

# Bottlenecks

# Bottlenecks - The slowest part of your program

A bottleneck is the slowest part of a process that limits overall speed.

In this course, we've analyzed algorithmic bottlenecks using Big O notation.

However, performance bottlenecks can arise from many sources beyond just algorithmic complexity:

- Training a model (high computational cost)
- Feature engineering (intensive transformations)
- Data loading (slow I/O operations)
- Networking (delays in fetching remote data)

 *Key takeaway:* Finding and addressing bottlenecks is critical for optimizing performance.

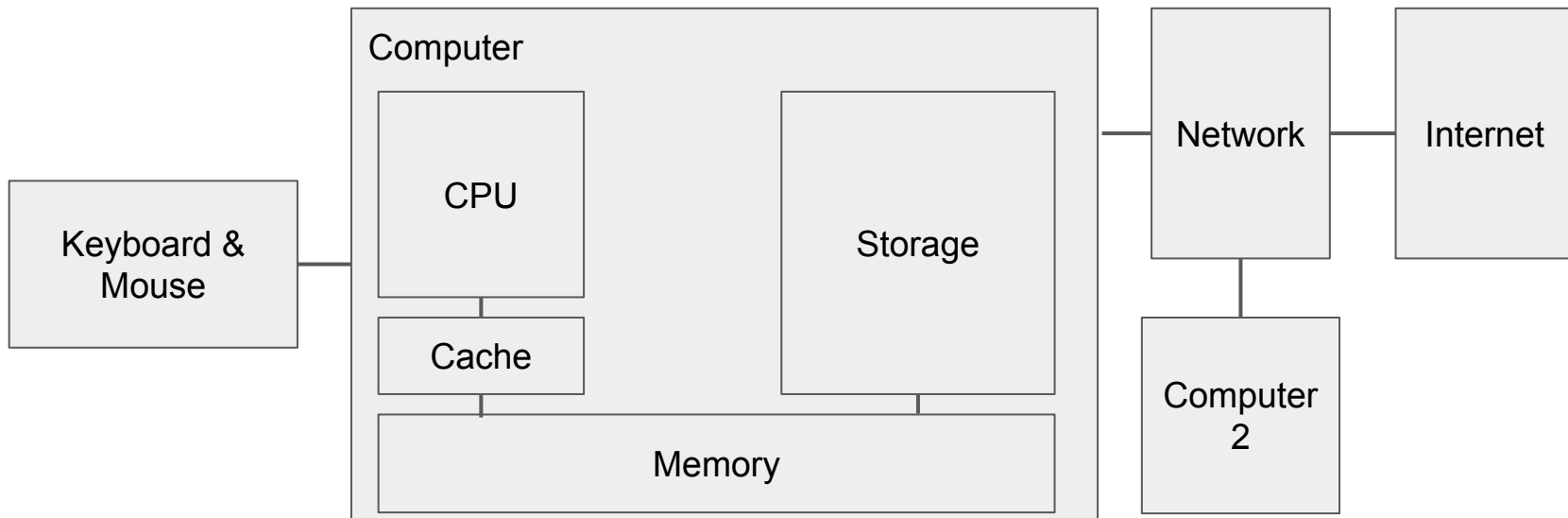
# Bottlenecks can exist elsewhere as well

Bottlenecks aren't always due to poor algorithms.

External factors, like hardware, network latency, or I/O constraints, often slow down programs.

Real-world scenario:

- A deep learning model may run efficiently on a GPU, but if training data is stored on a slow external drive, the training will still be slow.



# Cache & Memory - Fast

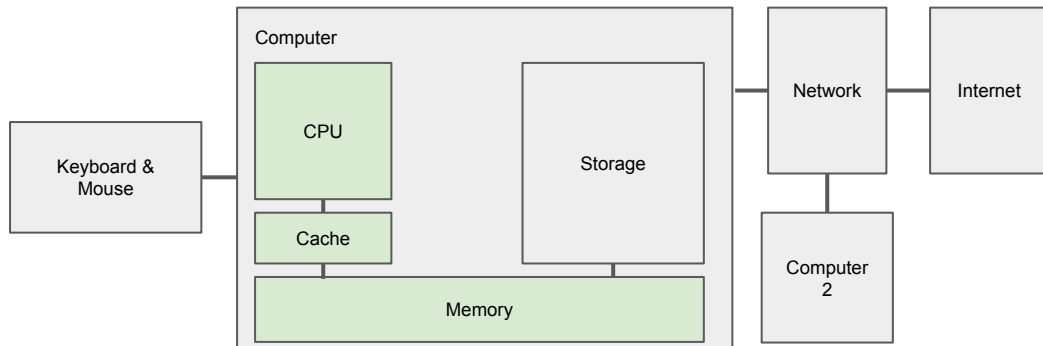
The CPU cache and RAM are the **fastest** places to store and retrieve data.

## Typical speeds:

- **L1/L2/L3 Cache:** ~100GB/s+
- **RAM:** ~20–100GB/s

## Implications:

- If data fits in memory, it's usually not a bottleneck.
- If data doesn't fit and must be read from disk, performance slows down.



📌 *Key takeaway:* Memory access is rarely the problem, unless working with massive datasets.

# Storage - Medium

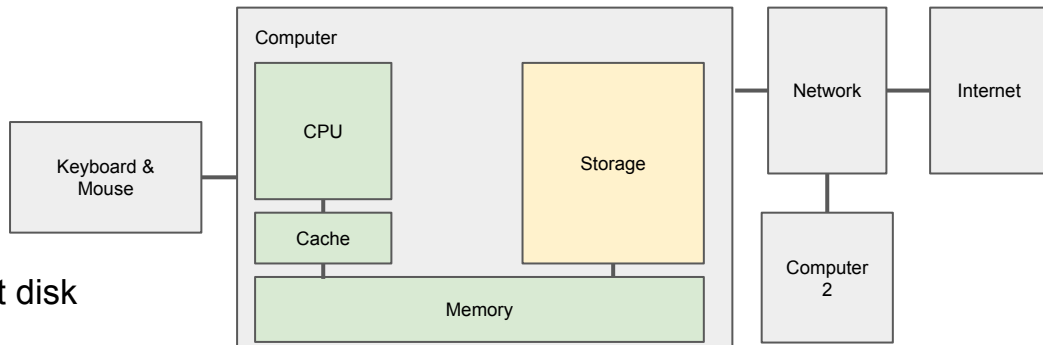
Storage is slower than memory but faster than networks.

- **HDD:** ~100MB/s (slowest)
- **SATA SSD:** ~500MB/s
- **NVMe SSD:** ~5,000MB/s

## Potential bottlenecks:

- If your dataset doesn't fit in RAM, frequent disk access slows performance.
- **Solution:** Optimize **data formats** (e.g., use **Parquet** instead of CSV), leverage memory-mapped files, or upgrade storage.

📌 *Key takeaway:* Disk I/O is a major bottleneck when working with large datasets.



# Networking - Potentially slow

## Networking speeds vary:

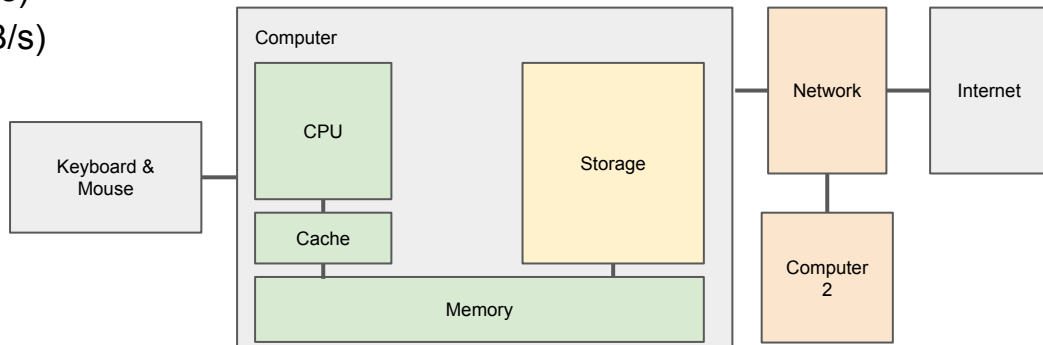
- **Local network (LAN):** ~1Gbps+ (125MB/s)
- **Internet speeds:** ~10–100Mbps (1–10MB/s)

## Bottlenecks arise when:

- Fetching data from remote servers.
- Making frequent API calls.

## Solutions:

- Download once and **cache** locally.
- Use **batch requests** instead of many small ones.
- Use **compression** when transferring data.

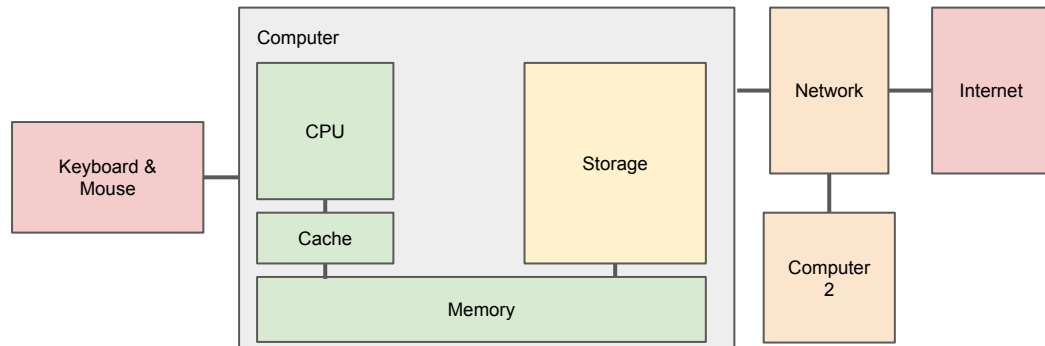


📌 *Key takeaway:* Always reduce dependence on remote data access when possible.

# Manual input = Extremely slow

## Example of slow manual tasks:

- Labeling a dataset by hand.
- Reviewing model predictions manually.



## Solutions:

- Automate wherever possible (e.g., using **active learning** or **semi-supervised learning**).
- Parallelize manual tasks across multiple people.

📌 *Key takeaway:* Human input is the ultimate bottleneck—**eliminate it whenever possible**.

# Bottlenecks in Machine Learning Pipelines

Bottlenecks in ML workflows can occur at different stages:

- **Data ingestion** (loading datasets, API calls)
- **Feature extraction** (expensive transformations)
- **Model training** (large models, slow optimizers)
- **Inference** (real-time predictions, deployment latency)

Example:

- Training a deep learning model on **CPU** vs. **GPU**—huge speed difference!
- Using **batch inference** instead of one-by-one predictions speeds up production systems.



*Key Takeaway:* Each stage of ML has its own performance bottlenecks—**profile before optimizing**



# Identifying Bottlenecks

# Identifying Bottlenecks

To fix bottlenecks, first **identify them**.

Common techniques:

- **Profiling:** Measure where time is spent.
- **Estimations:** Predict expected speeds.
- **A/B Testing:** Compare optimizations.

 *Key takeaway:* You can't optimize what you don't measure.

# Approach 1: Profiling

**Profiling tools help measure time taken per step:**


- Python: `cProfile`, `timeit`
- Linux: `perf`, `strace`

```
data = load_data_from_internet() # Takes 30 seconds
```

```
features = calculate_features(data) # Takes 1.5 seconds
```

**Analysis:**

- Loading data is the bottleneck.
- Potential solution: **Pre-cache the data locally.**


 *Key takeaway:* Profiling is the first step to optimization.

# Approach 2: Calculations

Predict performance using simple math.

## Example:

- **Data size:** 20GB
- **Speed calculations:**
  - **Internet (50MB/s) → 400s**
  - **HDD (150MB/s) → 133s**
  - **SSD (500MB/s) → 40s**
- **Optimization insight:** If model training takes 20s, disk I/O dominates, so **upgrading storage matters**.

 *Key takeaway:* Estimations help predict where bottlenecks will appear.

# Approach 3: A/B Testing

## Test Different Approaches to Find the Fastest Solution


- A/B testing compares two implementations to see which performs better.
- **Steps:**
  1. Implement two or more variations of a function or workflow.
  2. Measure performance under the same conditions.
  3. Choose the best-performing approach.

## Example: Optimizing Data Loading

- **Version A:** Load data from a CSV file.
- **Version B:** Load the same data from a Parquet file.
- **Test Result:** Parquet is **5x faster** due to better compression and indexing.

## When to Use A/B Testing?

- ✓ When multiple possible solutions exist.
- ✓ When performance gains are uncertain.
- ✓ When optimizing for real-world conditions.

 **Key Takeaway:** Always test optimizations—**assumptions can be wrong!**

How do we resolve bottlenecks

# The Trade Off Triangle - Speed, Cost, Accuracy

Optimizing bottlenecks often involves tradeoffs between:

- **Speed** (How fast does it run?)
- **Cost** (How expensive is it to compute/store?)
- **Accuracy** (Do we lose precision by optimizing?)

Example:



- Using **lower-precision models (e.g., float16 vs. float32)** speeds up inference but may reduce accuracy.
- Storing **compressed datasets** saves space but may slow down access.

 **Key Takeaway:** Optimization isn't just about speed—**balance trade offs based on your goals.**

# How do we resolve bottlenecks

## Memory Bottlenecks

- **Optimize algorithms** (reduce time complexity).

## Storage Bottlenecks


- **Upgrade hardware** (CPU, RAM, SSD).
- **Improve data management** (better formats, indexing).

## Network Bottlenecks

- **Optimize networking** (caching, compression).

## Manual Input Bottlenecks

- **Automate manual work** whenever possible.

 *Key takeaway:* The best solution depends on the type of bottleneck.




# Memory bottlenecks

## Causes:

- Inefficient data structures (e.g., linked lists instead of arrays).
- Algorithmic inefficiencies (e.g., **bubble sort** instead of **merge sort**).
- Primarily what we focused on in this module

## Solutions:

- Use **efficient algorithms**.
- Optimize **data representations** (e.g., NumPy over lists).

 *Key takeaway:* Optimize algorithms to reduce memory access bottlenecks.

# Storage bottlenecks

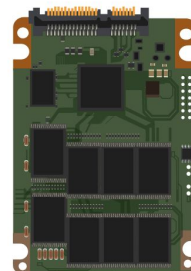
## Fixes:

- Upgrade to **NVMe SSD**.
- Use **RAM disks** for temporary storage.
- Use **indexed** file formats (Parquet, HDF5).
- **Benchmark read speeds** (e.g., Linux `dd` command).

### TYPES OF HARD DRIVES



SATA HDD



SATA SSD



M.2 NVMe



M.2 SATA

📌 *Key takeaway:* Faster storage speeds up I/O-bound processes.

# Optimizing Data Formats

Choosing the right data format can **eliminate bottlenecks**:

- **CSV** (Slow, large, no indexing)
- **Parquet** (Fast, compressed, indexed)
- **HDF5** (Optimized for large datasets)

```
import pandas as pd
```

```
df = pd.read_csv("data.csv") # Slow
```

```
df.to_parquet("data.parquet") # Convert to faster format
```

```
df = pd.read_parquet("data.parquet") # Much faster loading
```



*Key Takeaway:* **Optimized data formats** can improve speed **without extra computing power**.

# Network bottlenecks

## Fixes:

- Use **wired networks** over WiFi.
- Upgrade to **faster network cables (CAT7)**.
- **Cache frequently used data**.
- Use **parallel downloads** when possible.

 *Key takeaway:* Reduce reliance on slow external networks.

# Choosing the Right Compute Resources

Your hardware can be a **bottleneck** or a **boost**:

- **More CPU cores** → Faster parallel processing.
- **More RAM** → Load larger datasets without disk swapping.
- **GPUs/TPUs** → Accelerate deep learning.

Example:

- Training a model on **Colab CPU vs. Colab GPU**—huge speedup with GPU!



*Key Takeaway:* Hardware upgrades often provide the easiest performance gains.

# Manual input Bottlenecks

## Fixes:

- Automate processes.
- Run other tasks in parallel while waiting for input.
- Use **AI-assisted tools** (e.g., **auto-labeling** for datasets).

 *Key takeaway:* The best optimization for manual input is **eliminating it**.

# Real-World Case Studies in Bottlenecks

**Netflix:** Moved from on-premises data centers to the cloud to improve **scalability and streaming speed**.

**Google Search:** Optimized indexing algorithms to speed up **retrieval from trillions of documents**.

**Deep Learning:** Training GPT-3 required **thousands of GPUs** to avoid long training times.

 *Key Takeaway:* Industry leaders constantly optimize to remove bottlenecks—**so should you!**

# Final Thoughts

Bottlenecks appear **everywhere** in data science.

Optimizing one part might **expose another bottleneck**.

Always **profile first** before optimizing.

Smart **hardware choices** and **algorithmic efficiency** go hand-in-hand.



*Takeaway:* The slowest part of your system determines overall speed—find it, fix it, and repeat.