# Lecture 1: The Revolution is here!
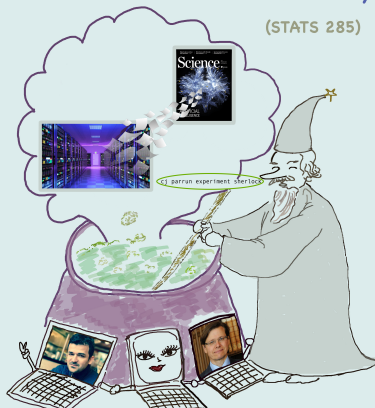
D Donoho/ H Monajemi
Stats 285 Stanford

20180925

# Stats 285 Fall 2018

# Outline

## Disclaimer

*This document contains images obtained by routine Google Images searches. Some of these images may perhaps be copyright. They are included here for educational noncommercial purposes and are considered to be covered by the doctrine of* **Fair Use**. *In any event they are easily available from Google Images.*

*It's not feasible to give full scholarly credit to the creators of these images. We hope they can be satisfied with the positive role they are playing in the educational process.*

**The Computing Discontinuity**
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

**Mobile is Eating the world**
Mobile Drives IT Revolution
AWS is Eating the World
New AWS Services are Proliferating

# The Mobile Revolution

**The Computing Discontinuity**
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

**Mobile is Eating the world**
Mobile Drives IT Revolution
AWS is Eating the World
New AWS Services are Proliferating

# Smartphones are Spreading Everywhere



**SMARTPHONE USERS: UP 800M**

The world in 2020
By 2020 80% of the adults on earth will have a smartphone

**The Computing Discontinuity**
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

Mobile is Eating the world
**Mobile Drives IT Revolution**
AWS is Eating the World
New AWS Services are Proliferating

# 24/7 Deluge Spawns Global Computational Services

**The Computing Discontinuity**
**The Revolution in Computational Science**
**Case Study: Deep Learning**
**Resistance**
**Painless Computational Experiments**

Mobile is Eating the world
**Mobile Drives IT Revolution**
AWS is Eating the World
New AWS Services are Proliferating

## Cloud Paradigm

Cloud Paradigm:

▶ Billions of smart devices each drive queries to cloud servers

▶ Millions of business relying on cloud for all needs

Symbiosis of cloud and economy is *lasting* and *disruptive*.

**The Computing Discontinuity**
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

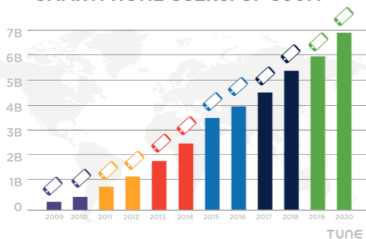Mobile is Eating the world
Mobile Drives IT Revolution
**AWS is Eating the World**
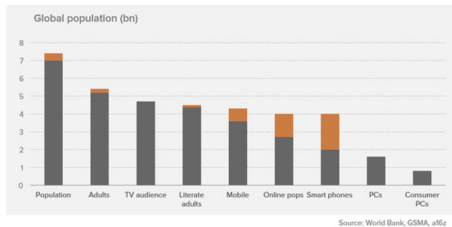New AWS Services are Proliferating

# AWS is Eating the world: Stock Market



World's
**RICHEST PERSON**
"JEFF BEZOS"
iTechHacks

amazon



TECH

TECH | MOBILE | SOCIAL MEDIA | ENTERPRISE | CYBERSECURITY | TECH G...

## Amazon shares soar after massive earnings beat

- Amazon reported its third quarter results Thursday after the bell.
- It was a huge beat across the board.
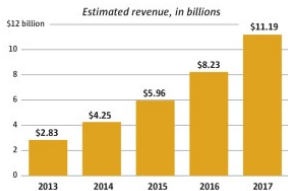- Amazon shares jumped over 7 percent in after hours trading.

Eugene Kim | @eugenekim222
Published 3:24 PM ET Thu, 26 Oct 2017 | Updated 6:55 PM ET Thu, 26 Oct 2017

CNBC

**The Computing Discontinuity**
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

Mobile is Eating the world
Mobile Drives IT Revolution
**AWS is Eating the World**
New AWS Services are Proliferating

# AWS is Eating the World, II



**Amazon Web Services sales**
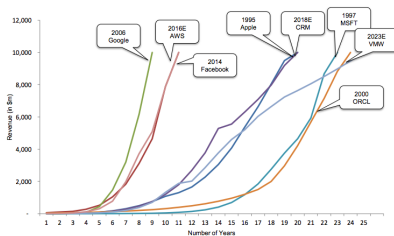
Amazon will break out specific sales data for AWS on Thursday for the first time. Here's Robert W. Baird & Co. analyst Colin Sebastian estimates.

*Estimated revenue, in billions*

Source: Robert W. Baird & Co.      KELLY SHEA / THE SEATTLE TIMES



Figure 9: AWS is the Fastest-Growing Enterprise Technology Company Ever

Source: Deutsche Bank Estimates, Public Company Filings

**The Computing Discontinuity**
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

Mobile is Eating the world
Mobile Drives IT Revolution
**AWS is Eating the World**
New AWS Services are Proliferating

# AWS is Eating the World: III

**The Computing Discontinuity**
**The Revolution in Computational Science**
**Case Study: Deep Learning**
**Resistance**
**Painless Computational Experiments**

Mobile is Eating the world
Mobile Drives IT Revolution
AWS is Eating the World
**New AWS Services are Proliferating**

# AWS Services Are Ubiquitous

**The Computing Discontinuity**
**The Revolution in Computational Science**
**Case Study: Deep Learning**
**Resistance**
**Painless Computational Experiments**

Mobile is Eating the world
Mobile Drives IT Revolution
AWS is Eating the World
**New AWS Services are Proliferating**

## AWS Services are Proliferating

### AWS Pace of Innovation

AWS has been continually expanding its services to support virtually any cloud workload, and it now has more than 90 services that range from compute, storage, networking, database, analytics, application services, deployment, management, developer, mobile, Internet of Things (IoT), Artificial Intelligence (AI), security, hybrid and enterprise applications. AWS has launched a total of 236 new features and/or services year to date* - for a total of 3,149 new features and/or services since inception in 2006.
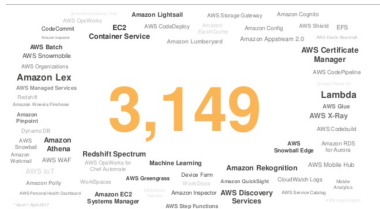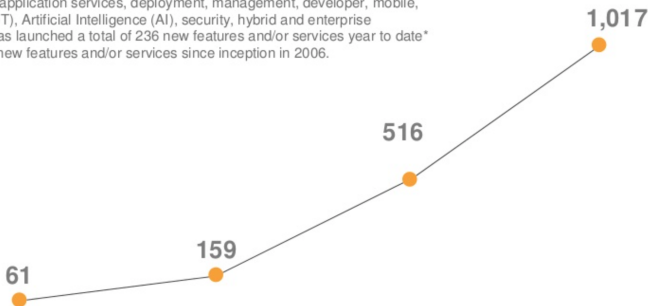
1,017

516

159

61

**The Computing Discontinuity**
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

Mobile is Eating the world
Mobile Drives IT Revolution
AWS is Eating the World
**New AWS Services are Proliferating**

## Stack Paradigm I

Stack Paradigm:

▶ Organizations combine software components from other providers in a `stack`

▶ Massive new capabilities emerge by hybridizing components

Examples:

▶ Uber (next slide)

▶ Netflix relies on AWS

▶ Snap, Dropbox etc. small teams

**The Computing Discontinuity**
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

Mobile is Eating the world
Mobile Drives IT Revolution
AWS is Eating the World
**New AWS Services are Proliferating**

# Stack Paradigm II



Uber doesn't own their cars. They also don't directly employ their own drivers. So, one might ask, what do they own exactly as a core asset? The core application and ecosystem around the Uber experience is their primary asset and differentiator. But to deliver that experience, they apply rigorous focus.

At the practical level, when you look at the technology components of Uber's world-renowned app, they decided to rely on other core platforms and technologies to power many of the key elements.

Jeetu Patel, *Software is still eating the world*, *TechCrunch*, Jan 2016

## Explosion of Computational Resources

Cloud Paradigm:

- ▶ Billions of smart devices each drive queries to cloud servers
- ▶ Millions of business relying on cloud for all needs

Symbiosis of cloud and economy is *lasting* and *disruptive*.

Cloud provides *any user* **same-day** delivery:

- ▶ Tens to hundreds of thousands of hours of CPU
- ▶ Pennies per CPU hour

Any user can consume *1 Million CPU hours* over a few days for a few \$10K's.

## Massive Computational Power Will Transform *Science*
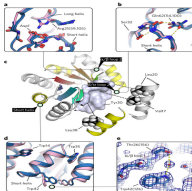
**Traditionally:**

- ▶ Deduction (in math)

- ▶ Induction (in physical sciences)

**Emerging new approach:**

- ▶ Massive computational experiments

# Massive Computations in Science

Traditionally computational fields



Protein Design

(Huang et al. 2016)

AI

(Alahi et al. 2016)

Oil Field Devel.

(Shirangi et al. 2015)

## Massive Computations in Science

Traditionally <span style="color:red">non-</span> computational field – Mathematics



Borwein/Bailey



Borwein/Devlin



Individual Articles

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

**The Sudden Emergence of Deep Learning**
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
Framework Wars

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

**The Sudden Emergence of Deep Learning**
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
Framework Wars

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

**The Sudden Emergence of Deep Learning**
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
Framework Wars

Andrej Karpathy
@karpathy

Follow

Came to visit first class of @cs231n at Stanford. 2015: 150 students, 2016: 350, this year: 750. #aiinterestsingularity

12:11 PM - 4 Apr 2017

155 Retweets  623 Likes

19    155    623

michael_nielsen @michael_nielsen · Apr 4
Replying to @karpathy @cs231n
Faster than Moore's Law. At this rate - doubling each year - in 24 years everyone on Earth will be enrolled :-)

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

**The Sudden Emergence of Deep Learning**
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
Framework Wars

## Synchronies, 1

Over same timeframe – 2010-2014

▶ Instagram, Snapchat emerge to global prominence

▶ Deep Learning catapults to global attention

Coincides with emergence of

▶ Smartphone photography

▶ Cloud computing

▶ Cloud storage of selfie/smartphone photography

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

**The Sudden Emergence of Deep Learning**
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
Framework Wars

## Synchronies, 2

*"Six decades into the computer revolution, four decades since the invention of the microprocessor, and two decades into the rise of the modern Internet, all of the technology required to transform industries through software finally works and can be widely delivered at global scale."*
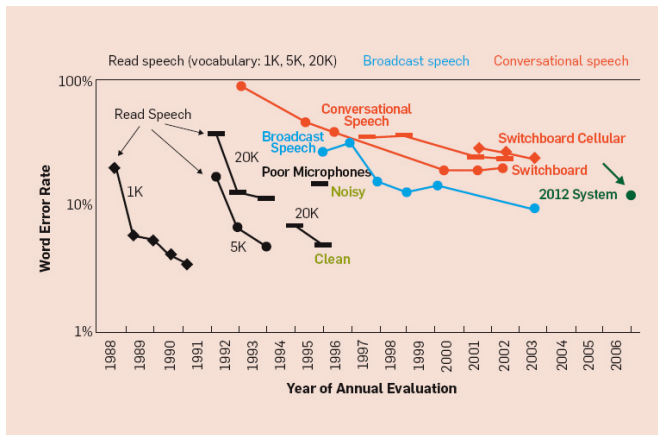
Marc Andreesen - WSJ - 2011

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
**The Slow Emergence of the Common Task Framework**
CTF Goes Mainstream
Lessons from Case Study
Framework Wars

## Common Task Framework (1980's)

Under CTF we have the following ingredients

(a) A **publicly available training dataset** involving, for each observation, a list of (possibly many) feature measurements, and a class label for that observation.

(b) A set of **enrolled competitors** whose **common task** is to **infer** a class **prediction rule from the training data**.

(c) A **scoring referee**, to which competitors can submit their prediction rule. The referee runs the prediction rule against a testing dataset which is sequestered behind a Chinese wall. The referee objectively and automatically reports the score achieved by the submitted rule.

See Mark Liberman's description (Liberman, 2009).

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
Framework Wars

# CTF *Really* Works!

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
**The Slow Emergence of the Common Task Framework**
CTF Goes Mainstream
Lessons from Case Study
Framework Wars
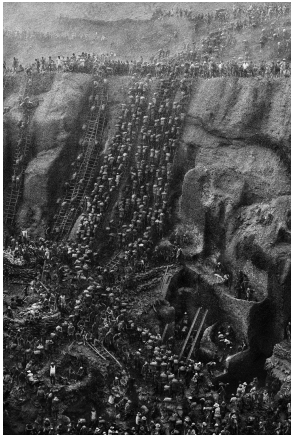
## CTF Lifestyle – 1

1. Researchers set up local copies of Challenge
   - ▶ Data – Training, Test carved out of public dataset
   - ▶ Scoring – same as challenge scoring rule
2. Researcher's job: *'tuning models'*
   - ▶ Think up a family of model variations – *'tweak's*
   - ▶ Run a full *'experiment'* – suite of tweaks – *'grid'*
   - ▶ Score each tweak
   - ▶ Submit best-scoring result to central authority
3. Successful researchers perpetually motivated by
   *Game-ification*: tweaking, scoring, winning.
4. Researchers who tweak more often, win more often!.
5. If easier to implement tweaks and faster to evaluate them,
   more likely to win!.

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
**The Slow Emergence of the Common Task Framework**
CTF Goes Mainstream
Lessons from Case Study
Framework Wars

## CTF Lifestyle – 2



Sebastiao Salgado, *Work*

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
The Slow Emergence of the Common Task Framework
**CTF Goes Mainstream**
Lessons from Case Study
Framework Wars

## CTF Goes Mainstream

1. Netflix Challenge (2009)
   $1 Million Prize

2. Kaggle (2010)
   1 Million'th competitor expected Sept. 2017

3. Fei-Fei Li masterminds ImageNet 2008-2010

4. Hinton's Deep Learning Team wins ImageNet 2012

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
Framework Wars

## ImageNet Classification Error (Top 5)

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
The Slow Emergence of the Common Task Framework
**CTF Goes Mainstream**
Lessons from Case Study
Framework Wars

Andrej Karpathy ✓
@karpathy

Follow

You can now understand state of the art AI with before high school math. You forward a neural net and repeat guess&check. works well enough.

12:53 PM - 14 Mar 2017

**50** Retweets  **207** Likes

💬 12          ⟲ 50          ♡ 207

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
The Slow Emergence of the Common Task Framework
**CTF Goes Mainstream**
Lessons from Case Study
Framework Wars

# Researchers Preparing for NIPS 2017



Sebastiao Salgado, *Work*

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
**Lessons from Case Study**
Framework Wars

# Lessons from Deep Learning Case Study

1. *Researchers who tweak more often, win more often!*
2. *If easier to implement tweaks and faster to evaluate them, more likely to win!*
3. Successful Research Environment
   - ▶ Easy to tweak models
   - ▶ Easy to score tweaks
   - ▶ Fast to score tweaks
4. Successful researchers perpetually motivated by *Game-ification*: tweaking, scoring, winning.
5. Easier to stay motivated when easier and more comfortable to play the game.
   - ▶ Elegant expression of tweaks
   - ▶ Rapid turn-around for scoring

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
**Framework Wars**

## Framework Wars – 1

Influential Frameworks for Deep Learning

▶ **Matlab**
– pre-framework

▶ **TensorFlow**
– open source (Originally by Google Brain)

▶ **Torch**
– scientific computing framework written in Lua

▶ **PyTorch**
– Python package for scientific computing (310 contributors)

▶ **Keras**
– A Python wrapper around TensorFlow, CNTK and Theano

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
Framework Wars

# Framework Wars – 2

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
**Framework Wars**

# Framework Wars – 3

The Computing Discontinuity
The Revolution in Computational Science
**Case Study: Deep Learning**
Resistance
Painless Computational Experiments

The Sudden Emergence of Deep Learning
The Slow Emergence of the Common Task Framework
CTF Goes Mainstream
Lessons from Case Study
**Framework Wars**

## Framework Wars - 4

The real action is all in frameworks

1. Dream up, test, and publish better ...
   - ▶ Types of models
   - ▶ Types of tweaks
   - ▶ Properties for evaluation

2. Implement better *frameworks* ...
   - ▶ More elegant expression of models, tweaks
   - ▶ Distributed Learning across clusters
   - ▶ Smoother collection and analysis of results

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
**Resistance**
Painless Computational Experiments

**Intellectual impoverishment**
Solution: The Great Enrichment

## Resistance – 1

*We are at a university!*

1. Q: *Where's the intellectual activity in tuning?*
2. Q: *I didn't come here to do hard manual labor!*
3. Q: *I didn't come here to compete as mindless drones!*

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
**Resistance**
Painless Computational Experiments

**Intellectual impoverishment**
Solution: The Great Enrichment

## Resistance – 2

*We are at a university!*

1. Q: *Where's the intellectual activity in tuning?*
2. Q: *I didn't come here to do hard manual labor!*
3. Q: *I didn't come here to compete as mindless drones!*

What we *see*:



Sebastiao Salgado, *Work*

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
**Resistance**
Painless Computational Experiments

**Intellectual impoverishment**
Solution: The Great Enrichment

## Resistance 3

*We are at a university!*

1. Q: *Where's the intellectual activity in tuning?*
2. Q: *I didn't come here to do hard manual labor!*
3. Q: *I didn't come here to compete as mindless drones!*

What we **imagine**:

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
**Resistance**
Painless Computational Experiments

**Intellectual impoverishment**
Solution: The Great Enrichment

# Metaphor: Computers as Slavery

Traditionally, 'using computers' involves interactively running programs (Excel, Point-and-click)

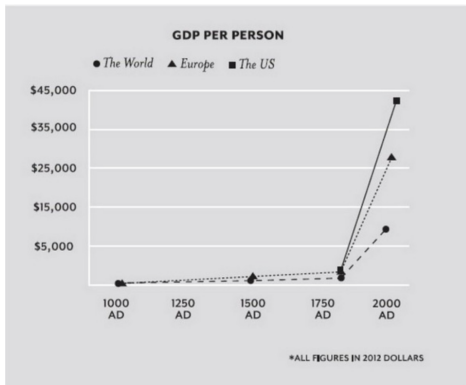Claerbout's Dictum: "... dependence on an interactive program can be a form of slavery"
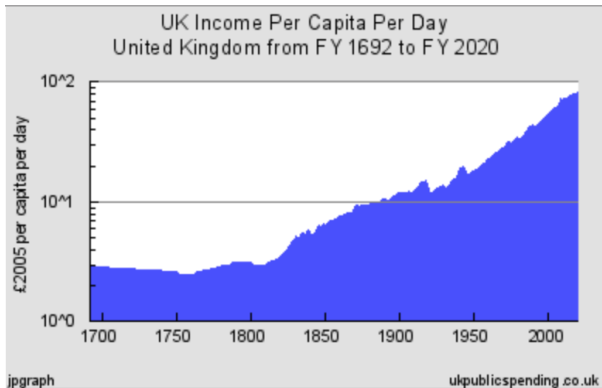
http://sepwww.stanford.edu/sep/jon/reproducible.html



Photo: Jon Claerbout    Cartoon: http://fritsAhlefeldt.com

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
**Resistance**
Painless Computational Experiments

Intellectual impoverishment
Solution: The Great Enrichment

# Digression: The Great Enrichment (Deidre McKloskey) 1

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
**Resistance**
Painless Computational Experiments

Intellectual impoverishment
Solution: The Great Enrichment

# Digression: The Great Enrichment (Deidre McKloskey) 2

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
**Resistance**
Painless Computational Experiments

Intellectual impoverishment
**Solution: The Great Enrichment**

# The Great IT Enrichment – 1

The Computing Discontinuity
The Revolution in Computational Science
Case Study: Deep Learning
**Resistance**
Painless Computational Experiments

Intellectual impoverishment
**Solution: The Great Enrichment**

## The Great IT Enrichment - 2

Our vision.

*The intellectual poverty of the old interactive 'Excel'-era paradigm was real, but will be transcended.*

*New and better and more powerful abstractions will lift us out of the mud and out of slavery.*

## Coming Soon to a Scientifc field near you

In the near future,

- ▶ Scientific research will be transformed

    - ▶ *1 million CPU Hours* behind research papers and theses
    - ▶ *Widespread acceptance* of empirical/simulation evidence

- ▶ 1 million-hour hurdle manageable through *new frameworks*.

- ▶ Frameworks offer Convenient and Efficient

    - ▶ ... definition of experiments
    - ▶ ... management of jobs
    - ▶ ... gathering of results
    - ▶ ... analysis and presentation

- ▶ Output:
    - ▶ Better science
    - ▶ Better math

# Course Focus: Frameworks for Massive Experiments, 1

- ▶ Traditional issues
    - ▶ Experiments implicitly defined by executing unorganized code
    - ▶ Hard to understand what the baseline is, what variations are
    - ▶ Code dependencies unclear
    - ▶ Ordeal to get all the jobs to run, maybe gave up early
    - ▶ Tedious to harvest all the data, maybe missing some data
    - ▶ Confusing manual compilation and reporting
- ▶ Modern Frameworks
    - ▶ Systematic structure to coding
    - ▶ Base experiment clearly defined
    - ▶ Tweaks clearly defined
    - ▶ Code dependencies explicit
    - ▶ Grid of Jobs run systematically
    - ▶ Automatic transparent access of (cluster, AWS,...)
    - ▶ Data Harvested automatically to central data repository
    - ▶ Data analyzed automatically using defined tools

# Course Focus: Frameworks for Massive Experiments, 2

► Example Frameworks

    ► By individual research teams:
        ► ClusterJob – Hatef Monajemi
        ► CodaLab – Percy Liang

    ► By startups:
        ► Databricks
        ► Civis Analytics
        ► Domino Data Labs

# A Look Ahead: `https://stats285.github.io`

Guest Lectures

Tue, 10/02/2018
Mark Piercy
Stanford (SRCC)

Tue, 10/16/2018
Gregory Kurtzer
Sylabs

Tue, 10/23/2018
Ali Zaidi
Microsoft

Tue, 11/13/2018
Riccardo Murri
University of Zurich

Tue, 11/20/2018
Wes McKinney
Ursa Labs

| | |
|---|---|
| Mark Piercy | SRCC |
| Gregory Kurtzer | Sylabs |
| Ali Zaidi | Microsoft |
| Riccardo Murri | University of Zuerich |
| Wes McKinney | Ursa Labs |

# Global Economy → Computing → Science