

UC San Diego

# **DSC 102**

# **Systems for Scalable Analytics**

Spring 2024

Haojian Jin

Now for the course logistics ...

# Prerequisites

- ❖ **DSC 100** (or equivalent) is necessary
- ❖ Transitively **DSC 80**; a mainstream ML algorithmics course is necessary
- ❖ Proficiency in Python programming
- ❖ For all other cases, email me with proper justification; a waiver can be considered

<https://haojian.github.io/DSC102SP24/>

# Components and Grading

- ❖ **3 Programming Assignments: 40%** (8% + 16% + 16%)
  - ❖ No late days! Plan your work well ahead.
  - ❖ **Plan your credit as well!**
- ❖ **Midterm Exam: 15%**
  - ❖ TBD; in-class only (50min)
- ❖ **Cumulative Final Exam: 35%**
  - ❖ 3hrs long but 4hrs limit
- ❖ **10 (of 12) Peer Instruction Activities: 10%**
- ❖ **Extra Credit Evaluation Activities: 2%** (likely)
- ❖ LMK ahead of time if you need makeup exam slot

<https://haojian.github.io/DSC102SP24/>

# Grading Scheme

Hybrid of relative and absolute; grade is better of the two

Grade	Relative Bin (Use strictest)	Absolute Cutoff (>=)
A+	Highest 5%	95
A	Next 10% (5-15)	90
A-	Next 15% (15-30)	85
B+	Next 15% (30-45)	80
B	Next 15% (45-60)	75
B-	Next 15% (60-75)	70
C+	Next 5% (75-80)	65
C	Next 5% (80-85)	60
C-	Next 5% (85-90)	55
D	Next 5% (90-95)	50

**Example:** Score 82 but 33%ile; Rel.: B-; Abs.: B+; so, B+

# Programming Assignments

- ❖ **PA0: Setting up AWS and Dask**
- ❖ **PA1: Data Exploration with Dask**
- ❖ **PA2: Feature Eng. and Model Selection with Spark**
- ❖ **Expectations on the PAs:**
  - ❖ Teams of 1-3; see webpage on academic integrity
  - ❖ I will cover the concepts and tools' tradeoffs in the lectures
  - ❖ TAs will explain and demo the tools; handle all Q&A
  - ❖ You are expected to put in the effort to learn the details of the tools' APIs using their documentation on your own!

<https://haojian.github.io/DSC102SP24/>

# Course Administtrivia

- ❖ **Lectures: MWF 3pm-3:50pm PT at Mandeville Center - B-202**
  - ❖ Attendance optional but encouraged; podcast available
  - ❖ No need for clickers.
- ❖ **Discussions:**
  - ❖ Only for talks on PAs by TAs, for pre-exam review by me
- ❖ **Instructor:** Haojian Jin; haojian@ucsd.edu
  - ❖ OHs: **Wednesday 4-5 pm PT at HDSI 341**
- ❖ **Slack** for all communications
- ❖ **Canvas** for PA submission, Peer Evaluation Activities, Final Exam

<https://haojian.github.io/DSC102SP24/>

# Office hours

- ❖ Haojian Jin's OHs: Wednesday 4:00 PM - 5:00 PM
- ❖ Course content.
- ❖ Tony Li's OHs: Thursday 3pm
- ❖ Ariane Yu's OHs: Tuesday 1pm
- ❖ Assignments, HDSI 3<sup>rd</sup> floor. Near conference rooms.
- ❖ Post questions to the ta-public channel.
- ❖ Avoid asking repetitive questions.

<https://haojian.github.io/DSC102SP24/>



# General Dos and Do NOTs

## **Do:**

- ❖ Follow all announcements on Piazza
- ❖ Try to join the lectures/discussions live
- ❖ Raise your hand before speaking
- ❖ View/review podcast videos asynchronously by yourself
- ❖ To contact me/TAs, use private Slack; if you really need to email, use “DSC 102:” as subject prefix

## **Do NOT:**

- ❖ Harass, intimidate, or intentionally talk over others
- ❖ Violate academic integrity on the PAs, exams, or other components; I am *very strict* on this matter!

# Reasonable person.

- (1) Everyone will be reasonable.
- (2) Everyone expects everyone else to be reasonable.
- (3) No one is special.
- (4) Do not be offended if someone suggests you are not being reasonable.

Now for the course structure ...

## DSC 102 will get you thinking about the fundamentals of systems for scalable analytics

1. **“Systems”**: What resources does a computer have? How to store and efficiently compute over large data? What is cloud?
2. **“Scalability”**: How to scale and parallelize data-intensive computations?
3. **For “Analytics”**:
  1. **Source**: Data acquisition & preparation for ML
  2. **Build**: Model selection & deep learning systems
  3. **Deploying** ML models
4. Hands-on experience with scalable analytics tools

# Data Systems Concerns in ML

## Key concerns in ML:

**Q:** How do “ML Systems” relate to ML?

## Runtime efficiency (sometimes)

## Additional key *practical* concerns in ML Systems:

ML Systems : ML :: Computer Systems : TCS  
Scalability (and efficiency at scale) | Long standing

# Usability

# Manageability

# Developability

*Long-standing  
concerns in the  
**DB systems**  
world!*

**Q. Q. What do you think will be the top 5 emerging technologies in the next 5 years?**

# Conceptual System Stack Analogy

## Relational DB Systems

## ML Systems

### Theory

First-Order Logic  
Complexity Theory

Learning Theory  
Optimization Theory

### Program Formalism

Relational Algebra

Tensor Algebra  
Gradient Descent

### Program Specification

SQL

TensorFlow?  
Scikit-learn?

### Program Modification

Query Optimization

???

### Execution Primitives

Parallel Relational  
Operator Dataflows

Depends on ML Algorithm

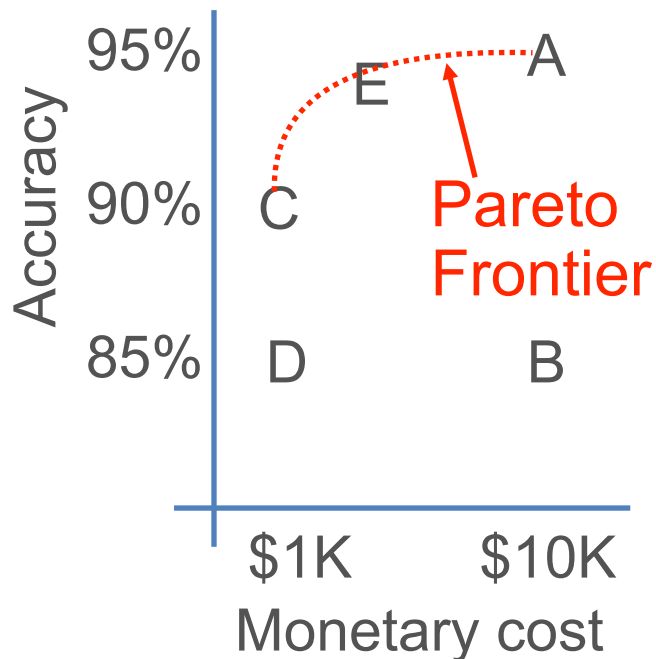
### Hardware

CPU, GPU, FPGA, NVM, RDMA, etc.

# Real-World ML: Pareto Surfaces

*Q: Suppose you are given ad click-through prediction models A, B, C, and D with accuracies of 95%, 85%, 90%, and 85%, respectively. Which one will you pick?*

*Q: What about now?*



- ❖ Real-world ML users must grapple with multi-dimensional *Pareto surfaces*: accuracy, monetary cost, training time, scalability, inference latency, tool availability, interpretability, fairness, etc.
- ❖ *Multi-objective optimization* criteria set by application needs / business policies.

# Learning Outcomes of this course

- ❖ **Explain** the basic principles of the memory hierarchy, parallelism paradigms, scalable data systems, and cloud computing.
- ❖ **Identify** the abstract data access patterns of, and opportunities for parallelism and efficiency gains in, data processing and ML algorithms at scale.
- ❖ **Outline** how to use cluster and cloud services, dataflow (“Big Data”) programming with MapReduce and Spark, and ML tools at scale.
- ❖ **Apply** the above programming skills to create end-to-end pipelines for data preparation, feature engineering, and model selection on large-scale datasets.
- ❖ **Reason** critically about practical tradeoffs between accuracy, runtimes, scalability, usability, and total cost.

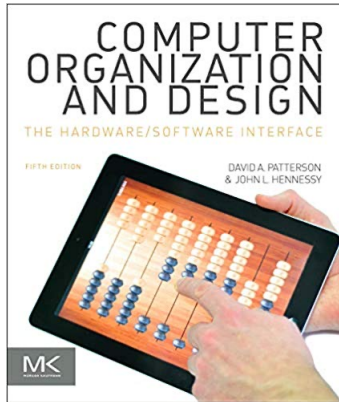


# Tentative Course Schedule

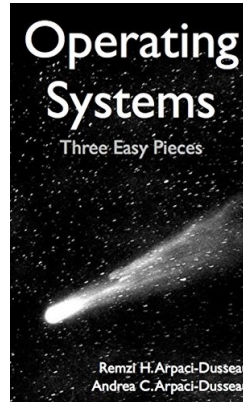
Week	Topic	
Systems Principles	Basics of Machine Resources: Computer Organization	
	Basics of Machine Resources: Operating Systems	
	4	Basics of Cloud Computing
4-5	Parallel and Scalable Data Processing: Parallelism Basics	
Scalability Principles	Midterm Exam on TBD	
	Parallel and Scalable Data Processing: Scalable Data Access	
	6-7	Parallel and Scalable Data Processing: Data Parallelism
7-8		
9	Scalable Analytics Systems	Dataflow Systems
10		ML Model Building Systems
11	Final Exam on Dec 15	

There will be 2 industry guest lectures (maybe 3)

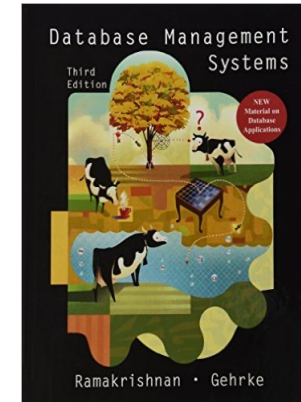
# Suggested Textbooks



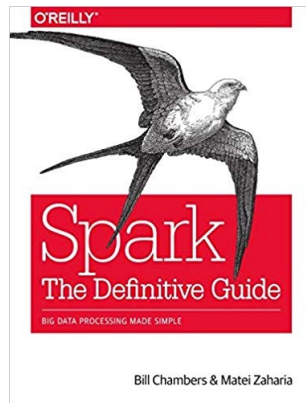
Aka “CompOrg Book”



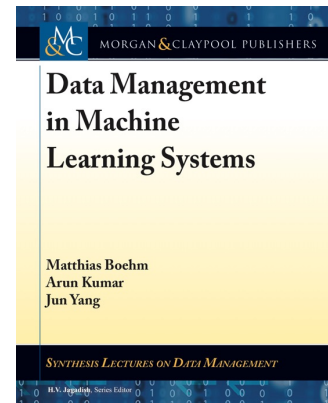
Aka “Comet Book”



Aka “Cow Book”



Aka “Spark Book”

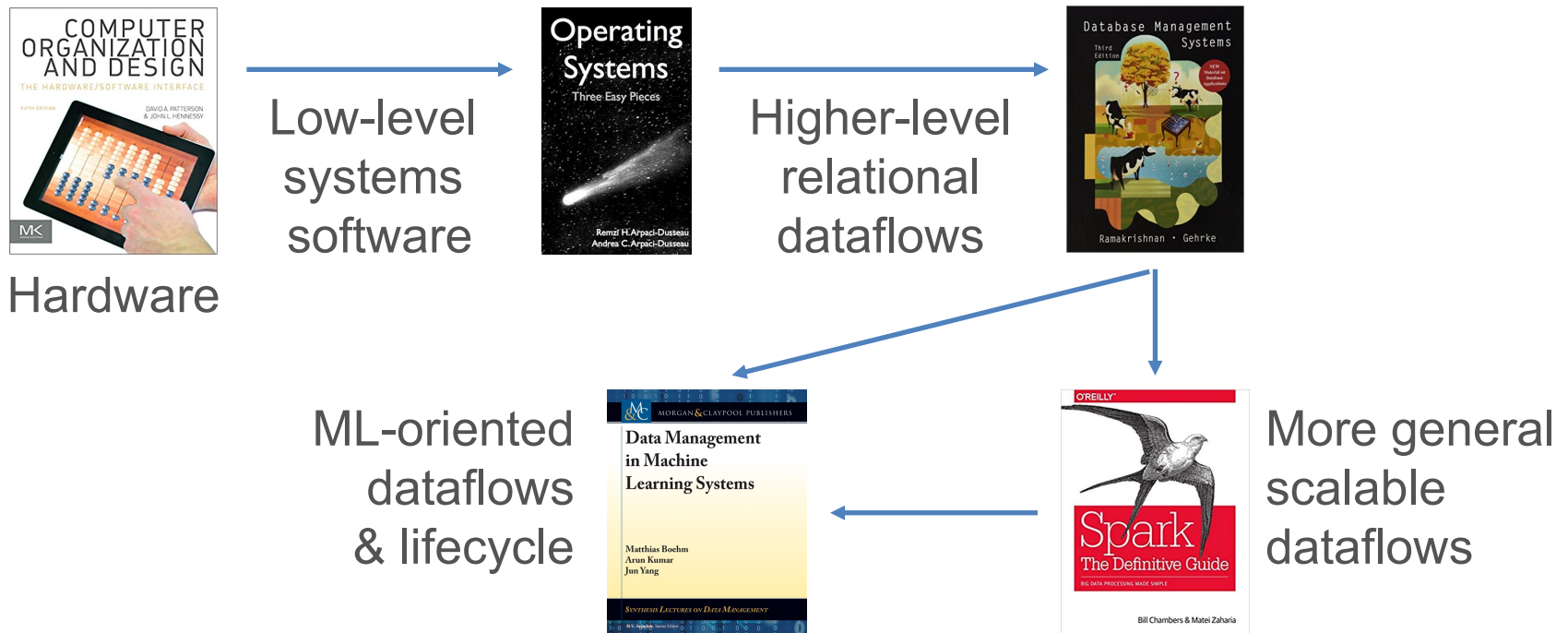


Aka “MLSys Book”

(Free PDFs available online; also check out our library)

# Why so many textbooks?!

1. Computer systems are about carefully layering *levels of abstraction*.



2. Analytics/ML Systems is a recent/emerging area of research.

3. Also, DSC 102 is the first UG course of its kind in the world!

# Tentative Course Schedule

Week	Topic
1-2	Basics of Machine Resources: Computer Organization
Systems Principles	Basics of Machine Resources: Operating Systems
	Basics of Cloud Computing
4-5	Parallel and Scalable Data Processing: Parallelism Basics
6	<b>Midterm Exam on TBD</b>
6-7	Parallel and Scalable Data Processing: Scalable Data Access
7-8	Parallel and Scalable Data Processing: Data Parallelism
9	Dataflow Systems
10	ML Model Building Systems
11	<b>Final Exam on Dec 15</b>

There will be 2 industry guest lectures (maybe 3)

# **DSC 102**

## **Systems for Scalable Analytics**

Topic 1: Basics of Machine Resources  
Part 1: Computer Organization

Ch. 1, 2.1-2.3, 2.12, 4.1, and 5.1-5.5 of CompOrg Book

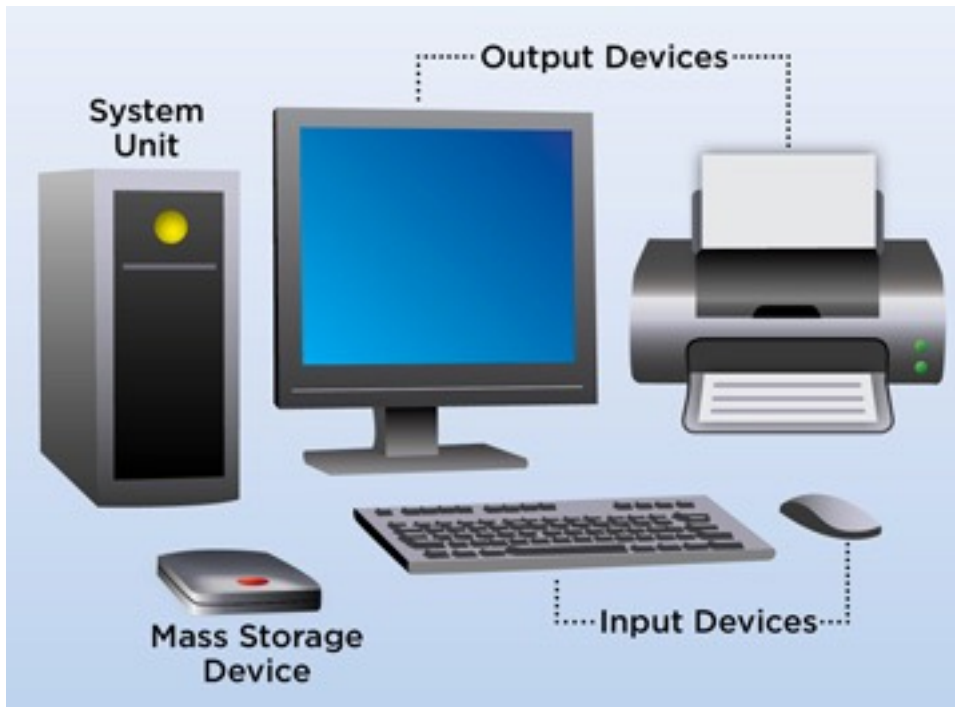
***Q: What is a computer?***

A programmable electronic device that can store, retrieve,  
and process digital data.

# Outline

- ➔ ❖ Basics of Computer Organization
  - ❖ Digital Representation of Data
  - ❖ Processors and Memory Hierarchy
- ❖ Basics of Operating Systems
  - ❖ Process Management: Virtualization; Concurrency
  - ❖ Filesystem and Data Files
  - ❖ Main Memory Management
- ❖ Persistent Data Storage

# Parts of a Computer



## **Hardware:**

The electronic machinery (wires, circuits, transistors, capacitors, devices, etc.)

## **Software:**

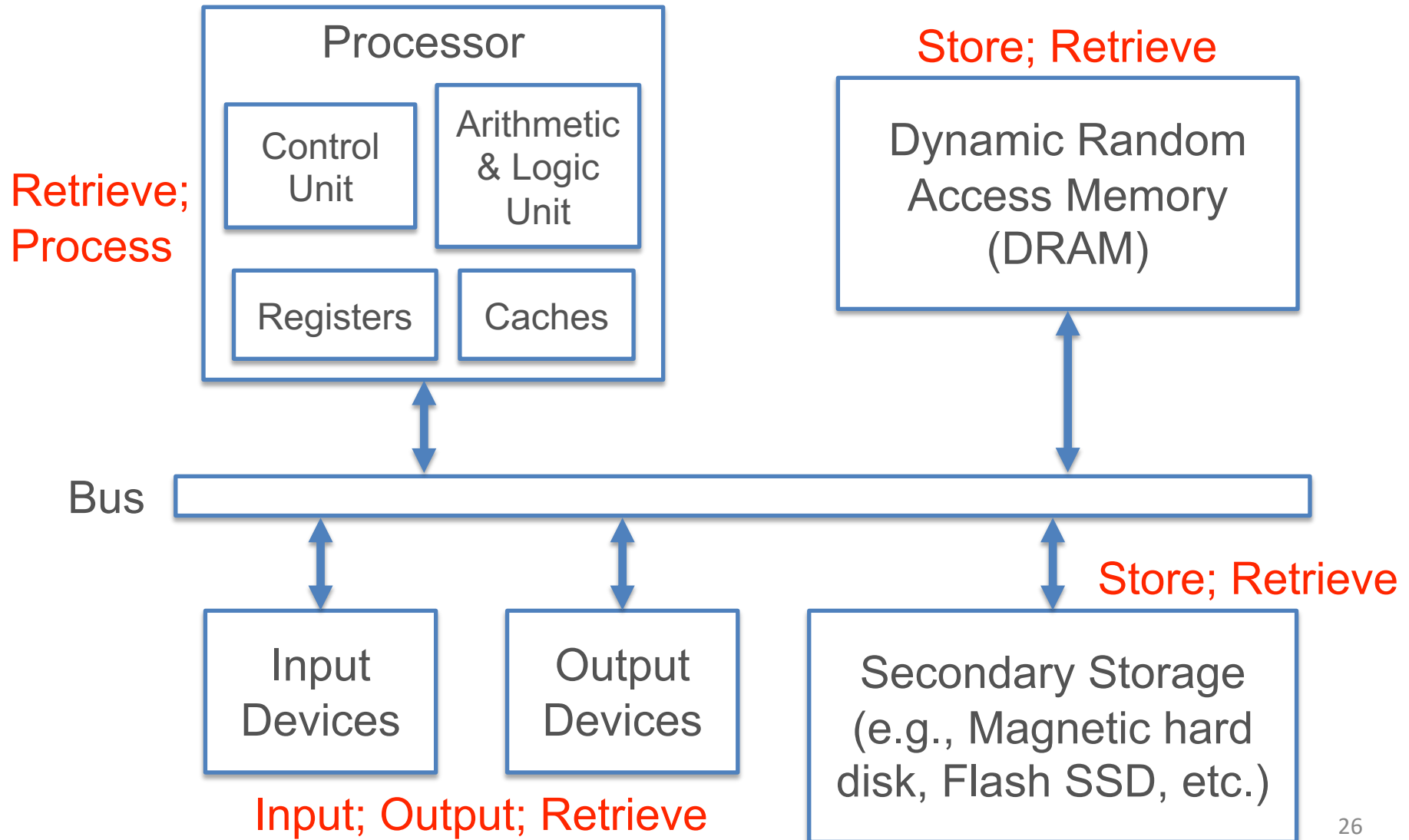
Programs (instructions) and data

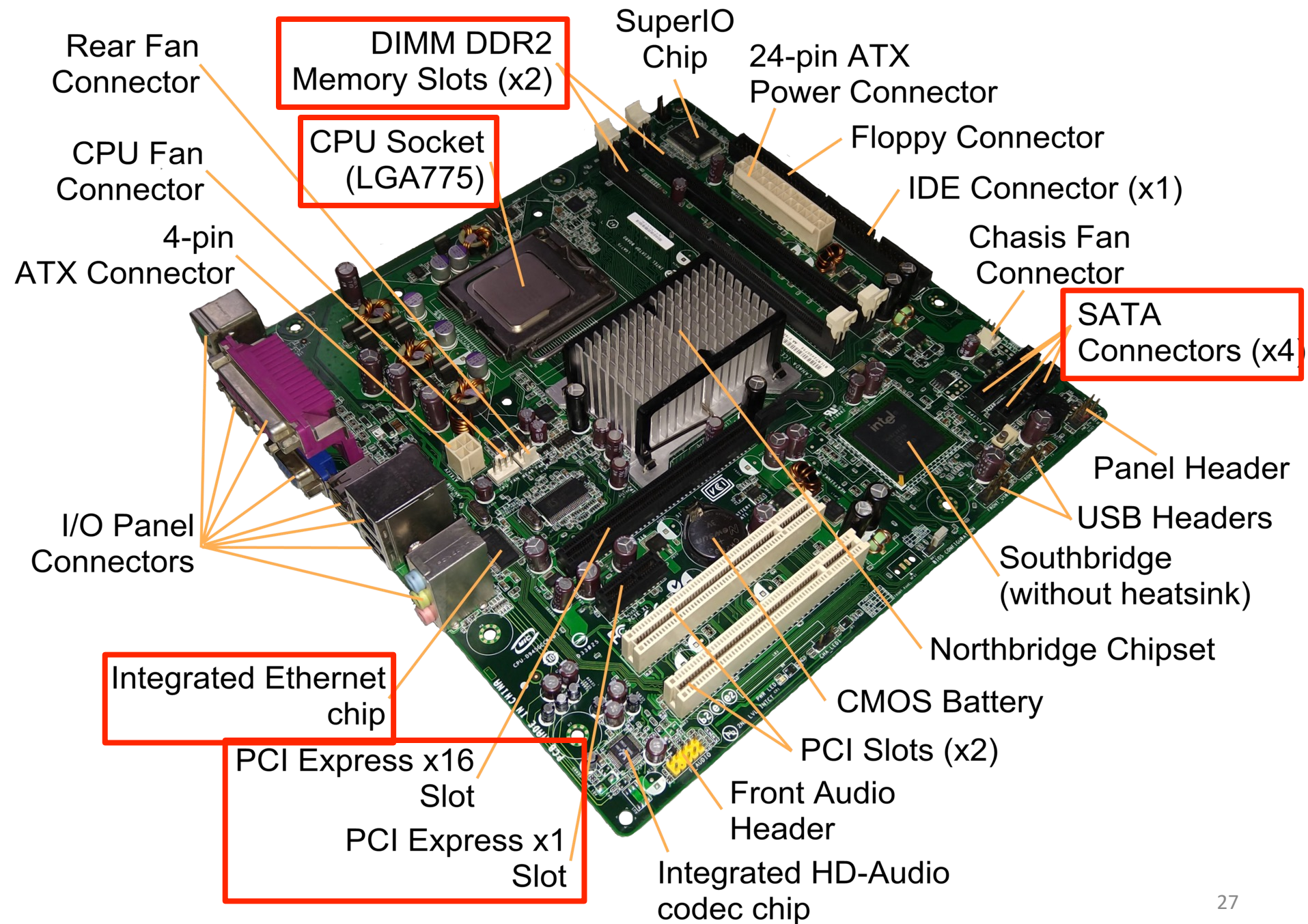


# Key Parts of Computer Hardware

- ❖ **Processor** (CPU, GPU, etc.)
  - ❖ Hardware to orchestrate and execute *instructions* to manipulate *data* as specified by a *program*
- ❖ **Main Memory** (aka Dynamic Random Access Memory)
  - ❖ Hardware to store *data* and *programs* that allows very fast location/retrieval; byte-level *addressing* scheme
- ❖ **Disk** (aka secondary/persistent storage)
  - ❖ Similar to memory but *persistent*, *slower*, and higher capacity / cost ratio; various addressing schemes
- ❖ **Network** interface controller (NIC)
  - ❖ Hardware to send data to / retrieve data over network of interconnected computers/devices

# Abstract Computer Parts and Data





# Key Aspects of Software

## ❖ Instruction

- ❖ A command understood by hardware; finite vocabulary for a processor: Instruction Set Architecture (ISA); bridge between hardware and software

## ❖ Program (aka code)

- ❖ A collection of instructions for hardware to execute

## ❖ Programming Language (PL)

- ❖ A human-readable *formal* language to write programs; at a much higher level of *abstraction* than ISA

## ❖ Application Programming Interface (API)

- ❖ A set of functions (“interface”) exposed by a program/set of programs for use by humans/other programs

## ❖ Data

- ❖ Digital representation of *information* that is stored, processed, displayed, retrieved, or sent by a program

# Main Kinds of Software

## ❖ Firmware

- ❖ Read-only programs “baked into” a device to offer basic hardware control functionalities

## ❖ Operating System (OS)

- ❖ Collection of interrelated programs that work as an intermediary platform/service to enable application software to use hardware more effectively/easily
- ❖ Examples: Linux, Windows, MacOS, etc.

## ❖ Application Software

- ❖ A program or a collection of interrelated programs to manipulate data, typically designed for human use
- ❖ Examples: Excel, Chrome, PostgreSQL, etc.

# Outline

- ❖ Basics of Computer Organization
  - ❖ Digital Representation of Data
  - ❖ Processors and Memory Hierarchy
- ❖ Basics of Operating Systems
  - ❖ Process Management: Virtualization; Concurrency
  - ❖ Filesystem and Data Files
  - ❖ Main Memory Management
- ❖ Persistent Data Storage

# Why bother with these in Data Science?

- ❖ Basics of Computer Organization

You will face myriad and new data types

- ❖ Digital Representation of Data

- ❖ Processors and Memory Hierarchy

Compute hardware is evolving fast

- ❖ Basics of Operating Systems

- ❖ Process Management: Virtualization; Concurrency

- ❖ Filesystem and Data Files

- ❖ Main Memory Management

You will need to use new methods on evolving data file formats on clusters / cloud

- ❖ Persistent Data Storage

Storage hardware is evolving fast