

# DSC 102

# Systems for Scalable Analytics

Winter 2022

Arun Kumar

# About Myself



2009: Bachelors in CSE from IIT Madras, India

Summers: 110F!



2009–16: MS and PhD in CS from UW-Madison  
PhD thesis area: Data systems for ML workloads

Winters: -40F!



2016-: Asst. Prof. at UC San Diego CSE  
2019-: + Asst. Prof. at UC San Diego HDSI  
2021: Assoc. Prof. at CSE & HDSI

Ahh! :)

# My Current Research

New abstractions, algorithms, and software systems  
to “**democratize**” ML-based data analytics from  
a data management/systems standpoint

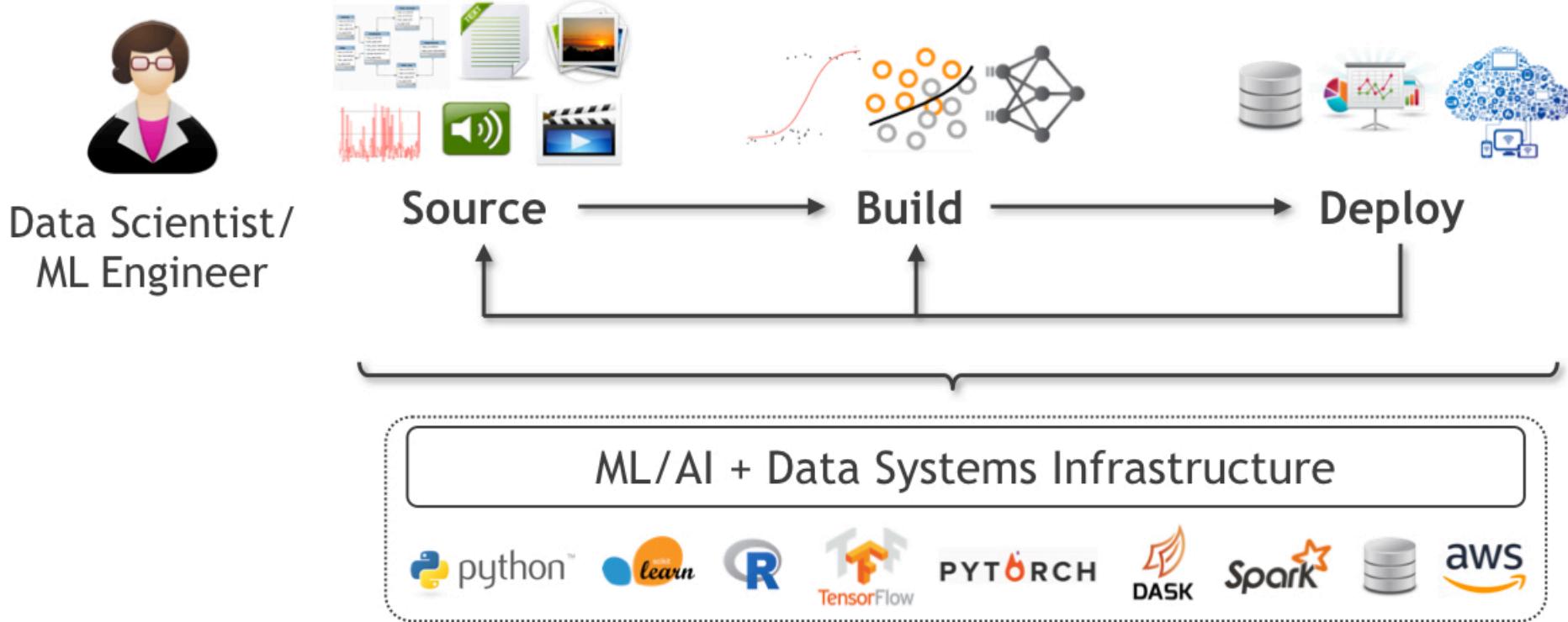
**Democratization =**      **System Efficiency**      +      **Human Efficiency**  
  (**Lower resource costs**)      +      (**Higher productivity**)

Practical and scalable data systems for ML analytics

Inspired by *relational database systems* principles

Exploit insights from *learning theory* and *optimization theory*

# My Current Research



**Research Approach :** *Abstract* key steps + *Formalize* computation + *Automate* grunt work + *Optimize* execution

What is this course about? Why take it?

# 1. Netflix's “spot-on” recommendations

## NETFLIX ORIGINAL **STRANGER THINGS**

95% Match 2017 2 Seasons 4K Ultra HD 5.1

When a young boy vanishes, a small town uncovers a mystery involving secret experiments, terrifying supernatural forces and one strange little girl.

*Winona Ryder, David Harbour, Matthew Modine*  
TV Shows, TV Sci-Fi & Fantasy, Teen TV Shows



### Popular on Netflix



### Recently Watched



# How does Netflix know that?

# Large datasets + Machine learning!

**Everything is a Recommendation**



**Over 80% of what people watch comes from our recommendations**

**Recommendations are driven by Machine Learning**

6

Log all user behavior (views, clicks, pauses, searches, etc.)  
Recommender systems apply ML to TBs of data from all users and movies to deliver a tailored experience

# 2. Structured data with search results

Google pradeep khosla

All News Images Videos Maps More Settings Tools

About 274,000 results (0.51 seconds)

[Pradeep Khosla - UC San Diego Office of the Chancellor - University ...](#)  
chancellor.ucsd.edu/chancellor-khosla ▾  
Pradeep K. Khosla became UC San Diego's eighth Chancellor on August 1, 2012. As UC San Diego's chief executive officer, he leads a campus with more than ...

[Pradeep K. Khosla - UC San Diego Office of the Chancellor](#)  
chancellor.ucsd.edu/chancellor-khosla/khosla-biography ▾  
Chancellor, University of California San Diego. Pradeep K. Khosla, an internationally renowned electrical and computer engineer, is the eighth Chancellor of the ...

[Pradeep Khosla - Wikipedia](#)  
[https://en.wikipedia.org/wiki/Pradeep\\_Khosla](https://en.wikipedia.org/wiki/Pradeep_Khosla) ▾  
Pradeep K. Khosla is an academic computer scientist and university administrator. He is the current chancellor of the University of California, San Diego. He was ...

[Pradeep Khosla | LinkedIn](#)  
<https://www.linkedin.com/in/pradeepkhosla> ▾  
Greater San Diego Area - Chancellor, UC San Diego - Avigilon  
View Pradeep Khosla's professional profile on LinkedIn. LinkedIn is the world's largest business network, helping professionals like Pradeep Khosla discover ...

[Robotics Institute: Pradeep Khosla](#)  
[www.ri.cmu.edu](http://www.ri.cmu.edu) > people ▾



More images

**Pradeep Khosla**

Chancellor of the University of California, San Diego

Pradeep K. Khosla is an academic computer scientist and university administrator. He is the current chancellor of the University of California, San Diego. [Wikipedia](#)

**Born:** March 13, 1957 (age 60 years), Mumbai, India

**Spouse:** Thespine Kavoulakis

**Education:** Indian Institute of Technology Kharagpur, Carnegie Mellon University

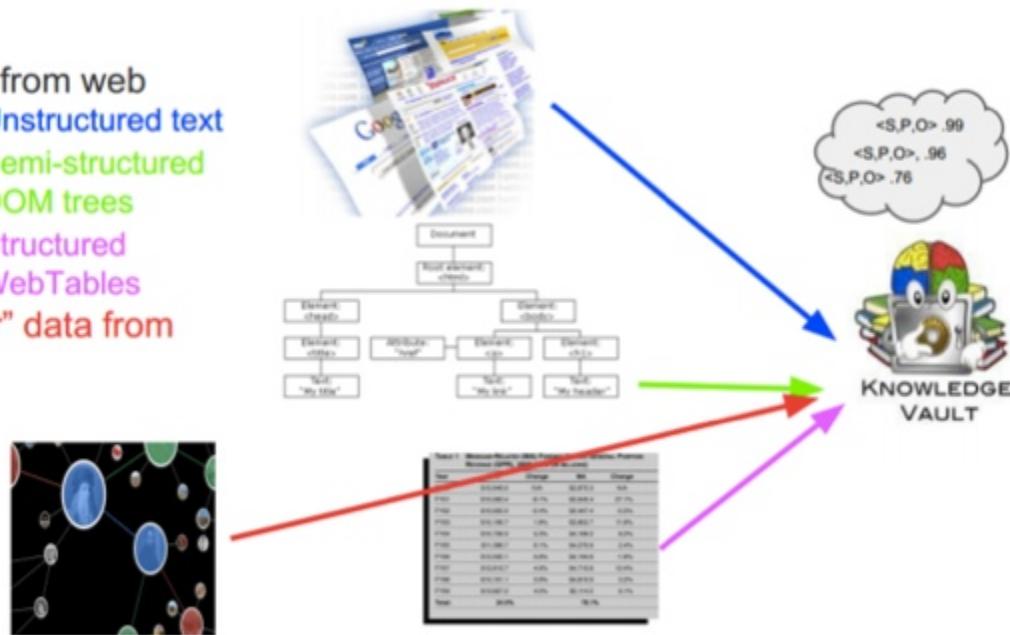
**Residence:** Audrey Geisel University House, La Jolla, CA

# How does Google know that?

# Large datasets + Machine learning!

Knowledge Vault\* fuses all these signals together

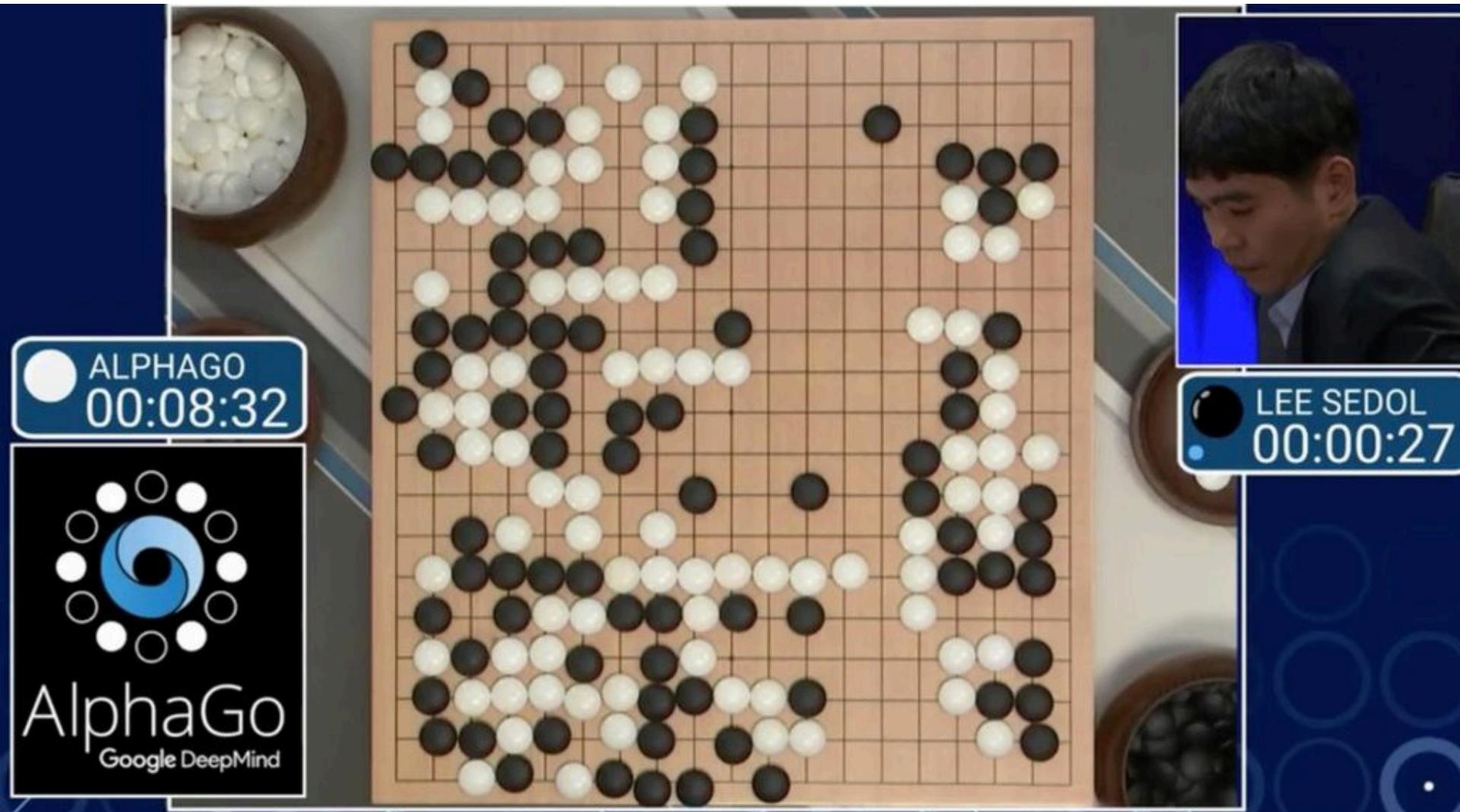
- Data from web
  - Unstructured text
  - Semi-structured DOM trees
  - Structured WebTables
- “Prior” data from FB



\* Details in a paper submitted to WWW'14 (Dong et al)

Knowledge Base Construction (KBC) process extracts tabular/relational data from large amounts of text data

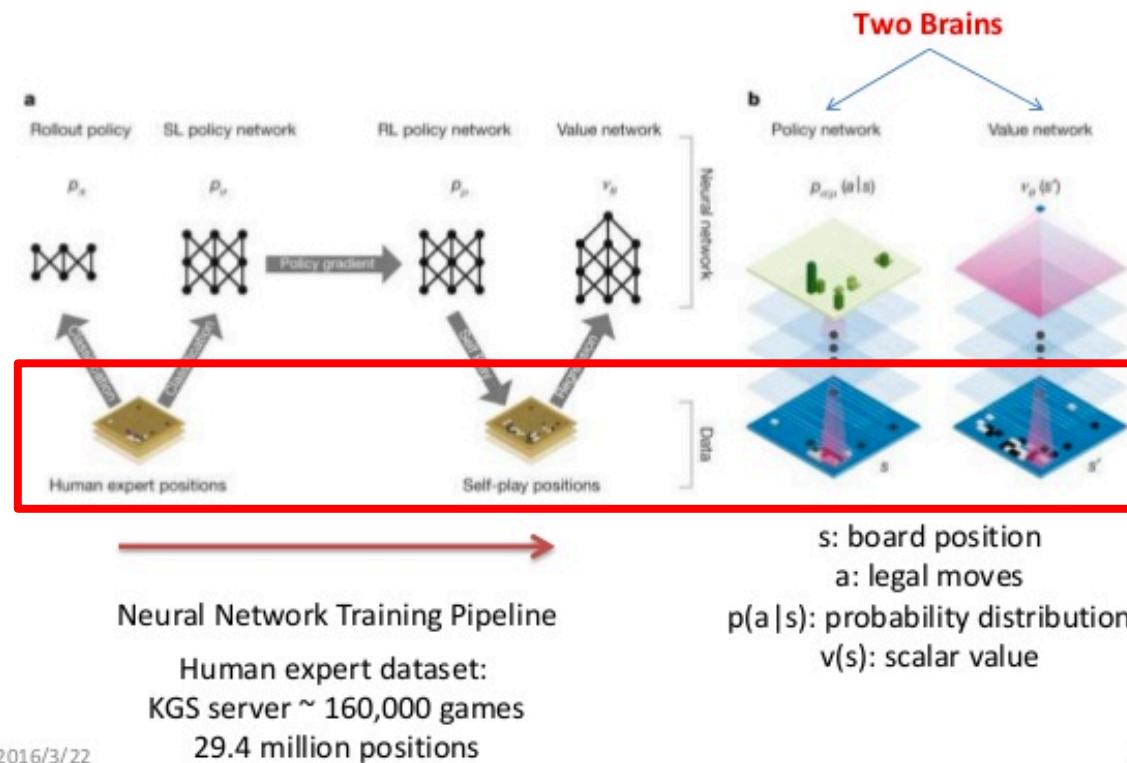
### 3. AlphaGo defeats human champion!



# How did AlphaGo achieve that?

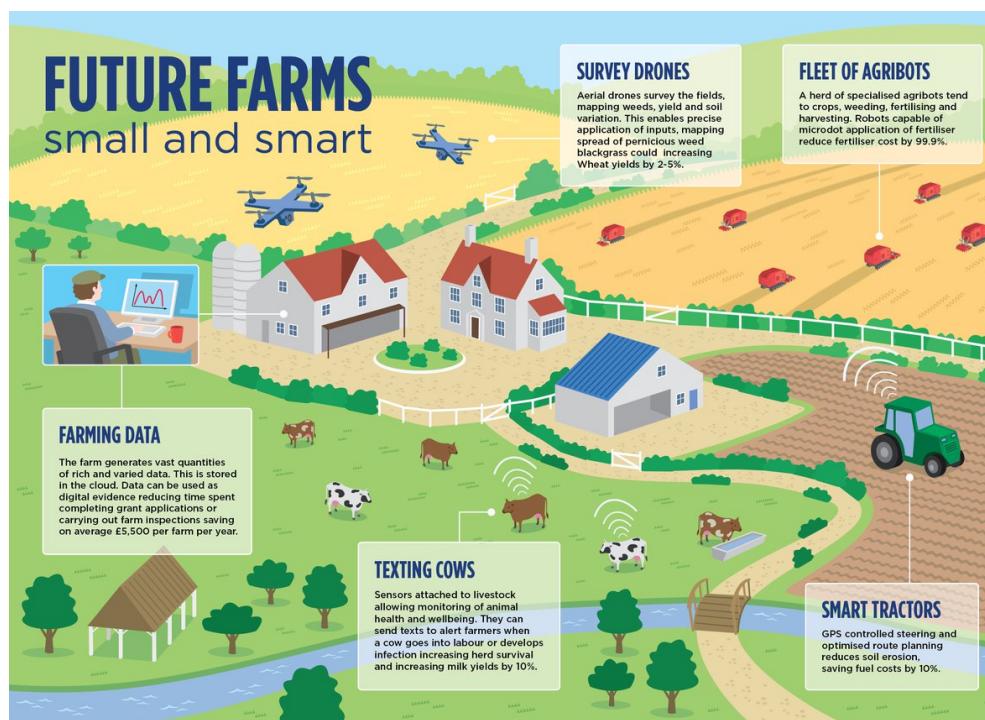
# Breakthrough powered by deep learning!

## Architecture of AlphaGo

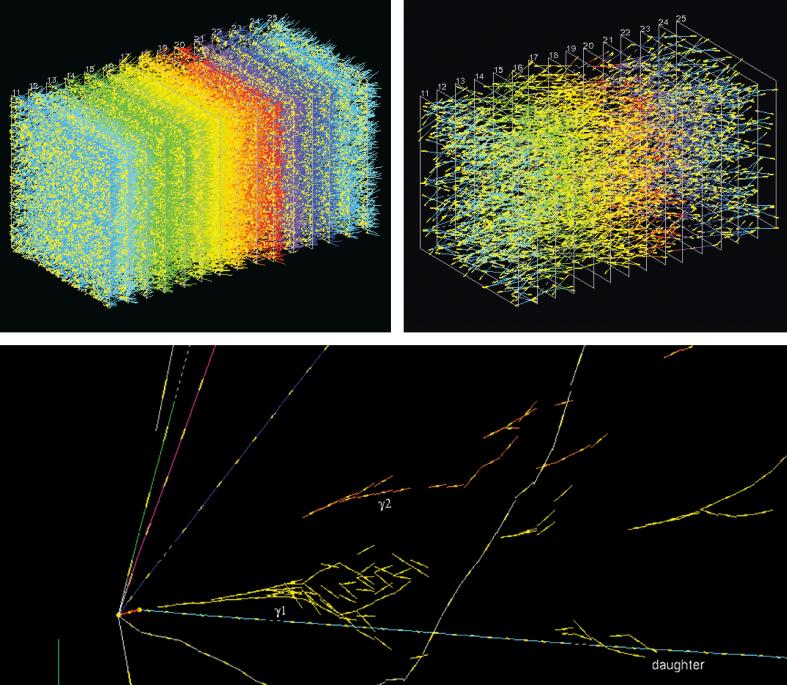


Deep CNNs to visually process board status in plays

# Innumerable “enterprise” applications



“Domain sciences” and healthcare tech  
are also becoming data+ML intensive

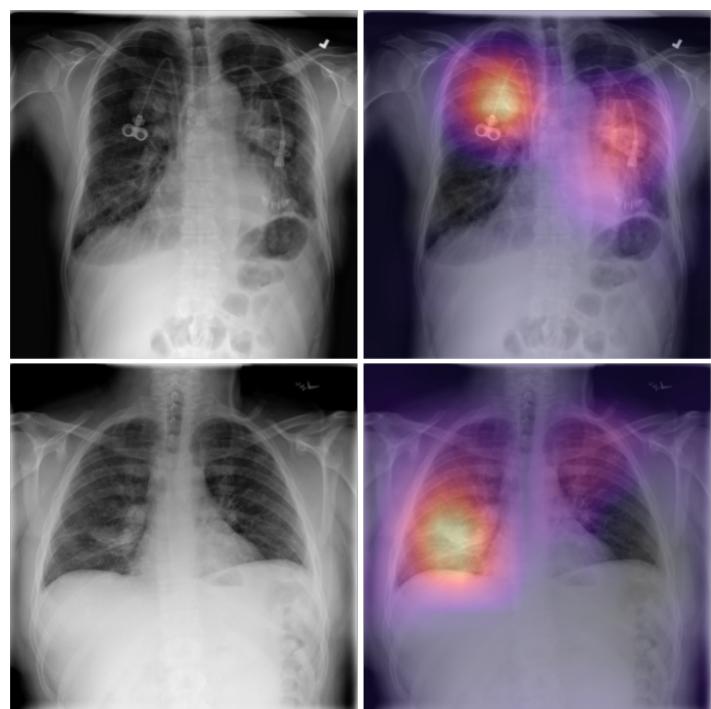


This is Data Release 16.

Data Surveys Instruments



The Sloan Digital Sky Survey: Mapping the Universe



**Software systems for data analytics and  
ML over large and complex datasets are  
now critical for digital applications in  
many domains**

# The Age of “Big Data”/“Data Science”

## The New York Times

SundayReview | NEWS ANALYSIS

The Age **Forbes** / Entrepreneurs

By STEVE LOHR F

MAR 25, 2015 @ 7:33 PM 4,407 VIEWS



Email



Share



Tweet



Save

Josh Steimle, CON

For roughly a decade, information about Big Data. The IDC industry will experience by 2018. What this

# Forbes

## Drowning In Big Data - Finding Insight In A Digital DATA Josh Steimle, CON by Thomas H. Davenport and D.J. Patil FROM THE OCTOBER 2012 ISSUE

SUMMARY   SAVE   SHARE   COMMENT 5   TEXT SIZE   PRINT   BUY COPIES \$8.95



## Harvard Business Review

**W**hen Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social experience. As one LinkedIn manager put it, "It was like arriving at a conference reception and realizing you don't know anyone. So you just stand in the corner sipping your drink—

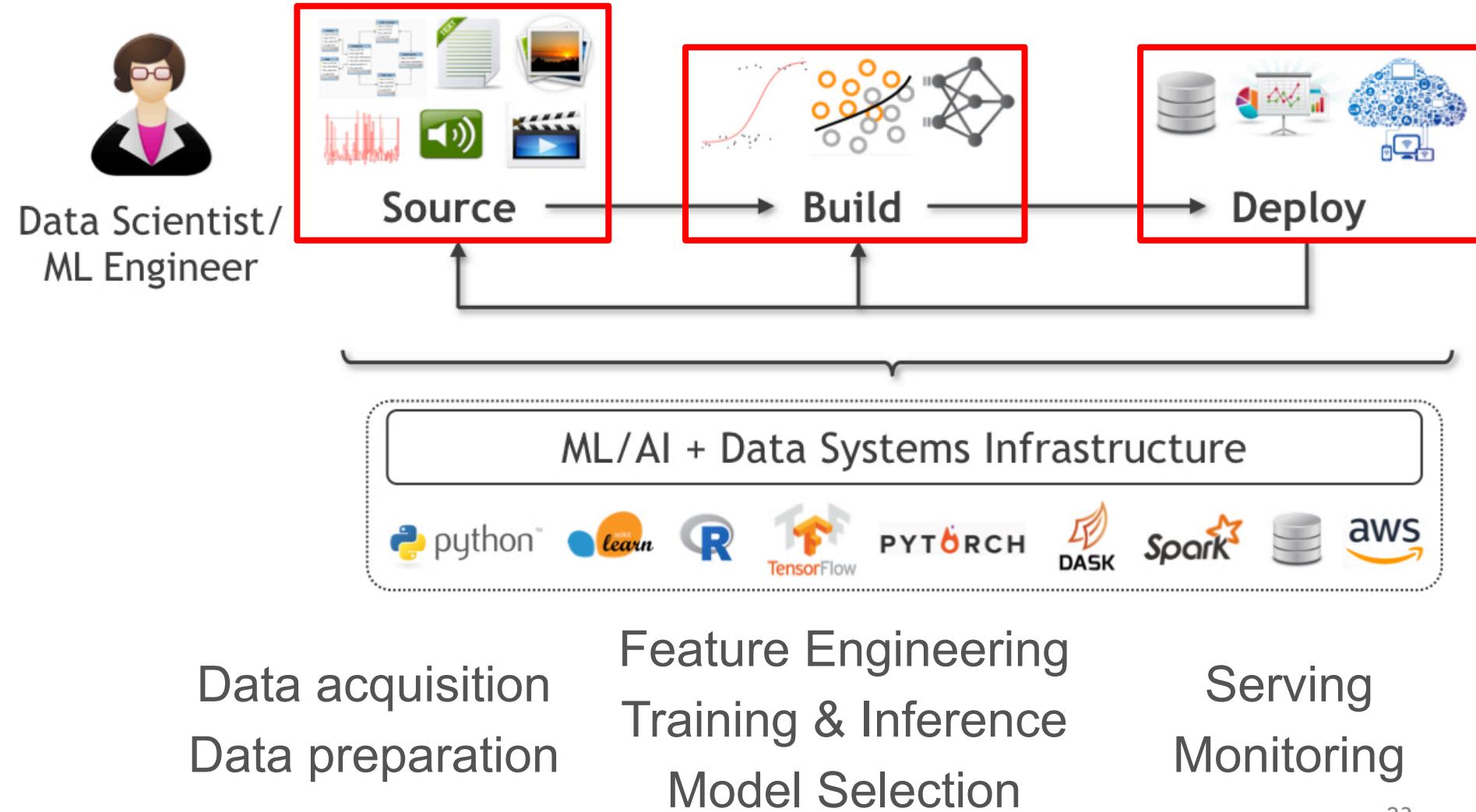
*Data data everywhere,  
All the wallets did shrink!*

*Data data everywhere,  
Nor any moment to think?*

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

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**”:
  - 3.1. **Source**: Data acquisition & preparation for ML
  - 3.2. **Build**: Model selection & deep learning systems
  - 3.3. **Deploying** ML models
4. Hands-on experience with scalable analytics tools

# The Lifecycle of ML-based Analytics



# ML Systems

*Q: What is a Machine Learning (ML) System?*

- ❖ A data processing system (aka *data system*) for mathematically advanced data analysis operations (inferential or predictive):
  - ❖ Statistical analysis; ML, deep learning (DL); data mining (domain-specific applied ML + feature eng.)
  - ❖ *High-level APIs* to express ML computations over (large) datasets
  - ❖ *Execution engine* to run ML computations efficiently

# Categorizing ML Systems

## ❖ Orthogonal Dimensions of Categorization:

1. **Scalability:** In-memory libraries vs Scalable ML system (works on larger-than-memory datasets)
2. **Target Workloads:** General ML library vs Decision tree-oriented vs Deep learning, etc.
3. **Implementation Reuse:** Layered on top of scalable data system vs Custom from-scratch framework

# Major Existing ML Systems

## General ML libraries:

In-memory:



Disk-based files:



Layered on RDBMS/Spark:



Cloud-native:



Azure Machine Learning



Amazon SageMaker

“AutoML” platforms:



DataRobot



Decision tree-oriented:



Microsoft  
LightGBM

Deep learning-oriented:



TensorFlow



# 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)

Usability

Manageability

Developability

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

Q: ~~Q: What are the difficulties in building big data systems for the R&D teams?~~

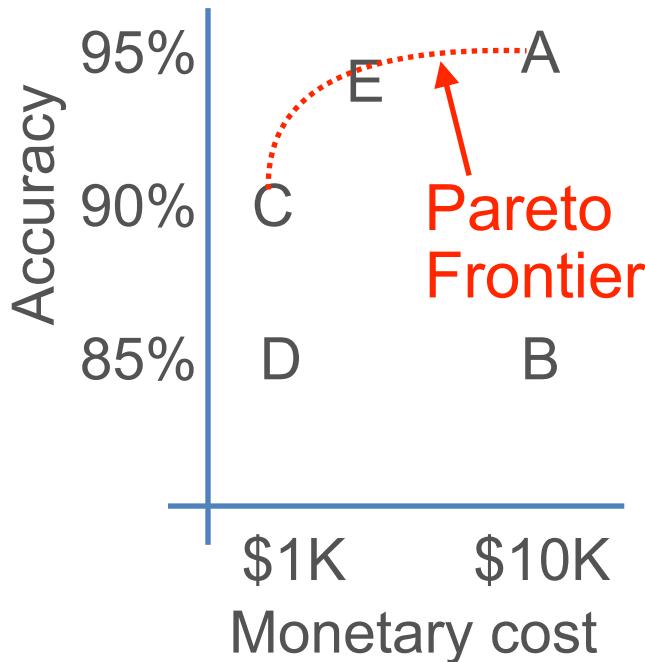
# Conceptual System Stack Analogy

	Relational DB Systems	ML Systems
Theory	First-Order Logic Complexity Theory	Learning Theory Optimization Theory
Program Formalism	Relational Algebra	Matrix 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, cloud computing, and containerization.
- ❖ *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 deep learning inference with TensorFlow and Keras.
- ❖ *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, efficiency, scalability, usability, and total cost.

# What this course is NOT about

- ❖ NOT a course on databases, relational model, or SQL
  - ❖ Take DSC 100 instead (pre-requisite)
- ❖ NOT a course on internal details of RDBMSs
  - ❖ Take CSE 132C instead
- ❖ NOT a training module for how to use Spark
- ❖ NOT a course on ML or data mining *algorithmics*;  
instead, we focus on ML *systems*

Now for the (boring) 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 the instructor with proper justification; a waiver can be considered

[http://cseweb.ucsd.edu/~arunkk/dsc102\\_winter22](http://cseweb.ucsd.edu/~arunkk/dsc102_winter22)

# Components and Grading

- ❖ **3 Programming Assignments: 35% (7% + 14% + 14%)**
  - ❖ No late days! Plan your work well ahead.
- ❖ **Best 5 of 6 Surprise Quizzes: 15%; in-class using iClicker**
- ❖ **Midterm Exam: 15%**
  - ❖ **Wed, Feb 9; in-class (50min)**
- ❖ **Cumulative Final Exam: 35%**
  - ❖ **Mon, Mar 14; 180min**
- ❖ All quizzes and exams are *in-person* only; plan accordingly
- ❖ LMK ahead of time if you need makeup exam slot

# 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
F	Lowest 5%	< 50

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

# Tentative Course Schedule

Week	Topic
1-3	Systems Principles; Basics of Computer Org. and Operating Systems
4	Basics of Cloud Computing
5-6	Parallel and Scalable Data Processing: Parallelism
Scalability Principles	
7-8	Midterm Exam on Wed, Feb 9 Parallel and Scalable Data Processing: Scalability
8	Dataflow Systems
9-10	ML Data Sourcing
10	ML Model Building Systems
10	ML Deployment
11	Final Exam on Wed, Mar 17

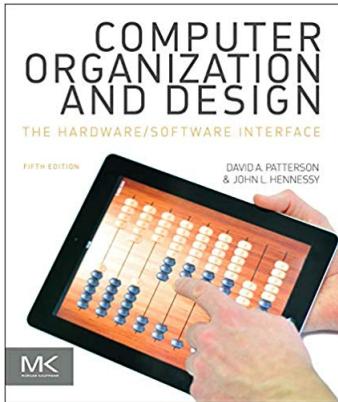
There will likely be 2 industry guest lectures; speakers TBD<sup>6</sup>

# Tentative Schedule for Prog. Assignments

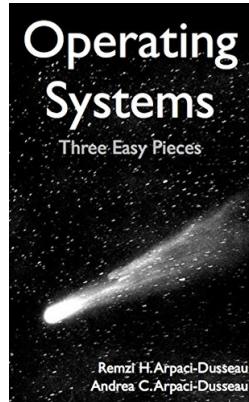
Date	Agenda
Jan 12	PA 0 released
Jan 14	Discussion on PA 0 by TA
Jan 25	<b>PA 0 due</b>
Jan 26	PA 1 released
Jan 28	Discussion on PA 1 by TA
Feb 16	<b>PA 1 due</b>
Feb 17	PA 2 released
Feb 18	Discussion on PA 2 by TA
Mar 8	<b>PA 2 due</b>

**No late days!** Plan your work upfront!  
Try to join the Discussion slot talks by the TAs.

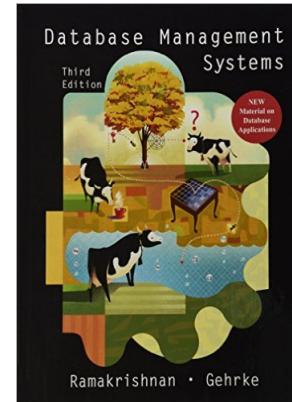
# Suggested Textbooks



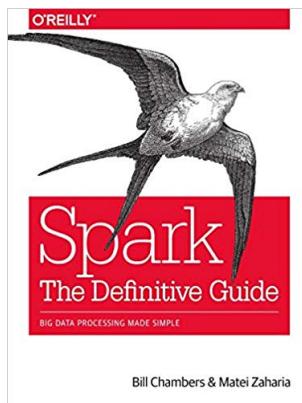
Aka “CompOrg Book”



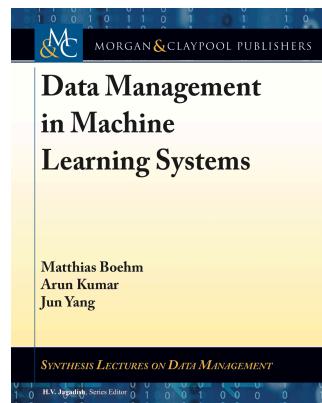
Aka “Comet Book”



Aka “Cow Book”



Aka “Spark Book”

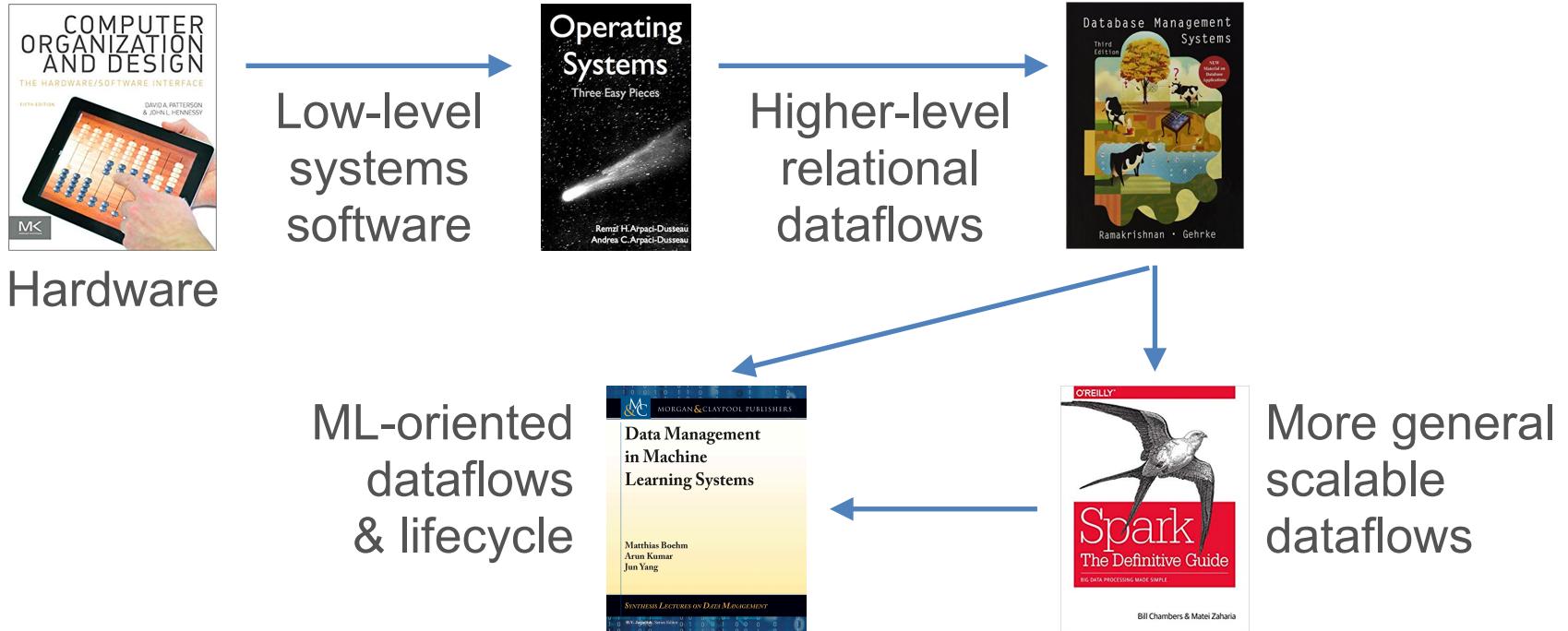


Aka “MLSys Book”

(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!

# Course Administrivia

- ❖ **Lectures:** MWF 1:00pm-1:50pm PT
  - ❖ Virtual-only for weeks 1 and 2; in-person only afterward
  - ❖ Will take live Q&A throughout
- ❖ **Instructor:** Arun Kumar; arunkk [at] eng.ucsd.edu
  - ❖ Office hours: Wed 2:00-3:00pm PT
- ❖ See **Piazza** (or Canvas) for all Zoom meeting links, logistical announcements; see Canvas for gradebook
- ❖ **TAs:** Umesh Singla, Vignesh Nandakumar, and Pradyumna Sridhara (see webpage for details)

[http://cseweb.ucsd.edu/~arunkk/dsc102\\_winter22](http://cseweb.ucsd.edu/~arunkk/dsc102_winter22)

# General Dos and Do NOTs

## ***Do:***

- ❖ Follow all announcements on Piazza/Canvas
- ❖ Try to join the lectures/discussions live
- ❖ View/review videos asynchronously by yourself
- ❖ Participate in discussions in class / on Piazza
- ❖ Raise your hand before speaking
- ❖ Use “DSC102:” as subject prefix for all emails to me/TA

## ***Do NOT:***

- ❖ Record lectures on your side without permission of instructor and other students
- ❖ Harass, intimidate, or intentionally talk over others
- ❖ Violate academic integrity on the graded quizzes, exams, or programming assignments; I am *very strict* on this matter!