



# Computational Structures in Data Science



UC Berkeley EECS  
Lecturer  
Michael Ball



UC Berkeley EECS  
Adj. Asst. Prof.  
Dr. Gerald Friedland

January 24, 2020

## Lecture #1: Welcome to CS88!



<http://cs88.org>



# CS88 Team - Michael

---

- **Michael Ball**
  - [ball@Berkeley.edu](mailto:ball@Berkeley.edu) – You're best off by using Piazza! ☺
  - 625 Soda Hall
  - <http://michaelball.co> – I don't update this much...
    - » It was great procrastination when I was a CS student.
  - Office hours: TBD @ 625 Soda
  - A few minutes after class
- **Things I do:**
  - Intro CS Research
    - » Tools, curriculum
  - Training TAs
  - Building Educational Software (Gradescope)
  - Tools for web accessibility



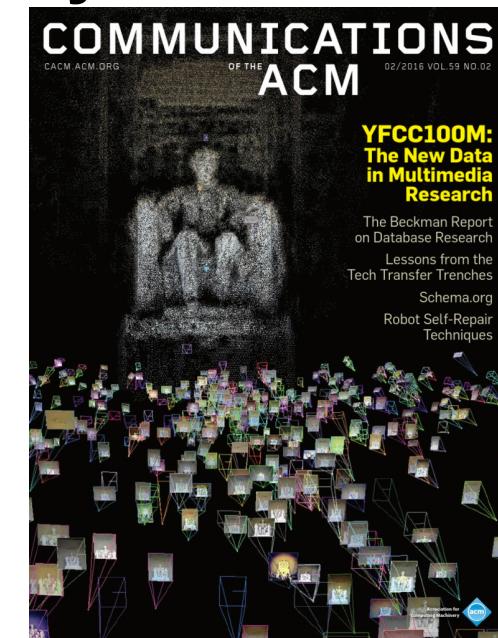


# CS88 Team - Gerald

- Dr. Gerald Friedland
  - [fractor@Berkeley.edu](mailto:fractor@Berkeley.edu) – You're best off by using Piazza! ☺
- Projects you might want to check out:  
<http://mmcommons.org>
  - Work with 100M images, 1M videos in your own Amazon instance.

<http://www.teachingprivacy.org>

Creating teaching materials informing about data over sharing.



 <Teaching Privacy>



# CS88 — Our Awesome Team!

## Head Teaching Assistants



Alex Kassil

Email: [alek.kassil@berkeley.edu](mailto:alek.kassil@berkeley.edu)



Brian Mi

Email: [bmi@berkeley.edu](mailto:bmi@berkeley.edu)



Julia Yu

Email: [juliayu@berkeley.edu](mailto:juliayu@berkeley.edu)

## Teaching Assistants



Alec Kan

Email: [alec.kan@berkeley.edu](mailto:alec.kan@berkeley.edu)



Cameron Malloy

Email: [cmalloy@berkeley.edu](mailto:cmalloy@berkeley.edu)



Shreya Kannan

Email: [shreyakannan@berkeley.edu](mailto:shreyakannan@berkeley.edu)



Sophia Qin

Email: [sophia.qin@berkeley.edu](mailto:sophia.qin@berkeley.edu)



Srinath Goli

Email: [srig@berkeley.edu](mailto:srig@berkeley.edu)



Vandana Ganesh

Email: [vandanag@berkeley.edu](mailto:vandanag@berkeley.edu)



---

# Welcome

- We are all here to learn:  
Knowledge (end) – Knowledge (start)



# Goals today

- Introduce you to
  - the field
  - the course
  - the team
- Answer your questions
- Big Ideas:
  - Abstraction
  - Data Types





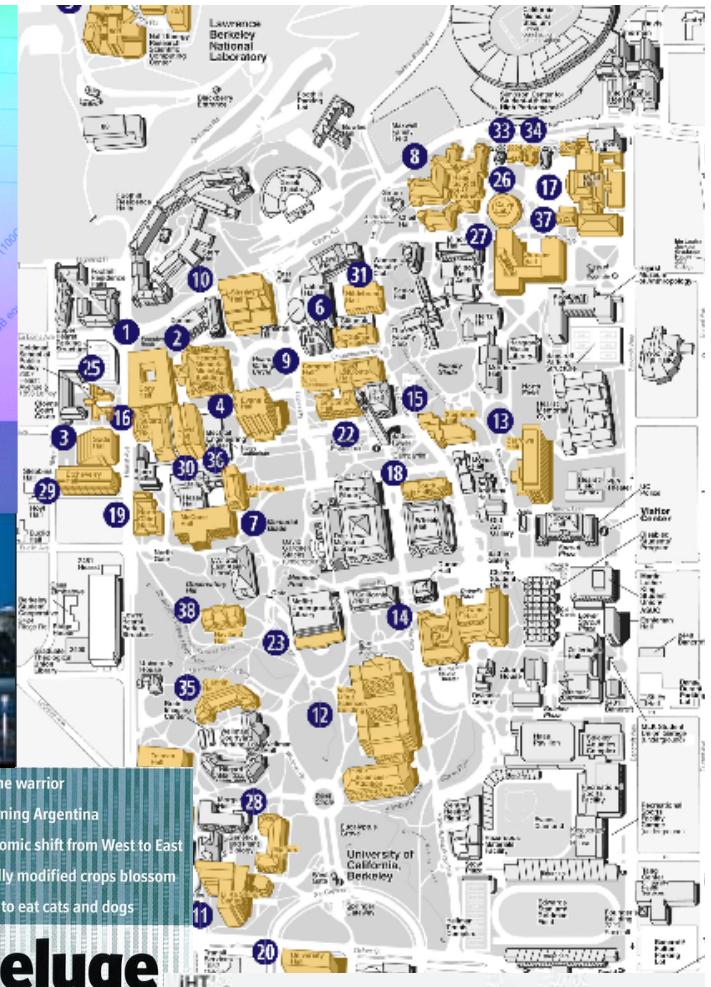
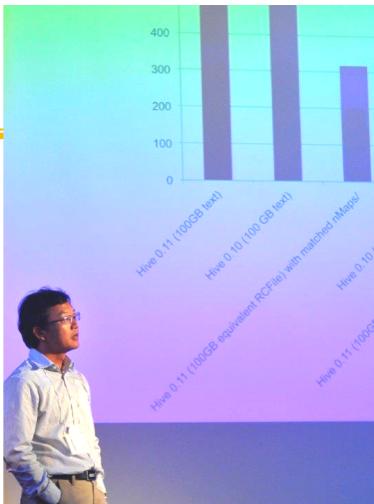
# Logistics

---

- **Labs are Mon & Tues**
  - Please see the signup form on Piazza
  - Go to any lab section next week
  - Attendance is tracked, but not required. ("Bonus" points.)
- **There are 2 sections on CalCentral, it does not matter which you are enrolled in.**
- **1-2pm lecture is webcast, 3pm is not.**

# Data Science

Nearly every field of discovery is transitioning from “data poor” to “data rich”



Berkeley  
UNIVERSITY OF CALIFORNIA

Data Science growing organically everywhere

**WIRED** Spark: Open Source Superstar Rewrites Future of Big Data



Reconstructing the movies in your mind



Bin Yu, Statistics  
Jack Gallant, Neuroscience



Earthquake Strong Shaking in 11 seconds  
Richard Allen Earth & Plan. Science Geospatial Lab



**KBase**  
PREDICTIVE BIOLOGY  
DOE Systems Biology Knowledgebase

Adam Arkin,  
Bioengineering



Fernando Perez,  
Brain Imaging Center  
iPython tools and community  
Charles Marshall  
Rosie Gillespie  
Integrative Biology  
Digitized Museum

The New York Times  
Incomes Flat in Recovery but Not for the 1%  
Feb 15, 2013  
Emmanuel Saez, Economics



**The Economist**

OBAMA  
The warrior

Misgoverning Argentina

The economic shift from West to East

Genetically modified crops blossom

The right to eat cats and dogs

**The data deluge**

AND HOW TO HANDLE IT: A 14-PAGE SPECIAL REPORT



**Analytics in Healthcare**

**Analytics: The Nervous System of IT-Enabled Healthcare**

The healthcare industry is moving from volume-based reimbursement to value-based reimbursement that is designed to achieve higher quality, lower costs, and better patient experience. To succeed, healthcare providers are forming accountable care organizations (ACOs) and restructuring their care delivery systems.



Berkeley  
UNIVERSITY OF CALIFORNIA

UCB CS88 Sp20 L1

# A National Challenge

Increasingly US jobs require data science and analytics skills. Can we meet the demand? The current shortage of skills in the national job pool demonstrates that business-as-usual strategies won't satisfy the growing need. If we are to unlock the promise and potential of data and all the technologies that depend on it, employers and educators will have to transform.

By 2021, **69% of employers** expect candidates with DSA skills to get preference for jobs in their organizations. Only **23% of college** and university leaders say their graduates will have those skills.

Report | McKinsey Global Institute

Big data: The next frontier for innovation, competition, and productivity

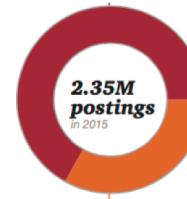
May 2011 | by James Manyika, Michael Chui, Brad Brown, Jacques Bughin, Richard Dobbs, Charles Roxburgh, Angela Hung Byers



pwc.com/us/dsa-skills

**Investing in America's data science and analytics talent**

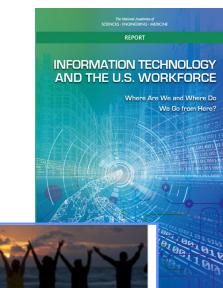
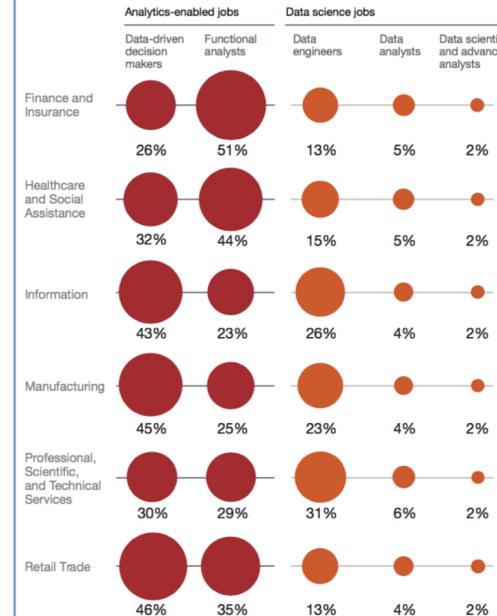
The case for action



April 2017

BHEF

Of 2.35 million job postings in the US.



Fourth Industrial Revolution

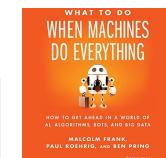
The fourth sector is a chance to build a new economic model for the benefit of all

Augmenting Human Intelligence



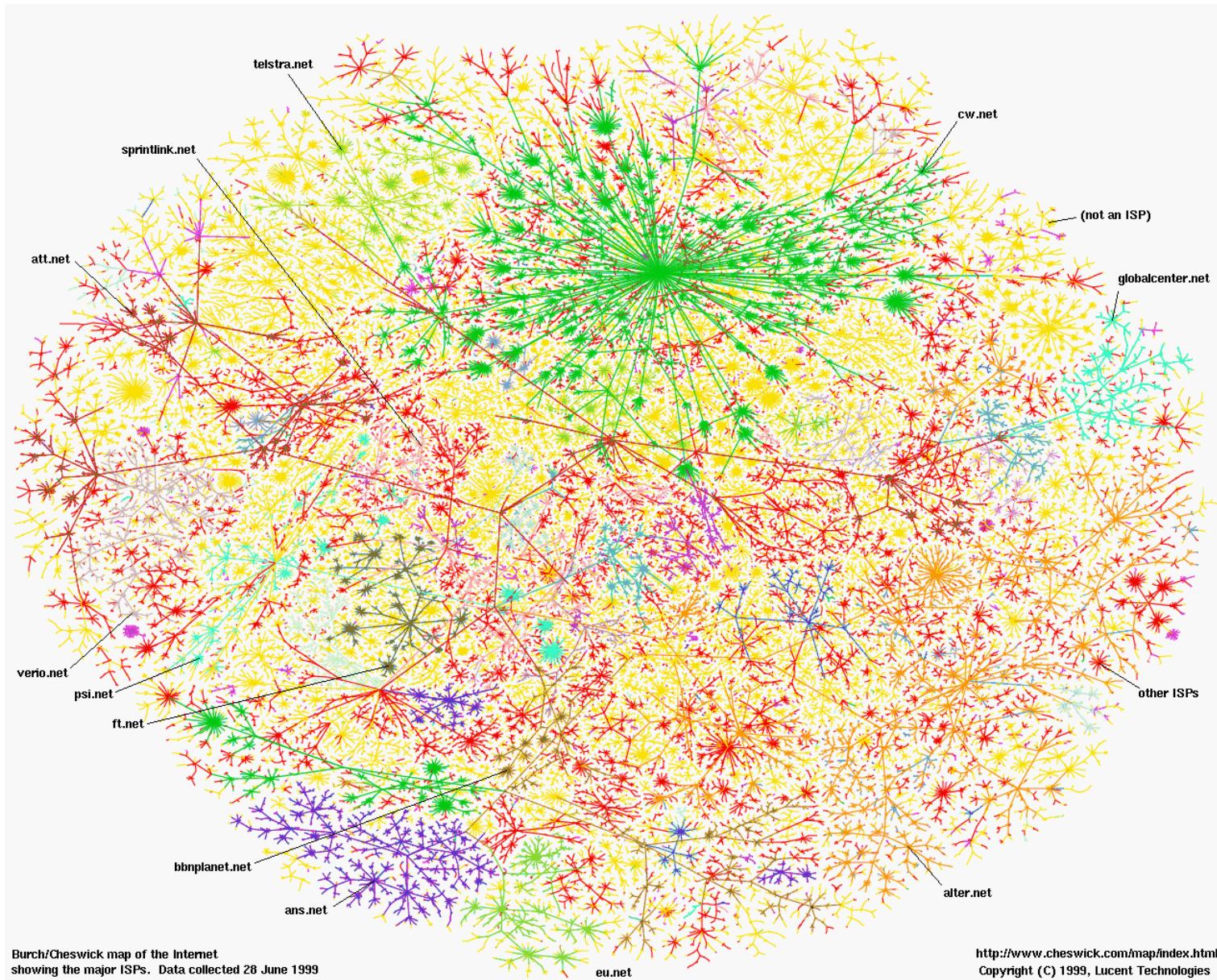
WORLD ECONOMIC FORUM

The Fourth Industrial Revolution:  
what it means, how to respond





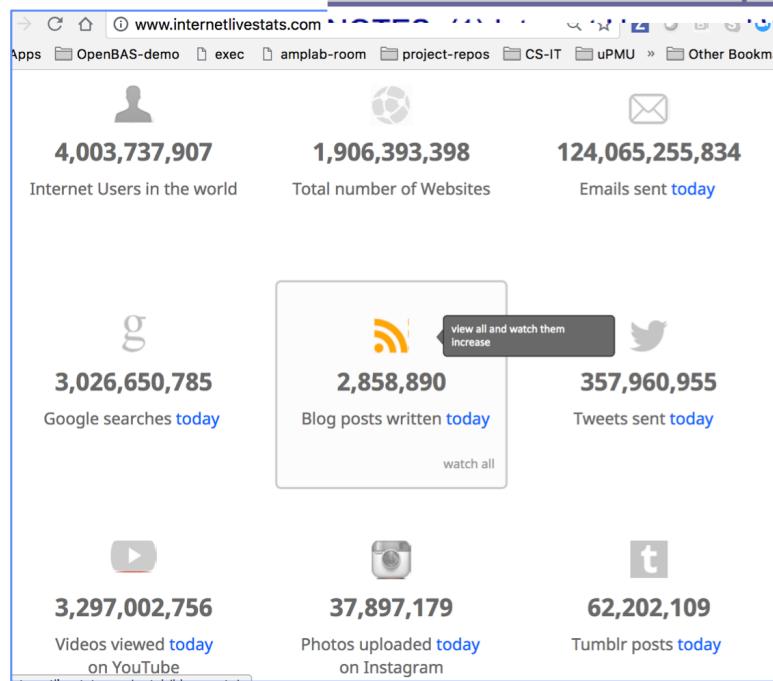
# Greatest Artifact of Human Civilization ...





## WORLD INTERNET USAGE AND POPULATION STATISTICS DEC 31, 2017 - Update

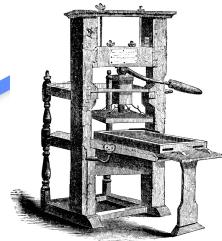
World Regions	Population (2018 Est.)	Population % of World	Internet Users 31 Dec 2017	Penetration Rate (% Pop.)	Growth 2000-2018
<a href="#">Africa</a>	<b>1,287,914,329</b>	16.9 %	<b>453,329,534</b>	35.2 %	9,941 %
<a href="#">Asia</a>	<b>4,207,588,157</b>	55.1 %	<b>2,023,630,194</b>	48.1 %	1,670 %
<a href="#">Europe</a>	<b>827,650,849</b>	10.8 %	<b>704,833,752</b>	85.2 %	570 %
<a href="#">Latin America / Caribbean</a>	<b>652,047,996</b>	8.5 %	<b>437,001,277</b>	67.0 %	2,318 %
<a href="#">Middle East</a>	<b>254,438,981</b>	3.3 %	<b>164,037,259</b>	64.5 %	4,893 %
<a href="#">North America</a>	<b>363,844,662</b>	4.8 %	<b>345,660,847</b>	95.0 %	219 %
<a href="#">Oceania / Australia</a>	<b>41,273,454</b>	0.6 %	<b>28,439,277</b>	68.9 %	273 %
<b>WORLD TOTAL</b>	<b>7,634,758,428</b>	<b>100.0 %</b>	<b>4,156,932,140</b>	<b>54.4 %</b>	<b>1,052 %</b>





# Era of Transformation

Age of Enlightenment



Industrial Revolution



Connected





# A Connected World of Data

---

- The world's knowledge at our finger tips
- *Digitization* of life, industry and society
- Intimately connected to billions of us, globally
- Explosion of observational instruments
  - Genomics, Microscopy, Astronomical, ...
- Vast Computational power to do analytics
- Synthetic design exploration thru simulation
- Machine reading of everything
- Statistical machine learning algorithms to “discover” structure



# What if I could ... ?

---



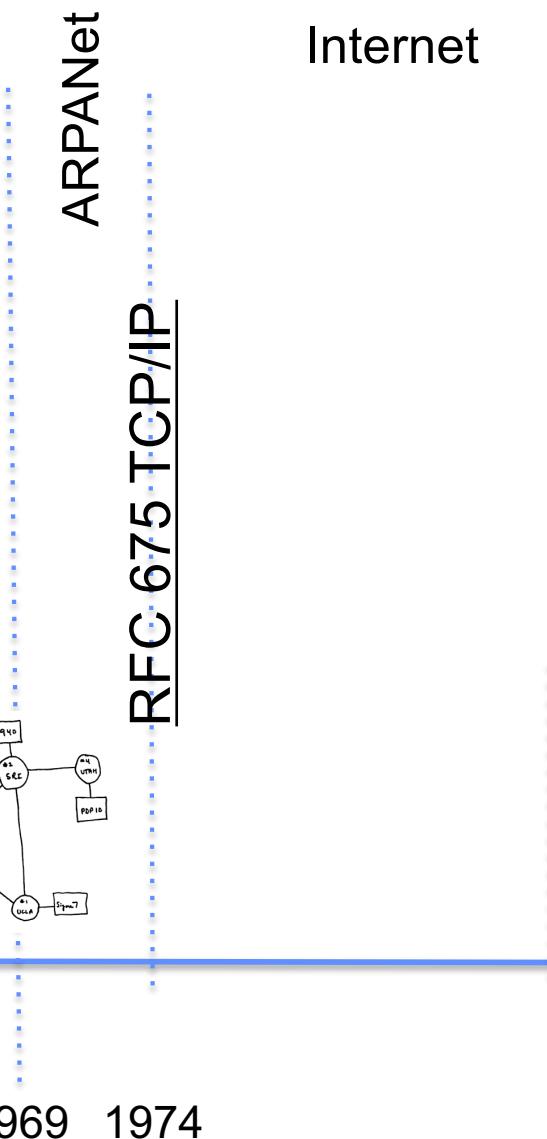
- See the world's digital footprints?
- Read everything that's ever been written?
- Take it all in and dive down anywhere as far as the science can take me?
- Learn the physical/chemical/biological /sociological/neurological... models from the data?
- Explore billions of designs and pick the one I want?
- ... ?



# A Connected World



3.0 B 11/15



Internet Users in the world

**3,293,151,639**

g

**2,652,887,737**

Google searches **today**

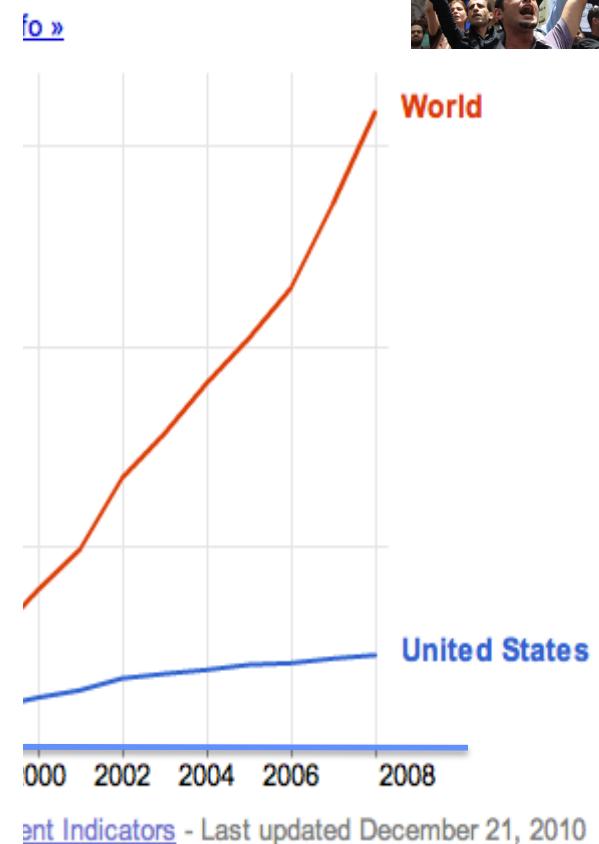


**5,835,884,253**

Videos viewed **today**  
on YouTube

1/24/20

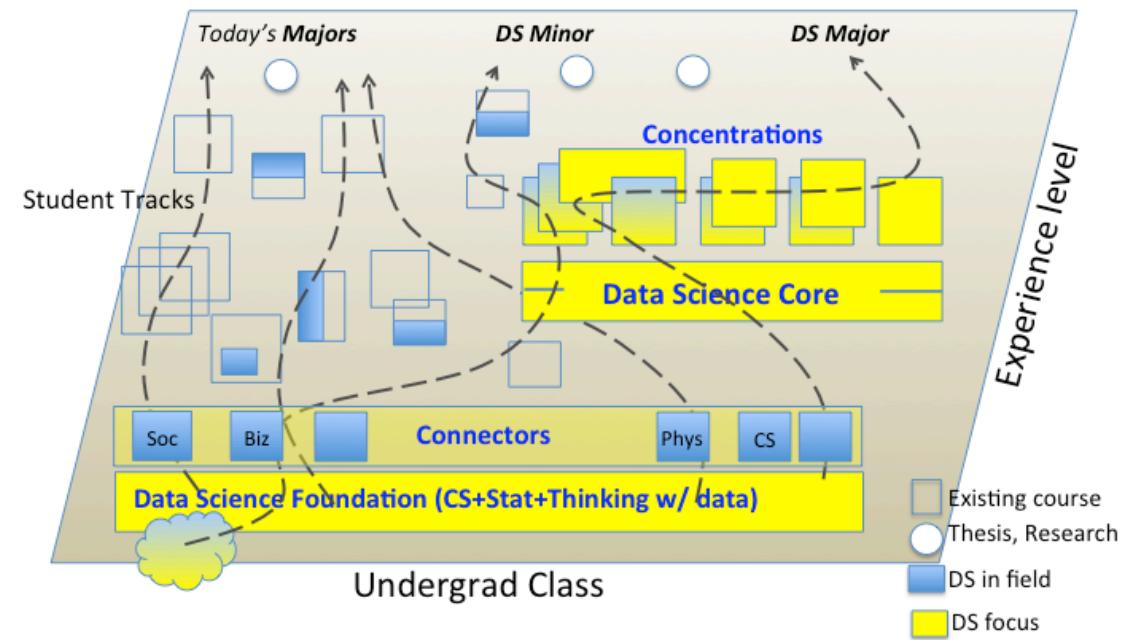
15





# Data 8 – Foundations of Data Science

- Computational Thinking + Inferential Thinking in the context of working with real world data
- Introduce you to several computational concepts in a simple data-centered setting
  - Authoring computational documents
  - Tables
  - Within Python3 and “SciPy”





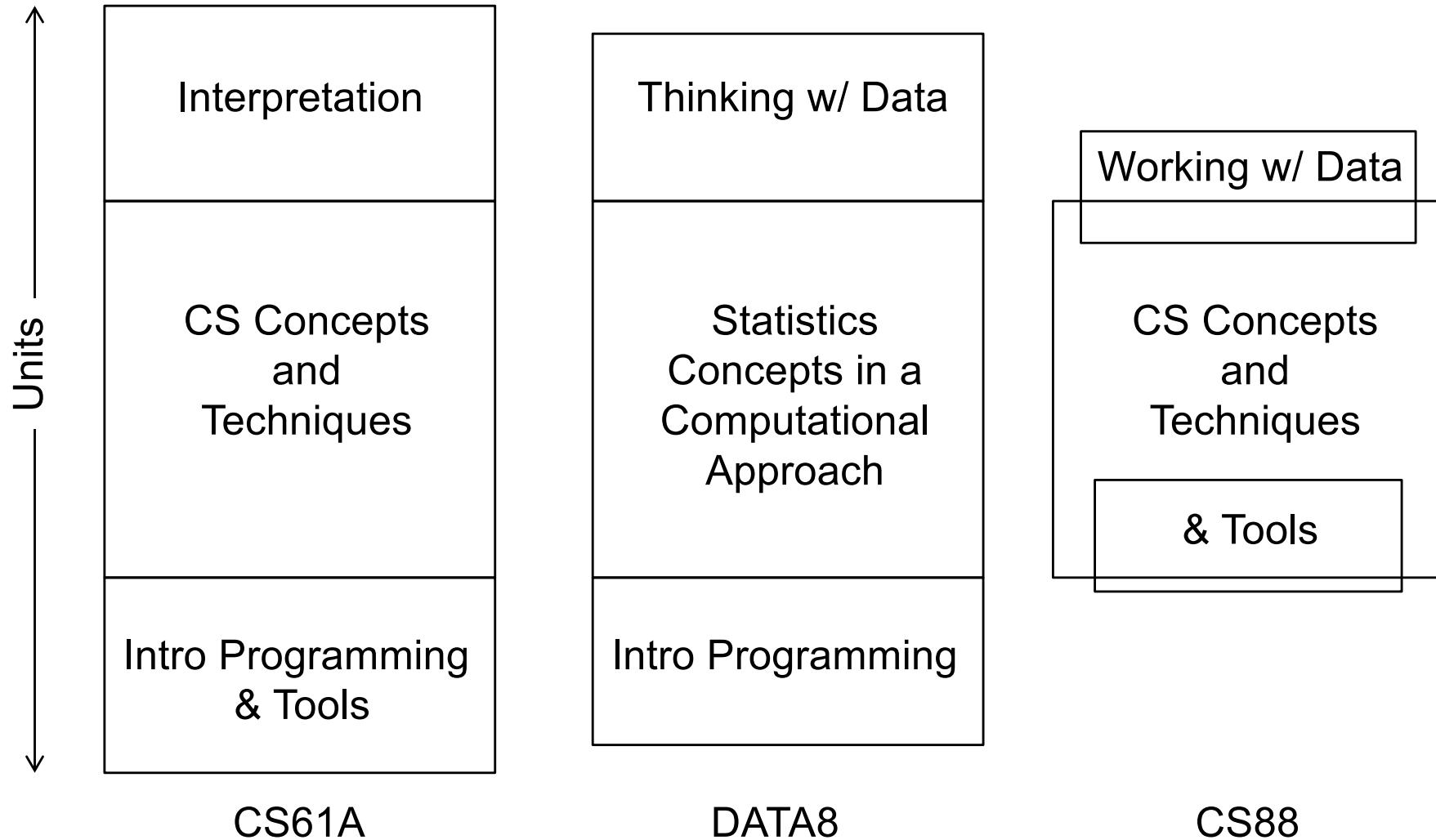
# CS88 – Computational Structures in Data Science

---

- **Deeper understanding of the computing concepts introduced in C8**
  - Hands-on experience => Foundational Concept
  - How would you create what you use in c8 ?
- **Extend your understanding of the structure of computation**
  - What is involved in interpreting the code you write ?
  - Deeper CS Concepts: Recursion, Objects, Classes, Higher-order Functions, Declarative programming, ...
  - Managing complexity in creating larger software systems through composition
- **Create complete (and fun) applications**
- **In a data-centric approach**



# How does CS88 relate to CS61A ?





# Opportunities for students

---

c8

c8 CS88

c8 CS88 CS61B

CS minor

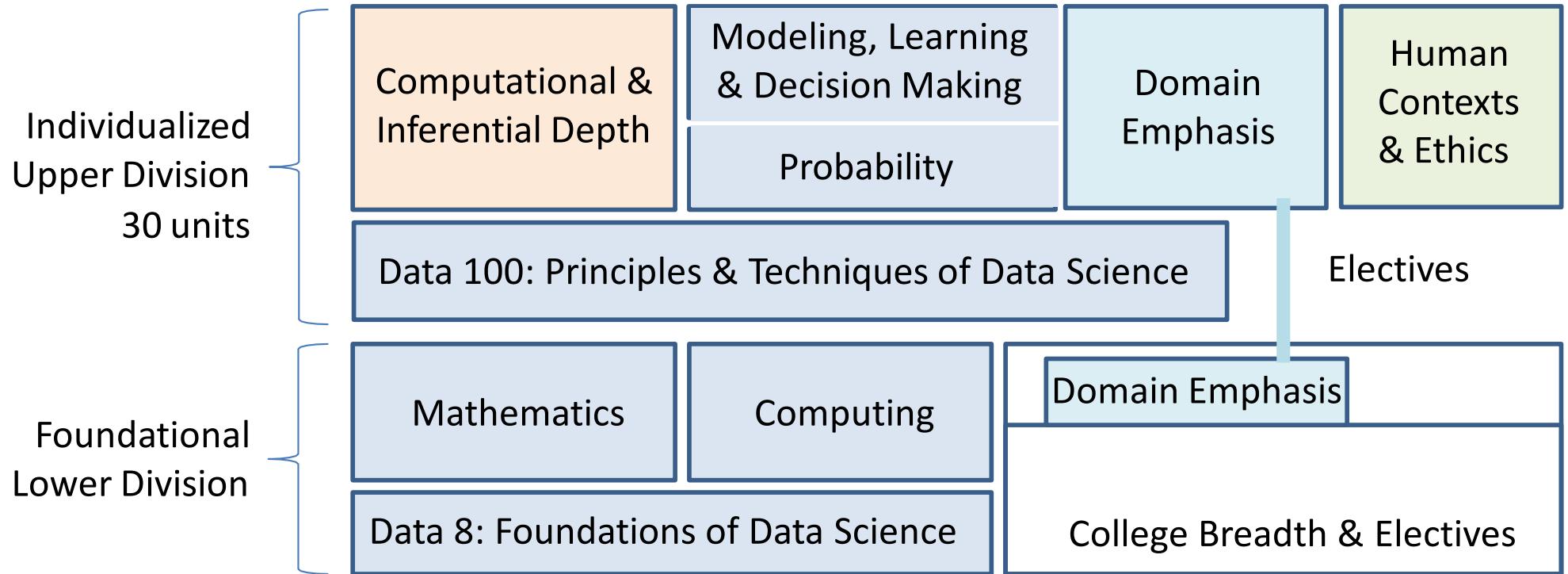
\*\*\*

CS major

c8 cs61a

cs61a

# The Data Science Major





# Course Structure

- **2 hours lecture, 2 hours lab, ~5 hours HW+study**
- **Lecture introduces concepts (quickly!), answers why questions.**
- **Lab provides concrete detail hands-on**
- **Homework (12) cements your understanding**
  - Out Tuesdays, Due Next Thurs (~9 days)
- **Projects (2) put your understanding to work in building complete applications**
  - Maps
  - Ants vs Some Bees
  - Maybe: open-ended final?

A screenshot of the Composing Programs website. The header includes a navigation bar with links like 'Apps', 'OpenBAS-demo', 'exec', 'amplab-room', 'project-repos', 'CS-IT', 'uPMU', 'Chair Viewer', 'DataSci', 'Confs', and 'DS8-88'. Below the header is a main menu with 'COMPOSING PROGRAMS' and links to 'TEXT', 'PROJECTS', 'TUTOR', and 'ABOUT'. The main content area features a welcome message about the site's focus on abstraction, programming paradigms, and Python 3, along with information about related sites (CS 61A Course, Version 1) and a call for instructors to fill out a survey.

- **Readings:** <http://composingprograms.com>
  - Same as CS61a



# Course Culture

- Learning
- Community
- Respect
- Collaboration
- Peer Instruction





# Piazza for {ask,answer}ing questions

Screenshot of the Piazza platform interface.

**Header:** piAZZZA CS 10 Questions - Statistics 35 | Search or ask a question... Add Question/Note | Dan Garcia Piazza Help

**Left Sidebar (QUESTION FEED):**

- This week:**
  - When are TA / professor office hours? (Sun) #instructor-question #admin
  - When can I meet up with a GSI or professor to get help with the course material? #admin
- Last week:**
  - So, I'm here... now how exactly does Pia (Mon) #logistics #welcome

**Question Detail View:**

**Question:** When are TA / professor office hours?  
When can I meet up with a GSI or professor to get help with the course material? #admin  
Last updated by Luke Segars 2 days ago

**Instructors' Response:** We haven't established our office hours yet, but we'll make that information available as soon as possible. Check back here for an update by the second week of classes.  
Last updated by Luke Segars 2 days ago

**Actions:** Good Question! | Good Answer! | Ask a Followup »

**Followup Discussions:** Still Confused? Ask New Followup

**Metrics:** AVERAGE RESPONSE TIME: N/A | SPECIAL MENTIONS: Luke Segars answered When are TA / ... in 1.1 hr. 2 days ago | USERS ONLINE THIS WEEK: 3 Online Now: 1

About Piazza | Privacy Policy | Copyright Policy | Terms of Use | Report a Bug!  
Copyright © 2013 Piazza Inc. All rights reserved.



# Where will we work?

---

- Your laptop
  - Using an editor and a terminal
- cs88.org
- datahub.berkeley.edu
  - “Jupyter Notebooks”. An industry standard tool.
  - Not as often, but an option



# iClicker Check In

---

- Are you enrolled in Data 8?
- A. I took it Fall 2018 or earlier
- B. I took it Spring 2019
- C. I'm taking it right now
- D. I am trying to enroll in Data 8
- E. I am not taking Data 8



# Pro-student Grading Policies

---

- **EPA**
  - Rewards good behavior
  - Effort
    - » E.g., Office hours, doing every single lab, hw, reading Piazza pages
  - Participation
    - » E.g., Raising hand in lec or discussion, asking questions on Piazza
  - Altruism
    - » E.g., helping other students in lab, answering questions on Piazza
- **You have 3 “Slip Days”**
  - Homework and Projects
  - You use them to extend due date, 1 slip day for 1 day extension
  - You can use them one at a time or all at once or in any combination
  - They follow you around when you pair up (you are counted individually)
    - » E.g., A has 2, B has 0. Project is late by 1 day. A uses 1, B is 1 day late



# Pro-student Grading Policies

---

- Turning it work will never harm you! Even if it's late or you can't get it completely correct.
- Labs are based primarily on effort.
- If you're struggling, come talk to us!
  - Especially your TAs. They have been there too.



# Abstraction

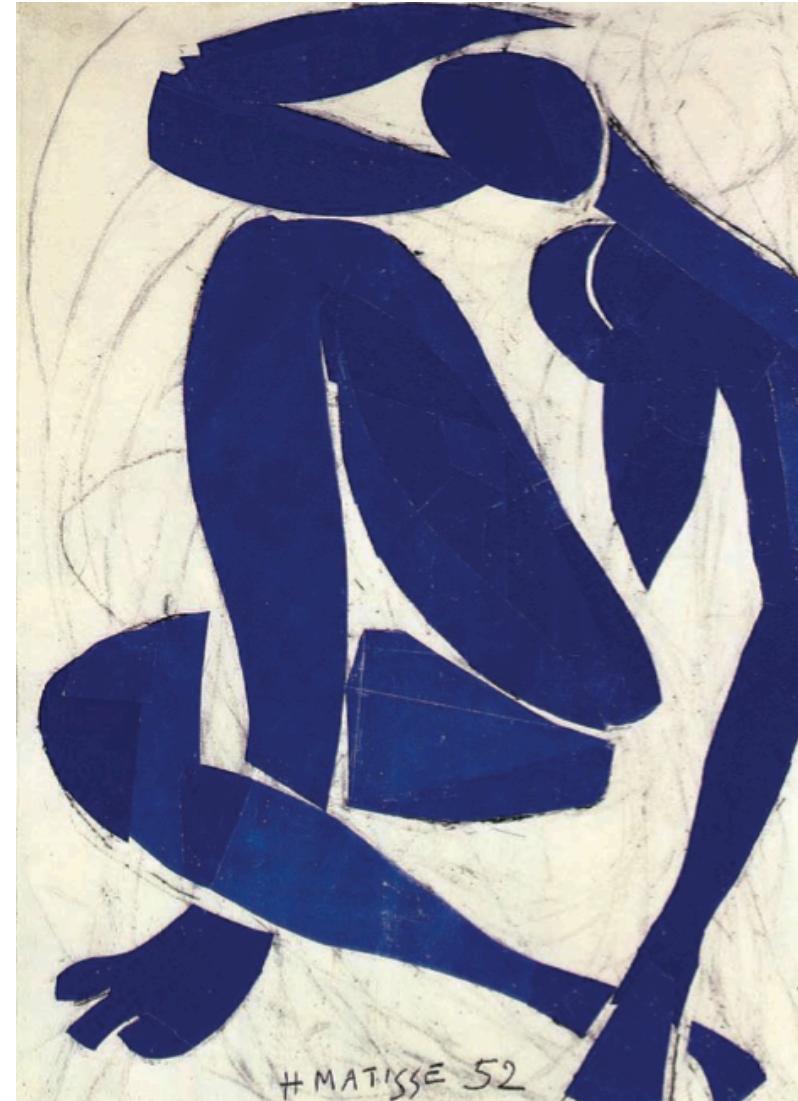
- Detail removal

**“The act of leaving out of consideration one or more properties of a complex object so as to attend to others.”**

- Generalization

**“The process of formulating general concepts by abstracting common properties of instances”**

- Technical terms:  
Compression, Quantization,  
Clustering, Unsupervised  
Learning



Henri Matisse “Naked Blue IV”



# Experiment

Standard Time Zones of the World

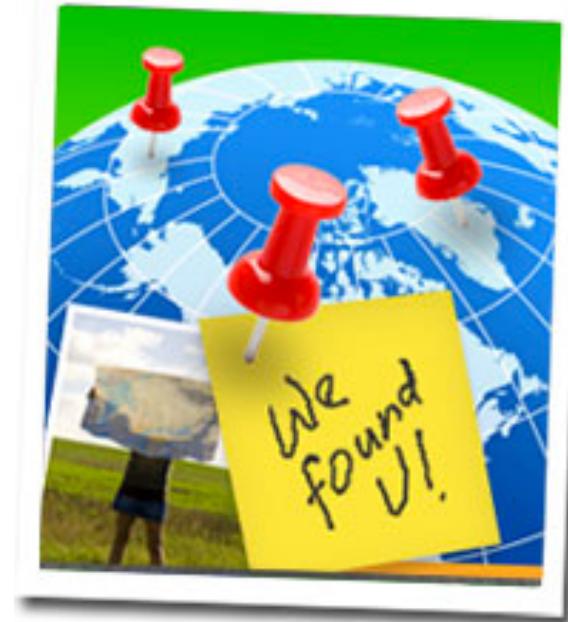




# Where are you from?

Possible Answers:

- Planet Earth
- Europe
- California
- The Bay Area
- San Mateo
- 1947 Center Street,  
Berkeley, CA
- $37.8693^\circ \text{ N}, 122.2696^\circ \text{ W}$



All correct but different levels of abstraction!



# Abstraction gone wrong!



## I Can Stalk U

Raising awareness about inadvertent information sharing

[Home](#)   [How](#)   [Why](#)   [About Us](#)   [Contact Us](#)

**What are people *really* saying in their tweets?**

 [denisluque](#): I am currently nearby <http://maps.google.com/?q=-23.6193333333,-46.5506666667>  
1 minute ago · [Map Location](#) · [View Tweet](#) · [View Picture](#) · [Reply to denisluque](#)

 [nikosofficial](#): I am currently nearby <http://maps.google.com/?q=48.8699833333,2.3282833333>  
5 minutes ago · [Map Location](#) · [View Tweet](#) · [View Picture](#) · [Reply to nikosofficial](#)

 [dilmanarede](#): I am currently nearby <http://maps.google.com/?q=-15.7878333333,-47.8291666667>  
7 minutes ago · [Map Location](#) · [View Tweet](#) · [View Picture](#) · [Reply to dilmanarede](#)

 [downtownvan](#): I am currently nearby <http://maps.google.com/?q=49.2833333333,-123.119833333>  
10 minutes ago · [Map Location](#) · [View Tweet](#) · [View Picture](#) · [Reply to downtownvan](#)

 [MommaGooseBC](#): I am currently nearby 15745 Weaver Lake Rd  
Maple Grove MN

**Links**

- Mayhemic Labs
- PaulDotCom
- SANS ISC
- Electronic Frontier Foundation
- Center for Democracy & Technology

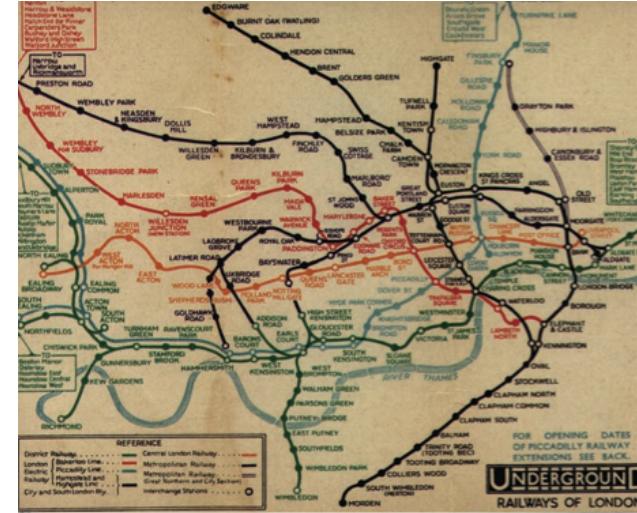
**How did you find me?**

Did you know that a lot of smart phones encode the location of where pictures are taken? Anyone who has a copy can access this information.



# Detail Removal (in Data Science)

- You'll want to look at only the interesting data, leave out the details, zoom in/out...
- Abstraction is the idea that you focus on the essence, the cleanest way to map the messy real world to one you can build
- Experts are often brought in to know what to remove and what to keep!



The London Underground 1928 Map & the 1933 map by Harry Beck.



# The Power of Abstraction, Everywhere!

- **Examples:**

- Functions (e.g.,  $\sin x$ )
- Hiring contractors
- Application Programming Interfaces (APIs)
- Technology (e.g., cars)

- **Amazing things are built when these layer**

- And the abstraction layers are getting deeper by the day!

*We only need to worry about the interface, or specification, or contract  
NOT how (or by whom) it's built*

**Above the abstraction line**

**Abstraction Barrier (Interface)**  
(the interface, or specification, or contract)

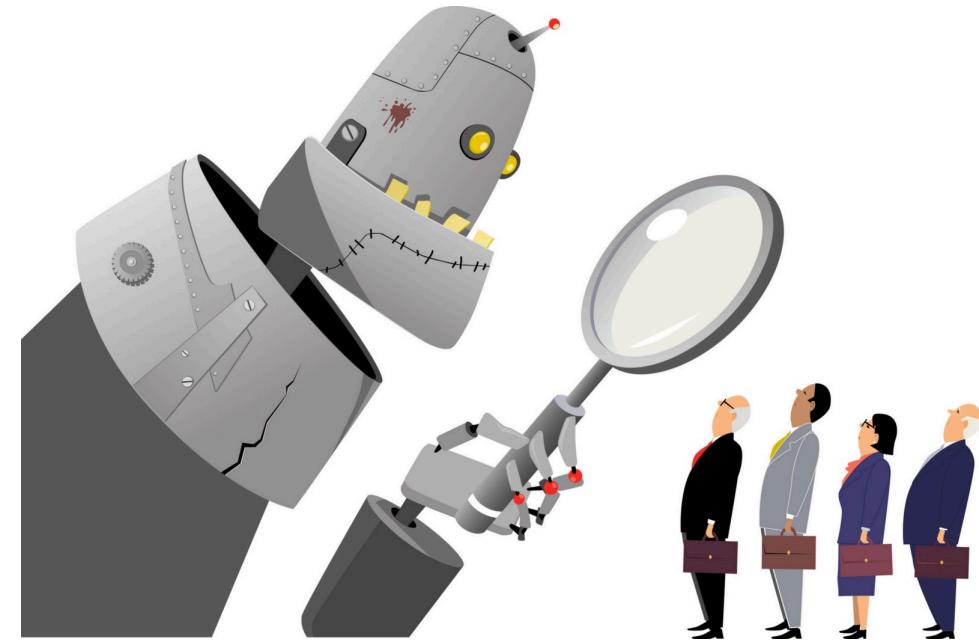
**Below the abstraction line**

*This is where / how / when / by whom it is actually built, which is done according to the interface, specification, or contract.*



# Abstraction: Pitfalls

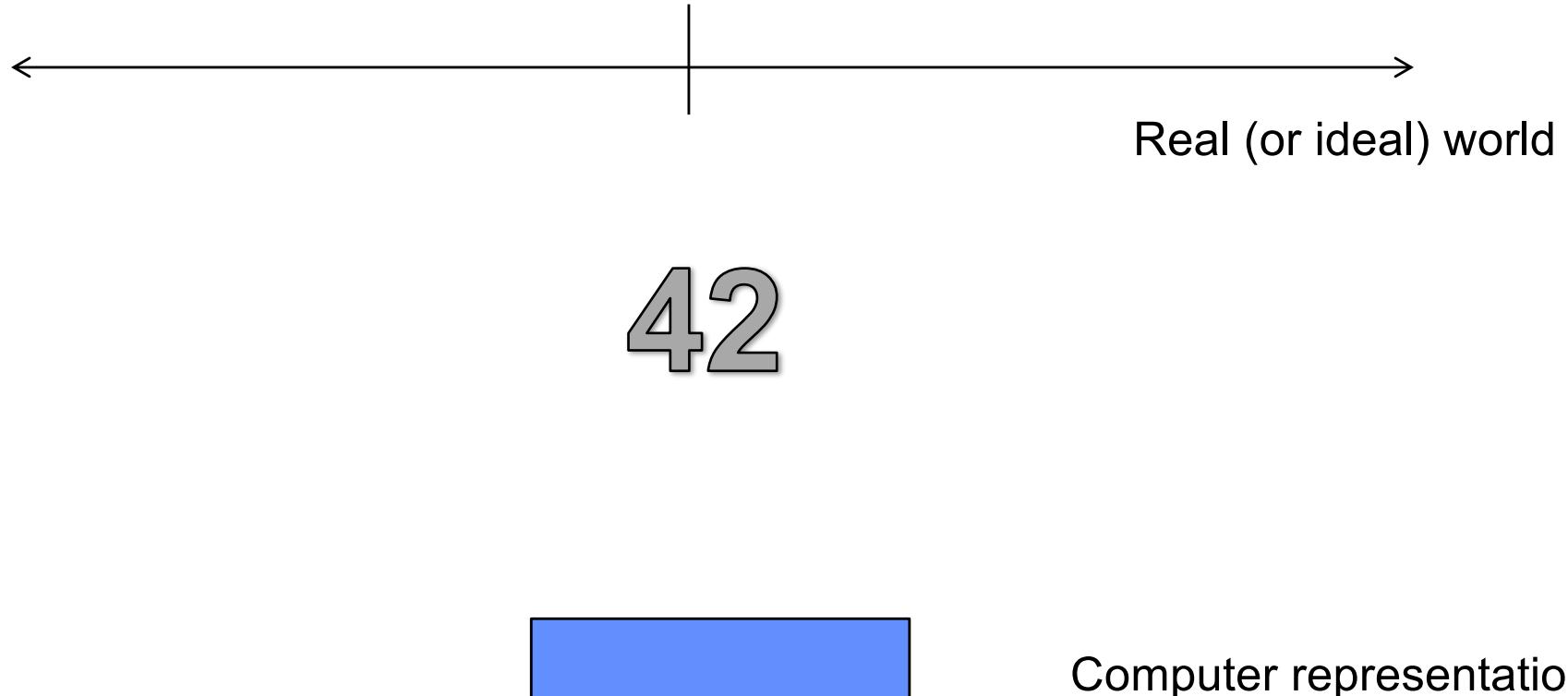
- Abstraction is not universal without loss of information (mathematically provable). This means, in the end, the complexity can only be “moved around”
- Abstraction makes us forget how things actually work and can therefore hide bias. Example: AI and hiring decisions.
- Abstraction makes things special and that creates dependencies. Dependencies grow longer and longer over time and can become unmanageable.





# Abstraction in CS: Data Type

- What's this?





# Data Types and Operations

---

- **Set of elements**
  - with some internal representation
  - E.g. Integers, Floats, Booleans, Strings, ...
- **Set of operations on elements of the type**
  - e.g. +, \*, -, /, %, //, \*\*
  - ==, <, >, <=, >=
- **Properties**
  - Commutative, Associative, ...
- **Expressions are valid well-defined sets of operations on elements that produce a value of a type**



# Lab and HW next week

---

- Lab will get you setup with Python locally
  - You'll use your own computer
  - Learn about text files and the terminal



# Thoughts for the Wandering Mind

---

A binary digit (bit) is a symbol from {0,1}.

- How many things can you represent with N bits?
- How many things can you represent with 1 digit (0-9)?
- 2 digits? 6 digits?