# A Berkeley View of Big Data

Ion Stoica
UC Berkeley

BEARS
February 17, 2011

# Big Data is Massive…

- Facebook:
  - 130TB/day: user logs
  - 200-400TB/day: 83 million pictures

- Google: > 25 PB/day processed data

- Data generated by LHC: 1 PB/sec

- Total data created in 2010:
  1 ZettaByte (1,000,000 PB)/year
  - ~60% increase every year

amplab

# …and Grows Bigger and Bigger!

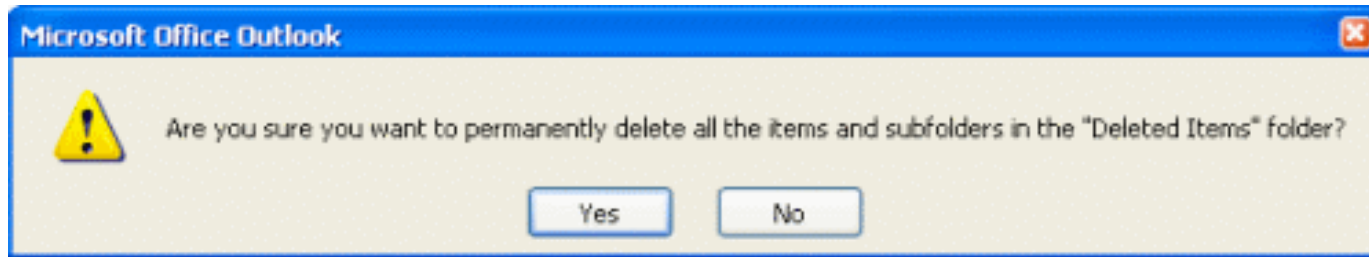- More and more devices

- More and more people

- Cheaper and cheaper storage
  - ~50% increase in GB/$ every year

 amplab

# …and Grows Bigger and Bigger!

- Log everything!
  - Don't always know what question you'll need to answer

- Hard to decide what to delete



  - Thankless decision: people know only when you are wrong!
  - "Climate Research Unit (CRU) scientists admit they threw away key data used in global warming calculations"

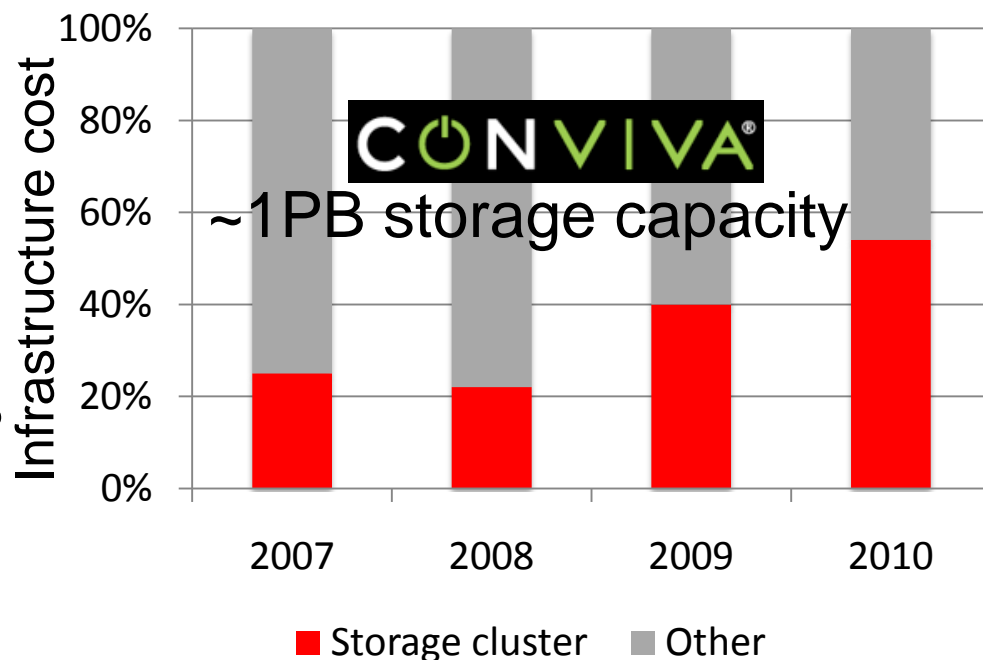- Stored data grows faster than GB/$

# What is Big Data?

Data that is <u>expensive</u> to manage, and hard to extract <u>value</u> from

- You don't need to be big to have big data problem!
  - Inadequate tools to analyze data
  - Data management may dominate infrastructure cost

amplab

# Big Data is not Cheap!

- Storing and managing 1PB data: $500K-$1M/ year
  - Facebook: 200 PB/year

- "Typical" cloud-based service startup (e.g., Conviva)
  - Log storage dominates infrastructure cost

~1PB storage capacity

# Hard to Extract Value from Data!

- Data is
  - Diverse, variety of sources
  - Uncurated, no schema, inconsistent semantics, syntax
  - Integration a huge challenge

- No easy way to get answers that are
  - High-quality
  - Timely

- Challenge: maximize value from data by getting best possible answers
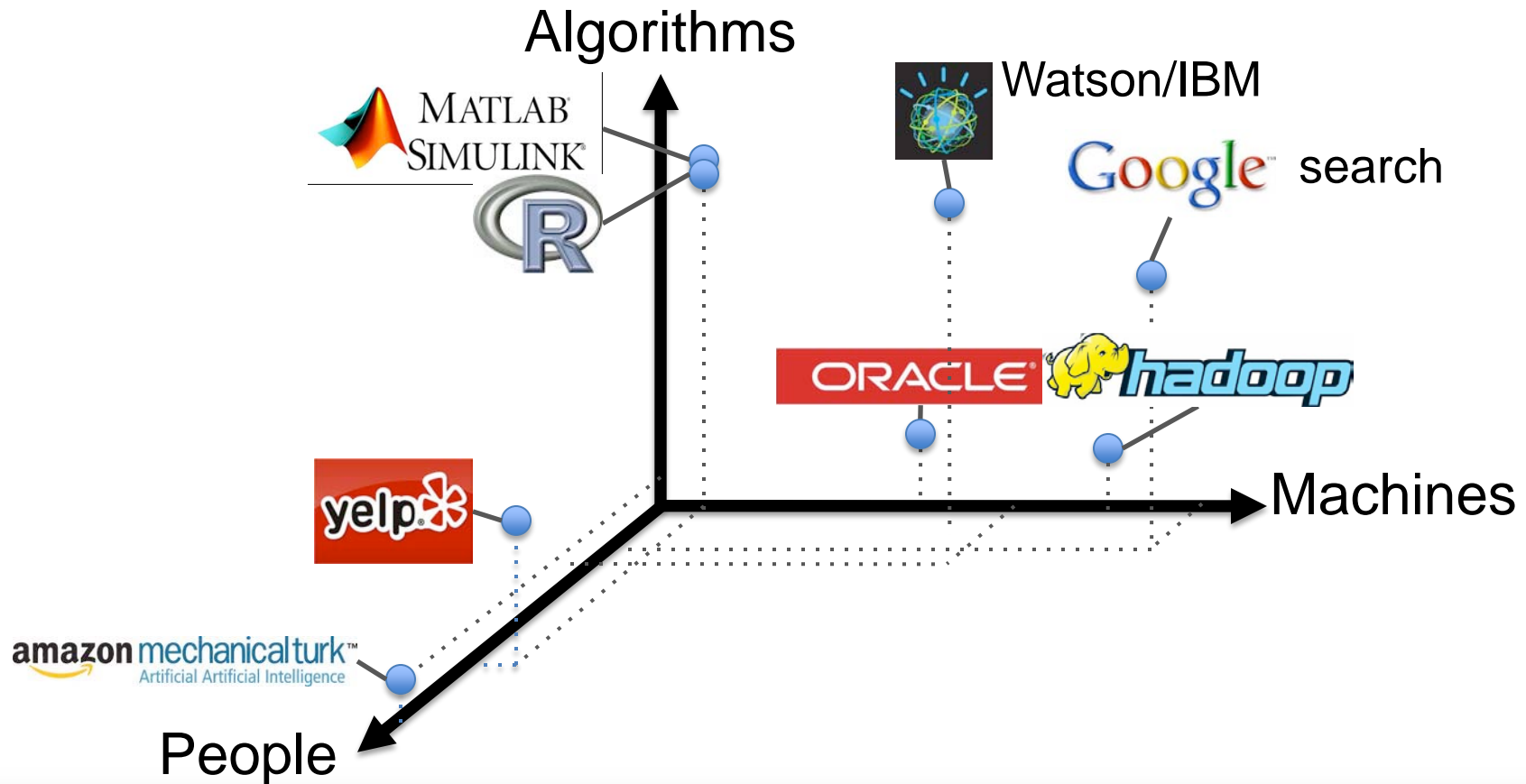
7

# Requires Multifaceted Approach

- Three dimensions to improve data analysis
  - Improving scale, efficiency, and quality of algorithms (**Algorithms)**
  - Scaling up datacenters (**Machines)**
  - Leverage human activity and intelligence (**People)**

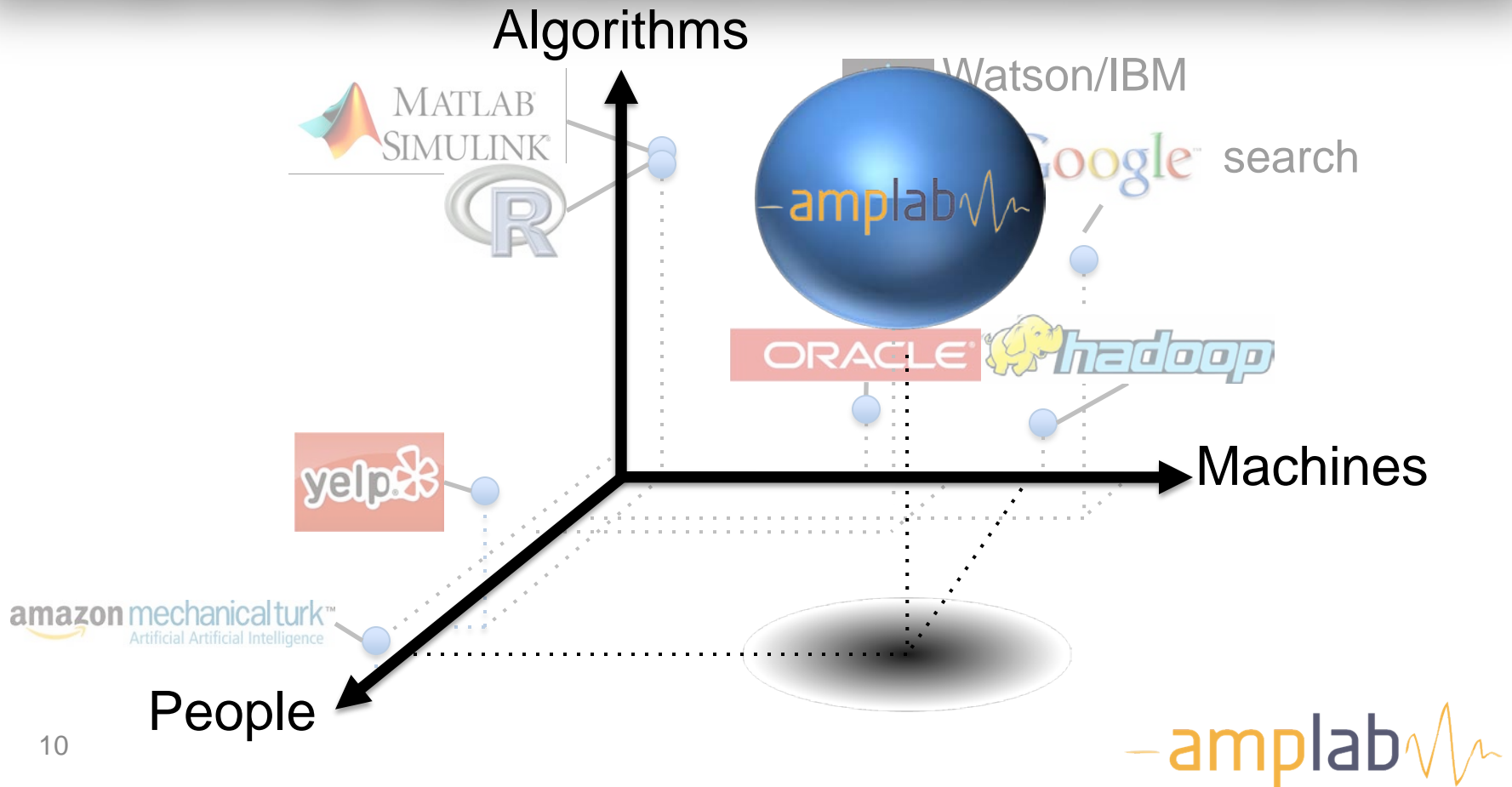- Need to adaptively and flexibly combine all three dimensions

amplab

- Today's apps: fixed point in solution space



Need techniques to dynamically pick best operating point

# The AMP Lab

## Make sense of data at scale by tightly integrating algorithms, machines, and people
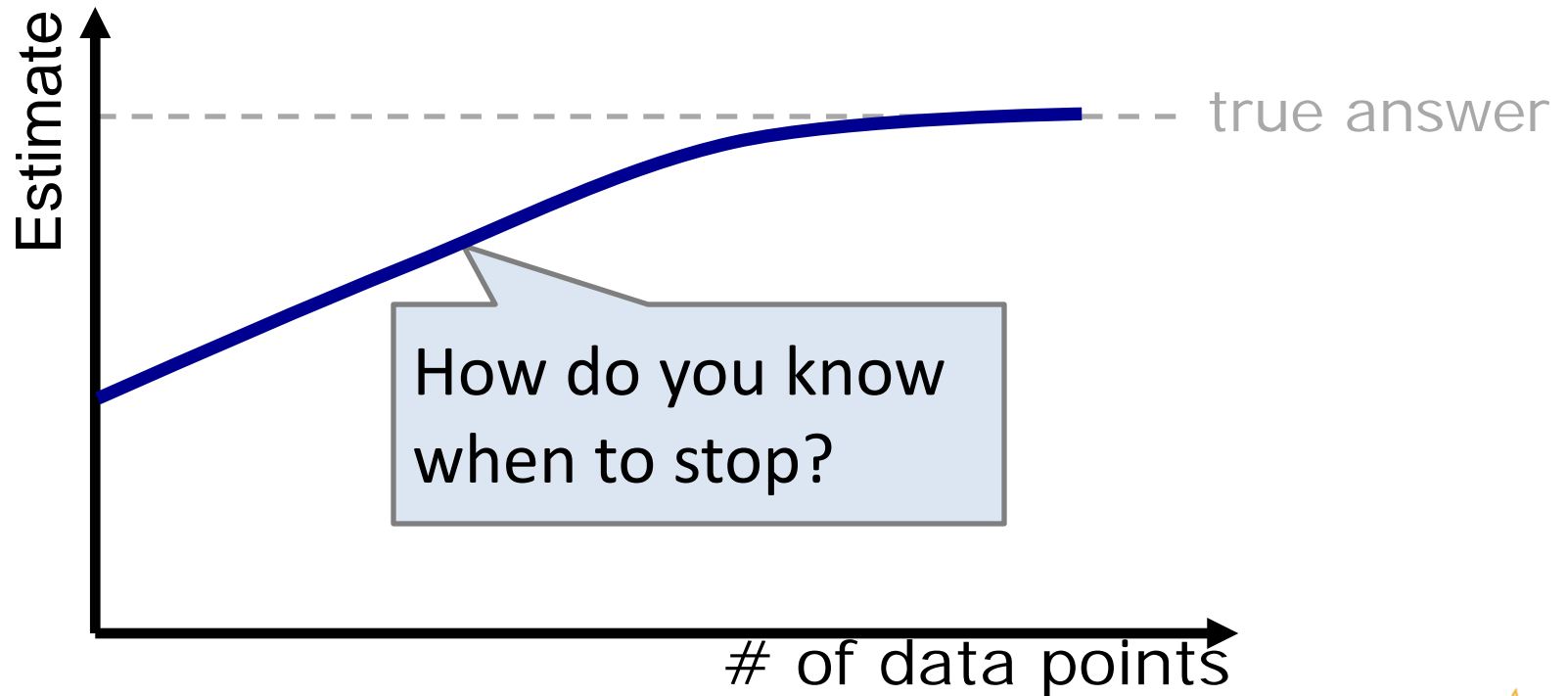
# AMP Faculty and Sponsors

- Faculty
  - Alex Bayen (mobile sensing platforms)
  - Armando Fox (systems)
  - Michael Franklin (databases): Director
  - Michael Jordan (machine learning): Co-director
  - Anthony Joseph (security & privacy)
  - Randy Katz (systems)
  - David Patterson (systems)
  - Ion Stoica (systems): Co-director
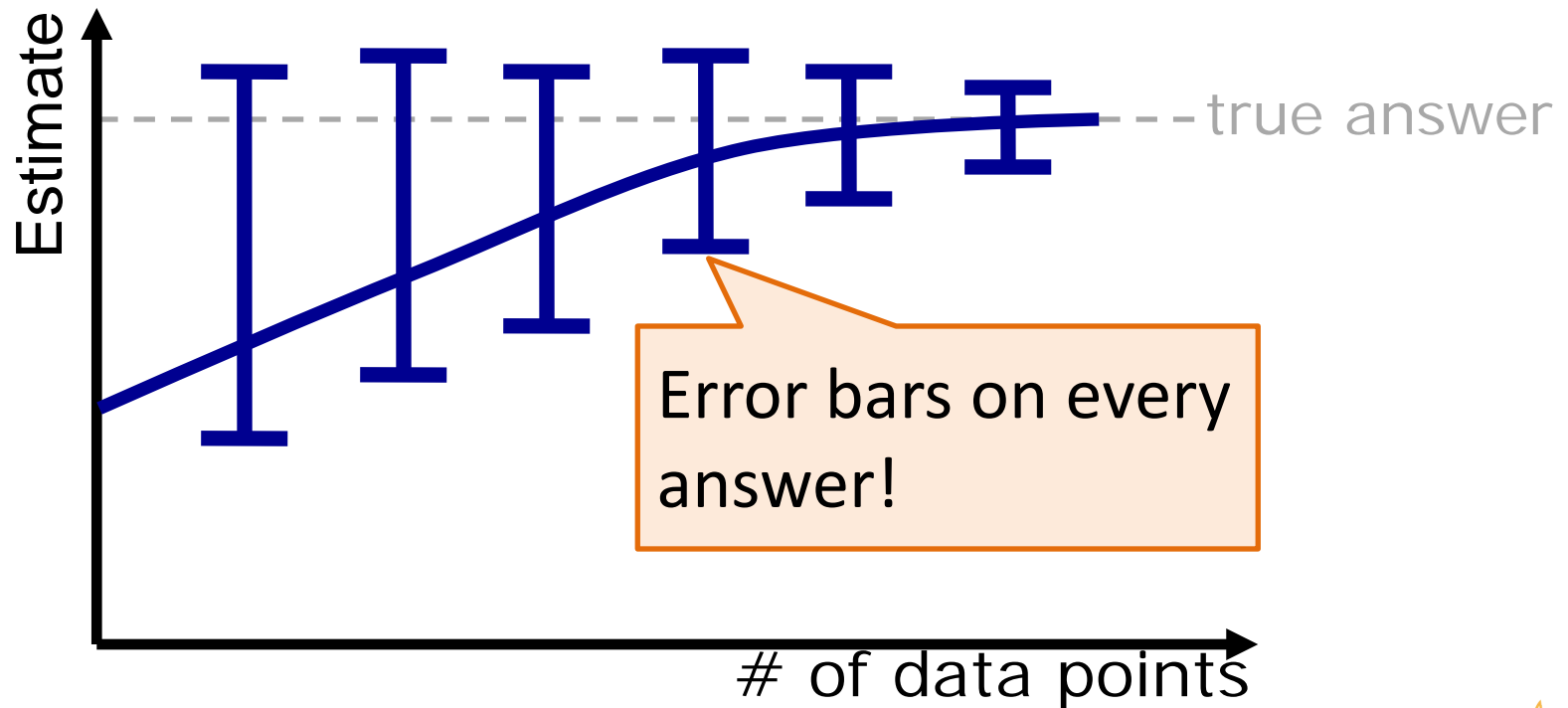  - Scott Shenker (networking)
- Sponsors:

- State-of-art Machine Learning (ML) algorithms do not scale
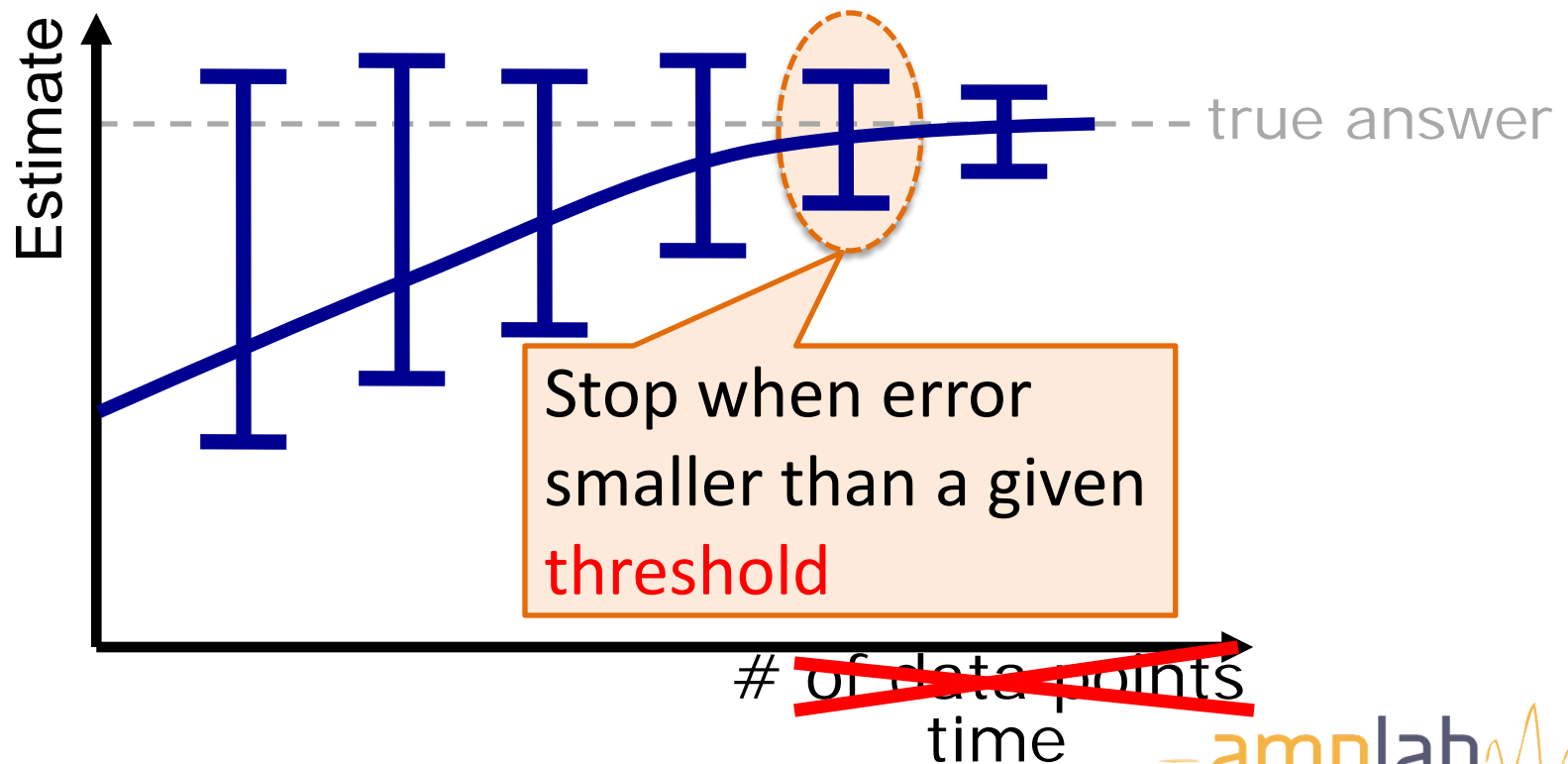  – Prohibitive to process all data points

- Given any problem, data and a budget
  - Immediate results with continuous improvement
  - Calibrate answer: provide error bars

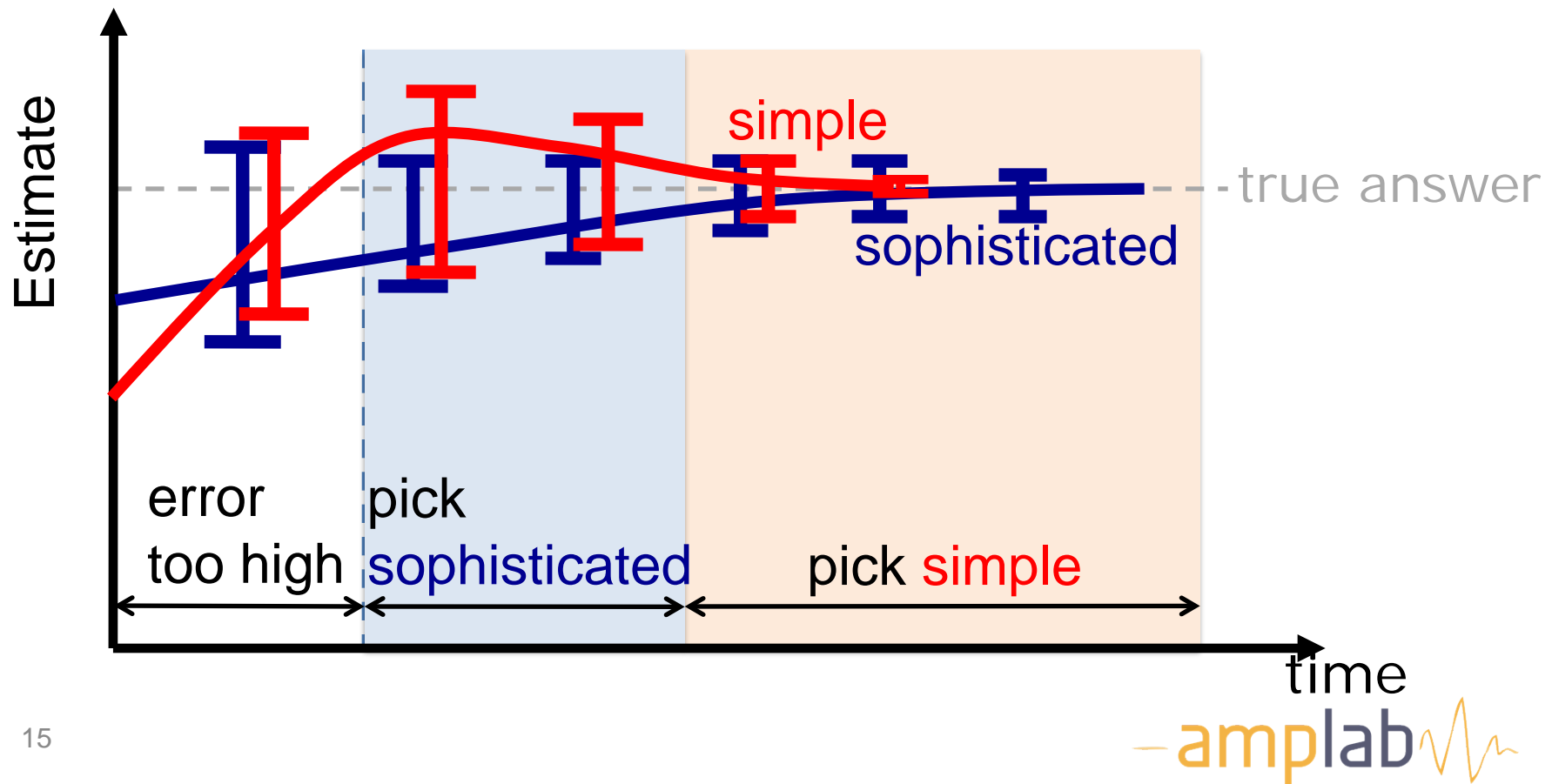- Given any problem, data and a time budget
  - Immediate results with continuous improvement
  - Calibrate answer: provide error bars



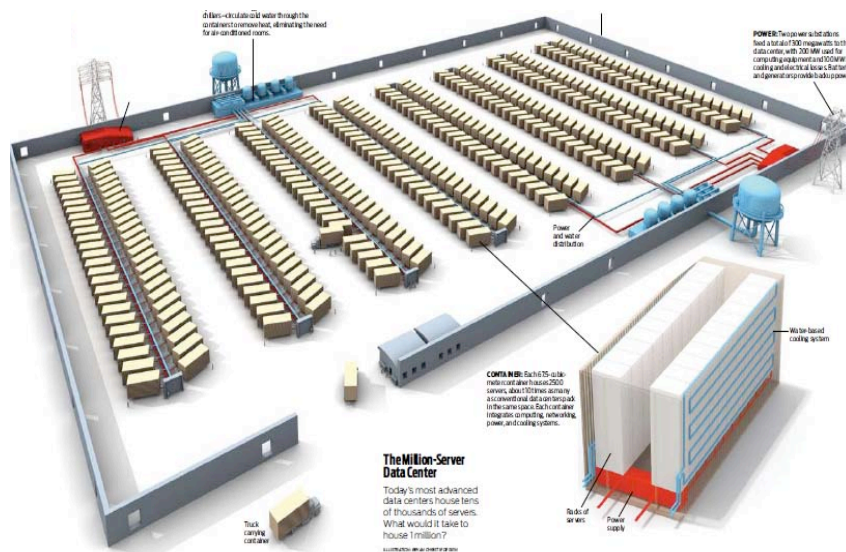Stop when error smaller than a given threshold

# of data points

time

amplab

- Given any problem, data and a time budget
  - Automatically pick the best algorithm

- "The datacenter as a computer" still in its infancy
    - Special purpose clusters, e.g., Hadoop cluster
    - Highly variable performance
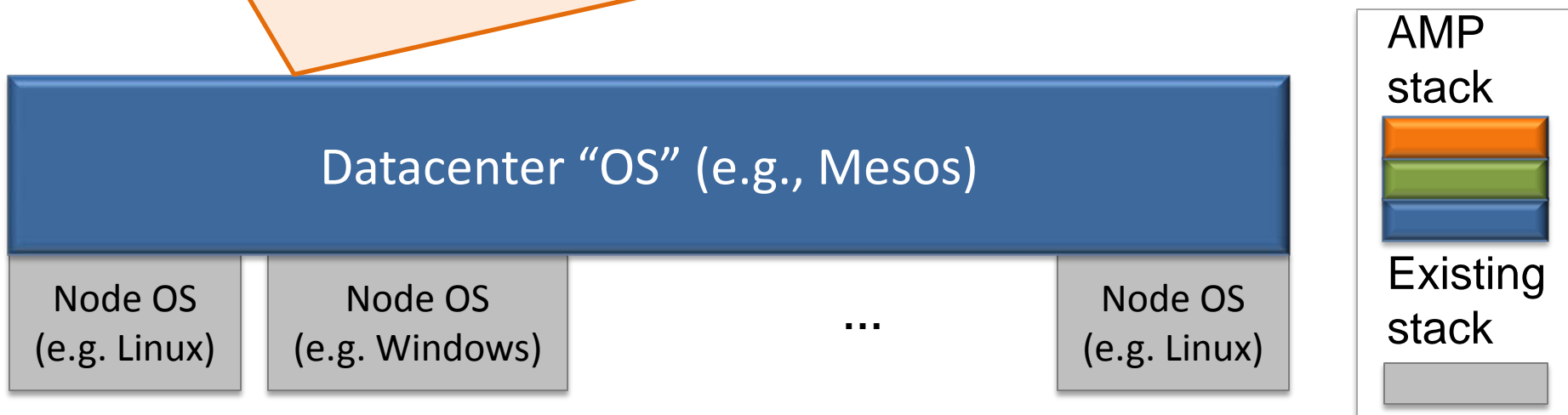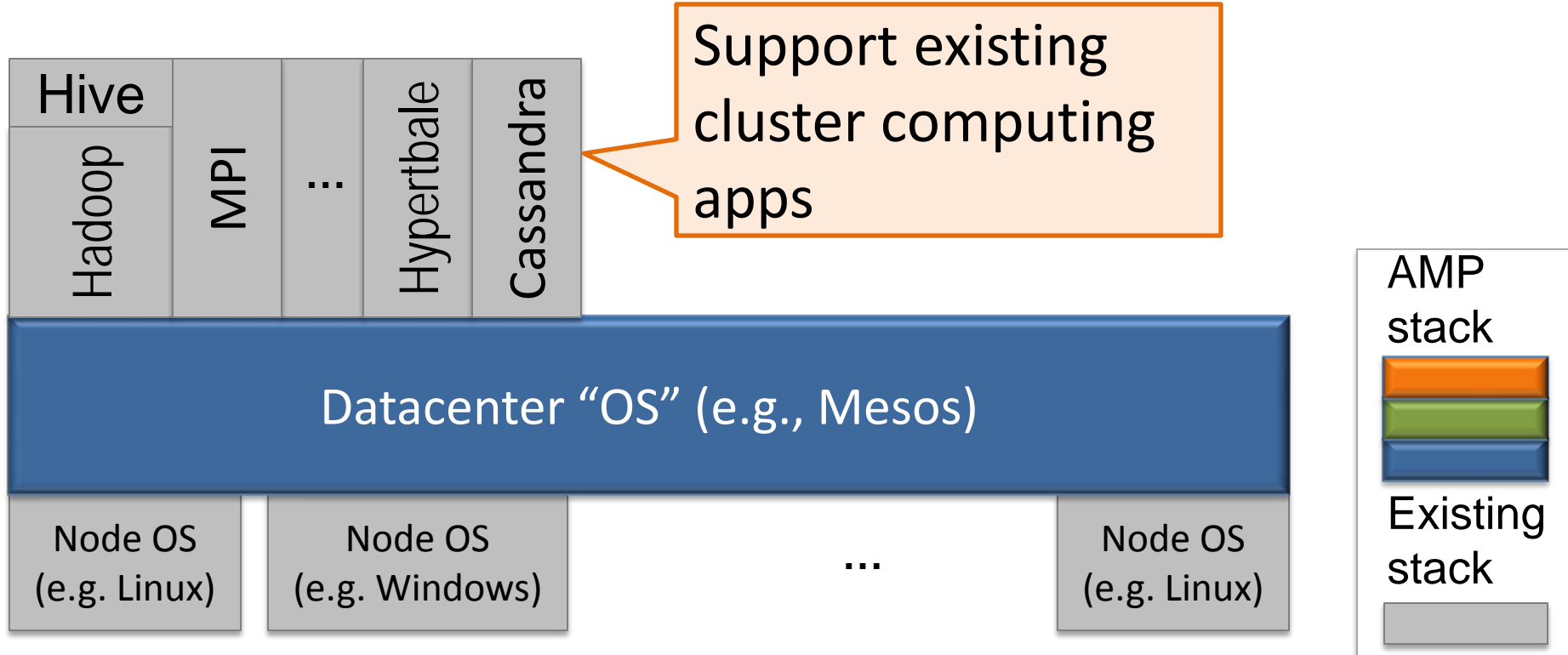    - Hard to program
    - Hard to debug



= ?

amplab

# Machines

- Make datacenter a real computer!

• Share datacenter between multiple cluster computing apps
• Provide new abstractions and services

Datacenter "OS" (e.g., Mesos)

| Node OS (e.g. Linux) | Node OS (e.g. Windows) | ... | Node OS (e.g. Linux) |
|---|---|---|---|

AMP stack

Existing stack

amplab

# Machines

- Make datacenter a real computer!



Hive

Hadoop | MPI | ... | Hypertbale | Cassandra

Support existing cluster computing apps

Datacenter "OS" (e.g., Mesos)

Node OS (e.g. Linux) | Node OS (e.g. Windows) | ... | Node OS (e.g. Linux)
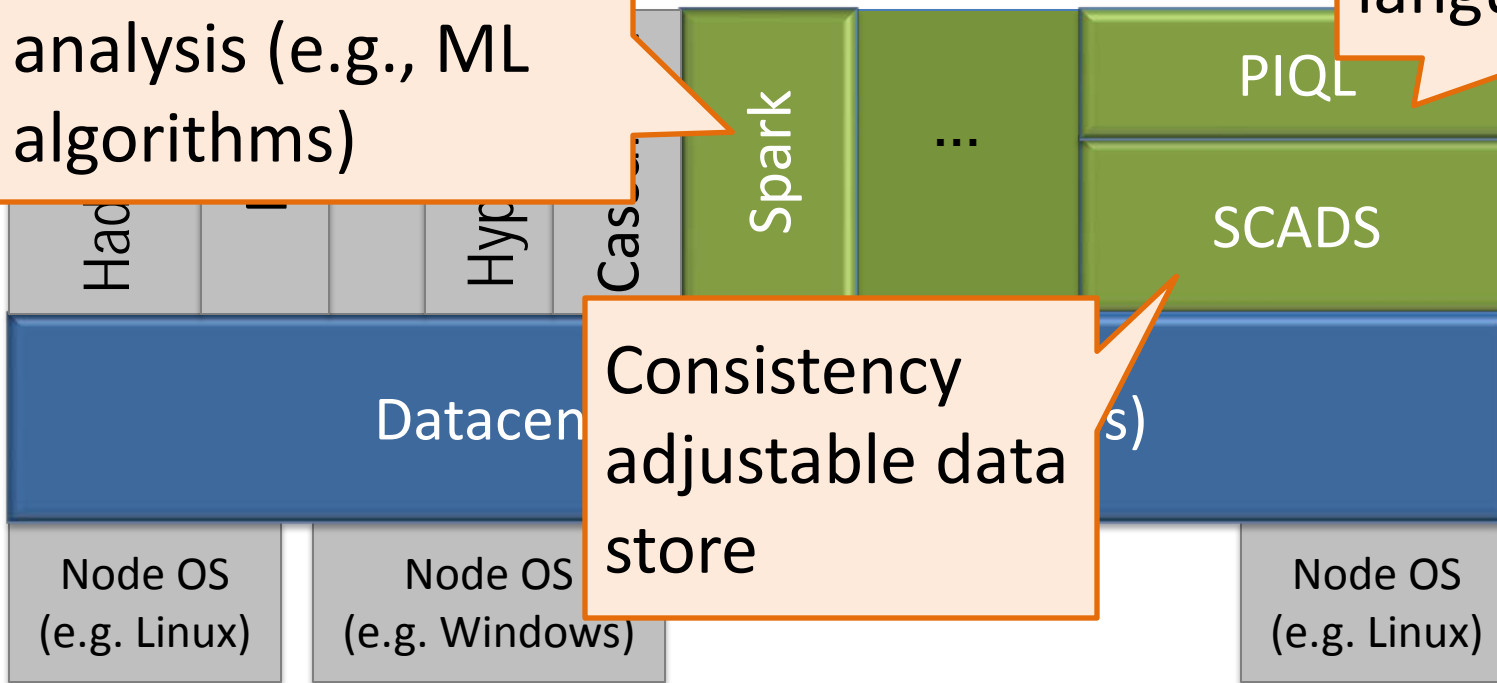
AMP stack

Existing stack

amplab

- Make datacenter a real computer!

Support interactive and iterative data analysis (e.g., ML algorithms)

Predictive & insightful query language
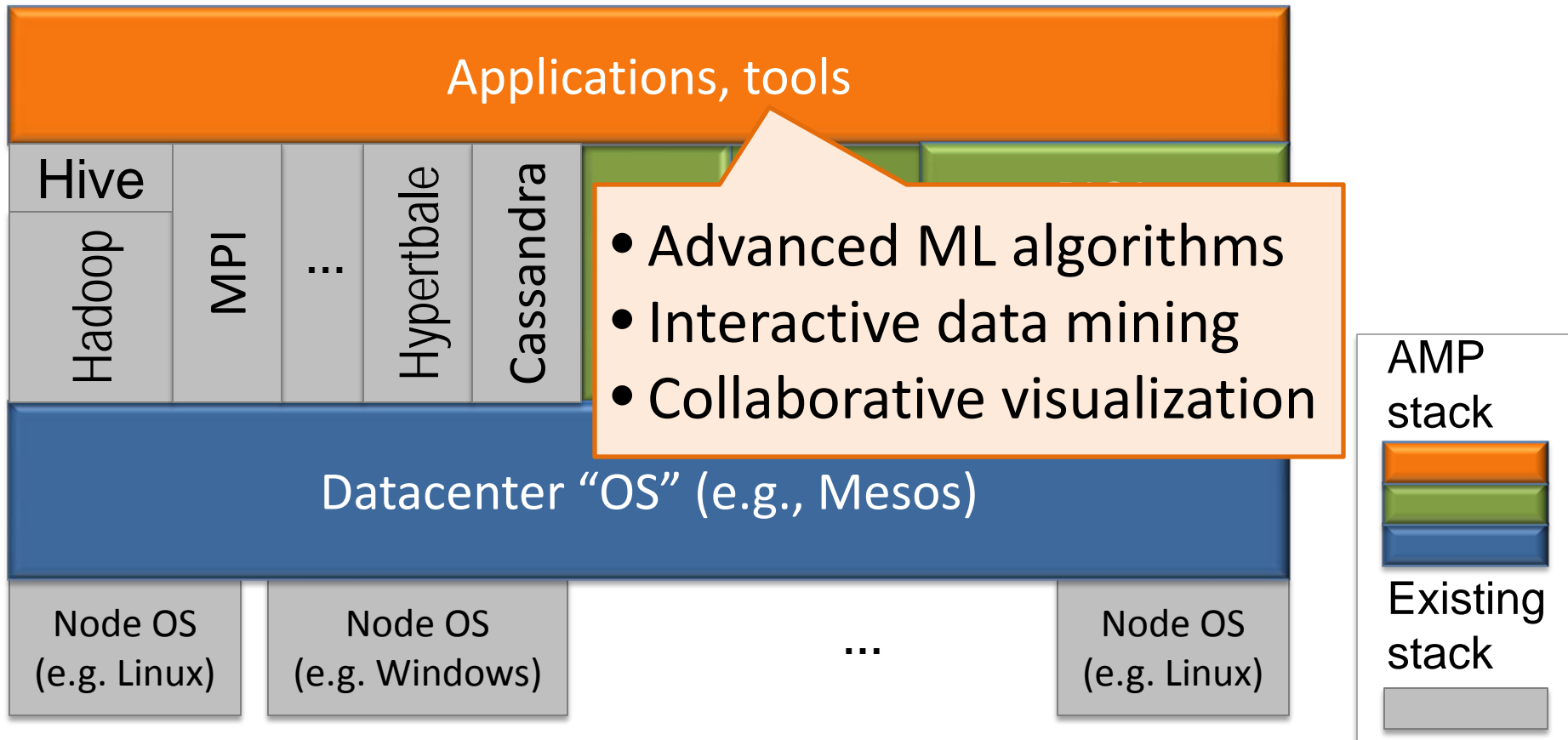
Consistency adjustable data store

Spark

...

PIQL

SCADS

Hadoop

Hypertable

Cassandra

Datacenter

Node OS (e.g. Linux)

Node OS (e.g. Windows)

Node OS (e.g. Linux)

AMP stack

Existing stack

19

amplab

# Machines

- Make datacenter a real computer!



Applications, tools

Hive

Hadoop

MPI

...

Hypertbale

Cassandra

- Advanced ML algorithms
- Interactive data mining
- Collaborative visualization

Datacenter "OS" (e.g., Mesos)

Node OS (e.g. Linux)

Node OS (e.g. Windows)

...

Node OS (e.g. Linux)

AMP stack

Existing stack

amplab

# People

- Humans can make sense of messy data!
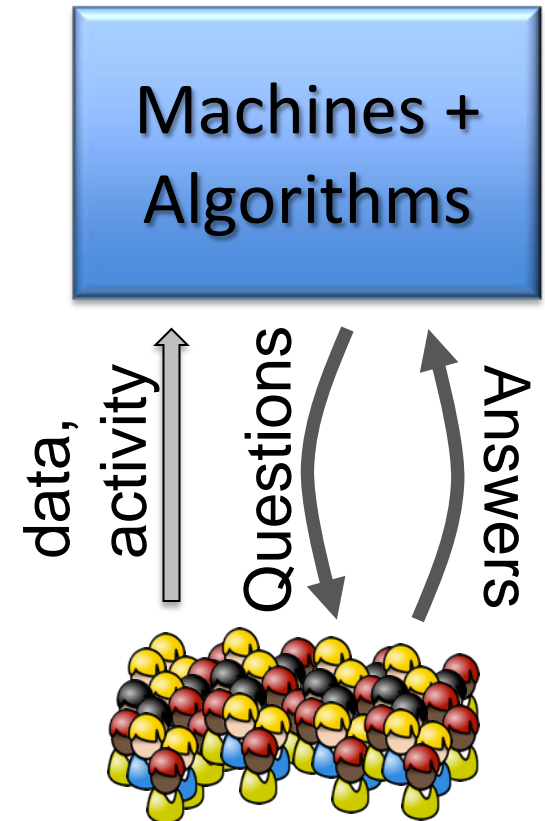
# <u>P</u>eople

- Make people an integrated part of the system!
  - Leverage human activity
  - Leverage human intelligence (crowdsourcing):
    - Curate and clean dirty data
    - Answer imprecise questions
    - Test and improve algorithms

- Challenge
  - Inconsistent answer quality in all dimensions (e.g., type of question, time, cost)

Machines + Algorithms

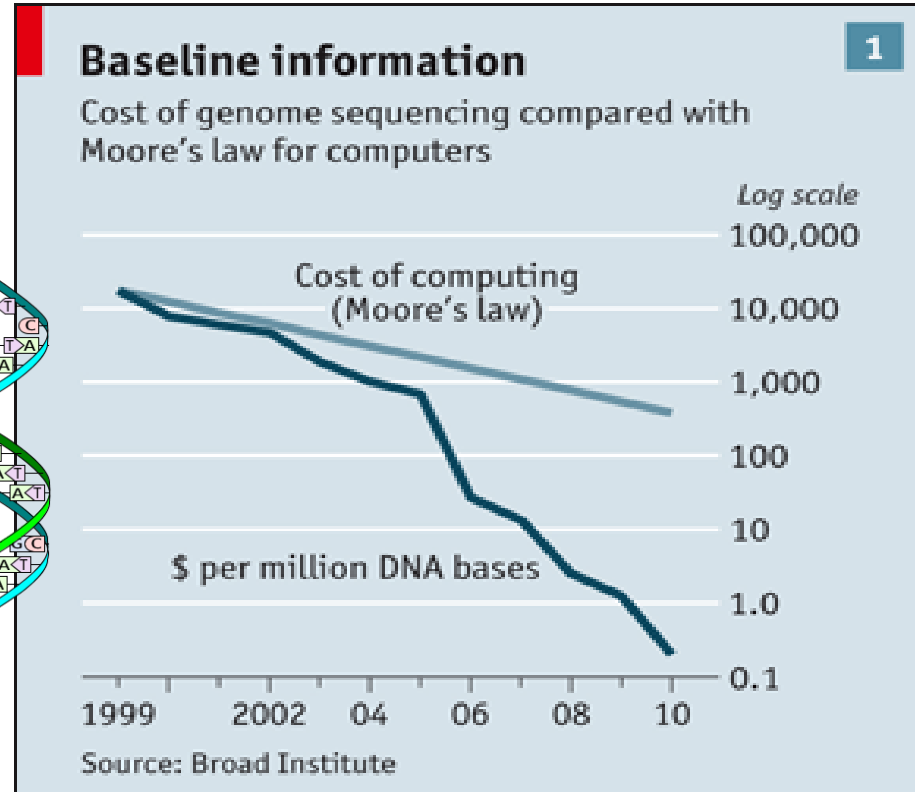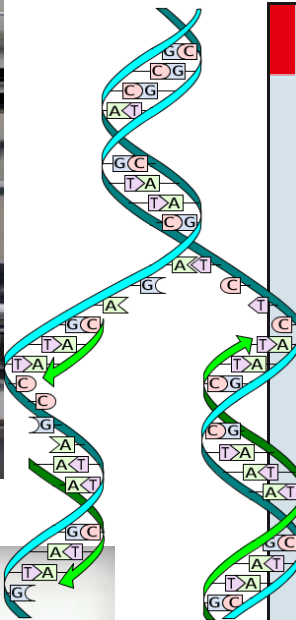data, activity

Questions

Answers

amplab

# Real Applications

- Mobile Millennium Project
  - *Alex Bayen*, Civil and Environment Engineering, UC Berkeley



- Microsimulation of urban development
  - *Paul Waddell*, College of Environment Design, UC Berkeley



- Crowd based opinion formation
  - *Ken Goldberg*, Industrial Engineering and Operations Research, UC Berkeley



- Personalized Sequencing
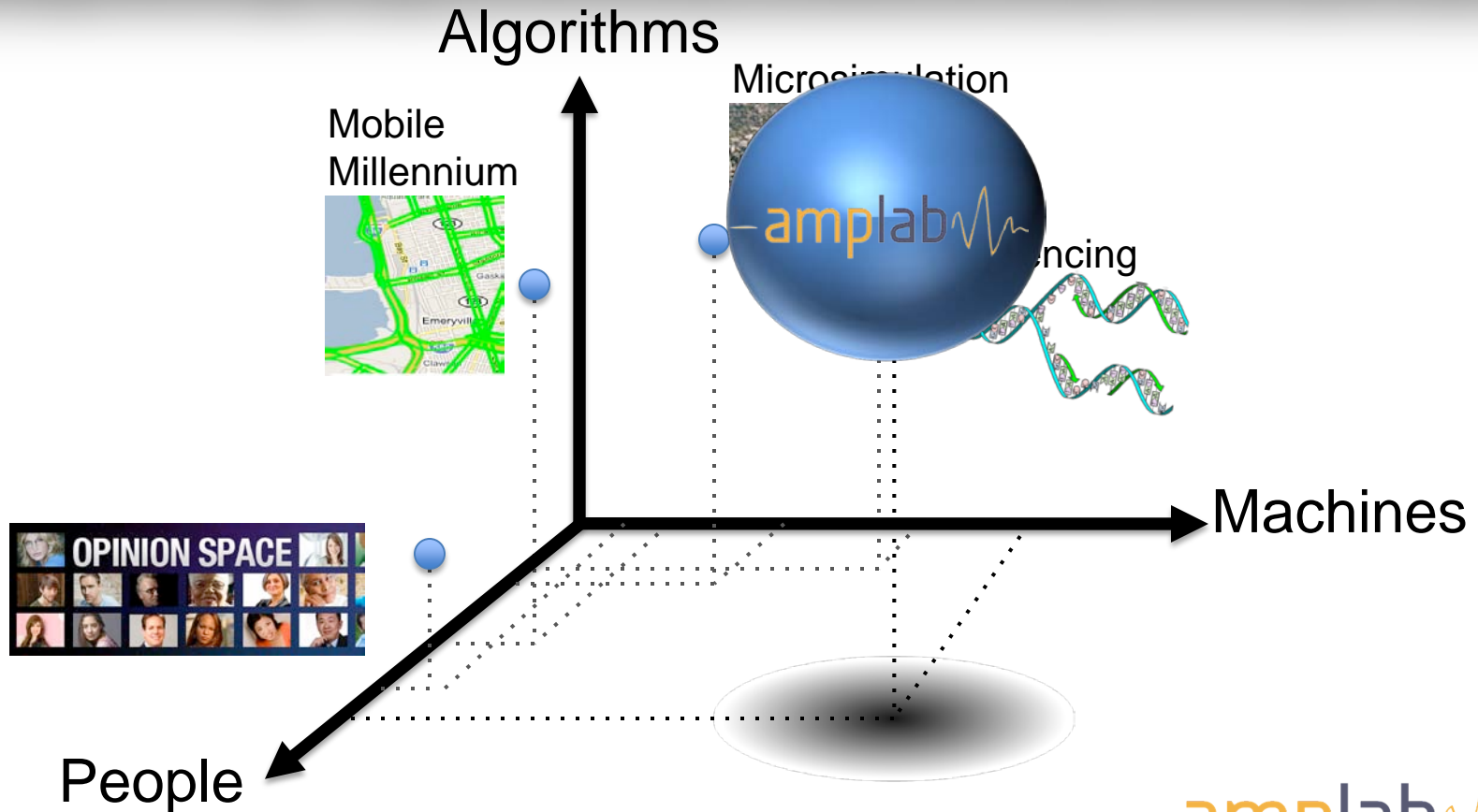  - *Taylor Sittler*, UCSF

amplab

# Personalized Sequencing



Baseline information

Cost of genome sequencing compared with Moore's law for computers

Source: Broad Institute

## Make sense of data at scale by tightly integrating algorithms, machines, and people

# Big Data in 2020

Almost Certainly:

- Create a new generation of big data scientist

- A real datacenter OS

- ML becoming an engineering discipline

- People deeply integrated in big data analysis pipeline

If We're Lucky:

- System will know what to throw away

- Generate new knowledge that an individual person cannot

_amplab_

# Summary

- Goal: Tame Big Data Problem
  - Get results with **right quality** at the **right time**
- Approach: Holistically integrate **A**lgorithms, **M**achines, and **P**eople
- Huge research issues across many domains