lógico.
- **Nível de Linguagem de Montagem:** É uma representação simbólica das instruções de máquina e fornece uma interface mais compreensível para programadores, substituindo os códigos binários por mnemônicos. esse é o nível onde as instruções são interpretadas e executadas pelo processador. Os programas desenvolvidos em linguagens de alto nível são traduzidos para uma linguagem de montagem ou Assembler, que apresenta um relacionamento direto com as instruções que o processador consegue executar.
- **Nível de Linguagem de Alto Nível:** Neste nível, os programadores utilizam linguagens como Python, C++, Java etc., que oferecem uma abstração mais alta em relação ao hardware.
- **Nível do Usuário:** É o nível no qual os usuários interagem diretamente com o software de aplicativos, como navegadores, editores de texto, jogos, entre outros.

Esses níveis representam uma progressão de abstração, indo dos detalhes físicos de transistores e portas lógicas até a interação amigável com o software em um nível de usuário.

### 1.1.1 Sistema de Computador Típico

Um sistema de computador moderno é um ecossistema sofisticado que integra diversos componentes que trabalham juntos para executar operações computacionais. Os principais elementos são:

- **Unidade Central de Processamento (CPU):** Responsável pela execução de operações e processamento de dados. Processadores modernos podem ter múltiplos núcleos (multi-core) e suportar a execução simultânea de várias tarefas através de tecnologias como Hyper-Threading e paralelismo.
- **Memória:** Armazena temporariamente dados e instruções enquanto o sistema está em funcionamento. A **Memória de Acesso Aleatório (RAM)** é volátil e de alta velocidade, crucial para o desempenho do sistema. Sistemas mais recentes utilizam DDR5 e **LPDDR5 (Low Power Double Data Rate 5)**, que é um tipo de memória RAM projetada para dispositivos que precisam de alto desempenho com baixo consumo de energia, como smartphones, tablets, laptops ultrafinos e dispositivos embarcados., oferecendo maior largura de banda e menor consumo de energia.
- **Dispositivos de Entrada/Saída (I/O):** Permitem a comunicação do computador com o ambiente externo. Dispositivos como **teclados, mouses, monitores, dispositivos de toque (touch screens)**, e **sensores** são comuns em sistemas modernos. A conectividade sem fio, como Bluetooth e Wi-Fi, também é amplamente usada para periféricos.
- **Barramentos de Sistema:** Facilitam a comunicação entre os componentes internos do sistema. Os barramentos modernos incluem o **PCIe (Peripheral Component Interconnect Express)**, que oferece alta velocidade para placas de vídeo, SSDs e outros dispositivos, e o **USB 4.0**, que conecta periféricos externos com altas taxas de transferência de dados.
- **Armazenamento Secundário:** Preserva dados de forma persistente, mesmo quando o sistema está desligado. Hoje, **unidades de estado sólido (SSD)** são o padrão, substituindo em grande parte os discos rígidos (HDD) por serem mais rápidos e mais confiáveis. **SSDs NVMe (Non-Volatile Memory Express)** conectados via PCIe oferecem altas velocidades de leitura e escrita, essenciais para desempenho superior em aplicações modernas.
- **Periféricos e Dispositivos Externos:** O computador pode ser expandido para incluir uma variedade de dispositivos, como **impressoras, scanners, dispositivos de realidade aumentada (AR) e realidade virtual (VR)**, câmeras, entre outros. A crescente demanda por conectividade com dispositivos inteligentes (IoT) também faz parte do ecossistema computacional.

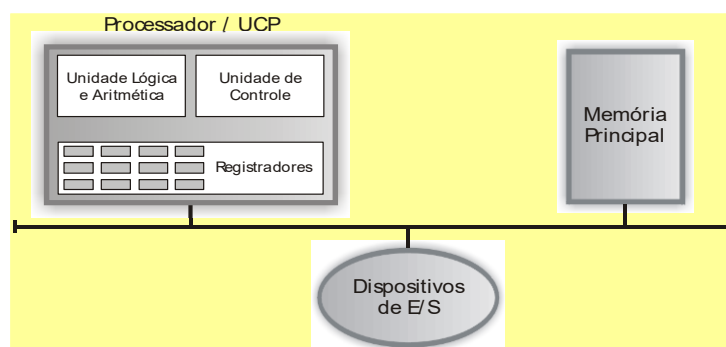


Figura 3 – Componentes de Hardware

As operações fundamentais realizadas por um computador correspondem às instruções que sua Unidade Lógica e Aritmética - ULA e outros componentes podem executar. Em um nível básico, podemos agrupar as operações de hardware em algumas categorias principais:

- **Operações Aritméticas:** incluem adição, subtração, multiplicação e divisão de números inteiros (a ULA geralmente suporta pelo menos soma e subtração diretamente, enquanto multiplicações e divisões podem ser executadas via circuitos dedicados ou em microcódigo). Em muitos processadores, há também suporte a operações aritméticas de ponto flutuante através de uma unidade de ponto flutuante (FPU). A ULA típica consegue determinar, além do resultado, propriedades como sinal (se o resultado é positivo ou negativo) e zero (resultado igual a zero). Operações de incremento/decremento (somar ou subtrair 1) também fazem parte do repertório aritmético.
- **Operações Lógicas/Booleanas:** englobam funções bit a bit como **AND** (E lógico), **OR** (OU lógico), **NOT** (negação, complemento de bits) e **XOR** (OU exclusivo). Essas operações tratam seus operandos como sequências de bits, efetuando a combinação lógica bit a bit. Por exemplo, uma instrução AND entre duas palavras de 32 bits realizará 32 operações AND em paralelo, um bit de cada operando por vez. Tais operações são úteis tanto para implementar operações aritméticas internas (AND e XOR são usadas em somadores binários) quanto para manipulação de máscaras e flags em nível de sistemas.
- **Operações de Deslocamento e Rotação de Bits:** são instruções que deslocam todos os bits de um operando para a esquerda ou direita, preenchendo com zero ou propagando o bit de sinal, ou fazendo rotação circular dos bits. *Shifts* são importantes por equivalerem a operações de multiplicação ou divisão por potências de 2 (um deslocamento para esquerda equivale a multiplicar por 2, se não houver overflow), além de permitirem isolamento de bits individuais ou campos dentro de palavras binárias. Processadores geralmente fornecem instruções de deslocamento lógico (preenche com zero) e aritmético (preserva o sinal no deslocamento para direita).
- **Operações de Comparação:** embora não sejam “operações” explícitas da ULA como uma unidade separada, a comparação de valores (igual, menor, maior) é realizada via operações aritméticas ou lógicas combinadas com análise de *flags*. Por exemplo, para saber se  $A = B$ , a ULA pode realizar  $A - B$  e verificar se o resultado foi zero (indicador de igualdade). Para  $A < B$ , verifica-se o sinal do resultado ou flags específicos. Muitas arquiteturas têm instruções de comparação que efetivamente realizam essas suboperações e ajustam registradores de status (flags) indicando o resultado da comparação (igual, menor, carry, overflow etc.).
- **Operações de Transferência de Dados:** embora não façam “cálculo” no sentido matemático, são operações fundamentais que movem dados de um lugar para outro. Incluem carregamento de registradores a partir da memória (load), armazenamento de registradores na memória (store) e movimentação de dados entre registradores (move). Essas instruções de movimentação não alteram o conteúdo binário além de copiá-lo, mas são essenciais para trazer dados para a ULA operar e enviar resultados de volta à memória ou para saída.
- **Operações de Controle de Fluxo:** novamente, não são operações de processamento de dados em si, mas afetam qual será a próxima instrução executada. Incluem desvios/jumps (condicionais ou incondicionais), chamadas de subrotina, retornos, etc. Essas operações dependem tipicamente do resultado de comparações (por exemplo: “se  $X == 0$ , salte para a instrução Y”) e permitem implementar estruturas de decisão e repetição nos programas.

### 1.1.2 Tipos de Computadores

Os tipos de computadores evoluíram significativamente ao longo do tempo, adaptando-se a diferentes necessidades de uso e avanço tecnológico. A seguir, são apresentados os principais tipos de computadores usados atualmente:

- **Desktops e Workstations:** Computadores de mesa tradicionais, projetados para uso geral em casa ou no escritório. As **workstations** são versões mais potentes, geralmente usadas por profissionais que necessitam de alto desempenho para tarefas como design gráfico, modelagem 3D, simulações científicas e edição de vídeos.
- **Laptops e Ultrabooks:** Dispositivos portáteis que combinam desempenho e mobilidade. Os **ultrabooks** são uma subcategoria que se destaca por serem mais finos, leves e energeticamente eficientes, com baterias de longa duração e uso de tecnologias avançadas, como **LPDDR5** e **SSDs NVMe** para desempenho otimizado.
- **Tablets e Convertíveis (2 em 1):** Dispositivos híbridos que podem atuar como laptop ou tablet. Geralmente equipados com telas sensíveis ao toque, os **2 em 1** possuem teclado destacável ou dobrável. São amplamente utilizados para produtividade e consumo de mídia, oferecendo uma combinação de mobilidade e funcionalidade.
- **Smartphones:** Computadores portáteis compactos com capacidades robustas de processamento e conectividade. Smartphones modernos possuem processadores de múltiplos núcleos, memória **LPDDR5** e armazenamento de alta velocidade, sendo capazes de executar tarefas avançadas, como edição de vídeos, jogos em alta definição e até tarefas de inteligência artificial.
- **Servidores:** Máquinas robustas que fornecem serviços e recursos para outros computadores em uma rede. **Servidores de data center** são projetados para manipular grandes volumes de dados, hospedagem de sites, e fornecer infraestrutura para **nuvem (cloud computing)**, oferecendo suporte a múltiplos usuários simultâneos e execução de processos complexos.
- **Mainframes:** Computadores de grande porte utilizados principalmente em corporações, bancos e governos para processar grandes volumes de transações simultâneas e dados críticos. Eles oferecem alta disponibilidade, escalabilidade e segurança.
- **Supercomputadores:** Máquinas extremamente poderosas, projetadas para executar cálculos complexos em alta velocidade. São usados em áreas como previsão do tempo, simulações científicas, inteligência artificial, e modelagem molecular. Os supercomputadores modernos utilizam arquiteturas paralelas e GPUs de alta performance para otimizar o processamento de grandes conjuntos de dados.
- **Sistemas Embarcados:** Computadores integrados em dispositivos específicos, como automóveis, eletrodomésticos, dispositivos médicos e equipamentos industriais. Esses sistemas são projetados para realizar tarefas específicas e podem ser encontrados em dispositivos da **Internet das Coisas (IoT)**, como smartwatches, câmeras inteligentes e sensores.
- **Computação em Nuvem (Cloud Computing):** Embora não seja um tipo de computador físico, a computação em nuvem permite o uso de recursos computacionais remotamente através da internet. As empresas e indivíduos podem acessar infraestrutura como serviço (IaaS), plataformas como serviço (PaaS), e software como serviço (SaaS) sem a necessidade de possuir hardware físico local, utilizando recursos de grandes data centers. IaaS (Infrastructure as a Service, ou Infraestrutura como Serviço) é um modelo de computação em nuvem que oferece recursos de infraestrutura de TI sob demanda por meio da internet. Em vez de comprar servidores físicos, redes, espaço de armazenamento ou data centers, uma organização pode alugar esses recursos de um provedor de nuvem. Exemplos: Amazon Web Services (AWS) – EC2; Microsoft Azure – Azure Virtual Machines; Google Cloud – Google Compute Engine. PaaS (Platform as a



6 Service), ou Plataforma como Serviço, é um modelo de computação em nuvem que fornece uma plataforma completa para desenvolvimento, execução e gerenciamento de aplicações, sem que o usuário precise lidar com a complexidade da infraestrutura subjacente (como servidores, armazenamento ou redes). O que inclui o PaaS?

- ✓ Sistema operacional
- ✓ Servidor web
- ✓ Banco de dados
- ✓ Ferramentas de desenvolvimento (como IDEs, bibliotecas, frameworks)
- ✓ Gerenciamento de aplicativos
- ✓ Serviços de monitoramento e segurança
- **Dispositivos de Realidade Aumentada e Virtual (AR/VR):** Computadores que oferecem experiências imersivas ao usuário, combinando ambientes virtuais ou sobreposições digitais ao mundo físico. Usados em setores como entretenimento, educação, simulações e até terapias médicas, esses dispositivos possuem alta capacidade de processamento gráfico e baixa latência para garantir uma experiência fluida.

## 1.2 Unidades Métricas nos Sistemas Computacionais

As unidades métricas são a linguagem numérica dos sistemas computacionais, utilizadas para quantificar capacidade de armazenamento, velocidade de transferência de dados e poder de processamento. Vão desde os bits individuais até os petabytes de dados, proporcionando uma escala precisa para descrever a capacidade e eficiência dos sistemas.

No nível de hardware, toda informação é representada de forma **binária**, utilizando dois símbolos básicos (0 e 1) que correspondem aos estados possíveis de circuitos eletrônicos (desligado/ligado, baixo/alto, falso/verdadeiro). O dígito binário – chamado *bit* – é a menor unidade de informação no computador. Por convenção, agrupa-se 8 bits para formar um **byte**, que é capaz de representar 256 valores distintos (de 0 a 255 em decimal). A capacidade das memórias e tamanhos de arquivos normalmente são medidos em bytes e seus múltiplos (KB, MB, GB, etc., onde 1 KB = 1024 bytes, 1 MB = 1024 KB, e assim por diante).

Sigla	Nome	Equivalência em Bytes	Potência de 2
<b>KB</b>	Kilobyte	1.024 bytes	$2^{10}$
<b>MB</b>	Megabyte	1.048.576 bytes	$2^{20}$
<b>GB</b>	Gigabyte	1.073.741.824 bytes	$2^{30}$
<b>TB</b>	Terabyte	1.099.511.627.776 bytes	$2^{40}$
<b>PB</b>	Petabyte	1.125.899.906.842.624 bytes	$2^{50}$
<b>EB</b>	Exabyte	1.152.921.504.606.846.976 bytes	$2^{60}$
<b>ZB</b>	Zettabyte	1.180.591.620.717.411.303.424 bytes	$2^{70}$
<b>YB</b>	Yottabyte	1.208.925.819.614.629.174.706.176 bytes	$2^{80}$
<b>RB</b>	Ronnabyte	1.237.940.039.285.380.274.899.124.224 bytes	$2^{90}$
<b>QB</b>	Quettabyte	1.267.650.600.228.229.401.496.703.205.376 bytes	$2^{100}$

Tabela 1 - Unidades de medida de armazenamento (em bytes)

Como mencionado, 8 bits formam um byte, que costuma ser a unidade básica para representação de um caractere de texto em muitos sistemas. Por exemplo, na codificação ASCII (um padrão antigo de codificação de caracteres), a letra "A" é representada pelo número 65 em decimal ( $01000001_2$  em binário), "B" por 66, etc. Atualmente, utiliza-se muito a codificação Unicode (UTF-8, UTF-16, etc.), que expande a representação para incluir caracteres acentuados, símbolos de diversas línguas e emojis, mas a ideia permanece: cada caractere ou símbolo é mapeado para um código numérico binário específico. Além de caracteres, qualquer outro tipo de dado (imagens, áudio, vídeos) é representado binariamente – por exemplo, imagens digitais são representadas por códigos para cada pixel (indicando cor), áudio por amostras numéricas de amplitude, e assim por diante. Tudo dentro do computador, seja código de programa ou dados do usuário, é reduzido a sequências de bits armazenadas na memória.

Em resumo, toda informação em um computador é representada por bits. A combinação de bits segue convenções para representar os diversos tipos de dados: inteiros em binário puro ou complemento de dois, reais em ponto flutuante, caracteres em códigos numéricos (ASCII/Unicode), dentre outros. O uso do sistema binário se dá porque é *simples e robusto* para implementação eletrônica – níveis de tensão podem representar 0 ou 1 de forma confiável, e portas lógicas podem manipular esses bits facilmente.

### 1.3 O Modelo de Von Neumann

O modelo de Von Neumann, proposto por John von Neumann, é a base conceitual da maioria dos computadores modernos. Ele introduziu uma abstração revolucionária, unificando as operações de processamento e a armazenagem de dados em um único conjunto de memória. Esse modelo permitiu que as instruções fossem executadas de forma sequencial, marcando uma inovação paradigmática na computação.

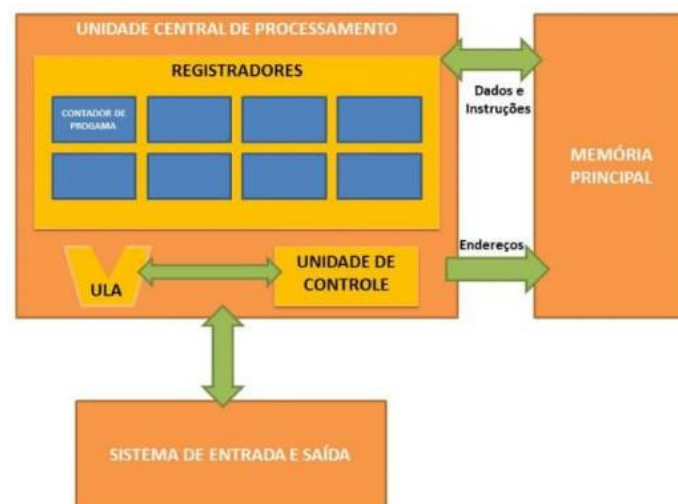


Figura 4 - Modelo de Von Neumann

A Arquitetura de Von Neumann é um modelo fundamental na concepção de computadores modernos. Foi proposta pelo matemático e físico John von Neumann nos anos 1940 e estabeleceu os princípios básicos que ainda são utilizados na maioria dos computadores hoje. Vamos detalhar os principais elementos desta arquitetura:

- Unidade Central de Processamento (CPU):

A CPU é o núcleo do sistema e é composta por:

- Unidade de Controle (UC): Responsável pela interpretação e execução de instruções. Controla o fluxo de dados e as operações.

- Unidade Lógica Aritmética (ULA): Executa operações matemáticas e lógicas como adição, subtração e operações de comparação.
- Memória: Na arquitetura de Von Neumann, tanto os dados quanto as instruções são armazenados na mesma memória endereçável. Isso significa que a CPU pode ler tanto dados como instruções da mesma área de armazenamento.
- Registrador de Instrução (IR): O IR mantém a instrução atualmente sendo executada pela CPU. Ele é atualizado conforme a CPU busca, decodifica e executa instruções.
- Contador de Programa (PC): O PC mantém o endereço da próxima instrução a ser buscada na memória. Após cada execução de instrução, o PC é incrementado para apontar para a próxima instrução.
- Arquitetura Sequencial: A execução de instruções em um computador de Von Neumann é sequencial e linear. Cada instrução é executada em ordem, uma após a outra, sem a capacidade de execução paralela de múltiplas instruções.
- Bus de Dados e bus de Endereços: A arquitetura de Von Neumann utiliza dois barramentos principais:

Bus de Dados: Transporta dados entre a CPU e a memória ou outros dispositivos.

Bus de Endereços: Transporta os endereços de memória, indicando a localização de leitura ou escrita.

- Características Principais:

Armazenamento de Programa: Programas e dados são armazenados na mesma memória, o que permite que um programa seja tratado como dados e vice-versa.

Acessibilidade Uniforme: Tanto dados quanto instruções são acessados da mesma forma, através de endereços.

Modificação de Programa: Um programa pode alterar a si mesmo, o que é uma característica essencial para execução de código dinâmico.

#### 1.4 Arquitetura de Harvard

A **arquitetura Harvard** separa fisicamente a memória de instruções da memória de dados. Na analogia da biblioteca, são duas estantes distintas: uma com livros de instruções de programa e outra com dados. A CPU possui dois conjuntos de barramentos paralelos – um para buscar instruções e outro para acessar dados – permitindo que, por exemplo, uma instrução seja obtida enquanto outra operação de dados é executada simultaneamente. Esse paralelismo elimina o gargalo de von Neumann, tornando o processamento mais rápido. Em contrapartida, a arquitetura Harvard é mais complexa e cara, pois exige duas memórias e dois barramentos separados. Ela é comum em microcontroladores e DSPs (Digital Signal Processor - é um tipo especializado de microprocessador projetado especificamente para realizar cálculos matemáticos complexos de forma rápida e eficiente sobre sinais digitais, como áudio, vídeo, comunicações e outros tipos de dados que exigem alta velocidade de processamento), onde a velocidade é crítica, enquanto a arquitetura de von Neumann prevalece em PCs e sistemas gerais devido à sua simplicidade.

### Quadro Comparativo: Von Neumann vs Harvard

Característica	Von Neumann	Harvard
<b>Memória</b>	Única memória compartilhada para dados e instruções.	Memórias separadas: uma só para dados e outra para instruções.
<b>Barramentos</b>	Um único conjunto de barramentos (dados, endereço e controle).	Dois conjuntos de barramentos independentes (um para dados, outro para instruções).
<b>Execução de instruções</b>	Em geral, uma instrução é buscada e depois executada sequencialmente.	Pode buscar instruções e dados em paralelo, permitindo execução mais rápida.
<b>Desempenho</b>	Sujeito ao <i>gargalo de von Neumann</i> : CPU espera alternadamente por dados/instruções.	Maior desempenho devido ao acesso simultâneo a memória de dados e de instruções.
<b>Custo e complexidade</b>	Menor custo, design mais simples.	Mais caro e complexo (duas memórias e buses).
<b>Uso típico</b>	Computadores gerais, PCs e servidores.	Sistemas embarcados, DSPs e microcontroladores (onde pipelining de instruções é importante).

#### 1.5 Unidade Central de Processamento (CPU)

A CPU é composta por vários elementos funcionais, incluindo registradores, decodificador de instruções, unidade de controle, unidade lógica aritmética (ULA) e registradores de estado. A ULA realiza operações como adição, subtração, comparação e operações lógicas.

A CPU pode ser vista como composta de:

- **Registradores Internos:** pequenos módulos de memória dentro da CPU, de altíssima velocidade, usados para armazenar temporariamente dados e endereços durante a execução. Existem *registradores de propósito geral* (que podem conter quaisquer dados, usados pelo programa) e *registradores especiais* ou *de controle* (usados pela CPU para finalidades específicas de controle). Entre os especiais, destacam-se:
- **PC (Program Counter ou Contador de Programa):** contém o endereço da próxima instrução a ser buscada na memória.
- **IR (Instruction Register ou Registrador de Instrução):** armazena a instrução atualmente em execução, após ser buscada da memória.
- **MAR (Memory Address Register ou Registrador de Endereço de Memória):** registra um endereço de memória a ser acessado.
- **MDR ou MBR (Memory Data/Buffer Register ou Registrador de Dados de Memória):** guarda temporariamente um dado lido da memória ou que será escrito na memória.

- **Registrador de Status/Flags:** conjunto de bits que refletem resultados de operações (por exemplo, flag Zero, Carry, Overflow, Sinal, etc.), utilizados para decisões de desvio e outras operações condicionais.
- Outros registradores específicos podem existir, como registros de índice, registrador de pilha (SP – Stack Pointer), registrador base de segmento (em arquiteturas segmentadas), etc., dependendo da arquitetura.
- **Unidade Lógica e Aritmética (ULA):** como visto, é responsável pelas operações de processamento de dados. A ULA conecta-se a vários desses registradores por meio de um *datapath* interno – isto é, linhas internas de dados que permitem transferir conteúdos de registradores para entrada da ULA e da ULA de volta para registradores de destino. Muitas arquiteturas possuem um conjunto de registradores de propósito geral que se comunicam com a ULA através de um barramento interno ou rede de interconexão.
- **Unidade de Controle (UC):** consiste em circuitos sequenciais e combinatórios (ou de um microprograma armazenado) que leem o IR, decodificam a instrução e emitem sinais de controle para todos os componentes citados (registradores, ULA, interfaces de memória/E/S). A UC geralmente inclui um elemento de sincronização (relógio) e um *decodificador de instruções*. No caso de UC microprogramada, inclui uma memória de controle onde está o microcódigo.
- **Interconexão Interna (Barramentos Internos):** dentro da CPU, há conexões que permitem a transferência de informações entre registradores, ULA e UC. Alguns designs usam um **barramento interno único**, onde apenas dois registradores por vez colocam dados para ULA operar e um terceiro recebe o resultado. Outros designs mais complexos têm múltiplos barramentos internos, permitindo várias transferências simultaneamente (por exemplo, leitura de dois registradores fonte simultaneamente em dois barramentos distintos para alimentar a ULA). De qualquer modo, a organização interna define como os dados fluem entre as partes do processador. A Figura 1 da seção anterior já ilustrava um diagrama simples da CPU dentro da arquitetura de von Neumann, com seus componentes interligados.

Em termos de **funcionamento geral**, podemos descrever o processo contínuo do processador assim: o PC fornece o endereço à memória para buscar uma instrução; a instrução entra no IR; a UC decodifica e aciona a ULA e outros componentes; a ULA possivelmente interage com registradores e memória para completar a instrução; o resultado pode atualizar registradores ou memória; então a UC incrementa ou ajusta o PC para a próxima instrução e repete. Essa é a chamada *sequência de instruções*. Quando uma instrução de desvio altera o PC para um endereço não sequencial, o fluxo de instruções muda de acordo (implementando loops, condicionais, chamadas de função, etc., no nível de máquina). Quando interrupções ocorrem, a UC salva o contexto (tipicamente empilhando PC e registradores importantes) e desvia para o vetor de interrupção apropriado – do ponto de vista do processador, é semelhante a um jump forçado para uma sub-rotina de serviço de interrupção, retornando depois.

Nos processadores modernos, essa estrutura básica é otimizada através de técnicas como **pipelining** (sobreposição de etapas de execução de múltiplas instruções) e **execução fora de ordem** (reordenar internamente instruções para melhor uso dos recursos). Por exemplo, um pipeline típico divide o ciclo de instrução em estágios: busca, decodificação, execução, acesso à memória, writeback, permitindo que até 5 instruções diferentes estejam em diferentes estágios simultaneamente. Contudo, essas otimizações não alteram a semântica básica da execução – elas servem para aumentar o throughput de instruções.

Portanto, o processador realiza constantemente as seguintes operações (Ciclo de instrução):

- Buscar instrução – o processador busca na memória a instrução a ser executada - uma instrução é uma ordem para que a UCP realize determinada operação (ex.: somar, subtrair, mover um dado de um local para outro, transferir um dado para um dispositivo de saída).

- Interpretar a instrução – a instrução é decodificada para determinar a ação que deve ser executada.
- Obter os dados – a execução da instrução pode necessitar a leitura de dados da memória ou dos dispositivos de entrada.
- Processar os dados – a execução da instrução pode necessitar de alguma operação aritmética ou lógica com os dados.
- Gravar os dados – a execução da instrução pode requerer a gravação dos dados na memória ou em um dispositivo de saída.

Para a máquina (computador) é necessário que cada instrução seja detalhada em pequenas etapas. Isso ocorre porque os computadores são projetados para entender e executar pequenas operações, ou seja, as operações mais básicas (ex.: soma, subtração). Essas pequenas etapas de uma instrução dependem do conjunto de instruções do computador. Assim, uma instrução de máquina pode ser definida pela formalização de uma operação básica que o hardware é capaz de realizar diretamente. Em outras palavras, consiste em transformar instruções mais complexas em uma sequência de instruções básicas e compreensíveis pelo processador.

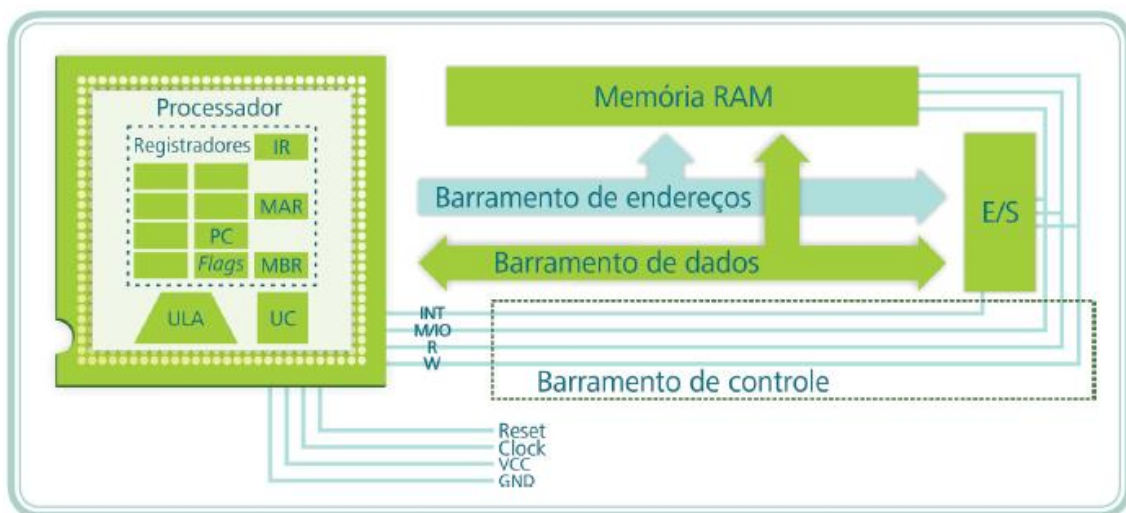


Figura 5 - Esquema de uma CPU

#### 1.5.1 A Unidade de Controle

É o "cérebro" da CPU. Ela é responsável por coordenar e controlar as operações de todos os outros componentes da CPU e determinar a sequência de operações a serem executadas. A UC interpreta as instruções da memória, decodifica-as e gera os sinais de controle necessários para coordenar as operações. Ainda, pela transferência de dados e instruções para dentro e para fora da CPU. Ela controla todas as ações a serem realizadas pelo computador, garantindo a correta manipulação dos dados e execução dos programas.

A Unidade de Controle, pode ser implementada de duas formas básicas: *UC cabeada* (hardwired) ou *UC microprogramada*. Numa UC cabeada, há circuitos lógicos sequenciais que, a partir do código de operação da instrução (e talvez do estado atual de execução), geram diretamente as sequências de sinais de controle necessários para cada etapa. Já numa UC microprogramada, as sequências de controle são armazenadas em uma memória especial (memória de controle) sob a forma de microinstruções – essencialmente um pequeno programa interno que, para cada instrução de máquina, executa um conjunto de passos elementares. As primeiras CPUs (décadas de 60-70) adotaram microprogramação para facilitar a implementação

de conjuntos de instruções complexos, enquanto designs mais modernos de alto desempenho tendem a usar controle cabeado para acelerar o ciclo de instruções. Em qualquer caso, para o estudioso iniciante, o importante é compreender o *papel* da UC: ela é quem “entende” a instrução e orquestra a ULA e demais componentes para realizá-la.

Vamos detalhar o **ciclo de instrução** sob controle da UC, já introduzido anteriormente, agora enumerando claramente as etapas:

1. **Busca (Fetch):** A UC coloca no barramento de endereços o conteúdo do **contador de programa (PC)** – que aponta para a próxima instrução – e emite um sinal de leitura de instrução. A memória responde fornecendo no barramento de dados o código binário da instrução localizada naquele endereço, o qual é carregado no **Registrador de Instrução (RI ou IR)** da CPU. (Durante esse processo, o PC é automaticamente incrementado para apontar para a próxima instrução sequencial, assumindo que não haja desvio.)
2. **Decodificação (Decode):** A UC analisa os bits da instrução no RI. Os campos da instrução são separados: identifica-se o opcode (código da operação) e os operandos. Caso algum operando resida na memória (por exemplo, instrução de acesso a variável), a UC calcula o **endereço efetivo** conforme o modo de endereçamento. Esse cálculo pode envolver somar registros (ex: registrador base + deslocamento) ou ler ponteiros em memória (modo indireto). A decodificação envolve lógica combinatória e eventualmente consulta a unidade de controle microprogramada (se existente) para definir a sequência de micro-operações necessárias.
3. **Busca de Operandos:** Se os operandos necessários não estão já em registradores internos, a UC aciona leituras de memória ou de registros de E/S para obtê-los. Por exemplo, se a instrução é ADD R1, X (sendo X um endereço de memória), nesta etapa a UC colocaria o endereço de X no barramento e leria o operando da memória para algum registrador interno (ou diretamente para a ULA). Operandos em registradores já estão prontos, e operandos imediatos vêm da própria instrução – nesses casos, esta etapa é trivial.
4. **Execução (Execute):** A UC agora configura a ULA para realizar a operação requerida. Ela envia à ULA um conjunto de sinais indicando qual função deve ser executada (por ex., somar), e quais registradores ou entradas usar (por exemplo, “pegue operando A do registrador R1 e operando B do registrador interno onde armazenamos X”). A ULA então realiza a operação – digamos, soma  $A + B$  – e o resultado é colocado no destino apropriado (um registrador ou talvez diretamente na memória, dependendo da instrução). A UC monitora sinais de status da ULA, como flags de overflow, zero ou carry, que podem ser armazenados em um **Registrador de Status** (também chamado registro de flags) para uso por instruções de desvio condicional.
5. **Escrita de Resultado (Writeback):** Caso a instrução produza um resultado que deve ser escrito na memória ou em registrador, essa é a etapa final em que isso é feito. Por exemplo, para um STORE R2, X, após ler a instrução e decodificá-la, a execução consiste em pegar o valor em R2 e colocar na memória no endereço X – a escrita efetiva acontece aqui. Para muitas instruções aritméticas, o resultado já foi colocado no registrador de destino durante a etapa de execução, então não há ação adicional. Contudo, se a operação envolveu E/S (por exemplo, enviar um dado para porta de saída), a UC cuidará de acionar o dispositivo correto durante essa etapa.

Terminada a execução da instrução corrente, o ciclo recomeça no passo 1 para buscar a próxima instrução (atual PC). Isso continua indefinidamente, a menos que ocorra alguma condição excepcional (como uma interrupção de hardware, uma exceção, ou o programa finalize).

Durante todas essas etapas, a Unidade de Controle aciona diversos sinais de controle: ela controla **multiplexadores internos** para roteamento de dados (por exemplo, selecionar se o operando da ULA vem de um registrador ou do barramento de memória), controla **ativação de escrita/leitura** em registradores e memória (sinaliza quando carregar um registrador ou escrever um valor), e coordena o **barramento** (assegurando que apenas um dispositivo esteja colocando dados no barramento de cada vez, evitando conflitos). Em sistemas microprogramados, cada

microinstrução da UC define um conjunto desses sinais para uma subetapa; em sistemas cabeados, redes lógicas de decodificação geram esses sinais temporizados pelo clock.

### 1.5.2 O decodificador de instruções

É um componente da Unidade de Controle que interpreta o código de operação (opcode) de uma instrução para determinar qual operação deve ser executada e quais operandos estão envolvidos.

Uma **instrução de máquina** é uma palavra binária que indica uma operação a ser realizada pela CPU e, geralmente, os operandos envolvidos nessa operação. As instruções constituem o vocabulário básico do processador, definindo tudo que ele pode fazer. Do ponto de vista do programador em linguagem de máquina/assembly, cada instrução possui um *mnemônico* (palavra simbólica) e campos de operandos. Por exemplo, uma instrução em assembly poderia ser ADD R1, R2, R3 (some os valores dos registradores R2 e R3 e armazene em R1) ou MOV AX, 1000h (mova o valor imediato 1000h para o registrador AX). Internamente, essas instruções são codificadas em binário de acordo com o formato definido pela arquitetura. Um formato comum é dividir a instrução em campos: um **código de operação (opcode)** que especifica qual operação executar (por ex., ADD, MOVE, JUMP) e campos que especificam os **operandos** – que podem ser registradores, endereços de memória ou valores imediatos. O tamanho da instrução em bits e o número de operandos variam conforme a arquitetura (arquiteturas RISC frequentemente têm instruções de tamanho fixo, por exemplo 32 bits, enquanto arquiteturas CISC possuem instruções de tamanho variável para permitir operandos mais complexos).

Uma instrução, em termos de arquitetura de computadores, é um comando ou operação específica que a CPU é capaz de executar. Ela é a unidade mais básica de ação que uma CPU pode compreender e executar. Cada instrução é codificada em linguagem de máquina, o que significa que é representada por um padrão de bits específico que a CPU é projetada para reconhecer e executar. Esses códigos binários são interpretados pela Unidade de Controle (UC) da CPU.

Uma instrução pode realizar várias operações, tais como:

- Transferência de Dados: Movimentar dados de um lugar para outro, por exemplo, entre registradores, memória e dispositivos de E/S.
- Operações Aritméticas e Lógicas: Realizar cálculos matemáticos ou operações lógicas, como adição, subtração, multiplicação, divisão, E lógico, OU lógico etc.
- Controle de Fluxo: Decidir o próximo passo do programa baseado em condições, como desvios condicionais (se-então-senão) e loops.
- Manipulação de Dados: Realizar operações de manipulação de bits, como deslocamento (shift), rotação etc.
- Controle de Interrupção: Gerenciar interrupções de hardware ou software.
- Operações de Entrada e Saída (E/S): Facilitar a comunicação com dispositivos periféricos, como leitura de teclado, gravação em disco etc.
- Operações de Transferência de Controle: Alterar o fluxo de execução do programa, como chamadas de sub-rotina e retornos de chamada.
- Operações de Controle de Sistema: Executar operações especiais de nível de sistema, como troca de contexto entre processos.

Cada instrução é composta por um código de operação (opcode) que indica a operação a ser realizada, bem como operandos, que podem ser valores, endereços de memória ou



registradores que participam da operação. Por exemplo, em uma instrução de adição de dois números, o opcode indicará que é uma operação de adição e os operandos serão os números a serem somados.

Assim, uma instrução é formada basicamente por dois campos:

- Código de operação (Opcode): um subgrupo de bits que identifica a operação a ser realizada pelo processador. É o campo da instrução cujo valor binário identifica a operação a ser realizada, conforme exemplo da abaixo. Esse valor é a entrada no Decodificador de Instruções na Unidade de Controle. Cada instrução deverá ter um código único que a identifique.
- Operando: um subgrupo de bits que identifica o endereço de memória onde está contido o dado que será manipulado, ou pode conter o endereço onde o resultado da operação será armazenado.

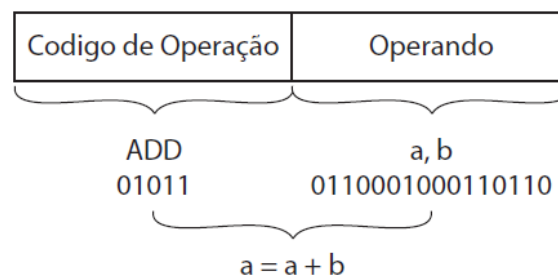


Figura 6 – Campos de uma instrução

Em resumo, uma instrução em arquitetura de computadores é um comando em linguagem de máquina que a CPU entende e executa para realizar uma operação específica. Cada instrução desempenha um papel crucial na execução de programas e no funcionamento geral do computador.

### 1.5.3 Unidade Lógica e Aritmética (ULA)

É o núcleo da CPU onde ocorrem operações matemáticas e lógicas. É aqui que adições, subtrações, multiplicações, divisões e operações lógicas como AND, OR e NOT são realizadas.

A **Unidade Lógica e Aritmética (ULA)** e a **Unidade de Controle (UC)** formam o núcleo funcional da CPU. A ULA é o subsistema responsável por realizar as operações de processamento: cálculos aritméticos (soma, subtração, etc.) e operações lógicas bit a bit. Já a UC é o subsistema que *dirige* todo o processo, coordenando as etapas de busca, decodificação e execução das instruções, emitindo sinais de controle apropriados para a ULA, registradores, memória e interface de E/S.

Em termos de implementação, a ULA é um circuito combinacional que recebe entradas (os operandos e possivelmente sinais de controle) e produz saídas (o resultado da operação e flags de status). Ela contém componentes como somadores (para adição/subtração), circuitos lógicos (AND, OR, etc.), e circuitos de deslocamento. A ULA típica pode realizar diversas operações dependendo do comando enviado pela UC. Por exemplo, quando a UC sinaliza uma operação de soma, a ULA conecta os operandos ao circuito somador interno; se sinaliza uma operação XOR, direciona os operandos para a rede lógica XOR, e assim por diante. Muitos processadores também incluem dentro da ULA comparadores para auxiliar em operações de comparação, e podem ter uma estrutura de cálculo pipeline para acelerar operações de múltiplos passos. Independentemente dos detalhes, do ponto de vista lógico podemos imaginar a ULA como uma “grande calculadora” que, a cada pulso de clock, realiza a operação solicitada sobre os dados que estiverem presentes em seus registradores de entrada.

#### 1.5.4 Registradores

São pequenas áreas de armazenamento de alta velocidade localizadas dentro da CPU. Eles armazenam dados temporários e resultados de operações, permitindo um acesso rápido e eficiente. Os registradores são vitais para o funcionamento eficaz da CPU.

#### 1.5.5 Memória Cache

A cache é uma memória de alta velocidade e pequena capacidade localizada dentro da CPU ou entre a RAM e a CPU. Ela armazena cópias de dados frequentemente usados da RAM, reduzindo a necessidade de acessos à memória principal, o que melhora significativamente o desempenho.

#### 1.5.6 Pipeline

É uma técnica que permite que a CPU processe múltiplas instruções simultaneamente, dividindo a execução em etapas. Cada estágio executa uma parte da instrução, permitindo um aumento significativo na taxa de execução de instruções.

### 1.6 Evolução dos Processadores

A evolução dos processadores tem sido uma jornada notável, marcada por avanços significativos em miniaturização, eficiência energética, e aumento do desempenho. Esses avanços permitiram que os processadores modernos se tornassem mais poderosos e versáteis, com impacto direto na computação cotidiana e em aplicações avançadas. A seguir estão as principais fases dessa evolução:

- **Miniaturização dos Transistores:** Desde a introdução dos transistores nos processadores, houve uma redução contínua em seu tamanho, conforme previsto pela **Lei de Moore**. Atualmente, processadores de última geração utilizam **nós de fabricação de 3 nm** (nanômetros), permitindo a inclusão de bilhões de transistores em um único chip, aumentando consideravelmente o desempenho e a eficiência energética.
- **Arquiteturas Superescalares:** As arquiteturas superescalares foram introduzidas para permitir que múltiplas instruções fossem processadas simultaneamente por diferentes unidades de execução dentro de um único núcleo de CPU, melhorando a taxa de execução de instruções e a eficiência geral.
- **Processadores Multi-Núcleo:** Uma das maiores inovações recentes foi a introdução de **processadores multi-core**, que possuem vários núcleos de processamento dentro de um único chip. Isso permite a execução paralela de múltiplas tarefas, sendo ideal para ambientes multitarefa e para aplicações que requerem grande capacidade de processamento, como jogos, modelagem 3D e inteligência artificial. Atualmente, processadores com **16, 24 ou mais núcleos** são comuns em CPUs de alto desempenho, tanto em desktops quanto em servidores.
- **Hyper-Threading e SMT (Simultaneous Multi-Threading):** Essas tecnologias permitem que cada núcleo físico de um processador execute múltiplos threads de software simultaneamente. Isso melhora a eficiência ao aproveitar ao máximo os recursos de cada núcleo, resultando em melhor desempenho para tarefas paralelizadas, como virtualização e servidores. **Vide adendo 1 e adendo 2.**

- **Unidades de Processamento Gráfico Integradas (iGPU):** Muitos processadores modernos agora vêm com **GPUs integradas**, capazes de realizar tarefas gráficas sem a necessidade de uma placa gráfica dedicada. As iGPUs são usadas para tarefas comuns, como reprodução de vídeo e gráficos 3D leves, enquanto sistemas mais avançados utilizam GPUs dedicadas para cargas de trabalho pesadas.
- **Aceleração por Inteligência Artificial (AI):** Processadores modernos, como os da série **Apple M1/M2** e os **Intel Core com IA integrada**, incluem **motores neurais** dedicados ou **unidades de processamento tensorial (TPU)**, projetados para acelerar tarefas de aprendizado de máquina e IA. Esses aceleradores são usados em aplicações como reconhecimento de imagem, processamento de linguagem natural e outros tipos de inferência de IA.
- **Redução no Consumo de Energia:** O foco atual dos fabricantes está tanto no desempenho quanto na eficiência energética. Processadores para dispositivos móveis e ultrabooks, como os **ARM-based chips** (Apple M1/M2, Qualcomm Snapdragon), são projetados para proporcionar alto desempenho com baixo consumo de energia, prolongando a vida útil da bateria sem comprometer a velocidade.
- **Tecnologias de Empacotamento (Chiplets):** Em vez de construir grandes monólitos de silício, algumas arquiteturas recentes, como as da **AMD com a série Ryzen e Epyc**, adotaram o uso de **chiplets**. Esses pequenos componentes são montados juntos para criar um processador completo, oferecendo maior flexibilidade de fabricação e melhorando o rendimento da produção. Esse design modular permite adicionar mais núcleos de processamento sem aumentar significativamente o custo ou o tamanho do chip.
- **Processadores Quânticos:** Embora ainda em desenvolvimento, os **computadores quânticos** e seus **qubits** estão começando a influenciar a pesquisa de processadores. As grandes empresas de tecnologia, como **Google** e **IBM**, estão liderando essa inovação, que promete resolver problemas complexos de forma exponencialmente mais rápida do que os processadores tradicionais.

## 1.7 Arquiteturas Modernas e Tendências

**RISC-V** é um ISA de código aberto criado na UC Berkeley em 2010. Baseado nos princípios RISC, ele é *livre de royalties*, permitindo que qualquer pessoa projete ou venda chips compatíveis sem custos de licença. O RISC-V foi concebido para alto desempenho com baixo consumo de energia, usando um conjunto “limpo e modular” que suporta larguras de 32, 64 e 128 bits. Essa modularidade torna fácil adicionar extensões (por exemplo, ponto flutuante ou vetores). Em razão de ser aberto, o RISC-V ganhou amplo suporte de comunidades e fabricantes (há diversos protótipos e sistemas comerciais já anunciados), o que acelera inovações em hardware. Por exemplo, distribuições Linux modernas já passaram a incluir suporte ao RISC-V (Linux 5.17, 2022), reforçando sua relevância atual.

### - Arquiteturas Paralelas (Multi-core e Multithreading)

O aumento da frequência de clock chegou a limites práticos, então as arquiteturas modernas usam **múltiplos núcleos** de processamento no mesmo chip para paralelizar o trabalho. Um *processador multicore* possui dois ou mais “cérebros” independentes (núcleos) em um único chip. Cada núcleo pode executar threads distintas ao mesmo tempo: em vez de um processador do tipo “engenheiro solo”, temos agora uma equipe que resolve várias partes do problema em paralelo. Por exemplo, num sistema dual-core, podem rodar duas instruções simultaneamente em threads diferentes. Além disso, técnicas como *hyper-threading* ou *SMT (Simultaneous Multithreading)* usam múltiplas threads por núcleo para aproveitar buracos de latência. Analogamente, imagine uma fábrica onde, em vez de um único robô fazendo todo o

trabalho, há várias estações de trabalho ligadas em série (**pipelining**) e em paralelo (multicore), acelerando enormemente a produção. Essa evolução permite que as CPUs modernas alcancem desempenho muito maior em aplicações multitarefa e paralelas.

### - Computação Neuromórfica

A **computação neuromórfica** busca inspirar-se no cérebro humano: usa hardware e software que simulam redes neurais biológicas. Nessa abordagem, dispositivos eletrônicos representam “neurônios” e “sinapses” que disparam sinais (spikes) de forma semelhante às células nervosas reais. O objetivo é criar chips que executem processamento paralelo massivo e adaptativo, com altíssima eficiência energética. Por exemplo, o chip TrueNorth da IBM tem milhões de *neurônios digitais* que processam informação de modo assíncrono, ideal para reconhecimento de padrões e IA avançada. Embora ainda emergente, a técnica é apontada como tendência, pois pode impulsionar a computação de alto desempenho e a inteligência artificial do futuro. Em analogia simples, enquanto as arquiteturas tradicionais seguem instruções sequenciais, as neuromórficas funcionam como um córtex eletrônico, onde muitos “neurônios” emitem sinais simultaneamente, permitindo aprendizagem e processamento paralelo inspirados pelo cérebro.

A computação neuromórfica se baseia em alguns pilares-chave:

- **Neurônios e sinapses artificiais:** inspirados nas células do sistema nervoso, são usados para processar e transmitir sinais.
- **Eventos assíncronos (spikes):** os sinais são transmitidos por pulsos discretos (semelhantes aos impulsos elétricos no cérebro), e não por ciclos de clock fixos como em CPUs convencionais.
- **Plasticidade sináptica:** as conexões (sinapses) podem ser fortalecidas ou enfraquecidas ao longo do tempo, simulando o aprendizado biológico (como o que ocorre no cérebro).
- **Processamento paralelo distribuído:** cada “neurônio” pode operar de forma independente e simultânea com os demais.

Como se compara às Arquiteturas Convencionais

Característica	Computação Convencional	Computação Neuromórfica
Unidade básica de operação	Unidade lógica (bit/instrução)	Neurônio artificial / sinal (spike)
Modelo de funcionamento	Sequencial, síncrono	Paralelo, assíncrono
Processamento	Determinístico, baseado em instruções	Inspirado em biologia, adaptativo
Eficiência energética	Alta exigência de energia	Extremamente eficiente energeticamente
Aprendizado	Explícito (software)	Implícito (plasticidade sináptica no hardware)

### Exemplos Reais:

#### - IBM TrueNorth

- Possui 1 milhão de neurônios artificiais e 256 milhões de sinapses.
- Consome apenas alguns watts, mesmo executando tarefas complexas de reconhecimento visual.

#### - Intel Loihi

- Chip neuromórfico com aprendizado embarcado.
- Pode adaptar-se a novos dados em tempo real, sem reprogramação.

#### Aplicações da Computação Neuromórfica

- Reconhecimento de padrões: voz, imagem, texto (como o cérebro faz com estímulos sensoriais).
- Robótica autônoma: movimentos adaptativos e respostas a estímulos não programados.
- IA embarcada de ultrabaixa potência: sensores inteligentes, dispositivos IoT, próteses neurais.
- Sistemas cognitivos: simulações do comportamento neural e da aprendizagem humana.

#### - Tendência Futura: Computação Quântica

A computação quântica explora fenômenos da física quântica como a superposição e o emaranhamento, utilizando qubits em vez de bits convencionais.

#### Fundamentos

- **Qubit:** Pode representar simultaneamente estados 0 e 1.
- **Emaranhamento:** Estados interligados de qubits que aumentam exponencialmente a capacidade computacional.

#### Arquitetura Básica

- Qubits
- Registradores quânticos
- Portas lógicas quânticas
- Sistemas avançados de controle e medição

#### Aplicações

- Criptografia avançada
- Simulação molecular precisa
- Inteligência artificial otimizada
- Problemas complexos de otimização

#### Desafios

- Decoerência (instabilidade dos estados quânticos)
- Correção de erros quânticos
- Escalabilidade

Característica	Computação Clássica	Computação Quântica
Unidade Básica	Bit (0 ou 1)	Qubit (superposição)
Processamento	Sequencial ou limitado	Paralelismo massivo
Desempenho	Limitado para problemas exponenciais	Alta eficiência em problemas específicos

## Perspectivas Futuras

Grandes empresas e centros de pesquisa têm avançado rapidamente na computação quântica, com perspectivas promissoras para aplicações práticas híbridas (quânticas-clássicas).

### - Quadro Comparativo Adicional

Além das tabelas anteriores, podemos contrastar características gerais das arquiteturas clássicas e das modernas:

Aspecto	Arquiteturas Clássicas	Arquiteturas Modernas
<b>Paradigma de Execução</b>	Sequencial (instrução por instrução)	Exploração de paralelismo: múltiplos cores e <i>pipelining</i> .
<b>Memória e Cache</b>	Memória única; pouca hierarquia de cache	Hierarquia profunda de memória (vários níveis de cache) e memórias separadas (modo Harvard modificado).
<b>Complexidade do ISA</b>	CISC predominante (instruções complexas)	RISC e ISA especializados (instruções simples e extensíveis)
<b>Componentes de CPU</b>	Um único núcleo forte; ULA e UC simples	Vários núcleos, ULA múltiplas instâncias, UC sofisticada (previsão de desvios, multithreading).
<b>Pipelines e Paralelismo</b>	Pipelines simples ou ausentes	Pipelines longos, execução superscalar, vetorização e multithreading.
<b>Frequência de Clock</b>	Buscando aumentar clocks absolutos	Teto de frequência, foco em mais núcleos e eficiência energética.
<b>Exemplos e Uso</b>	PCs antigos (x86), mainframes, microcontroladores sem cache.	Smartphones (ARM/RISC), servidores modernos, chips neurais e neuromórficos.

## 1.8 Avaliação de Desempenho

A medição do desempenho de uma CPU é um campo complexo, envolvendo várias métricas e técnicas de avaliação:

- IPC (Instruções Por Ciclo): Indica quantas instruções são executadas em média a cada ciclo de clock. É uma métrica chave para avaliar a eficiência de uma CPU.

- Frequência de Clock: Refere-se à taxa de operação da CPU, medida em ciclos por segundo. Uma frequência de clock mais alta geralmente resulta em um desempenho melhor.

- Cache Hits e Misses: Avaliam a eficiência do uso da memória cache. Hits ocorrem quando os dados estão presentes na cache e Misses quando estão requerendo uma busca na memória principal. "Misses" (falhas ou faltas), em um contexto relacionado à memória de computador, refere-se a situações em que a CPU (Unidade Central de Processamento) ou um componente de cache não encontra os dados ou instruções necessárias na memória cache, levando a uma busca na memória principal (RAM) ou em uma camada mais lenta de memória.

Essa operação de busca adicional em uma memória mais lenta resulta em um tempo de acesso maior em comparação com quando os dados são encontrados na memória cache. Portanto, minimizar o número de "misses" de cache é crucial para otimizar o desempenho de um

sistema computacional. Estratégias como a hierarquia de memória e algoritmos de gerenciamento de cache são utilizados para tentar reduzir a ocorrência de "misses"

- Benchmarking: Consiste em submeter a CPU a uma série de testes padronizados que simulam cargas de trabalho específicas. Os resultados do benchmarking são utilizados para comparar o desempenho relativo de diferentes CPUs.

### 1.8.1 Fatores de Desempenho

Os fatores que afetam o desempenho dos computadores são:

Algoritmo do programa; Linguagem de programação; Compilador; Conjunto de instruções do processador; Organização e tecnologia do processador; Clock do processador; Tamanho e frequência da memória cache; Frequência, latência e canais da memória RAM; Frequência do barramento do sistema; Taxa de transferência de dados dos discos. O desempenho de um computador pode variar significativamente de acordo com a combinação desses fatores.

A busca por aumento no desempenho dos processadores é constante. As principais evoluções nos processadores para obter um maior desempenho são através das seguintes técnicas:

- Aumento do clock (overclocking); Aumento no número de bits da CPU; Aumento na capacidade de endereçamento; Utilização de memória cache; Utilização de pipelines; Utilização de arquitetura vetorial - Hoje, a maioria das CPUs implementam arquiteturas que apresentam instruções para uma forma de processamento vetorial em vários conjuntos de dados, normalmente conhecidos como SIMD (Single Instruction Multiple Data);
- Utilização de arquitetura VLIW (Very Long Instruction Word) – tira proveito do paralelismo em nível de instrução, pois executa um grupo de instruções ao mesmo tempo;
- Utilização de Multithreading Simultâneo (SMT) – os bancos de registradores são replicados para que várias instruções possam compartilhar os recursos dos pipelines;
- Utilização de multicore – é a combinação de dois ou mais processadores num único chip; Incorporação da Unidade de Processamento Gráfico (GPU) na CPU.

### 1.9 Unidade Lógica Aritmética (ULA)

A ULA é o componente responsável por realizar operações aritméticas (como adição, subtração, multiplicação e divisão) e operações lógicas (como AND, OR, NOT) nos dados fornecidos pela Unidade de Controle.

### 1.10 O Modelo de Barramento do Sistema

Os barramentos do sistema são como as artérias e veias de um sistema computacional. O barramento de dados transporta os bits de informação entre os componentes, o barramento de endereço indica a localização específica dos dados na memória e o barramento de controle emite sinais para coordenar as operações complexas. Essa rede de comunicação é essencial para o funcionamento coordenado e eficaz do sistema.



Figura 7 - Slot barramento PCI Express – Conexão serial

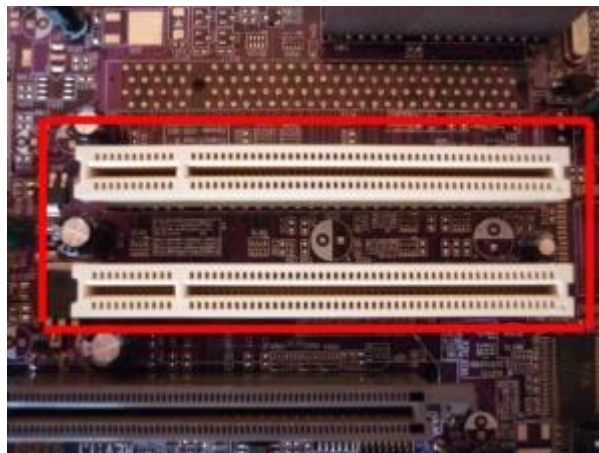


Figura 8 - Slot barramento PCI (Peripheral Component Interconnect) – conexão paralela

#### 1.10.1 Tipos de Barramentos

- Barramento de Dados: O barramento de dados é um conjunto de trilhas condutoras que permite a transferência de dados entre os componentes do computador, como a CPU, a memória e os dispositivos de entrada e saída. É bidirecional, o que significa que pode transportar dados em ambas as direções. A largura do barramento de dados, medida em bits, determina quantos bits de dados podem ser transmitidos simultaneamente. Por exemplo, um barramento de dados de 32 bits pode transmitir 32 bits de dados de uma só vez.
- Barramento de Endereços: O barramento de endereços é usado para selecionar uma localização específica na memória. Ele é utilizado pela CPU para indicar o endereço de memória a ser lido ou escrito. A largura do barramento de endereços determina o número máximo de endereços únicos que o sistema pode acessar. Por exemplo, um barramento de endereços de 16 bits pode endereçar até 65.536 locais de memória diferentes.
- Barramento de Controle: O barramento de controle transporta sinais de controle que coordenam as operações dos componentes do sistema, como a ativação de leitura ou escrita na memória, a ativação de um dispositivo de E/S, entre outros. Esses sinais de controle são essenciais para sincronizar as operações dentro do sistema e garantir que todos os componentes funcionem em harmonia.
- Barramento de Sincronização (Clock Bus): Em muitos sistemas, um barramento de sincronização, conhecido como clock bus, é utilizado para distribuir o sinal de clock para todos



os componentes do sistema. Esse sinal de clock é crucial para sincronizar as operações e garantir que todos os componentes funcionem em uníssono.

- **Barramento de Controle de Cache:** Em sistemas com hierarquia de memória, há um barramento dedicado para controlar as operações de leitura e escrita na cache. Esse barramento ajuda a coordenar a transferência de dados entre a cache e a memória principal.

- **Barramento de Controle de E/S (I/O Bus):** Em sistemas com múltiplos dispositivos de entrada e saída, pode haver um barramento específico para controlar as operações de E/S. Esse barramento facilita a comunicação entre a CPU e os dispositivos periféricos.

- **Barramento de Controle de DMA (Direct Memory Access):** Em sistemas que utilizam DMA para transferência de dados entre a memória e os dispositivos periféricos sem intervenção da CPU, existe um barramento dedicado para coordenar as operações de DMA.

- **Barramento de Controle de Interrupção:** Em sistemas que suportam interrupções, há um barramento de controle de interrupção para sinalizar à CPU a ocorrência de eventos que requerem sua atenção imediata.

Esses barramentos são vias de comunicação essenciais que permitem a interação eficiente entre os componentes de um sistema computacional. Comercialmente, podemos citar:

- **ATA (Advanced Technology Attachment):**

Era: Finais dos anos 1980 - Meados dos anos 2000

Largura de Banda (ATA-33): Até 33 MB/s

Principais Componentes Conectados: Unidades de disco IDE (Integrated Drive Electronics).

- **PCI Express (PCIe):**

Era: Início dos anos 2000 - Presente (ainda predominante em 2021)

Largura de Banda: Até 32 GT/s (Gen 4.0)

Taxa de Transferência Máxima (Gen 4.0 x16): 64 GB/s

O **PCIe 5.0** dobra a largura de banda do PCIe 4.0, sendo essencial para placas de vídeo de última geração e SSDs de alta performance.

Principais Componentes Conectados: Placas de vídeo, SSDs NVMe, placas de rede, entre outros.

É o barramento predominante em placas-mãe modernas. Barramento serial trabalhando no modo full-duplex. Os dados são transmitidos nesse barramento através de dois pares de fios chamados pista. Cada pista permite obter taxa de transferência máxima de 250 MB/s em cada direção. Podemos encontrar sistemas PCI Express com 1, 2, 4, 8, 16 e 32 pistas.

Por exemplo, a taxa de transferência de um sistema PCI Express com 8 pistas (x8) é de 2 GB/s (250 \* 8). O barramento PCI Express é hot plug, ou seja, é possível instalar e remover placas PCI Express mesmo com o micro ligado.

- **USB (Universal Serial Bus):**

Era: Meados dos anos 1990 - Presente

Largura de Banda: Varia (USB 2.0, 3.0, 3.1, 3.2, 4.0)

Taxa de Transferência Máxima (USB 4.0): Até 40 Gbps

Principais Componentes Conectados: Uma ampla variedade de periféricos, como teclados, mouses, impressoras, drives externos etc. Conector USB-C é o padrão mais recente — a

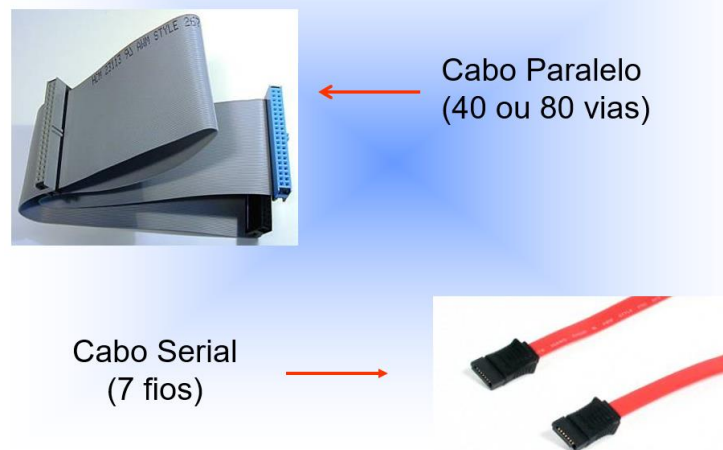
finalização de suas especificações aconteceu em agosto de 2014. O USB-C é compacto (tem 8,4 milímetros de largura por 2,6 milímetros de espessura) e reversível, ou seja, pode ser encaixado de qualquer lado.

- **SATA** (Serial Advanced Technology Attachment):

Era: Início dos anos 2000 - Presente

Largura de Banda (SATA 3.0): Até 6 Gbps

Principais Componentes Conectados: Unidades de armazenamento como HDDs e SSDs.



**Figura 9** - Cabos P(ATA) e SATA

- **NVMe** (Non-Volatile Memory Express):

Era: Meados dos anos 2010 - Presente

Largura de Banda (Gen 3.0): Até 32 Gbps

Largura de Banda (Gen 4.0): Até 64 Gbps

Principais Componentes Conectados: SSDs NVMe de alta performance.

### 1.11 Clock e Clock Cycle

O Clock é um sinal de pulso elétrico que regula a taxa de operação da CPU. Cada pulso de clock representa um ciclo de instrução. A Unidade de Contagem de Ciclos (Cycle Counter) mantém o controle do número de ciclos de clock necessários para executar uma instrução.

A frequência do clock é medida em Hertz (ciclos por segundo). Cada ciclo de clock representa uma unidade de tempo em que a CPU executa uma operação básica.

A execução de uma instrução pode levar vários ciclos de clock, dependendo da complexidade da instrução e da arquitetura da CPU.

Os processadores atuais operam com dois clocks, sendo um interno e um externo. O clock interno é sempre mais alto, e é usado para sincronizar as operações de processamento.

Quando falamos, por exemplo, sobre um “Intel i5 3 GHz”, estamos dizendo que o seu clock interno é de 3 GHz. O clock externo tem um valor menor, e é usado para sincronizar as operações de comunicação entre o processador, a memória e outros circuitos externos.

## 1.12 Conjuntos de Instruções (Instruction Set Architecture - ISA)

ISA é o conjunto de instruções, registradores, modos de endereçamento, tipos de dados, e convenções de chamada que definem a interface entre o hardware (CPU) e o software (compilador, sistema operacional, programas). Em outras palavras, a ISA é o “contrato” que define como o software se comunica com o processador. Ela especifica o que o processador pode fazer, mas não como essas funcionalidades são implementadas internamente — isso é papel da microarquitetura.

### Principais Elementos da ISA

Elemento	Descrição
<b>Conjunto de instruções</b>	Operações básicas que a CPU pode executar (ADD, SUB, LOAD, STORE, etc.).
<b>Registradores</b>	Pequenas áreas de armazenamento acessíveis diretamente pelas instruções.
<b>Modos de endereçamento</b>	Formas de localizar operandos (direto, indireto, indexado, imediato, etc.).
<b>Tipos de dados</b>	Tamanhos e formatos dos dados manipulados (inteiros, ponto flutuante, vetores).
<b>Formato das instruções</b>	Organização binária das instruções (campos de operação, registradores, constantes).
<b>Tratamento de interrupções</b>	Como a CPU responde a eventos externos.
<b>Instruções privilegiadas</b>	Instruções acessíveis apenas pelo sistema operacional.

### Exemplo simples de uma instrução ISA (RISC)

assembly

ADD R1, R2, R3

Essa instrução soma os valores contidos nos registradores R2 e R3, e armazena o resultado em R1. Essa instrução é definida na **ISA**. O compilador traduz o código-fonte para essa forma, e o processador a executa segundo sua microarquitetura.

### ISA vs Microarquitetura

Aspecto	ISA	Microarquitetura
O que é?	Interface visível ao programador	Implementação física da ISA
Visível ao software	Sim	Não
Exemplo	x86, ARM, RISC-V	Intel Core i7, AMD Ryzen, Apple M1

**Nota:** Vários processadores diferentes (ex: Intel e AMD) podem implementar a mesma ISA (x86), mas com microarquiteturas distintas.

### Tipos comuns de ISA

ISA	Características	Aplicações típicas
<b>x86</b>	CISC, muito difundida, complexa	PCs, notebooks, servidores
<b>ARM</b>	RISC, baixo consumo, eficiente	Smartphones, tablets, embarcados
<b>RISC-V</b>	RISC, aberta e extensível	Educação, pesquisa, IoT, servidores

ISA	Características	Aplicações típicas
<b>MIPS<sup>1</sup></b>	RISC, simples, muito usada academicamente	Embarcados, roteadores

A ISA é fundamental para garantir a compatibilidade entre software e hardware. Ela define como o processador “compreende” as instruções e como o compilador e o sistema operacional devem gerar código. Conhecer uma ISA permite programar em linguagem de montagem (assembly), desenvolver compiladores, criar emuladores ou projetar microprocessadores.

Quando um processador é desenvolvido, ele disponibiliza o conjunto de instruções (linguagem) que pode ser usado pelos programadores para escrever os programas. O projeto de um processador é centrado no conjunto de instruções de máquina que se deseja que ele execute, ou seja, do conjunto de operações primitivas que ele poderá executar. Quanto menor e mais simples for o conjunto de instruções, mais rápido é o ciclo de tempo do processador. De acordo com o tipo de instrução, um processador pode ser CISC ou RISC ou Híbrido.

Os computadores clássicos podem ter conjuntos complexos ou simples:

- **CISC (Complex Instruction Set Computer):** a CPU possui um grande número de instruções poderosas e complexas. Uma única instrução CISC pode executar várias operações de alto nível (como manipulação de memória ou de strings), reduzindo o número total de instruções de um programa. Isso facilita a geração de código de alto nível, mas torna o hardware mais complicado e lento por instrução (cada instrução pode ocupar vários ciclos de clock). Exemplos de arquiteturas CISC clássicas são o x86 (Intel/AMD), muito usado em desktops e servidores.
- **RISC (Reduced Instruction Set Computer):** segue a filosofia de instruções simples e executadas em um único ciclo de clock. A ideia é ter um conjunto pequeno e uniforme de instruções que executem operações básicas rapidamente. Em vez de operações complexas, tarefas mais sofisticadas são feitas por sequências simples de instruções. O design mais simples do hardware permite otimizar *pipelining* (execução sobreposta de instruções) e paralelismo interno. Exemplos notáveis de arquiteturas RISC incluem ARM - Advanced RISC Machine (presente em quase todos os smartphones) e RISC-V (criada na Universidade da Califórnia, em Berkeley, em 2010, com um diferencial fundamental: é totalmente aberta, livre de royalties e extensível). Na prática, RISC tende a oferecer maior eficiência energética e desempenho por watt, importante em dispositivos móveis.

---

<sup>1</sup> MIPS - Microprocessor without Interlocked Pipeline Stages, ou Microprocessador sem Estágios de Pipeline Intertravados. A arquitetura MIPS é um exemplo clássico de design RISC, reconhecida por sua clareza, uniformidade e eficiência. Embora tenha perdido espaço para ARM e RISC-V no mercado comercial, permanece relevante no ensino e na formação de profissionais da computação.

### Quadro Comparativo: CISC vs RISC

Característica	CISC (Complexo)	RISC (Reduzido)
Número de instruções	Muitas instruções complexas (variável em tamanho)	Conjunto pequeno de instruções simples e de formato fixo.
Ciclos por instrução	Várias instruções podem levar muitos ciclos cada (instrução complexa).	Instruções simples projetadas para 1 ciclo (pipeline eficiente).
Complexidade de hardware	Mais complexo: decodificação e execução das instruções é sofisticada.	Mais simples: implementação de cada instrução é direta, facilitando otimizações.
Pipelining e paralelismo	Menor eficiência no pipeline por instruções variadas; costuma exigir microcódigo.	Muito eficiente em <i>pipelining</i> : várias instruções podem ser processadas simultaneamente.
Registradores	Menos registradores internos (dependem de memória para muitos dados).	Em geral, muitos registradores de uso geral (reduz acesso à memória).
Tamanho do código	Código mais compacto (uma instrução CISC faz mais), mas mais complicado.	Código tende a ser mais longo (mais instruções simples), porém previsível.
Uso típico	Computadores pessoais, PCs e servidores legados (ex.: x86).	Dispositivos embarcados e móveis, servidores modernos, pois oferecem melhor desempenho por energia.

- Híbrido – Esta é uma abordagem mais moderna que busca aproveitar o melhor dos dois mundos. Ela pode envolver a implementação de uma arquitetura RISC com uma camada de tradução para executar instruções CISC legadas, ou pode combinar características de ambas as arquiteturas de forma a otimizar o desempenho para tarefas específicas.

Apesar dos fabricantes ainda venderem seus chips como sendo processadores RISC ou CISC, não existe praticamente nenhum processador atualmente que siga estritamente uma das duas filosofias, combinando características das duas arquiteturas, por questões de desempenho.

Originalmente todos os processadores Intel e AMD para os PCs eram puramente CISC enquanto os processadores dos Macintosh e vídeo games eram RISC.

## 2. Memória

### 2.1 Tipos de Memória

O ecossistema de memória em um sistema computacional é vasta e variada, composta por diversos tipos:

- RAM (Memória de Acesso Aleatório): É a arena onde dados são temporariamente armazenados para acesso rápido e eficiente durante a execução de programas. Ela é volátil, o que significa que os dados são perdidos quando o sistema é desligado.

- ROM (Memória Somente de Leitura): Esta é a biblioteca imutável de dados, onde informações essenciais, como o firmware do sistema, residem permanentemente. Ela é não volátil e contém instruções essenciais para inicialização e operação do sistema. Os Tipos:

- ✓ **Mask ROM:** gravada durante fabricação.
- ✓ **PROM:** programável uma única vez.
- ✓ **EPROM:** apagável por luz ultravioleta.
- ✓ **EEPROM:** apagável e programável eletricamente (ex: memória BIOS).

- Cache: Uma forma de memória de alta velocidade, localizada entre a RAM e a CPU, que guarda dados frequentemente acessados para otimizar o desempenho. Existem diferentes níveis de cache (L1, L2 etc.) com velocidades e capacidades variadas.

- Memória Virtual: Esta é uma extensão do espaço de endereçamento da memória RAM através do uso de armazenamento secundário (como discos rígidos). Permite que sistemas operacionais gerenciem mais dados do que poderiam caber na RAM física.

- Memória de Vídeo: Esta é uma área de memória dedicada a armazenar dados gráficos, como texturas e frames, que são utilizados pelo processador gráfico (GPU) para renderizar imagens em um monitor.

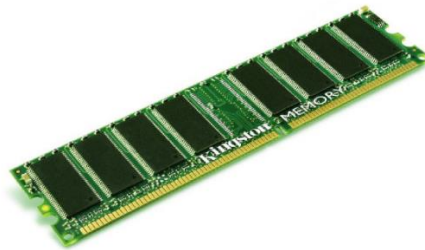


Figura 10 – Memória RAM

### 2.2 Características das Memórias

A memória desempenha um papel fundamental em qualquer sistema computacional, armazenando e fornecendo acesso rápido aos dados necessários para o processamento. A seguir, estão descritos os principais tipos de memória usados atualmente, bem como suas características e inovações recentes:

#### 2.2.1 RAM (Random Access Memory)

A **RAM** é um tipo de memória volátil, ou seja, os dados nela contidos são perdidos quando o sistema é desligado. A RAM é responsável pelo armazenamento temporário de dados e instruções durante a execução de programas. Existem vários tipos de RAM, cada uma com características específicas:

**- DRAM (Dynamic RAM):**

É o tipo mais comum de RAM usada em computadores e dispositivos móveis. Ela armazena dados em capacitores que precisam ser periodicamente recarregados para manter as informações. A **DDR4** é amplamente usada em desktops e laptops, enquanto a **DDR5**, mais recente, oferece velocidades mais rápidas, maior largura de banda e melhor eficiência energética. As memórias DDR5 atingem velocidades de até **6400 MT/s<sup>2</sup>** e são otimizadas para multitarefa e aplicativos intensivos em dados, como jogos e modelagem 3D.

**- LPDDR (Low Power DDR):** Uma variante de DRAM otimizada para dispositivos móveis, como smartphones e tablets. A mais recente **LPDDR5** combina alta eficiência energética com desempenho superior, permitindo maior duração da bateria e suporte para aplicativos exigentes, como streaming de vídeo em 4K, inteligência artificial e jogos.

**- SRAM (Static RAM):** A **SRAM** é uma forma de RAM não volátil que não precisa ser recarregada periodicamente. Ela é mais rápida e cara do que a DRAM, sendo usada principalmente em **memória cache** de alto desempenho, dentro de processadores.

### 2.2.2 Cache

A cache é uma memória de alta velocidade e baixa capacidade que armazena cópias dos dados mais frequentemente utilizados da RAM. Isso reduz o tempo de acesso à memória principal e melhora o desempenho do sistema. A memória cache (ou simplesmente "cache") é um tipo de memória localizada diretamente na CPU ou muito próxima dela. Sua função é armazenar temporariamente dados e instruções frequentemente utilizados para acelerar o acesso a essas informações pela CPU.

Vamos explorar detalhadamente os principais aspectos da memória cache:

- **Hierarquia de Memória:** A memória cache faz parte da hierarquia de memória de um sistema computacional. Essa hierarquia inclui diferentes níveis de memória, desde a memória cache até a memória RAM e, eventualmente, os dispositivos de armazenamento de longo prazo, como discos rígidos ou SSD.

- **Princípio da Localidade:** A eficiência da memória cache se baseia no princípio da localidade, que consiste em dois tipos:

- **Localidade Temporal:** Refere-se à tendência de um programa acessar as mesmas localizações de memória repetidamente em curtos intervalos de tempo. A memória cache explora isso armazenando dados recentemente utilizados.

---

<sup>2</sup> A unidade **MT/s (MegaTransfers por segundo)** é usada para medir a **taxa de transferência efetiva** da memória, e não a frequência do clock em si. **DDR (Double Data Rate):** 2 transferências por ciclo de clock (uma na borda de subida e outra na borda de descida). Por isso, a frequência efetiva em **MT/s é o dobro da frequência de clock em MHz**.

- Localidade Espacial: Refere-se à tendência de um programa acessar dados próximos em endereço de memória ao mesmo tempo. A memória cache também leva vantagem disso, pois armazena dados adjacentes em seus blocos.

- Blocos de Cache: A memória cache é organizada em blocos, onde cada bloco contém uma certa quantidade de dados. Quando a CPU acessa um dado da memória, ela na verdade traz um bloco inteiro para a memória cache. Isso é conhecido como princípio de transferência de bloco.

- Mapeamento de Cache: Existem diferentes métodos para mapear dados da memória principal para a memória cache. Alguns dos métodos mais comuns são:

- Mapeamento Direto: Cada bloco de memória tem uma localização específica na cache. Isso significa que uma localização específica da memória principal sempre mapeia para a mesma localização na cache.

- Associativo por Conjunto: Cada bloco de memória pode ser mapeado para um conjunto de locais na cache. Isso proporciona maior flexibilidade do que o mapeamento direto, pois um bloco de memória pode ser armazenado em qualquer local dentro do conjunto.

- Totalmente Associativo: Qualquer bloco de memória pode ser armazenado em qualquer local da cache. Isso oferece a máxima flexibilidade, mas é mais complexo e consome mais recursos.

- Políticas de Substituição: Quando a memória cache está cheia e precisa de espaço para um novo bloco, é necessário decidir qual bloco atual na cache será substituído. Existem diferentes políticas de substituição, como LRU (Least Recently Used), FIFO (First In, First Out) e aleatória.

- Níveis de Cache: Em sistemas modernos, é comum ter vários níveis de cache, como L1, L2 e, em alguns casos, até L3. Cada nível tem sua própria capacidade, latência e proximidade com a CPU. O L1 é o mais próximo da CPU e o L3 é o mais distante.

- Cache de Instruções e de Dados: Alguns processadores têm caches separadas para armazenar instruções e dados. Isso permite que a CPU busque instruções enquanto executa outras operações.

- Cache Unificada: Em contraste, uma cache unificada armazena tanto instruções quanto dados no mesmo espaço de cache.

- Tamanho, Associatividade e Latência: A eficácia da memória cache depende do tamanho (quantidade de dados que pode armazenar), da associatividade (como os dados são mapeados) e da latência (tempo de acesso) da cache.

A memória cache é um componente crucial para melhorar o desempenho do sistema, permitindo que a CPU acesse dados frequentemente utilizados de forma mais rápida do que se tivesse que buscá-los diretamente na memória principal (RAM). Sua eficiência depende da forma como é configurada e do padrão de acesso dos programas em execução.

### 2.2.3 ROM (Read-Only Memory)

A ROM é uma forma de memória não volátil que contém dados permanentes. O conteúdo é pré-gravado durante a fabricação e não pode ser alterado pelo usuário.

### 2.2.4 EPROM (Erasable Programmable Read-Only Memory)

A EPROM é uma forma de memória ROM que pode ser apagada e reprogramada utilizando radiação ultravioleta. É reutilizável, mas o processo é relativamente complexo e requer hardware especializado.



### 2.2.5 EEPROM (Electrically Erasable Programmable Read-Only Memory)

A EEPROM é uma forma de memória ROM que pode ser apagada e reprogramada eletricamente, sem a necessidade de radiação ultravioleta. É utilizada em dispositivos como BIOS de computadores e em cartões de memória.

### 2.2.6 Flash Memory

A Flash Memory é uma forma de memória não volátil que é eletricamente programável e apagável em blocos. É amplamente utilizada em dispositivos como pen drives, cartões de memória e SSD.

SSD significa "Solid State Drive" em inglês, e em português é traduzido como "Unidade de Estado Sólido". É um tipo de dispositivo de armazenamento de dados que utiliza memória flash para armazenar informações de forma permanente. O SSD funciona de maneira completamente diferente do HD. Ele não tem partes móveis e usa memória flash NAND, um tipo de memória eletrônica que permite acesso instantâneo aos dados. Vamos entender os principais componentes e conceitos do SSD:

#### - Memória NAND e Blocos de Armazenamento

O SSD é composto por células de memória NAND, que armazenam os dados de forma elétrica. Os dados não são organizados em setores e trilhas, mas sim em páginas e blocos.

- Página: A menor unidade de armazenamento no SSD, geralmente de 4 KB a 16 KB.
- Bloco: Conjunto de páginas, normalmente com 128 a 256 páginas.

Diferentemente dos discos rígidos tradicionais (HDD), que armazenam dados em discos magnéticos que giram a altas velocidades e precisa procurar fisicamente os setores, o SSD acessa diretamente a página desejada, tornando a leitura e a escrita **muito mais rápidas** os SSD não possuem partes móveis. Em vez disso, eles armazenam dados em chips de memória flash NAND, que mantêm os dados mesmo quando a energia é desligada. Diferentemente do HD.

A memória NAND pode ser dividida em diferentes tipos, cada uma com suas características:

Tipo	Bits por Célula	Velocidade	Durabilidade	Preço
SLC (célula de nível único)	1 pedaço	Muito rápido	Alta durabilidade	Muito cara
MLC (célula multinível)	2 bits	Rápida	Boa durabilidade	Preço intermediário
TLC (célula de nível triplo)	3 pedaços	Mais lento que MLC	Menor durabilidade	Mais barato
QLC (célula de quatro níveis)	4 bits	Lenta	Menor durabilidade	Mais barato ainda

Os SSDs mais baratos geralmente usam TLC ou QLC, enquanto os de alto desempenho usam MLC ou SLC.

#### - O Processo de Escrita e Exclusão de Dados no SSD

Os SSDs funcionam de maneira diferente do HD ao gravar e apagar arquivos. Para gravar um novo arquivo, o SSD escreve dados em células de memória NAND disponíveis. Porém, para apagar um arquivo, ele não pode apenas remover os dados individualmente. O SSD só consegue apagar blocos inteiros de dados. Isso significa que, mesmo que apenas uma página dentro de um bloco preciso seja alterada, todo o bloco deve ser desligado e reescrito, o que pode impactar o desempenho ao longo do tempo.

- Coleta de lixo e TRIM

Para evitar lentidão devido à reescrita de blocos inteiros, os SSDs utilizam duas técnicas importantes:

- Coleta de Lixo: O próprio SSD organiza os dados para otimizar o espaço livre e reduzir a necessidade de desligamento frequente.
- TRIM: Um comando do sistema operacional que ajuda o SSD a identificar e excluir blocos de dados não utilizados, melhorando a eficiência e prolongando a vida útil da unidade.

- Desgaste da Memória NAND e Nivelamento de Desgaste (Wear Leveling)

Cada célula NAND possui um número limitado de ciclos de escrita e desligamento.

Para evitar que algumas células se desgastem mais rapidamente do que outras, o SSD usa um sistema chamado wear leveling, que distribui os dados de forma equilibrada pelo dispositivo.

- Resumidamente, as principais características dos SSD são:

- Velocidade de Leitura e Escrita: Os SSD são notavelmente mais rápidos em comparação com os HDD. Eles oferecem tempos de acesso muito menores, o que significa que os dados podem ser lidos e gravados mais rapidamente.
- Durabilidade: Como não possuem partes móveis, os SSD são mais resistentes a choques físicos e vibrações do que os HDD. Isso os torna uma escolha popular em laptops e dispositivos móveis.
- Eficiência Energética: Os SSD consomem menos energia do que os HDD durante a operação, o que pode resultar em uma vida útil de bateria mais longa em laptops e dispositivos móveis.
- Durabilidade e Confiabilidade: A ausência de partes móveis também significa que há menos desgaste mecânico, o que pode resultar em uma vida útil mais longa em comparação com os HDD.



Figura 11 - M.2 NVMe (Non-Volatile Memory Express) e SATA - (Serial Advanced Technology Attachment)

- Tamanho Compacto: Os SSD são mais compactos e leves do que os HDD tradicionais, o que os torna ideais para dispositivos ultraportáteis e computadores de formato pequeno.
- Silenciosos: Por não terem partes móveis, os SSD não produzem ruídos durante a operação.

- **Desempenho Consistente:** Ao contrário dos HDD, onde o desempenho pode variar dependendo da localização dos dados no disco, os SSD oferecem um desempenho mais consistente, independentemente da localização dos dados. Devido a essas vantagens, os SSD se tornaram a escolha preferida para armazenamento em muitos dispositivos modernos, como laptops, desktops, tablets e smartphones. Eles também são amplamente usados em servidores e sistemas de armazenamento empresarial, onde o desempenho e a confiabilidade são críticos.

Existem vários tipos de SSD, cada um com suas próprias características e finalidades. Aqui estão os tipos mais comuns:

- **SSD SATA (Serial ATA):** O SSD SATA é um dos tipos mais comuns e amplamente utilizados. Ele se conecta à placa-mãe através da interface SATA, que é uma interface de barramento serial. Os SSD SATA oferecem uma melhoria significativa no desempenho em comparação com os discos rígidos tradicionais (HDD) e são compatíveis com a maioria dos sistemas.

- **SSD NVMe (Non-Volatile Memory Express):** Os SSD NVMe utilizam a interface PCIe (Peripheral Component Interconnect Express) para comunicação com a placa-mãe. Esta interface oferece taxas de transferência de dados significativamente mais rápidas do que o SATA, tornando os SSD NVMe ideais para sistemas que exigem desempenho de armazenamento extremamente alto, como em ambientes profissionais ou de jogos de alto desempenho. A **tecnologia NVMe Gen 5**, destaca-se pela maior velocidade e eficiência em relação à geração anterior, sendo essencial para aplicações que requerem altíssima performance.

- **SSD M.2:** Os SSD M.2 são um formato específico de SSD que utiliza a interface M.2. Eles são compactos e se conectam diretamente à placa-mãe, sem a necessidade de cabos. Os SSD M.2 podem ser tanto SATA quanto NVMe, dependendo do tipo de interface suportada pela placa-mãe.

- **SSD PCIe de Placa de Expansão:** Esses SSD são conectados à placa-mãe através de slots PCIe. Eles podem ser utilizados para adicionar armazenamento de alta velocidade a sistemas que têm slots PCIe disponíveis.

- **SSD Externo:** Alguns SSD são projetados para serem utilizados como armazenamento externo. Eles geralmente se conectam ao computador através de portas USB ou Thunderbolt (é uma tecnologia de interface de alta velocidade que permite a conexão de uma ampla variedade de dispositivos a um computador. Foi desenvolvida pela Intel em colaboração com a Apple), proporcionando uma opção de armazenamento portátil de alta velocidade.

Cada tipo de SSD oferece vantagens específicas em termos de desempenho, tamanho e compatibilidade. A escolha do SSD dependerá das necessidades do usuário, do sistema em questão e do orçamento disponível.

#### 2.2.7 Registradores de Memória

Além da RAM e ROM, a CPU também possui registradores de memória que são usados para armazenar temporariamente dados durante as operações.

A CPU possui diversos tipos de registradores, incluindo registradores de dados, de endereço, de controle e registradores especiais. Cada um tem uma função específica, como armazenar dados temporários, endereços de memória ou sinalizar condições especiais. Alguns exemplos de registradores:

- **Registrador de Dados (DR):** O Registrador de Dados é usado para armazenar temporariamente dados que estão sendo processados pela CPU. É também conhecido como Accumulator em algumas arquiteturas. Ele é frequentemente utilizado para realizar operações aritméticas e lógicas.

- Registrador de Endereço (AR): O Registrador de Endereço mantém o endereço de memória que a CPU irá acessar na próxima operação. Ele é usado para indicar à memória onde buscar ou armazenar dados. É fundamental para operações de leitura e escrita na memória.

- Registrador de Instrução (IR): O Registrador de Instrução mantém a instrução atualmente sendo executada pela CPU. Ele é lido pela Unidade de Controle para determinar a próxima operação a ser realizada. Ele pode ser dividido em campos que representam diferentes partes da instrução, como código de operação e operandos.

- Registradores de Estado (Flag Register): Os Registradores de Estado são usados para armazenar informações sobre o estado atual da CPU e das operações.

Exemplos comuns incluem:

Zero Flag: Sinaliza se o resultado de uma operação é zero.

Negative Flag: Indica se o resultado de uma operação é negativo.

Carry Flag: Usado em operações de aritmética para indicar se houve um “vai-um” na operação.

- Contador de Programa (PC): O Contador de Programa mantém o endereço de memória da próxima instrução a ser executada. Após cada execução de instrução, o PC é incrementado para apontar para a próxima instrução. É crucial para a execução sequencial de instruções em um programa.

- Registradores de Uso Geral (GPRs): Em algumas arquiteturas, existem registradores de uso geral que podem ser usados para armazenar dados temporários durante a execução de um programa.

- Registradores de Ponto Flutuante (FPU): Em CPUs que suportam operações de ponto flutuante, existem registradores especializados para realizar cálculos com números de ponto flutuante.

Esses registradores são locais de armazenamento de alta velocidade dentro da CPU que desempenham papéis cruciais durante a execução de programas. Cada registrador tem uma função específica e é essencial para a eficiência e a precisão das operações realizadas pela CPU.

#### 2.2.8 Virtual Memory

Virtual Memory é um sistema que permite que um sistema operacional utilize parte do disco rígido como uma extensão da memória RAM. Isso permite que programas utilizem mais memória do que está fisicamente disponível.

#### 2.3 Hierarquia de Memória

Muitos sistemas utilizam uma hierarquia de memória, onde a cache é a mais rápida e a mais próxima da CPU, seguida pela RAM e, finalmente, pelos discos rígidos ou SSD, que são mais lentos, mas têm maior capacidade.

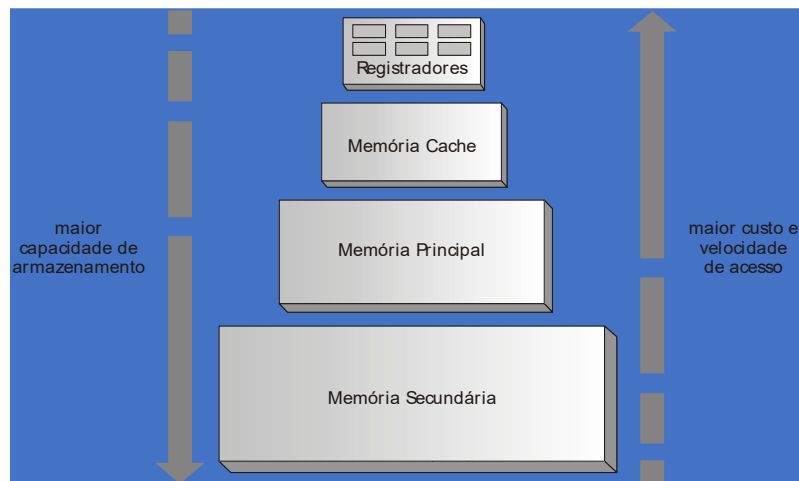


Figura 12 – Hierarquia de memórias

A hierarquia de memória é a magistral organização de camadas que compõem o espectro de memória de um sistema. Desde a memória cache ultrarrápida, que serve como um buffer de acesso imediato, até a memória secundária de alta capacidade, mas mais lenta, como os discos rígidos e SSD. A gestão eficaz dessa hierarquia é essencial para otimizar o desempenho do sistema.

#### 2.4 Encapsulamentos de Memória

A forma física dos módulos de memória influencia diretamente na capacidade e velocidade de operação. Os diferentes encapsulamentos, como DIMM (Dual Inline Memory Modules) e SIMM (Single Inline Memory Modules), Os módulos DIMM são capazes de transferir 64 bits de cada vez para o processador enquanto os SIMM transferem 32 bits.



Figura 13- Exemplo de modulo no formato SIMM de 72 vias

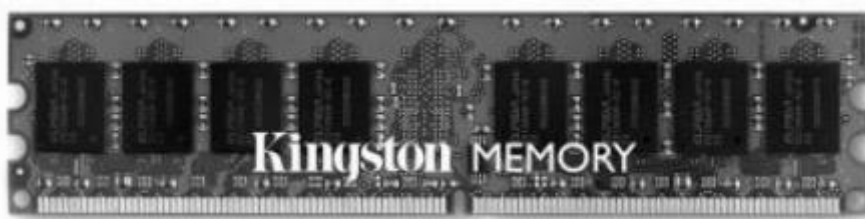


Figura 14- Exemplo de modulo no formato DIMM de 240 vias

### 3. Memória Secundária

#### 3.1 Discos Rígidos – (Hard Disk Drive – HDD)

Os discos rígidos são como os antiquários modernos de dados. Utilizam um conjunto de discos magnéticos rotativos para armazenar informações de forma persistente. Cada disco é revestido com uma camada de material magnético que pode ser magnetizada para representar dados binários. Uma cabeça de leitura/gravação lê e grava dados na superfície dos discos enquanto eles giram a altas velocidades. A capacidade de armazenamento de discos rígidos pode variar de gigabytes a vários terabytes, tornando-os a escolha predominante para armazenamento em massa em computadores de uso geral e servidores.

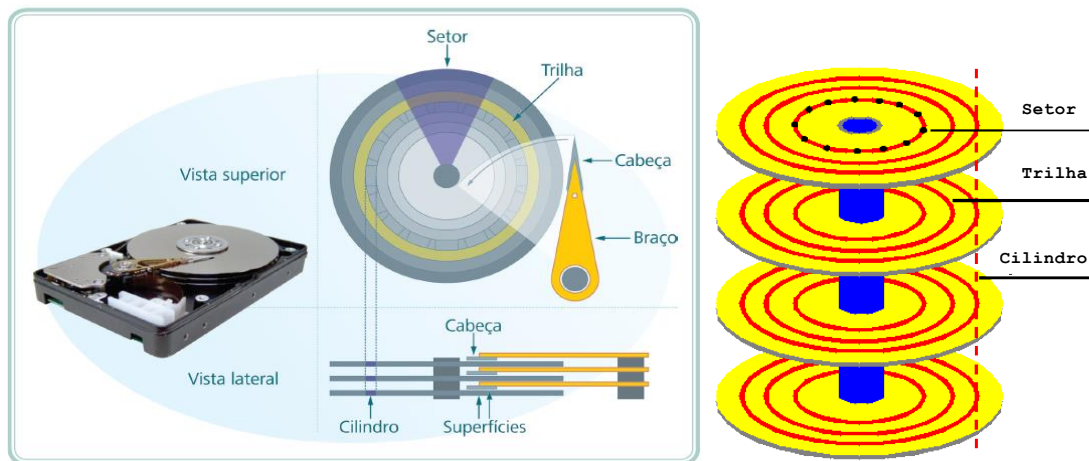


Figura 15 - HDD

Os HDD são organizados em:

- Trilha: é uma porção circular da superfície do disco.
- Setores: são divisões da trilha que contêm um número fixo de bytes, em que, para acessar os dados contidos na trilha, são indicados o número da superfície, o da trilha e o do setor em que os dados estão armazenados.
- Clusters: são números fixos de setores adjacentes tratados como uma unidade de armazenamento pelo sistema operacional.
- Cilindro: é a trilha em cada superfície que está sob a cabeça de leitura / gravação em determinada posição do braço de leitura / gravação.

#### 3.2 Discos Flexíveis (Disquetes)

Os disquetes, também conhecidos como discos flexíveis, foram uma das primeiras formas de mídia portátil para armazenamento de dados. Consistiam em discos magnéticos finos e flexíveis contidos em uma capa de plástico. Os disquetes foram populares nas décadas de 1980 e 1990, mas foram gradualmente substituídos por tecnologias de armazenamento mais avançadas devido à sua capacidade limitada e susceptibilidade a danos.

#### 3.3 Discos Óticos

Discos óticos, como CDs, DVDs e Blu-ray, representam a fronteira entre armazenamento físico e digital permanente. Eles funcionam gravando dados em uma camada reflexiva usando

um feixe de laser. CDs têm capacidade de armazenamento de cerca de 700 MB, DVDs podem armazenar até 9 GB e Blu-ray podem chegar a 50 GB. São frequentemente utilizados para distribuição de mídia, como filmes, jogos e software.

### 3.4 Drives de Estado Sólido (SSD)

Os SSD são a personificação da modernidade no armazenamento de dados. Utilizam memória flash para proporcionar velocidades de leitura e gravação que eclipsam os discos rígidos tradicionais. Uma das maiores vantagens dos SSD é a ausência de partes móveis, o que os torna mais resistentes a choques físicos e mais eficientes em termos de energia. São frequentemente utilizados em laptops e desktops para melhorar a responsividade e o desempenho do sistema.

### 3.5 Pen Drives

Os pendrives, também conhecidos como unidades flash USB, são dispositivos portáteis de armazenamento que utilizam memória flash para armazenar dados. São compactos, duráveis e oferecem uma maneira conveniente de transportar e transferir arquivos entre diferentes dispositivos. Pen drives estão disponíveis em uma ampla gama de capacidades, desde alguns gigabytes até vários terabytes.

### 3.6 Cartões de Memória

Os cartões de memória são pequenos dispositivos de armazenamento removíveis, frequentemente utilizados em câmeras, smartphones e outros dispositivos eletrônicos. Eles vêm em diferentes formatos, como SD (Secure Digital), microSD e CompactFlash, e oferecem uma maneira conveniente de expandir a capacidade de armazenamento desses dispositivos.



Figura 16 – Tipos de Memória Secundária



## 4. Placa-mãe

### 4.1 Componentes

A placa-mãe é o epicentro da interconexão de todos os componentes vitais de um sistema computacional. Ela possui uma complexa teia de elementos essenciais:

- Soquete do Processador: É a interface onde o processador é encaixado. Existem diferentes tipos de soquetes para acomodar diferentes famílias de processadores.
- Slots de Memória: São os encaixes onde os módulos de memória RAM são instalados. A placa-mãe suporta um certo tipo e quantidade de módulos de memória, o que influencia a capacidade total de RAM que o sistema pode ter.
- Conectores de Periféricos (USB, áudio etc.): Esses conectores permitem a conexão de periféricos externos como teclado, mouse, alto-falantes, entre outros. Cada placa-mãe tem uma variedade de portas USB e conectores de áudio para atender às necessidades do usuário.
- Chipsets: Os chipsets são conjuntos de chips que gerenciam a comunicação entre o processador, memória, periféricos e outros componentes do sistema. Existem dois tipos principais: o Northbridge, que lida com componentes de alta velocidade como a RAM e placa gráfica, e o Southbridge, que gerencia componentes de menor velocidade e funcionalidades de entrada/saída.
- BIOS (Sistema Básico de Entrada e Saída): A BIOS é um firmware embutido na placa-mãe que fornece as instruções básicas para inicializar o sistema. Ela é responsável por testar e inicializar os componentes essenciais durante o processo de inicialização.
- Slots de Expansão (PCIe etc.): Esses slots permitem a adição de placas de expansão, como placas gráficas, placas de som, placas de rede, entre outras. O padrão mais comum atualmente é o PCIe (Peripheral Component Interconnect Express).

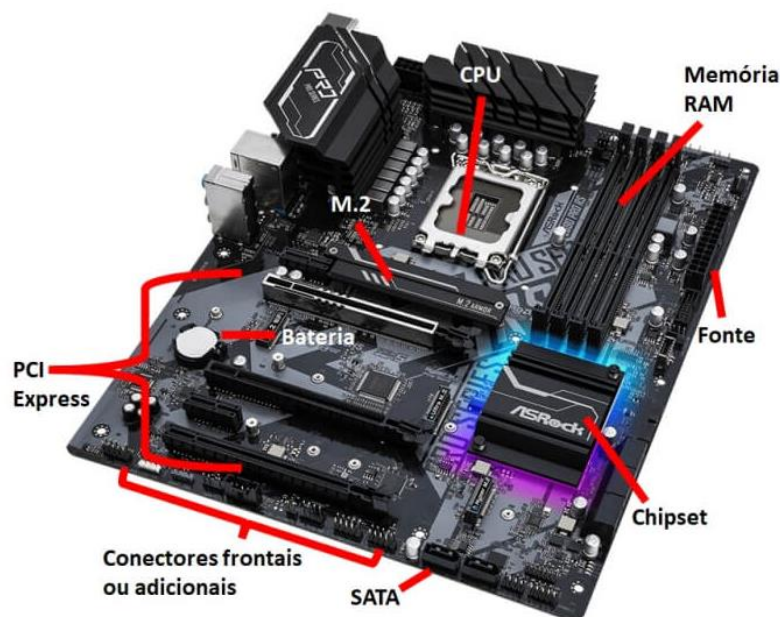


Figura 17 – Placa Mãe



## 4.2 Detalhes da Placa-mãe

A placa-mãe é meticulosamente projetada para suportar um tipo específico de processador, com slots de memória limitados por capacidade e velocidade. Ela também apresenta uma gama de conectores para periféricos, como USB, áudio, vídeo, entre outros. Além disso, a placa-mãe é equipada com slots de expansão, permitindo a adição de placas gráficas, placas de som e outros dispositivos para personalizar e aprimorar as capacidades do sistema.

Também conhecida como motherboard, mainboard ou placa principal, é uma das peças mais essenciais de um computador. Ela serve como a plataforma central que conecta todos os componentes e permite que eles funcionem em conjunto de maneira eficiente. Vamos explorar em detalhes os principais aspectos da placa-mãe.

### 4.2.1 Conexões e Slots

A placa-mãe possui uma variedade de conexões e slots onde os componentes essenciais e periféricos do computador são conectados. Isso inclui slots para a CPU, slots de memória RAM, slots de expansão, como o (Peripheral Component Interconnect Express - PCIe), conectores SATA (Serial Advanced Technology Attachment) para dispositivos de armazenamento, conectores USB, portas de áudio, entre outros.

### 4.2.2 Chipset

O chipset é um conjunto de chips localizados na placa-mãe que gerencia a comunicação entre a CPU, a memória, os dispositivos de armazenamento, as interfaces de E/S e outros componentes do sistema. Geralmente, é dividido em dois chips principais: Northbridge e Southbridge.

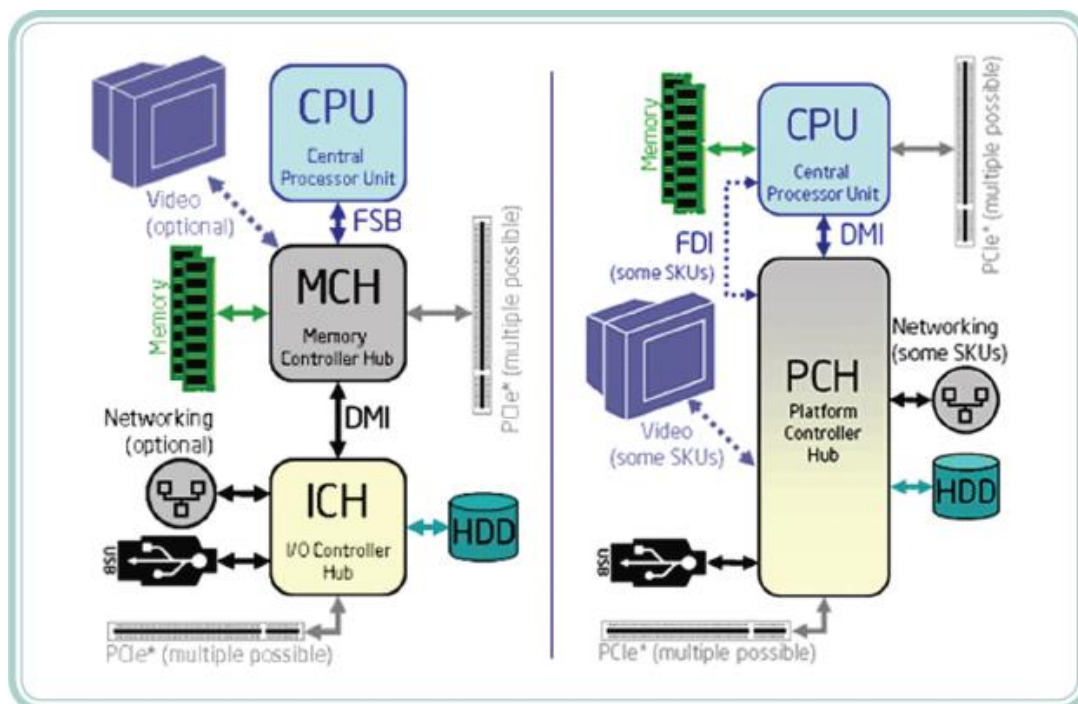


Figura 18 - Chipset ponte norte e sul (à esquerda) e único (à direita)

O chipset desempenha um papel fundamental na comunicação e na coordenação entre os diferentes componentes do sistema. Vamos detalhar mais sobre o chipset:

- **Função Principal:** O chipset age como uma ponte de comunicação entre a CPU, a memória, os slots de expansão, os dispositivos de E/S (entrada e saída) e outros componentes da placa-mãe. Ele gerencia o tráfego de dados e sinais entre esses elementos.

- Componentes do Chipset: O chipset consiste em dois principais chips:
  - Northbridge: Tradicionalmente, o Northbridge era responsável por gerenciar a comunicação de alta velocidade entre a CPU, a memória RAM e as placas de vídeo. No entanto, com o avanço da tecnologia, algumas das funções do Northbridge foram integradas diretamente no processador ou movidas para o Southbridge.
  - Southbridge: O Southbridge é encarregado de operações de E/S de menor velocidade, como a comunicação com dispositivos de armazenamento (como discos rígidos e SSDs), portas USB, interfaces de áudio, entre outros.
- Funções Específicas: O chipset executa várias funções específicas, como:
  - Controlar a interface entre a CPU e a RAM, determinando as capacidades de memória -suportadas e as velocidades de transferência.
  - Fornecer suporte para dispositivos de E/S, como portas USB, portas de áudio, interfaces de rede, entre outros.
  - Gerenciar os slots de expansão, como PCIe (Peripheral Component Interconnect Express), que são usados para conectar placas de vídeo, placas de som, e outros periféricos.
- Integração com Processadores Modernos: Em muitos sistemas modernos, algumas funções do chipset foram integradas diretamente no processador, especialmente as funções tradicionalmente associadas ao Northbridge. Isso ajuda a melhorar a eficiência e reduzir a latência nas comunicações.
- Arquiteturas de Chipset: Há diferentes arquiteturas de chipset, como o tradicional FSB (Front-Side Bus) e a arquitetura DMI (Direct Media Interface) usada pela Intel. Cada uma tem suas próprias especificidades em termos de comunicação entre a CPU e o chipset.
- Compatibilidade e Suporte: A escolha e a configuração do chipset são fatores críticos na determinação da compatibilidade de componentes em um sistema. Por exemplo, a escolha do chipset influencia quais tipos de processadores, memória e dispositivos de armazenamento podem ser utilizados.
- Evolução Tecnológica: O chipset é uma parte vital da placa-mãe e sua evolução ao longo dos anos tem sido fundamental para o aumento de desempenho e a introdução de novas tecnologias nos computadores.

Em resumo, o chipset é um componente crucial na arquitetura de um computador, facilitando a comunicação e a coordenação entre a CPU, a memória, os dispositivos de E/S e outros componentes da placa-mãe. Ele desempenha um papel fundamental no funcionamento eficiente do sistema como um todo.

#### 4.2.3 CPU Socket

O CPU socket é onde a CPU é instalada na placa-mãe. Cada socket é projetado para um tipo específico de processador, portanto, a compatibilidade entre a CPU e a placa-mãe é crucial.

#### 4.2.4 Slots de Memória RAM

A placa-mãe possui slots de memória onde os módulos de RAM são instalados. Eles são projetados para suportar diferentes tipos e velocidades de memória.

#### 4.2.5 Slots de Expansão

Os slots de expansão são conectores físicos presentes na placa-mãe de um computador que permitem a adição de componentes e funcionalidades ao sistema. Eles funcionam como portas de encaixe onde diferentes tipos de placas adicionais podem ser instaladas, expandindo as capacidades do computador. Esses slots são fundamentais para a flexibilidade e atualização dos sistemas de hardware. Esses slots permitem a adição de placas de expansão, como placas de vídeo, placas de som, placas de rede, entre outras. O padrão mais comum é o PCIe (Peripheral Component Interconnect Express). Também, as GPUs modernas, como as arquiteturas NVIDIA Ada Lovelace e AMD RDNA3, são usadas não apenas para gráficos avançados, mas também para aplicações de aprendizado de máquina e inteligência artificial, oferecendo processamento paralelo de alto desempenho. Essas GPUs são essenciais para tarefas que exigem grande capacidade computacional, como modelagem 3D, simulações científicas, e análise de dados complexos. Esses slots são fundamentais para a flexibilidade e atualização dos sistemas de hardware. Abaixo, seguem mais detalhes:

### Principais Tipos de Slots de Expansão

1. **PCI (Peripheral Component Interconnect):** Esse é um padrão mais antigo, usado para conectar placas de rede, placas de som, e outros dispositivos periféricos ao computador. Embora ainda seja encontrado em alguns sistemas mais antigos, está sendo substituído por padrões mais rápidos e eficientes.
2. **PCI Express (PCIe):** Este é o tipo mais moderno e utilizado atualmente. Existem diferentes versões (como PCIe 3.0, PCIe 4.0, e PCIe 5.0), cada uma oferecendo uma largura de banda maior para suportar dispositivos mais rápidos. O PCIe é utilizado principalmente para:
  - **Placas de vídeo (GPU):** Placas que são responsáveis pelo processamento gráfico. É essencial para jogos, modelagem 3D, e, mais recentemente, para aprendizado de máquina.
  - **Placas de rede:** Placas que adicionam ou melhoram a conectividade de rede, incluindo conexões Ethernet de alta velocidade.
  - **Placas de som:** Para proporcionar áudio de maior qualidade, especialmente em sistemas voltados para edição de áudio ou gamers.
  - **SSDs NVMe:** Unidades de armazenamento de estado sólido que se conectam via PCIe, proporcionando altas velocidades de leitura e escrita.
  - **AGP (Accelerated Graphics Port):** Padrão utilizado antigamente exclusivamente para placas de vídeo. Esse tipo de slot foi popular antes da introdução do PCIe, mas atualmente está obsoleto.

#### 4.2.6 Conectores de Armazenamento

A placa-mãe possui conectores para dispositivos de armazenamento, como discos rígidos (Hard Disk Drive - HDD) e unidades de estado sólido (SSD). Os tipos comuns incluem conectores SATA e, mais recentemente, M.2 para SSD e NVMe. SSD M.2 NVMe (Non-Volatile Memory Express) são uma evolução dos tradicionais SSD (Solid State Drives). Eles utilizam o protocolo NVMe para se comunicar com a CPU, o que proporciona uma significativa melhora no desempenho em comparação com os SSD SATA convencionais. A principal diferença entre os dois está na forma como se conectam ao computador. Os SSD NVMe utilizam slots M.2 (M.2 se refere a um formato SSD que se parece com um chiclete) ou placas de expansão PCIe (Peripheral Component Interconnect Express) conectada diretamente na placa, enquanto os SSD SATA se conectam através de cabos SATA.

Devido à sua alta taxa de transferência e baixa latência, os SSD NVMe são amplamente utilizados em sistemas onde o desempenho é crucial, como em laptops gamers, estações de trabalho e servidores. Enquanto o SATA I suportava até 150 MB/s de transferência de dados, o SATA III quadruplicou essa capacidade, chegando aos 600 MB/s. Já no caso da conexão PCI Express, a geração 2 já partia de 500 MB/s e atualmente, na geração 4, a taxa de transferência vai a 2000 MB/s – mais do que três vezes o máximo suportado pelo SATA III.

Por conta das suas características, há diferenças físicas entre os SSD SATA e NVMe. Os do tipo SATA são mais próximos dos HDD, ainda que mais finos. Já os SSD do tipo NVMe mais parecem com pentes de memória do que com unidades de armazenamento. Como eles são encaixados diretamente em um slot da placa-mãe, faz mais sentido que eles tenham esse formato.

#### 4.2.7 Conectores de E/S

Esses incluem portas USB, portas de áudio, portas Ethernet (para conexão de rede com fio), portas de vídeo, entre outros. Eles facilitam a conexão de periféricos e dispositivos externos. Os barramentos seguem padrões determinados por órgãos de normalização (IEEE). Os padrões surgiram para que os fabricantes de componentes não precisassem criar componentes para cada tipo de computador. Os barramentos de expansão de E/S apresentam diversos tipos com desempenhos variados. Dentre eles, podem-se citar os já ultrapassados: ISA, EISA, VLB e AGP; e os atuais PCI, PCI Express, SATA e USB-C.

#### 4.2.8 BIOS e UEFI

A BIOS (Basic Input/Output System) é um software de baixo nível armazenado em um chip na placa-mãe. Ela fornece instruções fundamentais para inicializar o computador e controlar hardware básico.

A Unified Extensible Firmware Interface (UEFI) é uma evolução significativa em relação à tradicional BIOS. A UEFI é parte integrante do firmware do sistema, responsável por fornecer uma interface entre o hardware do computador, o sistema operacional e os aplicativos. Ela é carregada durante o processo de inicialização (boot) e permanece residente na memória durante o funcionamento do sistema, fornecendo serviços e funcionalidades necessárias para operação do sistema.

Aqui estão algumas comparações entre UEFI e BIOS:

- Interface Gráfica e Recursos Avançados:

UEFI oferece uma interface gráfica de usuário, enquanto a BIOS normalmente oferece apenas uma interface de linha de comando ou uma interface de texto muito básica. UEFI suporta recursos avançados, como mouse e teclado USB, enquanto a BIOS pode ser limitada a dispositivos mais antigos.

- Capacidade de Armazenamento:

UEFI pode inicializar a partir de discos rígidos maiores que 2,2 TB e suporta partições GUID (Globally Unique Identifier) e o esquema de partição GPT (GUID Partition Table), enquanto a BIOS tem limitações em relação ao tamanho do disco rígido e está limitada ao esquema de partição MBR (Master Boot Record), com até 4 partições primárias, ou até 3 partições primárias e 1 partição estendida (com múltiplas partições lógicas dentro dela). O GPT, teoricamente até 128 partições primárias.

O MBR é um antigo esquema de partição usado pela BIOS para inicializar dispositivos de armazenamento, como discos rígidos. Ele usa uma tabela de partição MBR localizada no

primeiro setor (setor de inicialização) do disco. O MBR tem algumas limitações, como suporte apenas a discos com até 2,2 TB de capacidade e um máximo de quatro entradas de partição primária.

GPT é uma alternativa mais moderna ao MBR, usada pelo UEFI para inicializar dispositivos de armazenamento. Em vez de usar um setor de inicialização com uma tabela de partição como o MBR, o GPT armazena a tabela de partição em várias cópias ao longo do disco. GPT suporta discos com mais de 2,2 TB de capacidade e permite um número muito maior de partições (teoricamente, até 128 partições primárias). Em teoria, o GPT suporta discos com até 9,4 ZB (zettabytes) de capacidade. Isso é muitíssimo além da capacidade de armazenamento dos dispositivos atualmente disponíveis no mercado. Além disso, GPT oferece recursos de integridade de dados para proteção contra corrupção de dados.

GUID é um identificador único globalmente utilizado para identificar objetos, incluindo partições de disco no contexto do GPT. Um Identificador Único Universal (do inglês Universally Unique Identifier - UUID) é um número de 128 bits usado para identificar informações em sistemas de computação. O termo Identificador Único Global (Globally Unique Identifier - GUID) também é utilizado, sobretudo nos sistemas da Microsoft. Cada partição GPT é identificada por um GUID exclusivo. Esses identificadores únicos ajudam a evitar conflitos e garantem que cada partição seja única em um sistema.

Em resumo, enquanto o MBR é o antigo padrão usado pela BIOS com limitações significativas em termos de capacidade de armazenamento e número de partições, o GPT é uma alternativa mais moderna e flexível usada pelo UEFI. O GPT supera as limitações do MBR, suportando discos maiores, mais partições e oferecendo recursos avançados de integridade de dados.

- Segurança:

UEFI inclui recursos de segurança avançados, como Secure Boot, que ajuda a proteger o sistema contra a execução de código malicioso durante a inicialização. A BIOS não possui esses recursos de segurança.

- Suporte a Drivers e Protocolos:

UEFI suporta uma ampla gama de drivers e protocolos de rede, o que pode facilitar a inicialização remota e a manutenção do sistema. A BIOS pode ser limitada em termos de suporte a novos dispositivos e protocolos.

- Extensibilidade e Modularidade:

UEFI é mais modular e extensível do que a BIOS, permitindo que os fabricantes de hardware adicionem facilmente novos recursos e funcionalidades por meio de módulos de firmware. A BIOS é mais rígida e pode ser difícil de atualizar com novos recursos.

- Inicialização Mais Rápida:

UEFI geralmente oferece inicializações mais rápidas do sistema em comparação com a BIOS devido à sua arquitetura mais moderna e eficiente.

Em resumo, o UEFI representa uma evolução significativa em relação à BIOS, oferecendo uma interface mais amigável, suporte a hardware mais moderno, recursos de segurança aprimorados e maior flexibilidade para os fabricantes de hardware.

#### 4.2.9 Alimentação (Power Connectors)

A placa-mãe possui conectores de alimentação para receber energia do fornecimento de energia (fonte de alimentação) do computador. Geralmente, há um conector principal (ATX) e outros conectores para fornecer energia a componentes específicos, como a CPU.

#### 4.2.10 Área de Montagem e Layout

O layout da placa-mãe, incluindo a disposição dos componentes, conectores e slots, é projetado para acomodar diferentes configurações de gabinete e facilitar a instalação dos componentes.

#### 4.3 Padrões de Formatos de Placa-mãe

Os padrões de formato de placa-mãe, como ATX, MicroATX e Mini-ITX, determinam o tamanho e as características físicas. Esses padrões influenciam a compatibilidade com gabinetes e periféricos, garantindo um ambiente de sistema coeso. Por exemplo, as placas ATX são maiores e oferecem mais slots de expansão, enquanto as Mini-ITX são compactas e destinadas a sistemas compactos e eficientes.

#### 4.4 Sistema de Entrada e Saída (E/S)

O Sistema de Entrada e Saída (E/S) é crucial para a interação de um computador com o ambiente externo. Ele permite a comunicação com dispositivos periféricos, como teclado, mouse, disco rígido, monitores. Os Controladores de E/S são circuitos especializados ou componentes de hardware que gerenciam as operações de entrada e saída para dispositivos periféricos. Eles atuam como intermediários entre a CPU e os dispositivos de E/S. Dispositivos de E/S são mapeados em um espaço de endereços separado da memória principal. Isso significa que a CPU utiliza instruções específicas para acessar dispositivos de E/S, como leitura de dados de um teclado ou escrita em um disco rígido.

## 5. **ADENDO 1** - “threads” de processamento

Quando se adquire um computador comercial — seja um notebook, desktop ou estação de trabalho — é comum que nas especificações técnicas do processador apareça o número de “threads” de processamento. Esse número está diretamente ligado à capacidade de execução paralela do processador e influencia diretamente no desempenho em tarefas simultâneas. Mas, o que isso significa? Neste adendo vamos explicar

### 1. O que é uma *thread*?

Uma thread (em português, *linha de execução* ou *tarefa leve*) representa uma unidade mínima de execução de um processo. Um processo pode conter uma ou várias threads que compartilham os mesmos recursos (memória, arquivos abertos, etc.), mas são executadas de forma paralela ou concorrente.

No nível de hardware, uma thread pode ser vista como uma sequência independente de instruções que o processador é capaz de gerenciar simultaneamente com outras threads.

### 2. Relação entre núcleos físicos e threads

Um processador moderno pode ter múltiplos núcleos físicos (cores), e cada núcleo pode suportar uma ou mais threads por meio de uma tecnologia chamada Simultaneous Multithreading (SMT) — sendo a mais conhecida o Hyper-Threading da Intel.

Exemplo:

- Um processador com 4 núcleos físicos e Hyper-Threading ativado pode executar 8 threads simultaneamente (2 por núcleo).
- Já um processador com 6 núcleos e sem SMT terá apenas 6 threads.

### 3. Por que os fabricantes destacam o número de threads?

O número de threads dá uma estimativa da capacidade de paralelismo lógico do processador. Quanto mais threads simultâneas forem possíveis, maior a capacidade do processador de lidar com múltiplas tarefas ao mesmo tempo, como:

- Executar vários aplicativos sem travamentos.
- Rodar sistemas multitarefa (ex: abrir navegador, editores, planilhas e música simultaneamente).
- Executar aplicações multithreaded (ex: editores de vídeo, jogos modernos, softwares de engenharia, máquinas virtuais, entre outros).

### 4. Threads *lógicas* x *físicas*

É importante entender que o número de threads divulgado **não representa o número de núcleos físicos**, mas sim **o total de threads simultâneas que o processador pode gerenciar**.

**Resumo:**

Termo	Significado
Núcleo físico	Unidade real de processamento. Possui ALU, registradores, cache local etc.
Thread	Linha de execução lógica. Pode ser mapeada a um núcleo ou ser dividida via SMT.
SMT (Hyper-Threading)	Técnica para permitir múltiplas threads por núcleo, aumentando a eficiência.

## 5. Exemplo prático

Considere um **Intel Core i7-12700H**, muito comum em notebooks modernos:

- **Cores:** 14 (6 de desempenho + 8 de eficiência)
- **Threads:** 20

Isso significa que, apesar de ter 14 núcleos físicos, o processador consegue lidar com até 20 fluxos de execução distintos simultaneamente, o que melhora consideravelmente o desempenho em multitarefa e aplicações otimizadas para paralelismo.

## 6. Considerações finais

- **Mais threads nem sempre significa mais desempenho**, pois o ganho depende do tipo de aplicação e do suporte do sistema operacional.
- Para tarefas básicas (navegar, editar textos), 4 a 8 threads já são suficientes.
- Para aplicações pesadas (modelagem 3D, compilação, big data), mais threads oferecem ganhos substanciais se o software for bem projetado.

## Conclusão

O número de **threads de processamento** especificado ao comprar um computador indica a **quantidade de tarefas independentes que o processador pode executar de forma simultânea**, seja por meio de múltiplos núcleos físicos ou via técnicas como o SMT. Este número é um **indicador direto da capacidade multitarefa e de paralelismo lógico** do sistema, sendo um dos fatores-chave na escolha de um computador adequado para diferentes perfis de uso.



## 6. ADENDO 2 - Comparativo de Processadores e Threads de Processamento

Esta tabela apresenta um comparativo entre diversos processadores comerciais, indicando a quantidade de núcleos físicos, threads, a tecnologia de multithreading utilizada e os perfis de uso mais adequados. É um guia útil para escolha consciente de hardware, com foco em desempenho e paralelismo lógico.

Modelo do Processador	Fabricante	Núcleos Físicos	Threads	Tecnologia de SMT	Aplicação Típica
Intel Core i3-1215U	Intel	6 (2P + 4E)	8	Sim (Hyper-Threading nos P)	Notebooks básicos e intermediários
Intel Core i5-12400	Intel	6 (6P)	12	Sim	PCs domésticos multitarefa
Intel Core i7-12700H	Intel	14 (6P + 8E)	20	Sim (nos núcleos P)	Notebooks gamers e profissionais
Intel Core i9-13900K	Intel	24 (8P + 16E)	32	Sim	Desktops de alto desempenho
AMD Ryzen 5 5600G	AMD	6	12	Sim	PCs domésticos e home office
AMD Ryzen 7 5800X	AMD	8	16	Sim	Jogos, edição de vídeo, streaming
AMD Ryzen 9 7950X	AMD	16	32	Sim	Workstations, renderização
Apple M2	Apple	8 (4P + 4E)	8	Sim (gerenciado pelo sistema)	MacBooks e ambientes criativos
Intel Xeon W-2255	Intel	10	20	Sim	Servidores e estações de trabalho
AMD EPYC 9654	AMD	96	192	Sim	Servidores de alta performance

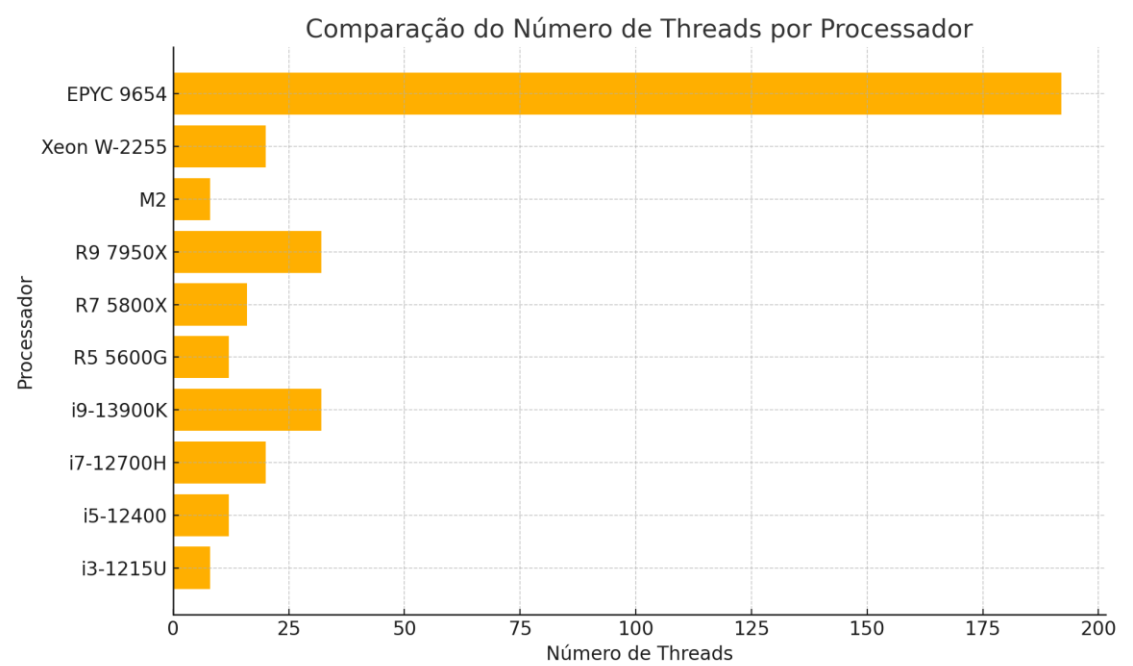
### Notas Explicativas

- P = Performance Core (Núcleo de Alto Desempenho)
- E = Efficient Core (Núcleo de Eficiência Energética) (\* vide nota a seguir)
- SMT (Simultaneous Multithreading) permite a execução de 2 threads por núcleo físico, melhorando o aproveitamento interno dos recursos da CPU.
- Processadores ARM (como o M2 da Apple) gerenciam paralelismo de forma diferenciada e nem sempre expõem claramente o número de threads.

Como Interpretar:

- Processadores com 8 a 12 threads são suficientes para multitarefas leves e moderadas.
- Modelos com 16 ou mais threads atendem bem a tarefas profissionais e jogos de alto desempenho.
- Workstations e servidores com 32+ threads são recomendados para tarefas altamente paralelizáveis, como renderização, análise de dados e virtualização.

O gráfico abaixo compara visualmente o número de threads entre os processadores listados na tabela anterior.



Nota: O termo **Efficient Core**, ou simplesmente **E-Core**, é uma denominação utilizada por fabricantes como a **Intel** (nas arquiteturas Alder Lake, Raptor Lake, entre outras) para descrever um tipo específico de núcleo de CPU projetado para maximizar a eficiência energética, ou seja, executar tarefas com o menor consumo possível de energia, mesmo que isso implique em desempenho reduzido em relação aos núcleos de alto desempenho (*Performance Cores* ou *P-Cores*).

Diferente dos P-Cores, que geralmente possuem Hyper-Threading (duas threads por núcleo), os E-Cores executam apenas uma thread por núcleo. Isso reduz a complexidade do projeto e o consumo energético. Eles são ideais para tarefas em segundo plano (como atualizações, notificações, gerenciamento do sistema operacional, antivírus, etc.) ou trabalhos em paralelo altamente escaláveis (como processamento em lote de dados pequenos).

**Comparativo Técnico: E-Cores vs P-Cores**

Característica	Performance Core (P-Core)	Efficient Core (E-Core)
Objetivo	Máximo desempenho para tarefas exigentes	Eficiência energética e paralelismo leve
Frequência de Clock	Alta	Média a baixa
Tamanho físico	Maior	Menor

Característica	Performance Core (P-Core)	Efficient Core (E-Core)
Capacidade de Threads	2 (via Hyper-Threading)	1 (sem SMT)
Ideal para	Jogos, renderização, simulações	Navegação, tarefas em segundo plano
Presente em	Núcleos principais	Núcleos auxiliares (arquitetura híbrida)

### Exemplo Prático

Considere o **Intel Core i7-12700H**:

- Possui **6 P-Cores** (com Hyper-Threading → 12 threads)
- E **8 E-Cores** (1 thread por núcleo → 8 threads)

◆ Total: **14 núcleos físicos e 20 threads**

Neste design híbrido:

- Os **P-Cores** executam jogos, edição de vídeo, tarefas críticas.
- Os **E-Cores** assumem funções como antivírus, streaming em segundo plano, cálculos simples de scripts, etc.