

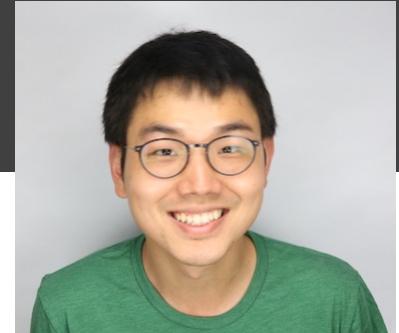
DSC 102

Systems for Scalable Analytics

Spring 2024

Haojian Jin

About Myself



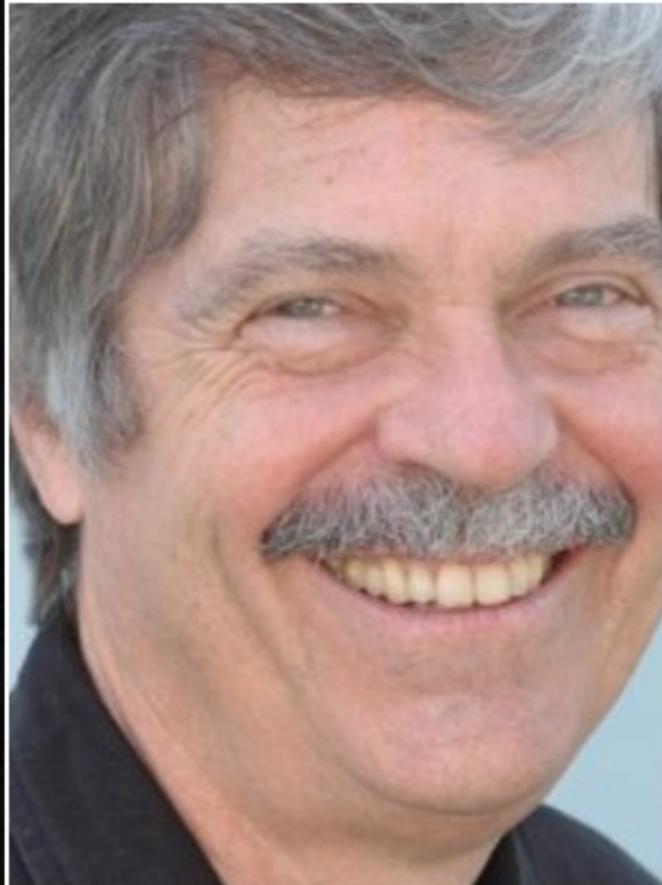
Haojian Jin (<http://haojianj.in/>)
Asst. Prof @ UCSD-HDSI

Data Smith Lab:

*We study the **security** and **privacy** of data systems by researching the **people** who design, implement, and use these systems.*

Ph.D. from CMU Human-Computer Interaction Institute
Ph.D. Thesis: Modular Privacy Flow

Before Ph.D.: worked at Yahoo Research, ran a startup
HCI, Software Engineering, Mobile Computing, AI.

A close-up portrait of Alan Kay, an elderly man with grey hair and a prominent grey mustache, smiling warmly at the camera.

The best way to predict the future is
to invent it.

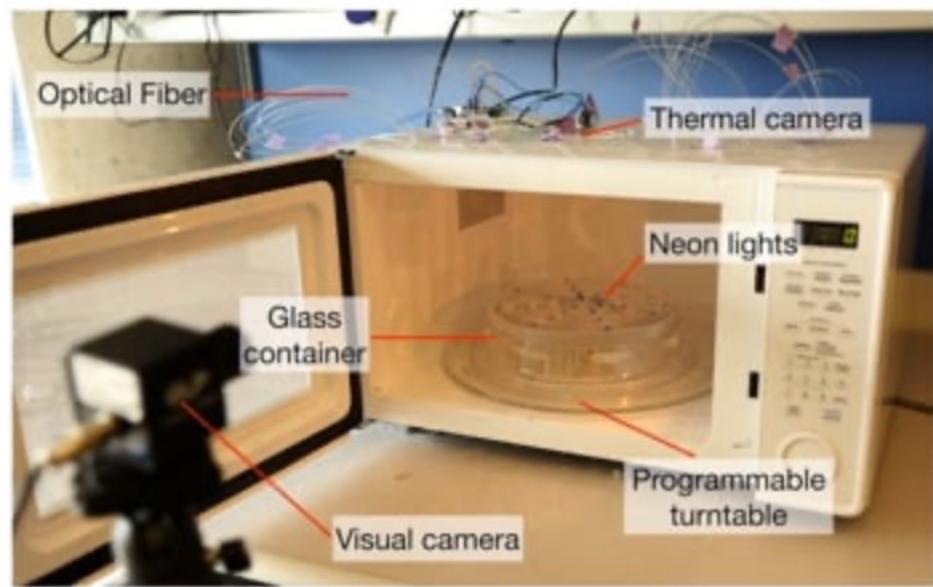
— *Alan Kay* —

Working Code Trumps All Hype!

Software Defined Cooking (SDC) using a microwave oven

Haojian Jin
Jingxian Wang
Swarun Kumar
Jason Hong

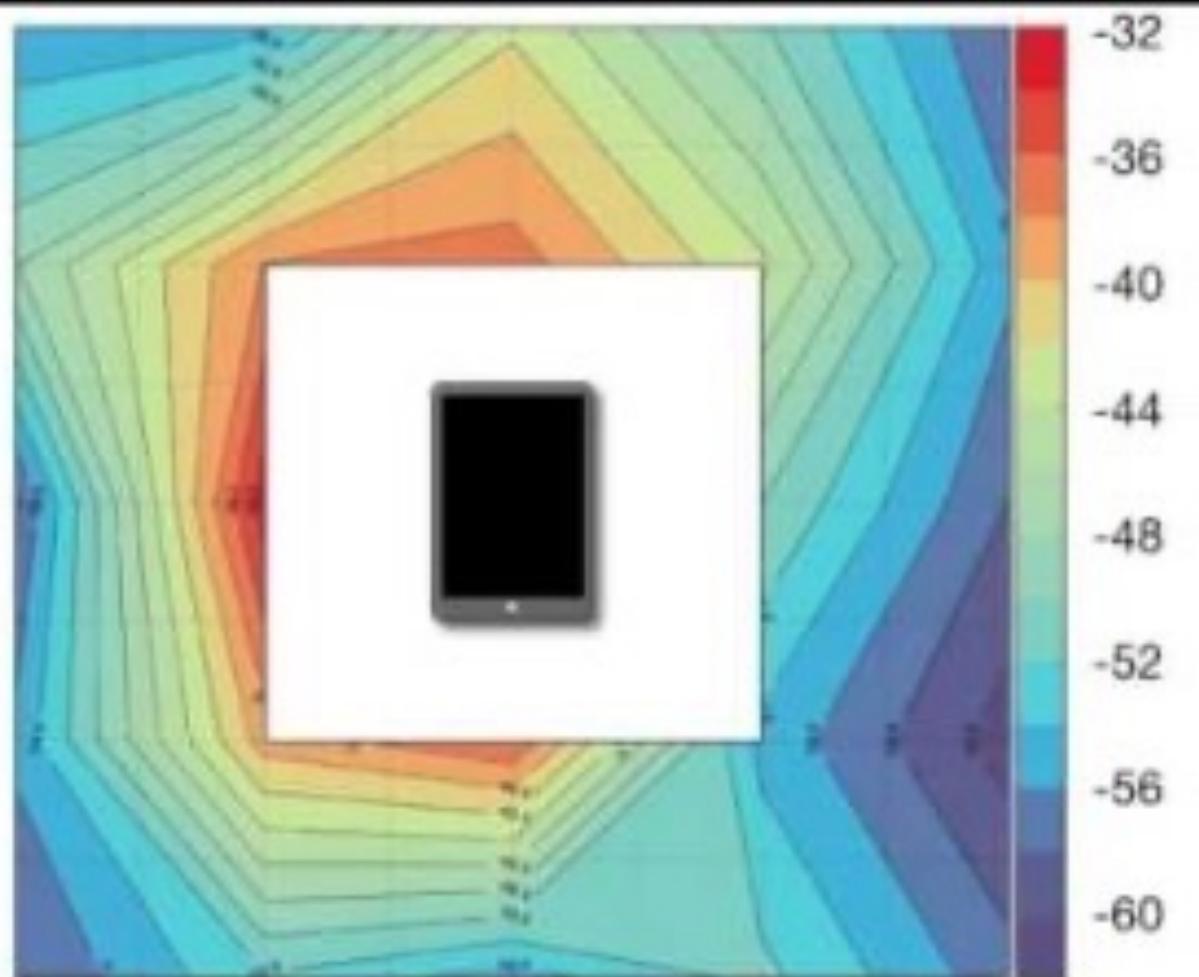
**Carnegie
Mellon
University**



Working Code Trumps All Hype!



Working Code Trumps All Hype!



And mapped the asymmetric RSSI distribution around the iPad.

Working Code Trumps All Hype!

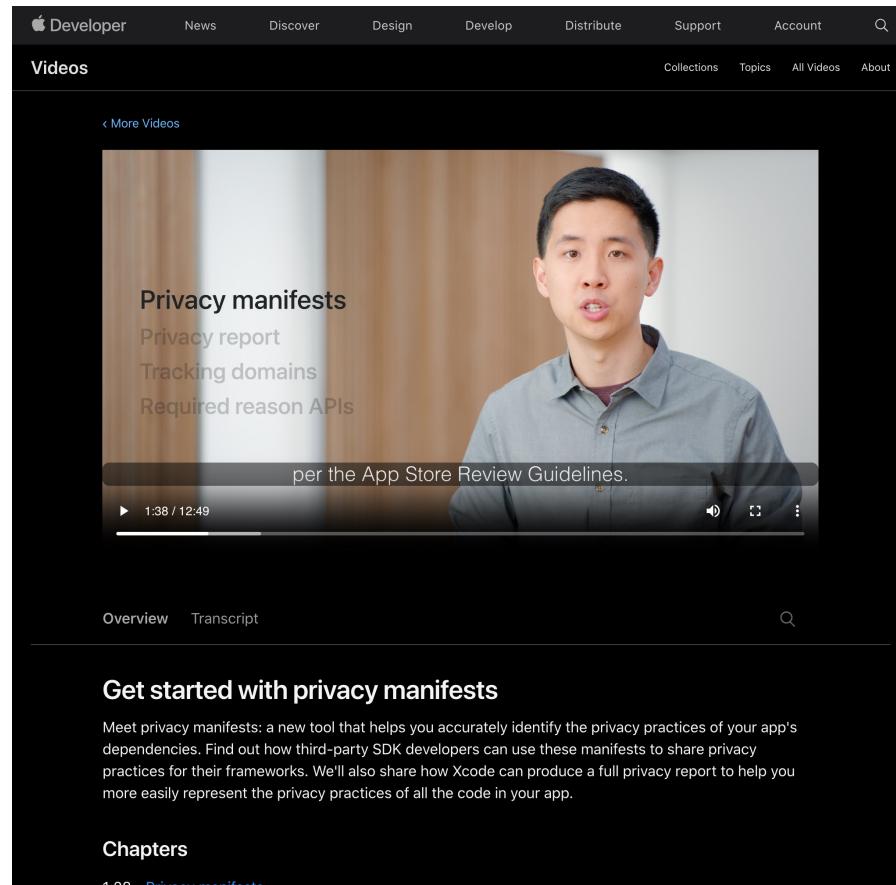
**“Uber” Would Like to Use
Your Location.**

Uber picks you up exactly where you are.
To start riding, choose “Allow” so the
app can find your location.

Don't Allow

OK

Some of them are getting adopted.



The screenshot shows the Apple Developer website's video player interface. At the top, there's a navigation bar with links for Apple Developer, News, Discover, Design, Develop, Distribute, Support, Account, and a search icon. Below that is a sub-navigation bar for 'Videos' with links for Collections, Topics, All Videos, and About. The main content area features a video player with a play button, a progress bar showing 1:38 / 12:49, and a video thumbnail of a man speaking. To the left of the video, there's a sidebar with text: 'Privacy manifests', 'Privacy report', 'Tracking domains', and 'Required reason APIs'. Below the video player, there's a transcript section with 'Overview' and 'Transcript' buttons, and a search icon. At the bottom, there's a link to '1:28 - Privacy manifests'.

- ✓ Categorized Purpose string (2017 -> 2022)
- ✓ Declared manifests (2020->2024)
- ✓ Operator-based API (??)

What is this course about? Why take it?



Reddit · r/UCSD

10+ comments · 2 years ago · :

DSC 102 in a nutshell. : r/UCSD



IVEBEENGRAPED · 3y ago

This class was honestly the most useful class I took at UCSD. I'm in my first job out of school, and half of what I do here is messing around with AWS and Spark like we did in that class. Would highly recommend, even to non-DS majors.



19



Reply

Share

...



atvrider512 · 3y ago

yeah lol this was me, I feel like this material is sooooo useful but the class was so disorganized I didn't get to fully learn and process it

2



Reply

Share

...



phatfat · 3y ago

this class hurt me

2



Reply

Share

...

https://www.reddit.com/r/UCSD/comments/npqcdm/dsc_102_in_a_nutshell/



statistician

Location



Statistician Salaries United States ▾

Overview

Salaries

Interviews

Insights

Career Path

How much does a Statistician make?

Updated Jan 4, 2022

Industry

Employer Size

Experience

All industries

All company sizes

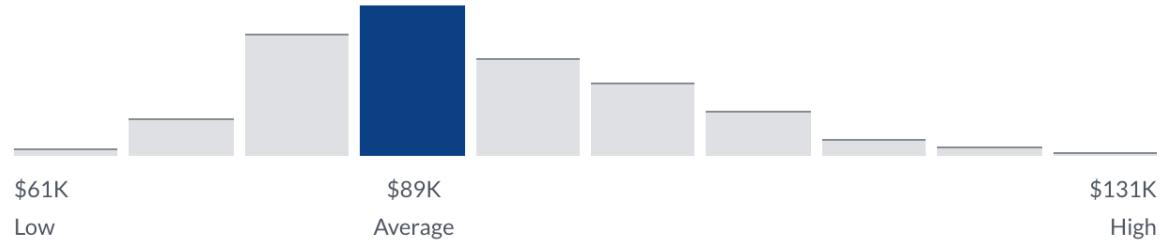
All years of Experience

Very High Confidence

\$88,989 /yr

Average Base Pay

2,398 salaries





Data Scientist Salaries United States ▾

Overview

Salaries

Interviews

Insights

Career Path

How much does a Data Scientist make?

Updated Jan 4, 2022

Industry

▼

Employer Size

▼

Experience

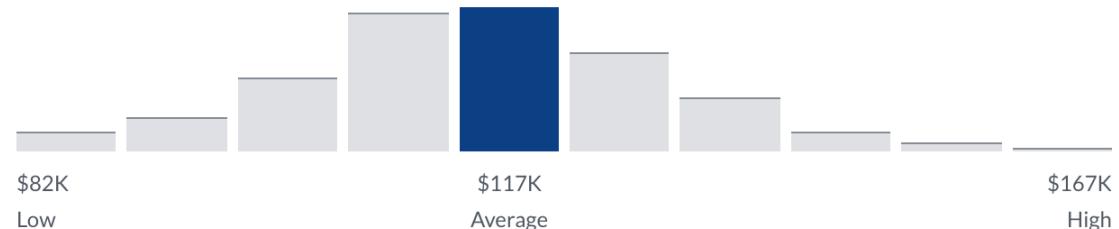
▼**i** To filter salaries for Data Scientist, [Sign In](#) or [Register](#).

Very High Confidence

\$117,212 /yr

Average Base Pay

18,354 salaries

**— 88,989****= 28,223!**

Professional Vision

DISCURSIVE PRACTICES are used by members of a profession to shape events in the domains subject to their professional scrutiny. The shaping process creates the objects of knowledge that become the insignia of a profession's craft: the theories, artifacts, and bodies of expertise that distinguish it from other professions. Analysis of the methods used by members of a community to build and contest the events that structure their lifeworld contributes to the development of a practice-based theory of knowledge and action.¹ In this article, I examine two contexts of professional activity: archaeological field excavation and legal argumentation. In each of these contexts, I investigate three practices: (1) *coding*, which transforms phenomena observed in a specific setting into the objects of knowledge that animate the discourse of a profession; (2) *highlighting*, which makes specific phenomena in a complex perceptual field salient by marking them in some fashion; and (3) *producing and articulating material representations*. By applying such practices to phenomena in the domain of scrutiny, participants build and contest *professional vision*, which consists of socially ordered seeing and understanding events that are answerable to the distinctive needs of a particular social group.

Coding
Highlighting
Reconstructing



What will happen if I click the second link?



mat

Yoga Web Search

Search

Refine results for mat:

[Asanas \(poses\)](#)

[Anatomy](#)

[Classes](#)

[Teachers](#)

[Yoga Mats :: The Yoga.com Store](#)

Yoga.com's yoga **mat**s and related products are great for living a life in balance. Buy our yoga **mat**s or browse any of yoga.com's related products.

www.yoga.com/store/supercategory.asp?Category_ID=295 - 38k - [Cached](#)

Labeled [Stores](#)



[Barefoot Yoga Co. || Yoga Mats » Yoga Clothing » Yoga Supplies](#)

Quality yoga products at great prices - yoga **mat**s, yoga clothing, yoga **mat** bags, yoga props, yoga rugs, yoga videos, books and more.

www.barefootyoga.com/ - 17k - [Cached](#)

[Yoga :: Yoga Books :: Yoga Videos :: Yoga Mats :: Yoga Clothing](#)

For all your yoga videos, books, clothing and **mat**s, as well as pilates equipment, Yoga.com is the best source on the web for fitness and wellness products.

www.yoga.com/ - 31k - [Cached](#)

Labeled [Stores](#)

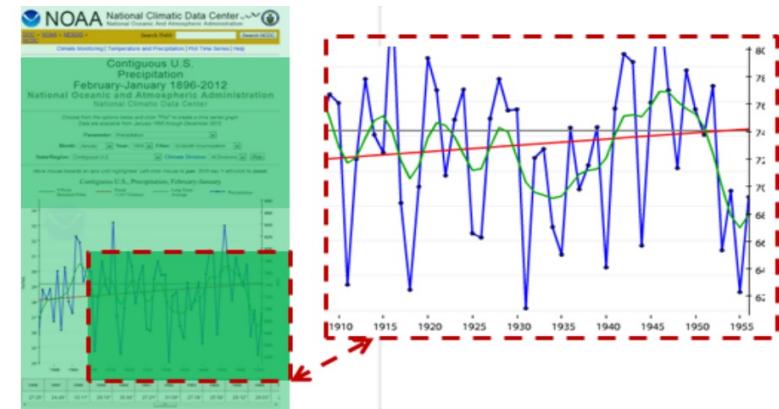
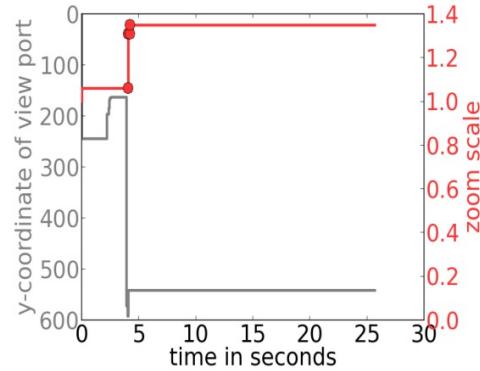
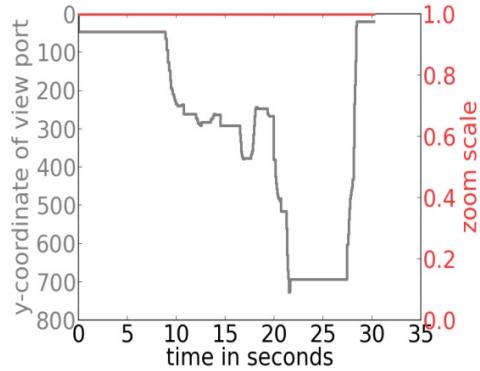
[lululemon athletica | products](#)

We recommend cleaning your **mat** with natural based disinfectants such as warm ...

Remember to wipe off your **mat** with a dry towel before you roll it up. ...

www.lululemon.com/products/ accessories/yoga/eco_ultra_mat - 11k - [Cached](#)

Labeled [Stores](#)



Mining Touch Interaction Data on Mobile Devices to Predict Web Search Result Relevance, SIGIR 2013

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



TV Schedule



FALL TV GRID 2022

New series are listed in **RED**

SUN	8:00	8:30	9:00	9:30	10:00	10:30
ABC	7 pm AFV CELEBRITY JEOPARDY		CELEBRITY WHEEL OF FORTUNE		THE ROOKIE	
CBS	7 pm 60 MIN THE EQUALIZER		EAST NEW YORK		NCIS: LOS ANGELES	
FOX	To be announced					
NBC	SUNDAY NIGHT FOOTBALL					
CW	FAMILY LAW (acquired)		CORONER (acquired)			
MON	8:00	8:30	9:00	9:30	10:00	10:30
ABC	BACHELOR IN PARADISE				THE GOOD DOCTOR	
CBS	THE NEIGHBORHOOD	BOB HEARTS ABISHOLA	NCIS		NCIS: HAWAII	
FOX	To be announced					
NBC	THE VOICE				QUANTUM LEAP	
CW	ALL AMERICAN		ALL AMERICAN: HOMECOMING			
TUE	8:00	8:30	9:00	9:30	10:00	10:30
ABC	BACHELOR IN PARADISE				THE ROOKIE: FEDS	
CBS	FBI		FBI: INTERNATIONAL		FBI: MOST WANTED	
FOX	To be announced					
NBC	THE VOICE		LA BREA		NEW AMSTERDAM	
CW	THE WINCHESTERS		PROFESSIONALS (acquired)			
WED	8:00	8:30	9:00	9:30	10:00	10:30
ABC	THE CONNERS	THE GOLDBERGS	ABBOTT ELEMENTARY	HOME ECONOMICS	BIG SKY	
CBS	SURVIVOR		THE AMAZING RACE		THE REAL LOVE BOAT	
FOX	To be announced					
NBC	CHICAGO MED		CHICAGO FIRE		CHICAGO P.D.	
CW	DC'S STARGIRL		KUNG FU			
THU	8:00	8:30	9:00	9:30	10:00	10:30
ABC	STATION 19		GREY'S ANATOMY		ALASKA	
CBS	YOUNG SHELDON	GHOSTS	SO HELP ME TODD		CSI: VEGAS	
FOX	To be announced					
NBC	LAW & ORDER		LAW & ORDER: SVU		LAW & ORDER: ORGANIZED CRIME	
CW	WALKER		WALKER INDEPENDENCE			
FRI	8:00	8:30	9:00	9:30	10:00	10:30
ABC	SHARK TANK		20/20			
CBS	S.W.A.T.		FIRE COUNTRY		BLUE BLOODS	
FOX	FRIDAY NIGHT SMACKDOWN (presumably!)					
NBC	COLLEGE BOWL (until November) LOPEZ VS LOPEZ YOUNG ROCK		DATELINE NBC			
CW	PENN & TELLER: FOOL US		WHOSE LINE IS IT ANYWAY? X2			
SAT	8:00	8:30	9:00	9:30	10:00	10:30
ABC	COLLEGE FOOTBALL					
CBS	Drama encores		Drama encores		48 HOURS	
FOX	To be announced					
NBC	Drama encores		DATELINE WEEKEND MYSTERY		SNL VINTAGE	11:30 SNL
CW	MAGIC WITH THE STARS		WORLD'S FUNNIEST ANIMALS x2			

How does Netflix know that?

Large datasets + Machine learning!

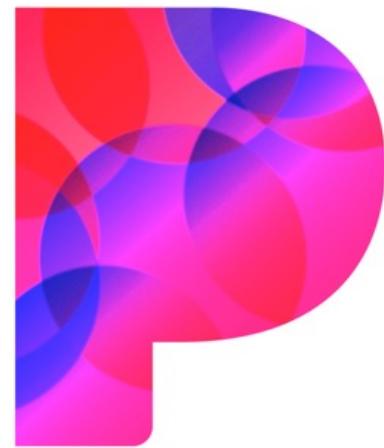


Log all user behavior (views, clicks, pauses, searches, etc.)
Recommender systems apply ML to TBs of data from all users and movies to deliver a tailored experience

Pandora v.s. Spotify

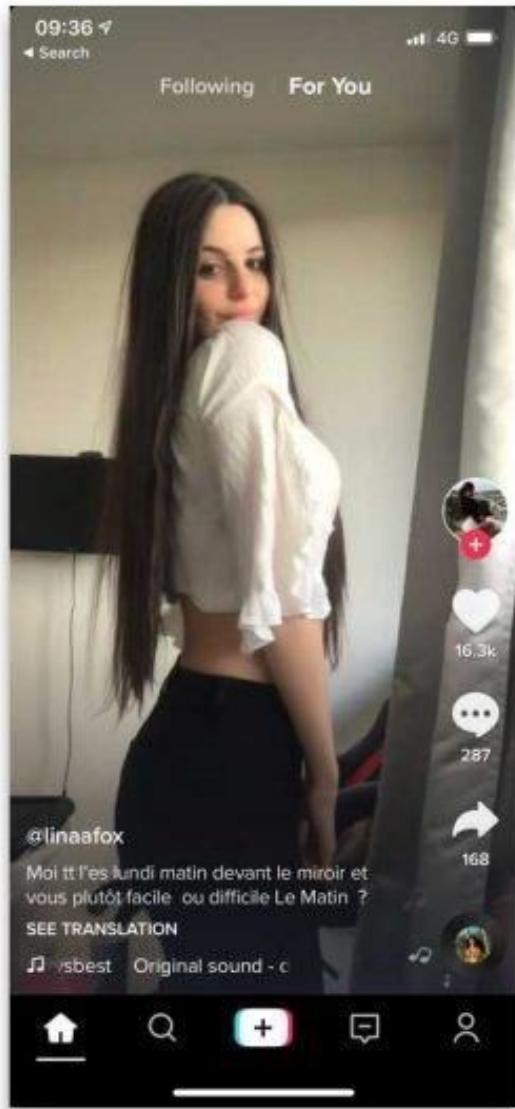


Music Genome Project

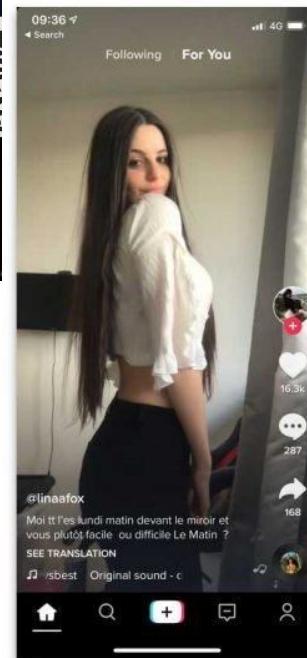


Collaborative filters (social)

Tiktok



From Netflix to TikTok: More than algorithms



2. Structured data with search results

Google X 🔍 ⚙️

All Images Books News Videos More Tools

About 13,000,000 results (0.52 seconds)

 **Alan Turing**
Mathematician ⋮

Overview Education Books Videos

     More images

[https://en.wikipedia.org › wiki › Alan_Turing](https://en.wikipedia.org/wiki/Alan_Turing) ⋮

Alan Turing - Wikipedia

Alan Mathison Turing OBE FRS (/'tjøərɪŋ/; 23 June 1912 – 7 June 1954) was an English mathematician, computer scientist, logician, cryptanalyst, philosopher, ...

Partner(s): [Joan Clarke](#); (engaged in 194... Known for: [Cryptanalysis of the Enigm...](#)

Awards: Smith's Prize (1936) Resting place: Ashes scattered in gard...

[The Enigma](#) · [Alan Turing law](#) · [Legacy of Alan Turing](#) · [Alan Turing Year](#)

About

Alan Mathison Turing OBE FRS was an English mathematician, computer scientist, logician, cryptanalyst, philosopher, and theoretical biologist.

[Wikipedia](#)

Born: June 23, 1912, [Maida Vale, London, United Kingdom](#)

Died: June 7, 1954, [Wilmslow, United Kingdom](#)

Academic advisor: [Alonzo Church](#)

Education: Princeton University (1936–1938), [MORE](#)

Influenced by: [Alonzo Church](#), [Kurt Gödel](#), [Ludwig Wittgenstein](#), [Max Newman](#)

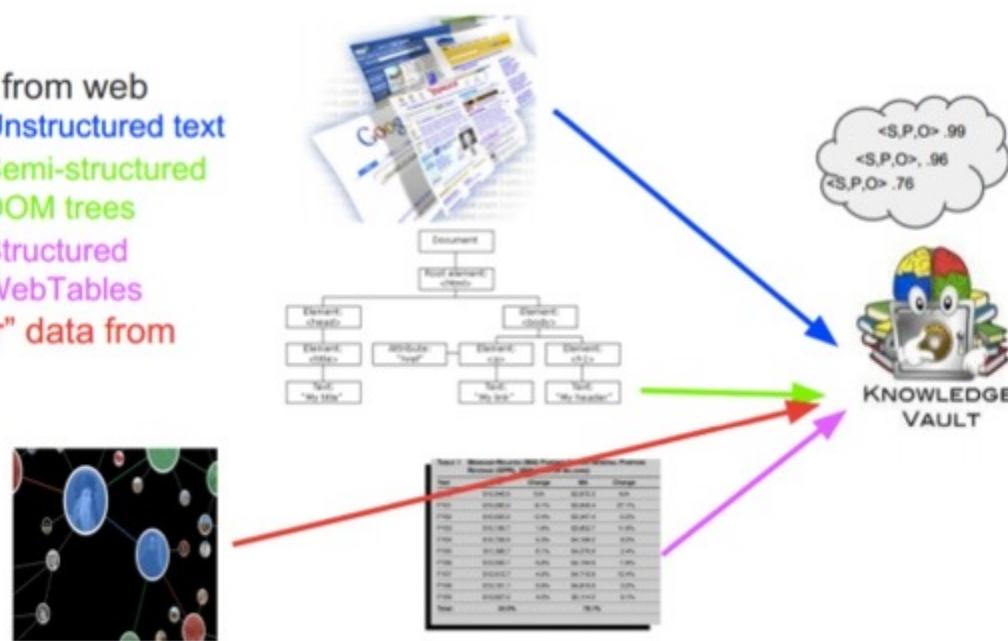
Notable students: [Robin Gandy](#), [Beatrice Worsley](#)

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

Innumerable “enterprise” applications

How Target Figured Out A Teen Girl Was Pregnant Before Her Father Did

se



307 comments, 167 called-out

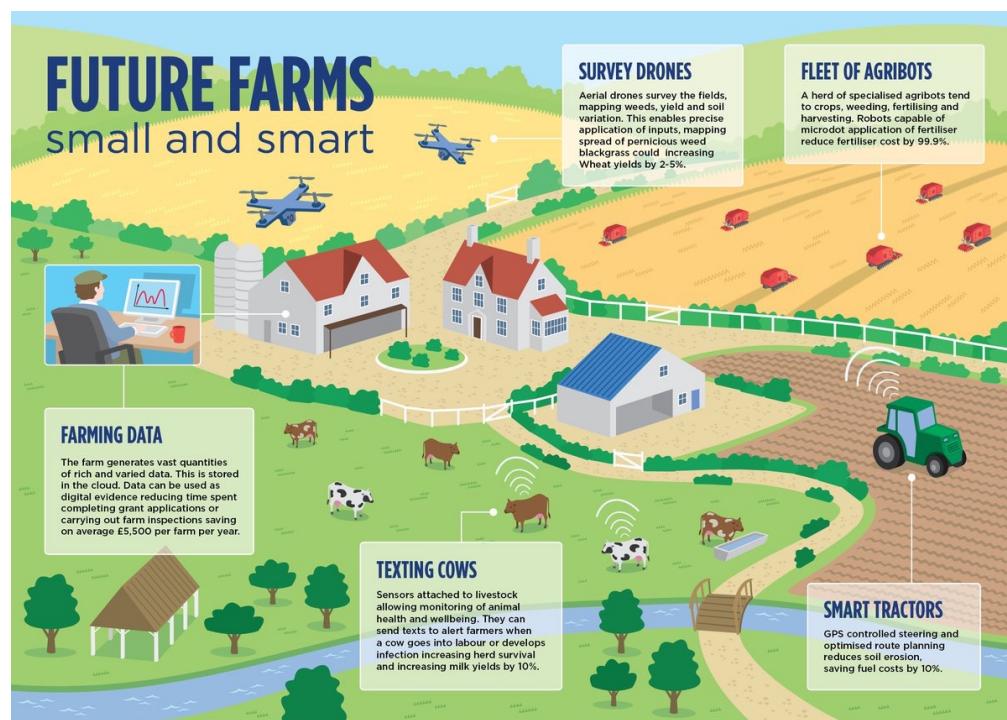
+ Comment Now + Follow Comments

Every time you go shopping, you share intimate details about your consumption patterns with retailers. And many of those retailers are studying those details to figure out what you like, what you need, and which coupons are most likely to make you happy. Target, for example, has figured out how to data-mine its way into your womb, to figure out whether you have a baby on the way long before you need to start buying diapers.

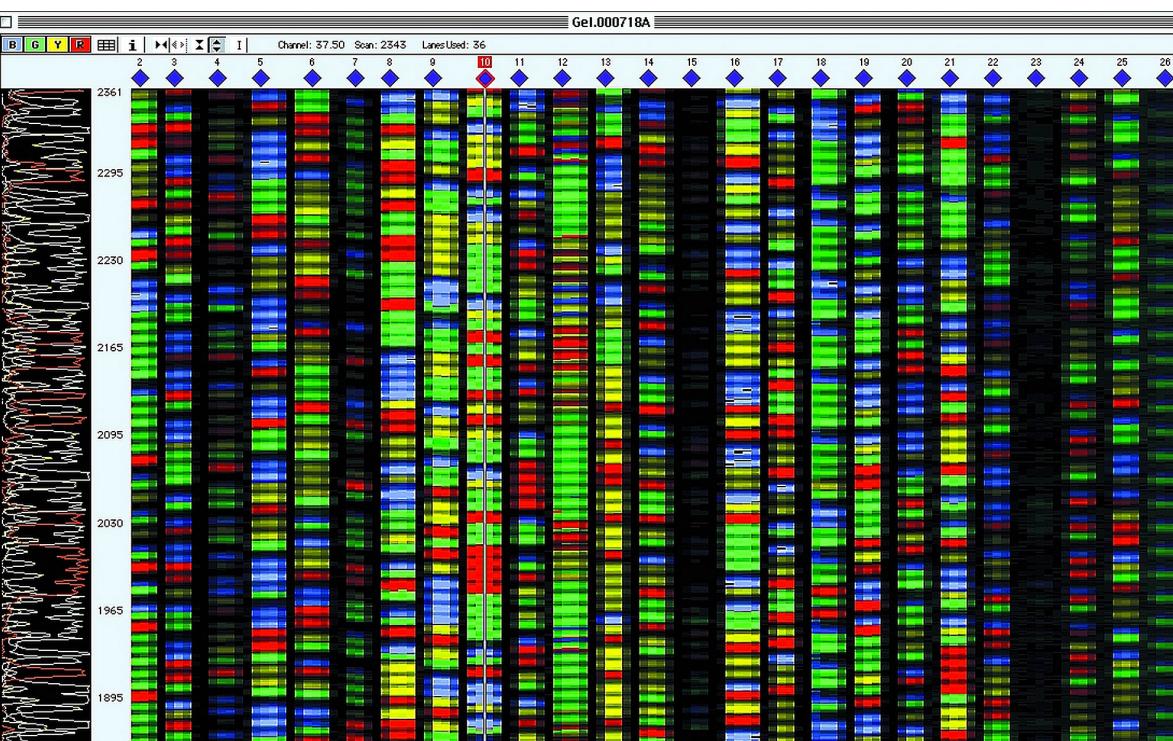
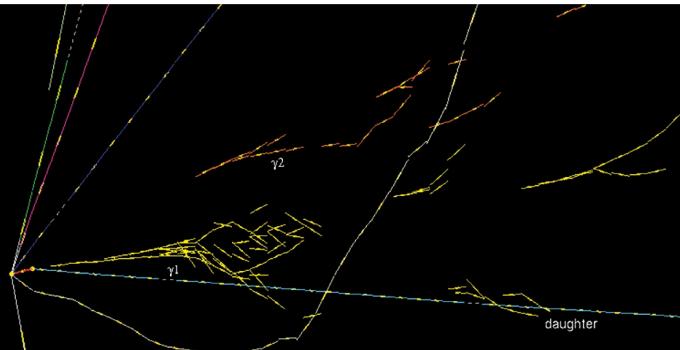
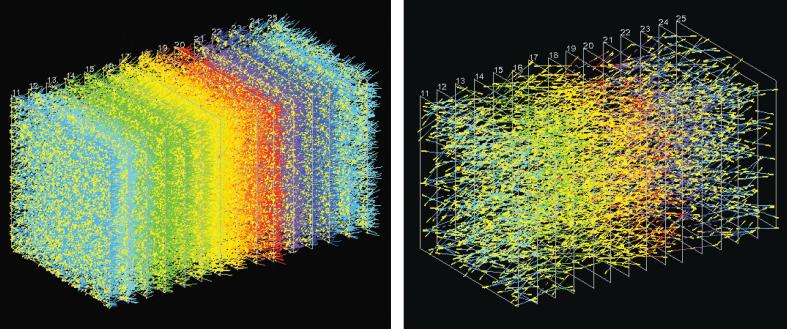
Charles Duhigg outlines in the New York Times how Target tries to hook parents-to-be at that crucial moment before they turn into



Target has got you in its aim



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

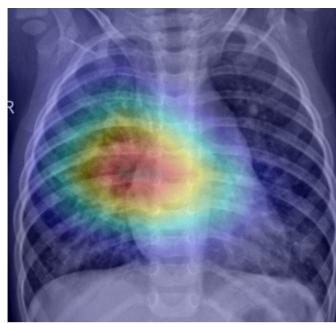


This is Data Release 16.

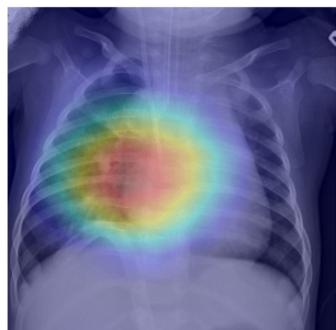
Data Surveys Instruments



(a)



(b)



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

Forbes

Drowning In Big Data - Finding Insight In A Digital DATA Josh Steimle, CON Digital Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

[SUMMARY](#) [SAVE](#) [SHARE](#) [COMMENT](#) [TEXT SIZE](#) [PRINT](#) [\\$8.95](#)

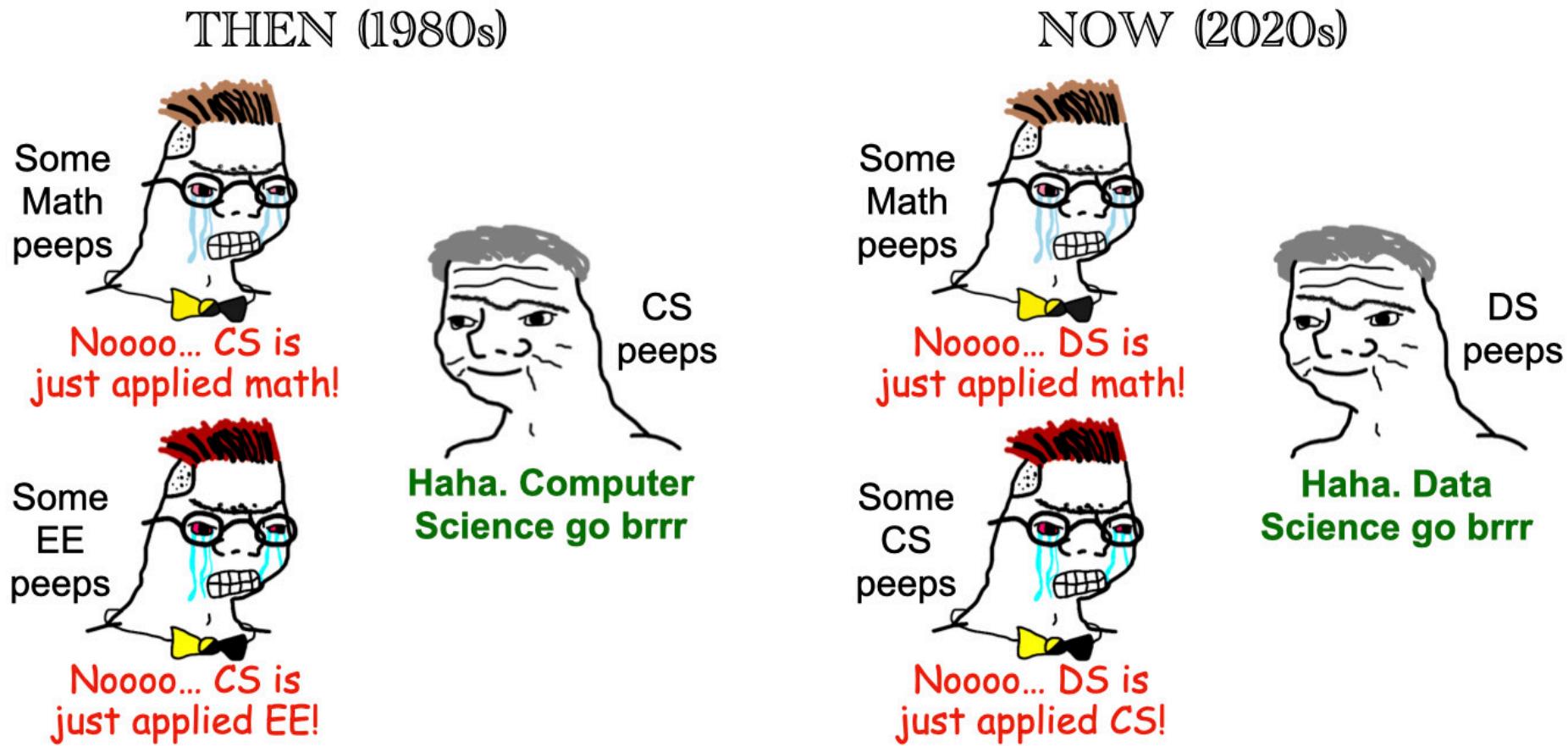


Harvard Business Review

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

When 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—

Meme from Previous DSC 102



The Discipline of Machine Learning

Tom M. Mitchell

July 2006

CMU-ML-06-108

School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213

*Machine Learning Department

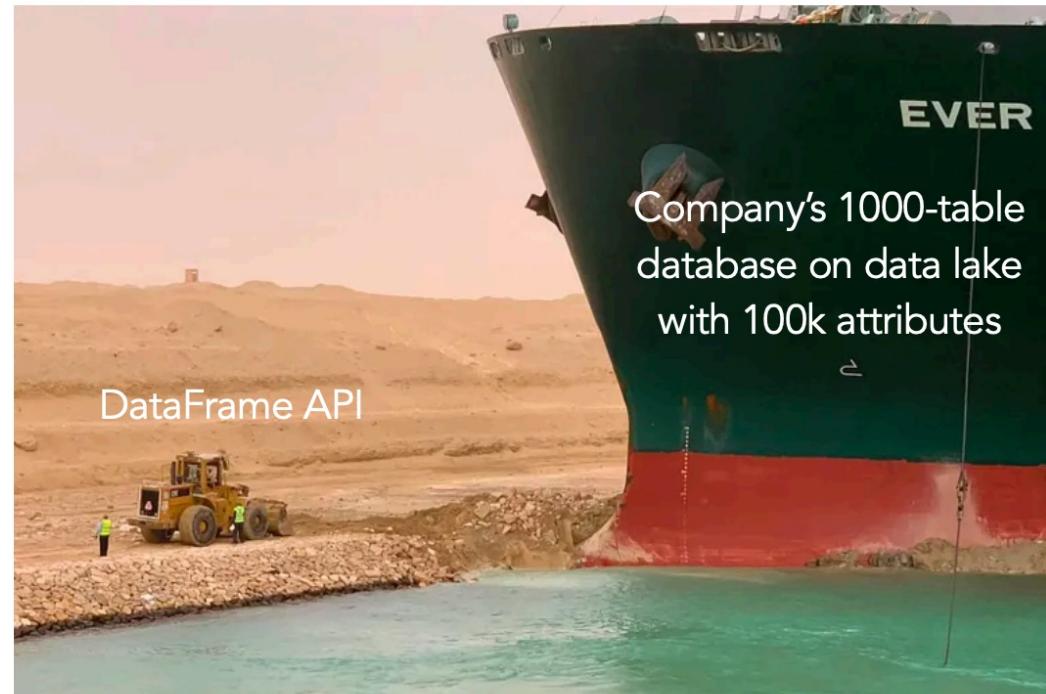
†School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA

Abstract

Over the past 50 years the study of Machine Learning has grown from the efforts of a handful of computer engineers exploring whether computers could learn to play games, and a field of Statistics that largely ignored computational considerations, to a broad discipline that has produced fundamental statistical-computational theories of learning processes, has designed learning algorithms that are routinely used in commercial systems for speech recognition, computer vision, and a variety of other tasks, and has spun off an industry in data mining to discover hidden regularities in the growing volumes of online data. This document provides a brief and personal view of the discipline that has emerged as Machine Learning, the fundamental questions it addresses, its relationship to other sciences and society, and where it might be headed.

DSC 204a Scalable Data Systems

- Haojian Jin



Vision

Data science professionals ought to be familiarized with data systems from a user's standpoint, as opposed to the conventional approach of a system implementer.

15-213/15-513/14-513 Introduction to Computer Systems (ICS)

Fall 2023

- 15-213 Pittsburgh: Tue, Thu 12:30 PM–01:50 PM, GHC 4401, [Brian Railing](#) and [Phillip Gibbons](#)
- 14-513 Pittsburgh: Tue, Thu 12:30 PM–01:50 PM, CIC 1202, [David Varodayan](#)

12 units

The ICS course provides a programmer's view of how computer systems execute programs, store information, and communicate. It enables students to become more effective programmers, especially in dealing with issues of performance, portability and robustness. It also serves as a foundation for courses on compilers, networks, operating systems, and computer architecture, where a deeper understanding of systems-level issues is required. Topics covered include: machine-level code and its generation by optimizing compilers, performance evaluation and optimization, computer arithmetic, memory organization and management, networking technology and protocols, and supporting concurrent computation.

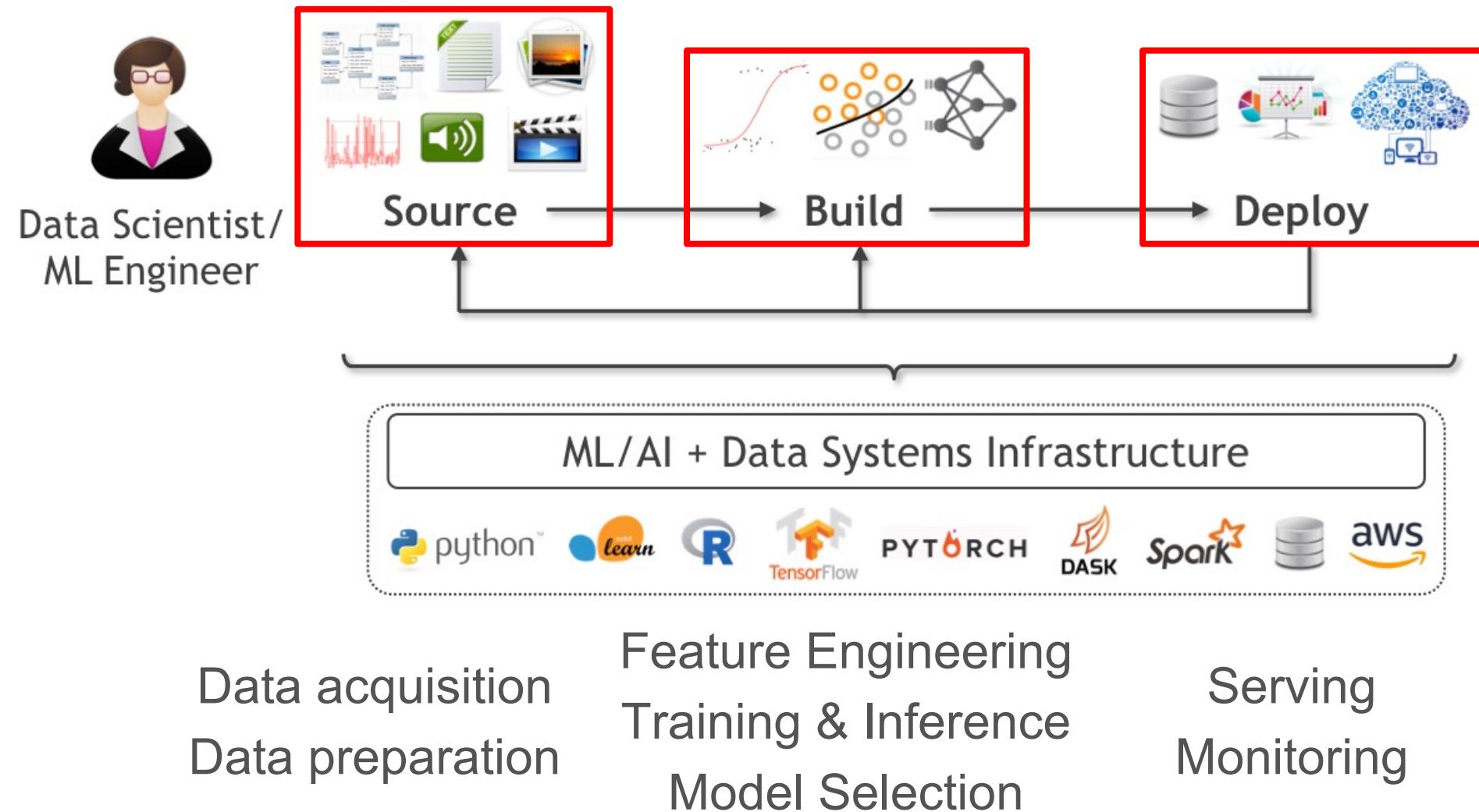
[Course Syllabus](#)

Prerequisites: 15-122

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

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 v. Scalable ML system (works on larger-than-memory datasets)
- 2. Target Workloads:** General ML library v. Decision tree-oriented v. Deep learning, etc.
- 3. Implementation Reuse:** Layered on top of scalable data system v. 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 if I didn't have the discipline to take my ideas from the PPT to code?

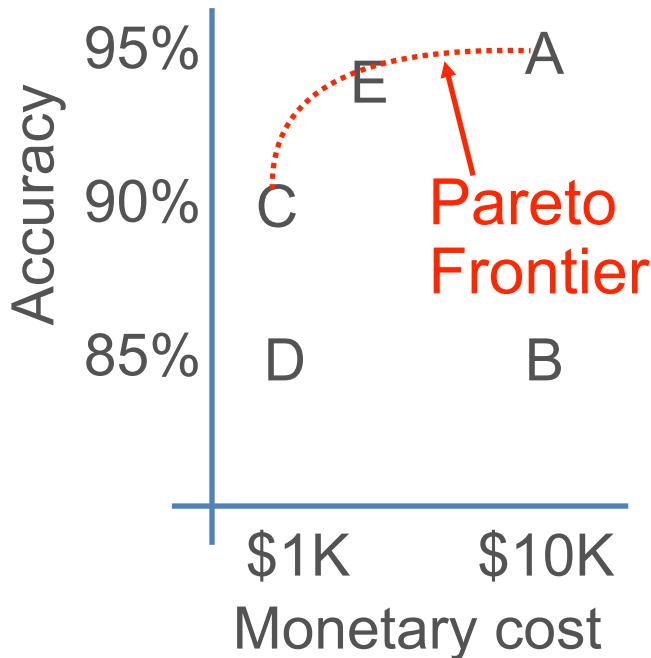
Conceptual System Stack Analogy

	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.

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*