# Genome 540 discussion

#### January 23st, 2025 Joe Min



Agenda

Homework 2 questions?

GPUs

Climate implications of AI and machine learning

#### Homework 2 questions?

#### GPUs

# But first, central processing units (CPUs)

Major units:

- Storage unit
  - RAM
  - Cache
- Control unit
- Arithmetic Logic Unit (ALU)



- CPU (Central Processing Unit)
- Executes general-purpose tasks (arithmetic, logic).
- Handles complex instructions and task switching.
- CPU has a lower core count (4-8 cores) than others, focusing on delivering higher performance for tasks that rely on a single core.

# CPUs: storage unit



https://spear-itn.eu/wp-content/uploads/2024/01/pyramid\_upgrade-1024x603.png https://images.slideplayer.com/27/9213216/slides/slide\_9.jpg "The holder"

CPUs operate on their internal memory (fast)

CPUs utilize Random Access Memory (RAM) as well as other caches to reduce trips to disk (slow)

# CPUs: control unit

"The manager"

- Translates instructions into arithmetic operations
- Routes instructions and data to the correct executing units

Sends data back as needed



# **CPUs: ALUs**

"The workhorse"

Actually carries out the instructions on the data

Simple binary logic builds into complex programs



https://study.com/cimages/videopreview/videopreview-full/arithmetic-logic-unit-alu-definition-design-and-function\_111389.jpg https://res.cloudinary.com/witspry/image/upload/witscad/public/content/courses/computer-architecture/primitive-alu-supporting-and-or-add-function.png

# Graphics processing units (GPUs)

- Same major units!
- But optimized for specific, parallelizable operations
- 100s-1000s of ALUs

Usually larger storage unit





- GPU (Graphics Processing Unit)
- Used for parallel processing tasks (graphics, AI).
- Has many cores (100s or 1000s of cores) for handling multiple operations simultaneously.
- GPUs use a specialized type of DRAM known as VRAM (Video RAM), specifically designed to handle graphical data and textures.

#### Side-by-side comparison



# Working together



https://www.cgdirector.com/wp-content/uploads/media/2022/06/How-GPU-Acceleration-works.jpg

#### Climate implications of AI and machine learning

### What energy problem?

In 2022, data centers used ~2% of global demand for energy

"By 2026, ... data centres' energy consumption will have increased by between 35% and 128%"

• Equivalent to the energy usage of Sweden or Germany

nature > outlook > article

OUTLOOK | 17 October 2024 | Correction 25 November 2024

#### Fixing AI's energy crisis

Hardware that consumes less power will reduce artificial intelligence's appetite for energy. But transparency about its carbon footprint is still needed.

By <u>Katherine Bourzac</u>

https://www.nature.com/articles/d41586-024-03408-z

https://iea.blob.core.windows.net/assets/6b2fd954-2017-408e-bf08-952fdd62118a/Electricity2024-Analysis and for exast to 2026.pdf and the set of the set

#### AI is exacerbating the issue

A ChatGPT prompt response takes roughly 10x the energy of a Google search

Google estimates its carbon emissions have increased by 48% since 2019

Microsoft, too, at 30% since 2020

### What uses energy?

"Any time that electrons move through chips, some energy is dissipated as heat."

Reliance on huge datasets, and moving it between processing units and long term storage, is the biggest factor

• 90% of energy usage is spent accessing memory

#### What uses water?

#### Making AI Less "Thirsty": Uncovering and Addressing the Secret <u>Water</u> Footprint of AI Models

Pengfei Li UC Riverside Jianyi Yang UC Riverside Mohammad A. Islam UT Arlington Shaolei Ren<sup>1</sup> UC Riverside

Paper that originally made the claim that every ChatGPT response uses, on average, 500mL of water

#### What uses water?

On-site water for data center cooling (scope 1)

Off-site water for electricity generation (scope 2)

Evaporation from cooling towers is the calculated quantity "wasted water"





Good work on homework 1!

Reminder:

Homework 2 is due this Sunday (Jan 26) at 11:59pm!