



Welcome to CS88

David E. Culler

CS8 – Computational Structures in Data Science

<http://inst.eecs.berkeley.edu/~cs88>

Lecture 1

August 27, 2018



Welcome

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

CS88 Team





S88 Team - uGSIs



Ting Ding
tingding96@berkeley.edu

Jessica Gao
gaojessicaping@berkeley.edu

Alex Kassil
alexkassil@berkeley.edu



Amir Shahatit
ashahatit@berkeley.edu

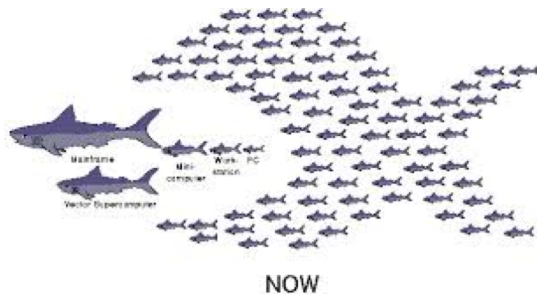
Andrew Tan
andrewtan@berkeley.edu

John Yang
john.yang20@berkeley.edu



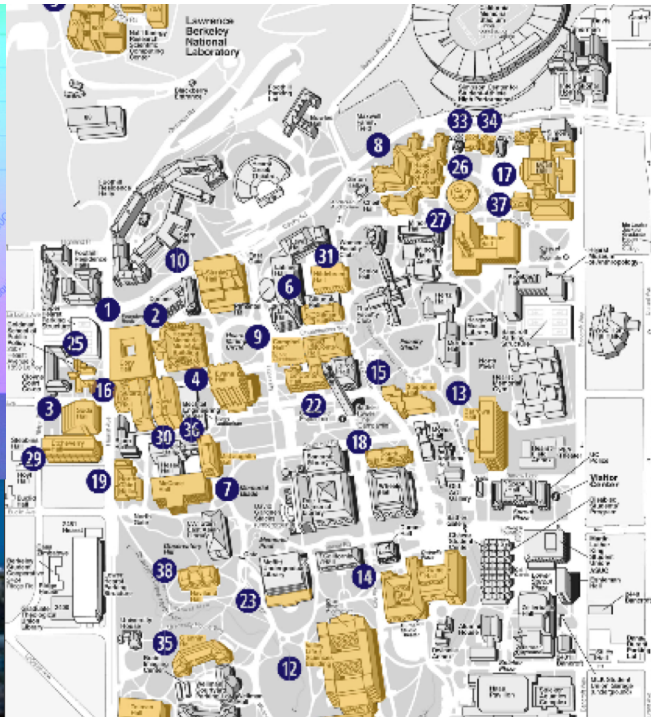
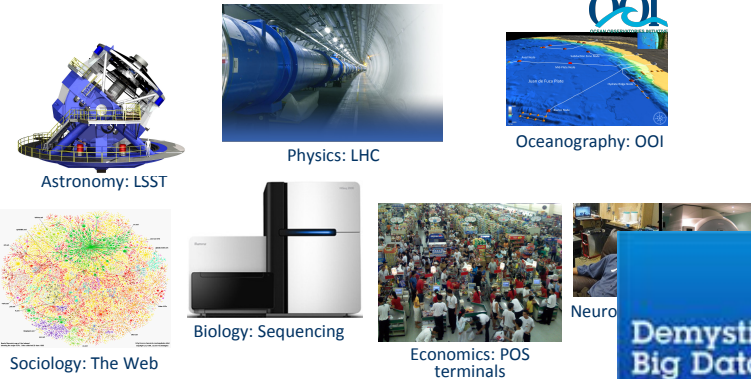
CS88 Team - me

- **David Culler (culler@berkeley.edu)**
 - Hearst Field Annex / 465 Soda Hall (amplab)
 - <http://www.cs.berkeley.edu/~culler>
 - Office hours: Mon 3-4 + TBD
- **Build things**
 - Cray Time Sharing System
 - OS386, OS286
 - Active Messages
 - Massive High Performance Clusters
 - TinyOS / Berkeley Motes, ...
 - LoCal, BOSS, ...



Data Science

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



Demystifying Big Data in Government

A practical guide to transforming the business of government

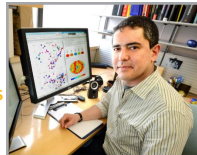


Data Science growing organically everywhere

WIRED Spark: Open Source Superstar Rewrites Future of Big Data
BY CADE METZ 08.19.13 6:30 AM



AMP Lab
Ion Stoica, CS
Michael Franklin, CS



Fernando Perez,
Brain Imaging Center
iPython tools and community

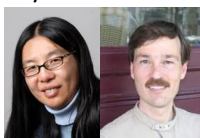
KBase
PREDICTIVE BIOLOGY
DOE Systems Biology Knowledgebase

Adam Arkin,
Bioengineering



Charles Marshall
Rosie Gillespie
Integrative Biology
Digitized Museum

Reconstructing the movies in your mind



Bin Yu, Statistics
Jack Gallant, Neuroscience



Richard Allen
Earth & Plan. Science
Geospatial Lab



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. To succeed, healthcare providers are leveraging accountable care organizations (ACOs) and restructuring their care delivery systems.

Collecting the Data	Clinical Intelligence (CI)	Business Intelligence (BI)	Performance Evaluation
80% of electronic health information	30% of US hospitals	33% of healthcare organizations use BI tools	YEAR 2015

A National Challenge

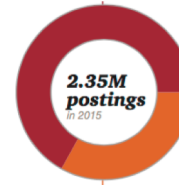
Report | McKinsey Global Institute
Big data: The next frontier for innovation, competition, and productivity

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



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.



April 2017

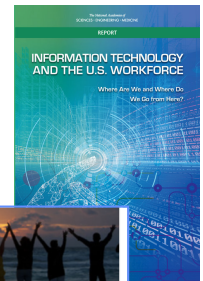
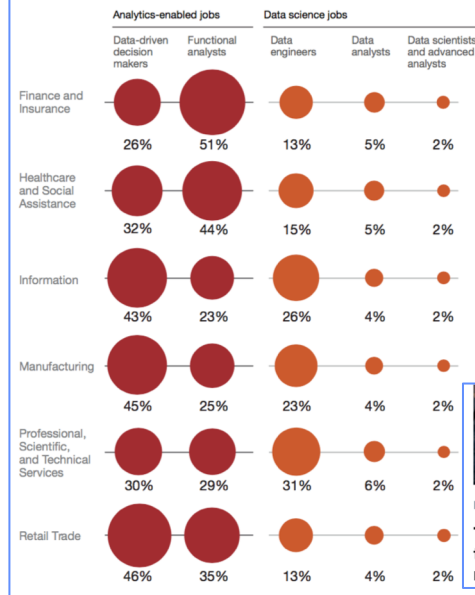


Investing in America's data science and analytics talent

The case for action



Of 2.35 million job postings in the US.



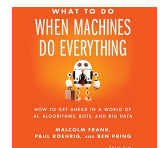
Fourth Industrial Revolution

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

Augmenting Human Intelligence



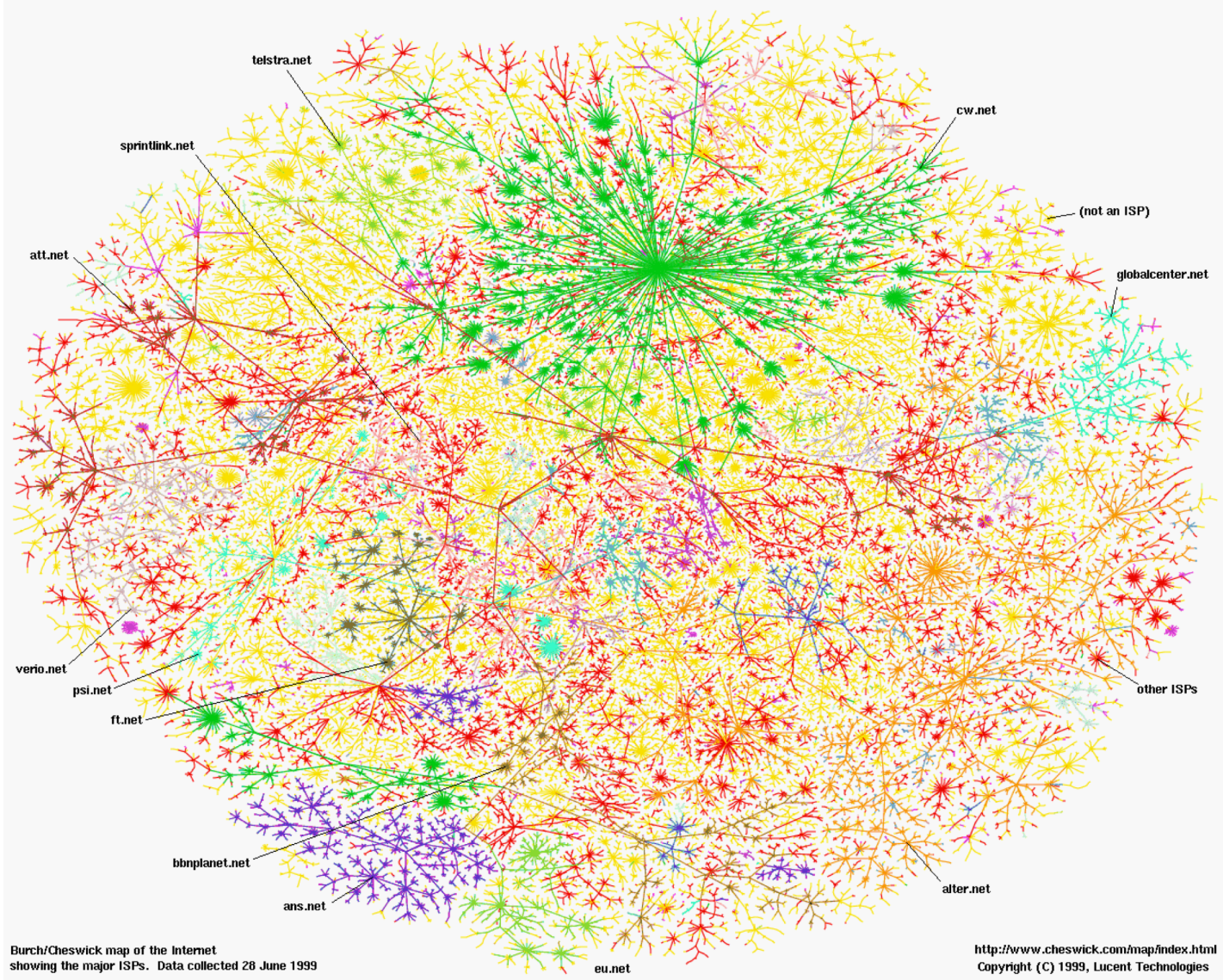
The Fourth Industrial Revolution: what it means, how to respond



pwc.com/us/dsa-skills



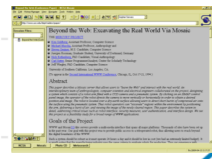
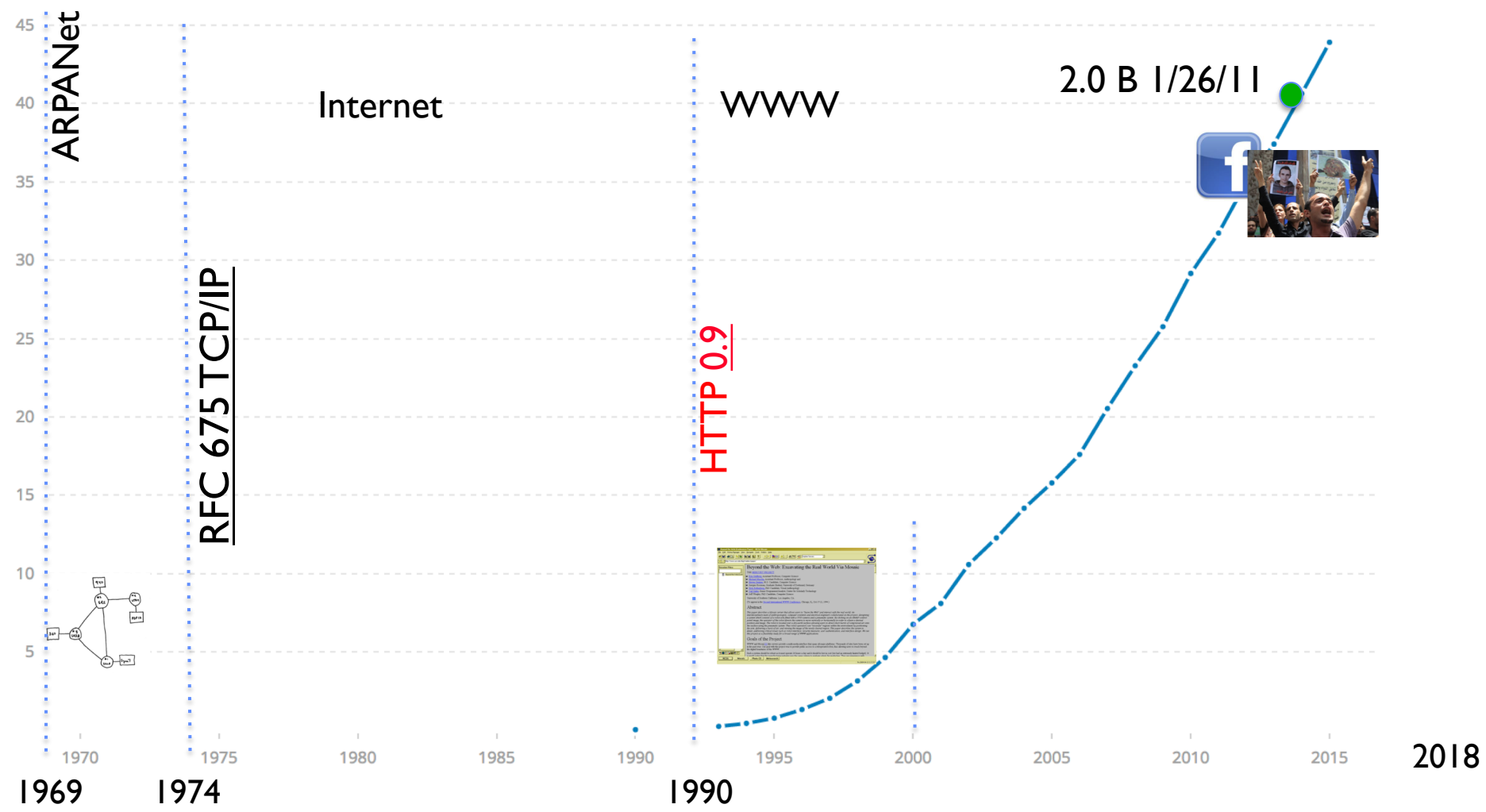
Greatest Artifact of Human Civilization ...



Burch/Cheswick map of the Internet showing the major ISPs. Data collected 28 June 1999

<http://www.cheswick.com/map/index.html>
Copyright (C) 1999, Lucent Technologies

The Global Village





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
Africa	1,287,914,329	16.9 %	453,329,534	35.2 %	9,941 %
Asia	4,207,588,157	55.1 %	2,023,630,194	48.1 %	1,670 %
Europe	827,650,849	10.8 %	704,833,752	85.2 %	570 %
Latin America / Caribbean	652,047,996	8.5 %	437,001,277	67.0 %	2,318 %
Middle East	254,438,981	3.3 %	164,037,259	64.5 %	4,893 %
North America	363,844,662	4.8 %	345,660,847	95.0 %	219 %
Oceania / Australia	41,273,454	0.6 %	28,439,277	68.9 %	273 %
WORLD TOTAL	7,634,758,428	100.0 %	4,156,932,140	54.4 %	1,052 %

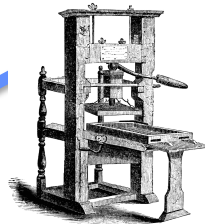
The screenshot shows the following statistics on the website:

- Internet Users in the world: 4,003,737,907
- Total number of Websites: 1,906,393,398
- Emails sent today: 124,065,255,834
- Google searches today: 3,026,650,785
- Blog posts written today: 2,858,890
- Tweets sent today: 357,960,955
- Videos viewed today on YouTube: 3,297,002,756
- Photos uploaded today on Instagram: 37,897,179
- Tumblr posts today: 62,202,109

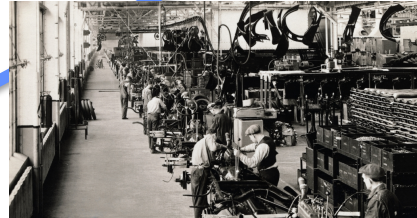


Era of Transformation

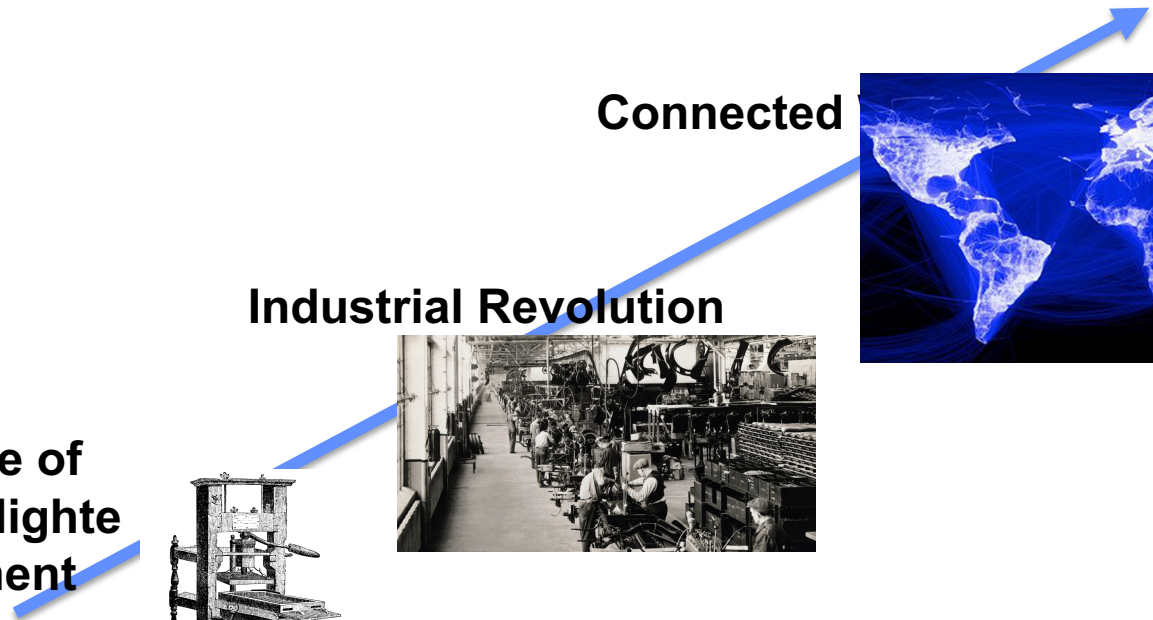
Age of Enlightenment



Industrial Revolution



Connected





A Connected World of Data

- The world's knowledge at our finger tips
- *Digitalization* 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 ... ?

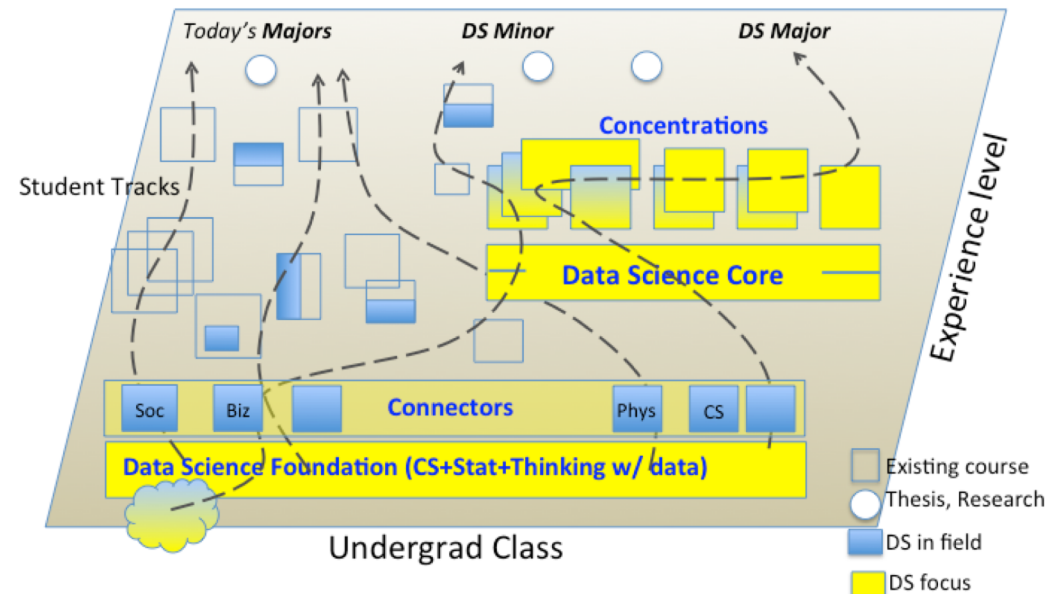


- See the world's digital footprints?
- Read everything that's ever been written?
- Take it all in and dive down anywhere as far as the science can take me?
- Learn the physical/chemical/biological /sociological/neurological... models from the data?
- Explore billions of designs and pick the one I want?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



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”



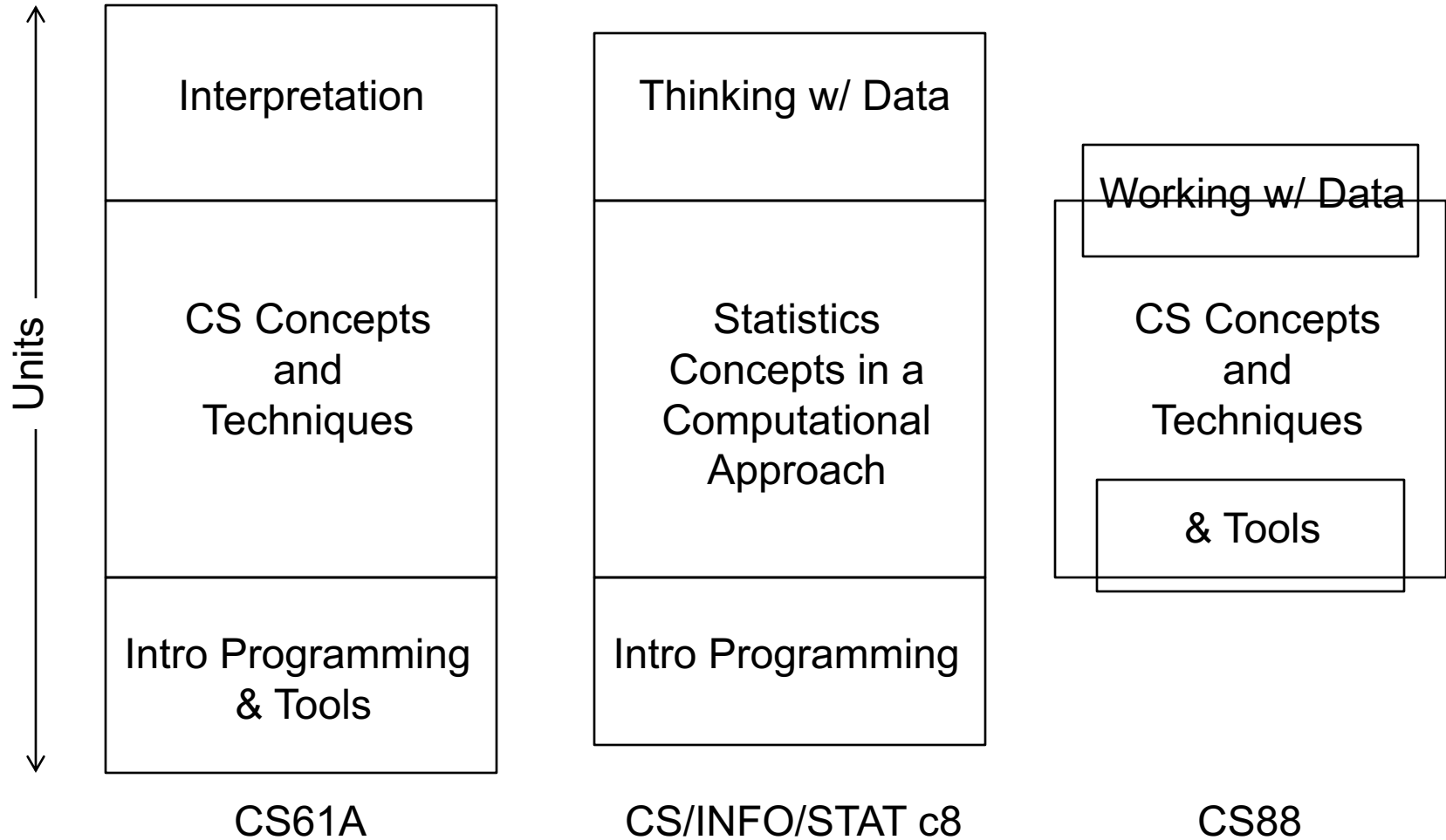
CS88 – Computational Structures in Data Science



- **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

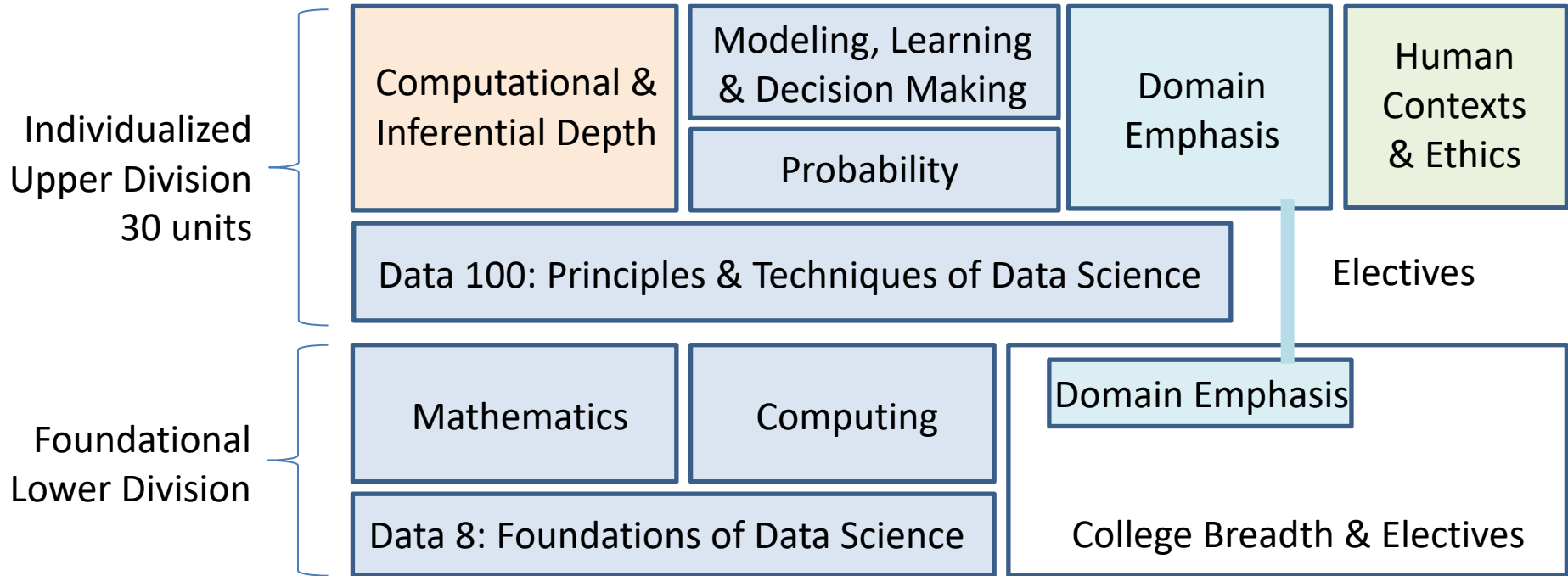
CS minor

CS major

c8 cs61a

cs61a

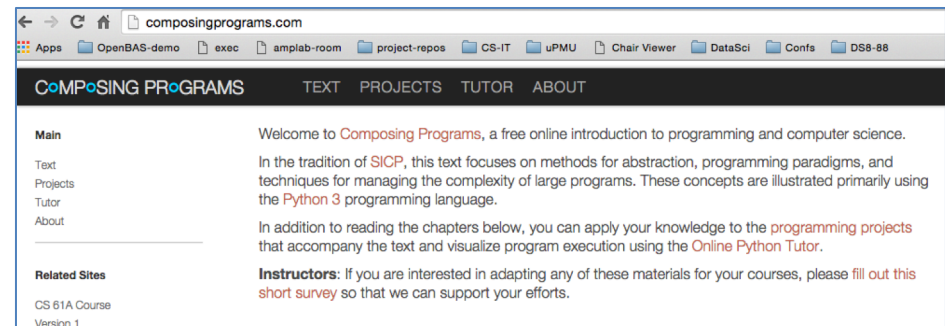
A New Data Science Major soon





Course Structure

- **Monday Lecture + Friday Lab/Discussion**
- **Lecture introduces concepts (quickly)**
- **Lab provides concrete detail hands-on**
- **Homework cements your understanding**
 - Out Friday, Due Thursday
- **Projects (3) put your understanding to work in building complete applications**
- **Readings: composingprograms.com**
 - Same as cs61a





Course Culture

- Learning
- Community
- Respect
- Collaboration
- Peer Instruction

Collaboration

Asking questions is highly encouraged

- Discuss all questions with each other (except exams)
- Submit lab assignments individually (graded on completeness)
 - If you come to lab, you can collaborate liberally
 - If you choose not to come to lab, you must work alone
- Submit homework individually and list collaborators
- Submit projects in pairs; find a partner in your lab

The Limits of collaboration

- Don't share solutions with each other (except project partners)
- Copying solutions will result in failing the course



Piazza for {ask,answer}ing questions

PIAZZA CS 10 Questions · Statistics **35** Search or ask a question... Add Question/Note Dan Garcia Piazza Help

Popular tags: #instructor-question #admin #logistics #welcome

QUESTION FEED FILTERS

▼ This week

When are TA / professor office hours? Sun
When can I meet up with a GSI or professor to get help with the course material? #admin
#instructor-question #admin

▼ Last week

So, I'm here... now how exactly does Pia: Mon
(No question details)
#logistics #welcome

question. 3 Views, 1 Follows Actions

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

Good Question!

instructors' response. Actions

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

Good Answer! **Ask a Followup** »

Start off a Students' Response

followup discussions.

Still Confused? Ask New Followup

AVERAGE RESPONSE TIME SPECIAL MENTIONS **USERS ONLINE THIS WEEK**

N/A Luke Segars answered **When are TA / ...** in 1.1 hr. 2 days ago **3**
Online Now: 1

About Piazza · Privacy Policy · Copyright Policy · Terms of Use · Report a Bug!



Where will we work?

- **datahub.Berkeley.edu**
- **The computer you carry around**
- **inst.eecs.Berkeley.edu**



Lab Sections Assignments

- **We will collect availability on Wednesday**
- **Attend any lab section on Friday.**
- **Assignments effective following Friday.**

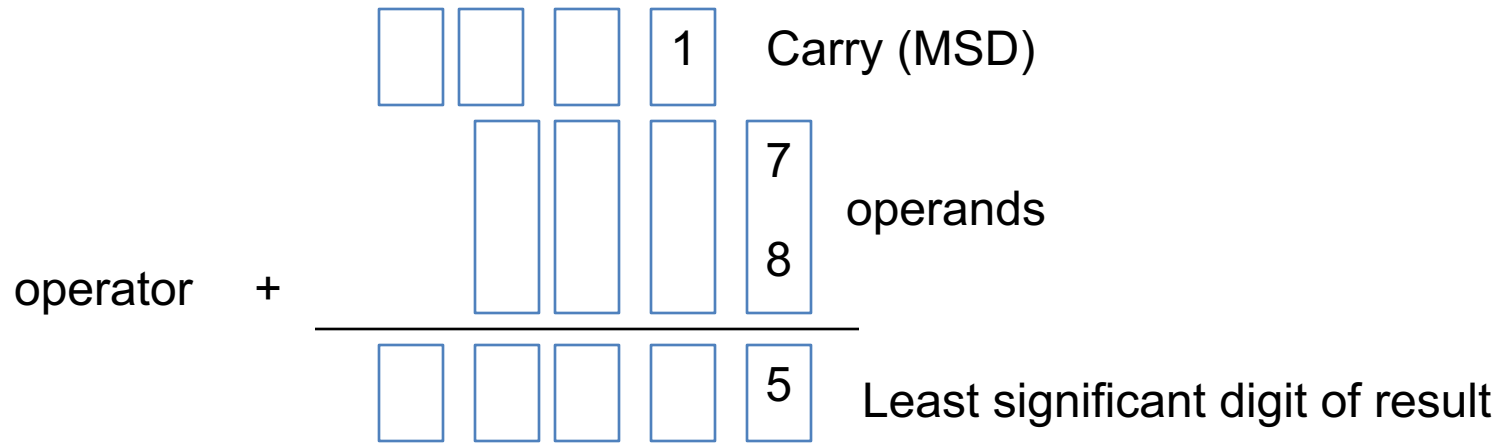


Algorithm

- **An algorithm (pronounced AL-go-rith-um) is a procedure or formula for solving a problem.**
- **In mathematics and computer science, an algorithm is a self-contained step-by-step set of operations to be performed.**
- **An algorithm is an effective method that can be expressed within a finite amount of space and time and in a well-defined formal language for calculating a function.**

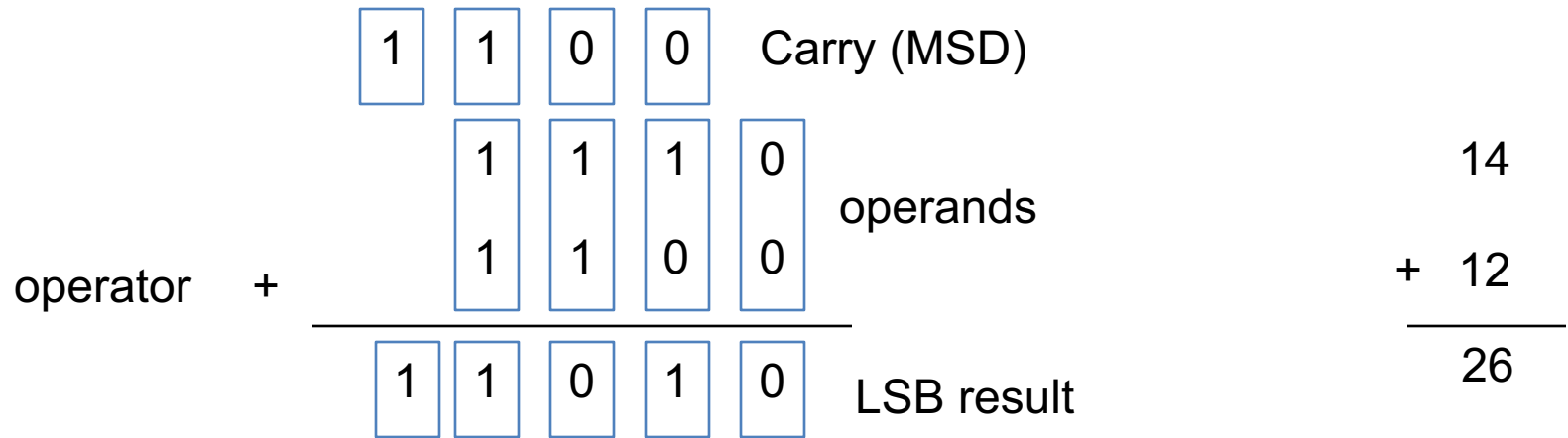


Algorithms early in life





Algorithms early in life (in binary)





A Simple Algorithm in Class

- **Count the number of students**



More interesting one, ...

- Betcha people in here share a birthday?

https://en.wikipedia.org/wiki/List_of_Presidents_of_the_United_States_by_date_of_birth

Presidents?

Abstraction

- **Detail removal**
 - “The act or process 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”



Henri Matisse "Naked Blue IV"



Experiment

Standard Time Zones of the World



Where are you from?

Possible Answers:

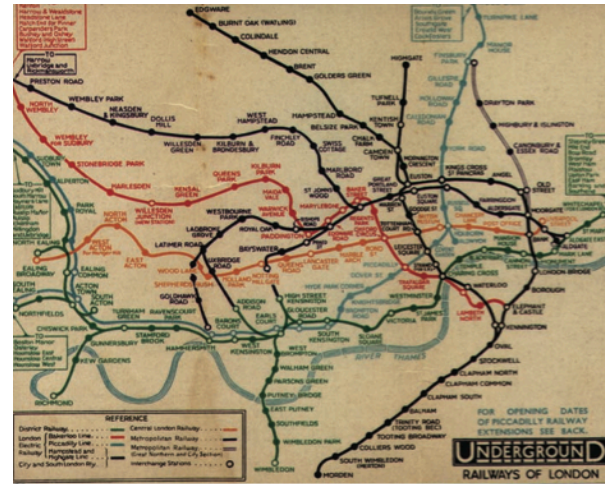
- China
- California
- The Bay Area
- San Mateo
- 1947 Center Street, Berkeley, CA
- 37.8693° N, 122.2696° W



All correct but different levels of abstraction!

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)

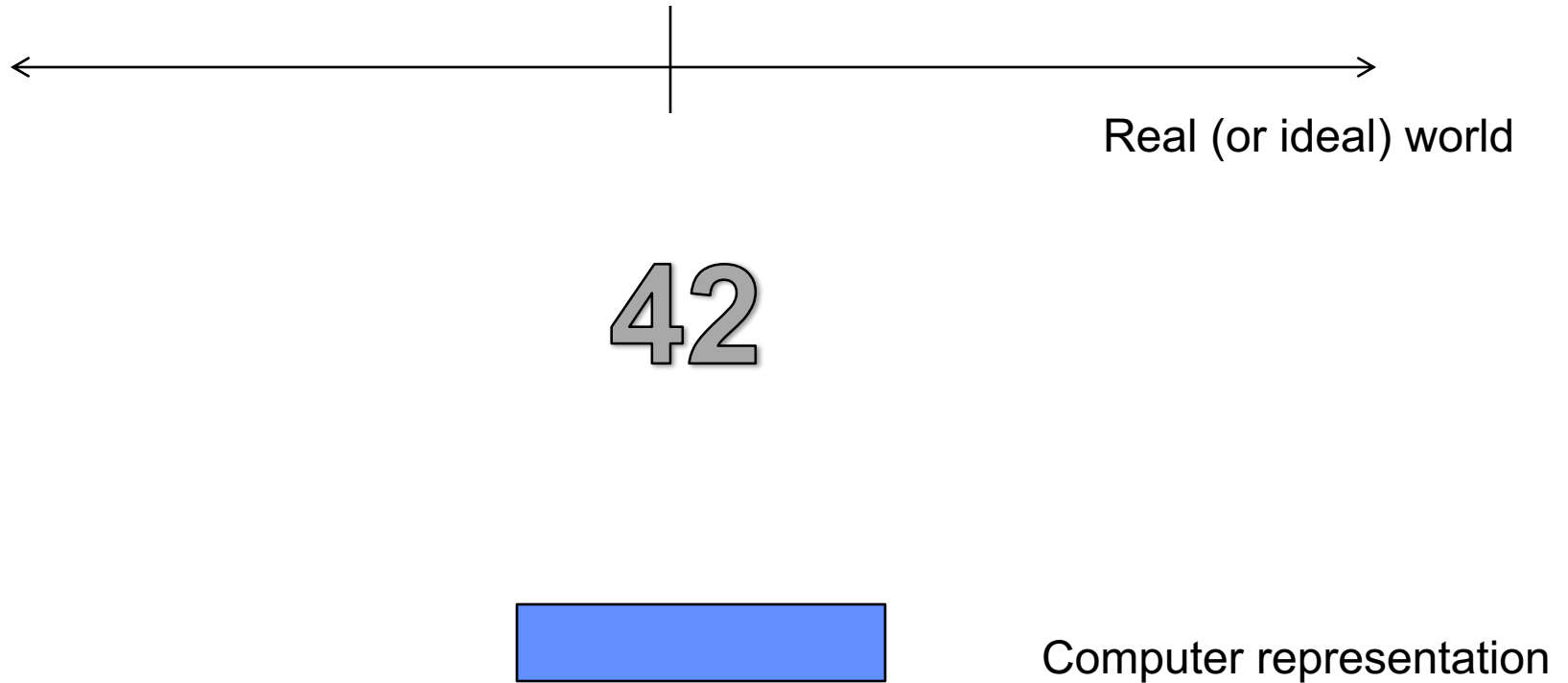
Below the abstraction line

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



Abstraction 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, ... , Closure (???)
- **Expressions are valid well-defined sets of operations on elements that produce a value of a type**



Questions

- What's the difference between '==' and '=' ?



Lab and HW this week

- Lab will get you to where you have a *program development environment*
 - Even on your computer
- HW will give practice and explain subtleties of **types, operators, and expressions**
 - In a program development environment



Question of the week

- How many “things” can you represent with **N** bits