

# CS226

# Big-Data Management

Instructor:  
Ahmed Eldawy

# Welcome (back) to UCR! Happy New Year!



# Class information

- Classes: M/W 2:00 – 3:20 PM
- Materials Science and Engineering 103
- Instructor: Ahmed Eldawy
- TAs: Bo Huang
- Office hours: T/Th 4:00 – 4:50 PM  
on Zoom
- Email: [eldawy@ucr.edu](mailto:eldawy@ucr.edu)
  - Subject: “[CS226]”
  - All caps, with brackets and no spaces
  - Email if you would like to meet in-person

# Code of Conduct

- Instructor and TA provide:
  - Clear and detailed instruction
  - Help and support
  - Appropriate feedback
- Students should:
  - Follow the instructions responsibly
  - Ask for help when needed
  - Complete the work on their own
- Both should be:
  - Respectful, understanding, and honest
- Refer to the course syllabus for more details

# Course work



Project (50%)



Course activities (15%)



Assignments (20%)



Exams (15%)

# Project

- Groups of 4-5 students
- Milestones
  - Group Selection
  - Project proposal (6%)
  - Project proposal presentation (8%)
  - Literature survey (7%)
  - Report outline (5%)
  - Final report and deliverables (10%)
  - Final presentation (10%)
  - Presentation questions (4%)
- This year, we focus on Digital Twins
- PhD Students?

# Use of Generative AI

- Use of generative AI, ChatGPT and the likes, is allowed
- You are required to disclose which parts were generated
- ChatGPT is not always right
- You are ultimately responsible for everything you deliver
- Exams will be on the test center

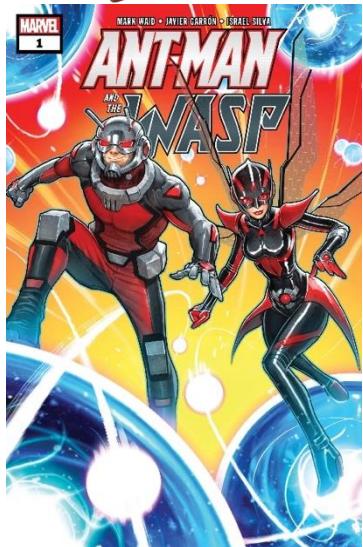
# Learning Outcomes

- Understand the characteristics and aspects of big data.
- Learn how big data is stored in the distributed file system.
- Get familiar with big-data programming models and query execution engines.
- Build real applications using big-data systems, e.g., Spark and AsterixDB.
- Know the architectural differences between big-data systems.

# Prerequisites

- Database Management Systems
- Relational data model
- Relational algebra
- SQL
- Indexing
- Query optimization
- Data normalization
- Transactions
- (Required task) Take the prerequisite test on Canvas

# Super Hero



# Ant-Man/Wasp



Get smaller to understand how ants work and what they are capable of.

Use this knowledge to control thousands of ants and do amazing things!

Optional task: Watch Ant-Man and the Wasp this weekend 😊



# Big-data Expert

- Understand how the big-data platforms really work
- Control those thousands of processors efficiently to carry out your task

# Syllabus

- Overview of big data
- Big-data storage
- Unstructured big-data processing
- (Semi-)Structured big-data processing
- Big-data layout and formats
- High-level applications on big-data



# Introduction

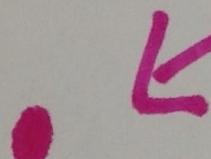
**Big Data**

Straight Ahead



All of the  
information

Information  
you  
need!



# The Market of Big Data

≡    Forbes / Tech / #BigData

Your roadmap for distributed care is here. [Read the eBook.](#) 

JAN 20, 2017 @ 09:27 AM 72,834  The Little Black Book of Billionaire Secrets

## 6 Predictions For The \$203 Billion Big Data Analytics Market

 **Gil Press, CONTRIBUTOR**  
I write about technology, entrepreneurs and innovation [FULL BIO](#) ▾  
Opinions expressed by Forbes Contributors are their own.  
  
Shutterstock

The creation and consumption of data continues to grow by leaps and bounds and with it the investment in big data

**RELATED KEYWORDS**

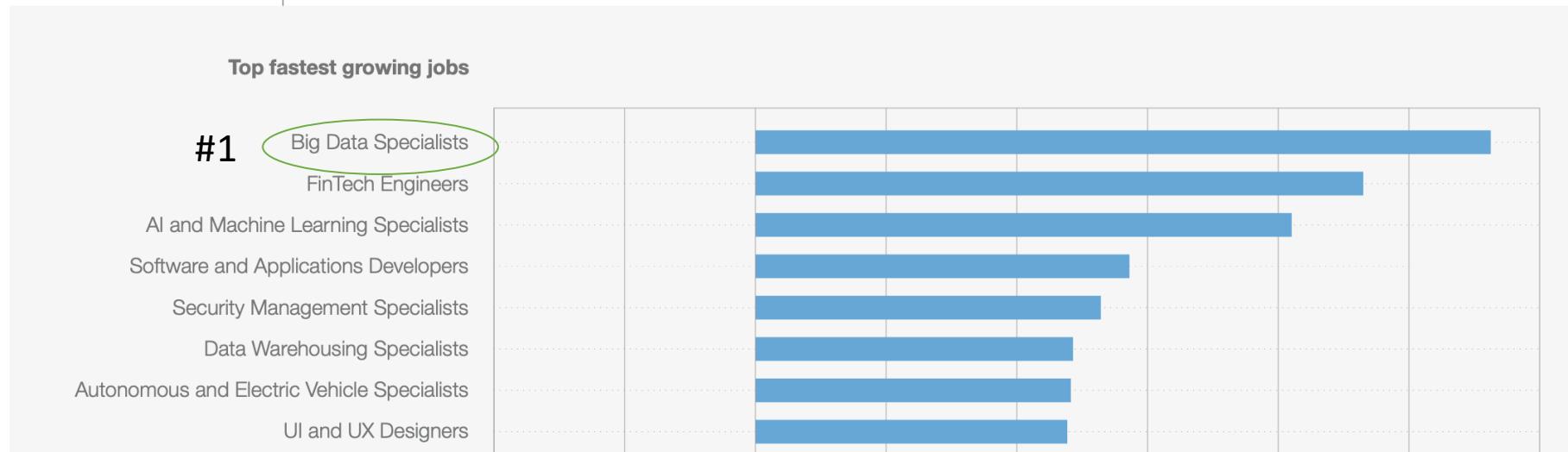
- 1. [BIG DATA ANALYTIC TOOLS](#) >
- 2. [BIG DATA ANALYTIC TRENDS](#) >
- 3. [BIG DATA TRENDS FOR 2018](#) >
- 4. [BIG DATA FOR BUSINESS](#) >
- 5. [NEW BIG DATA SOLUTIONS](#) >
- 6. [DATA MANAGEMENT PLATFORM](#) >
- 7. [DATA ENTRY SERVICES](#) >
- 8. [DATA ANALYTICS TRAINING](#) >
- 9. [BIG DATA COURSES](#) >
- 10. [GEOSPATIAL DATA MANAGEMENT](#) >

# Fastest Growing Jobs

FIGURE 2.2

## Fastest-growing and fastest-declining jobs, 2025-2030

Top jobs by fastest net growth and net decline, projected by surveyed employers



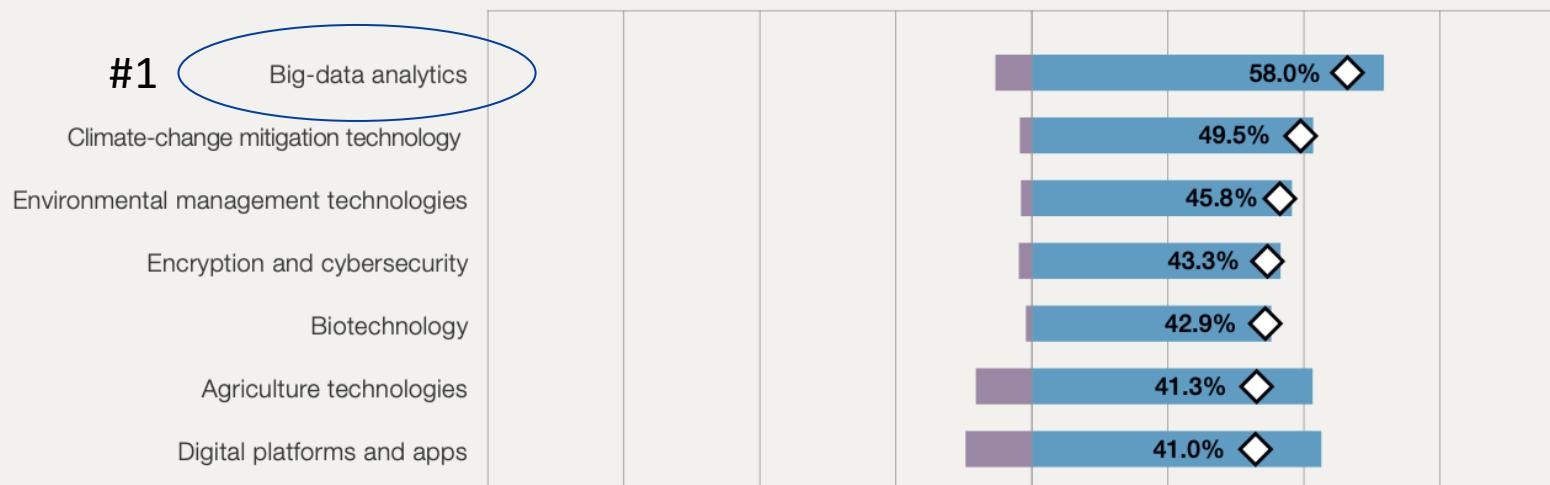
# Technology Impact on Jobs

FIGURE 2.5

## Expected impact of technology adoption on jobs, 2023–2027

Share of organizations surveyed that expect each technology to create or displace jobs, ordered by the job creation net effect.

The shares of organizations which expect the impact of adopting these technologies to be neutral are not plotted.



Source

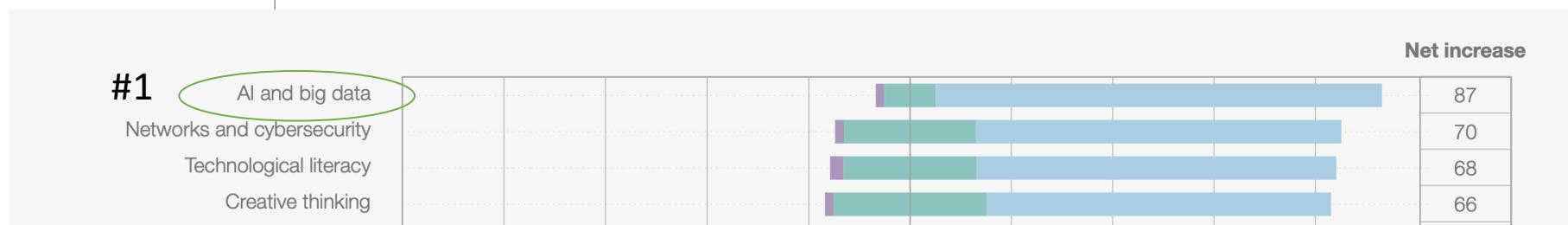
World Economic Forum, Future of Jobs Survey 2023.

# Skills on the Rise

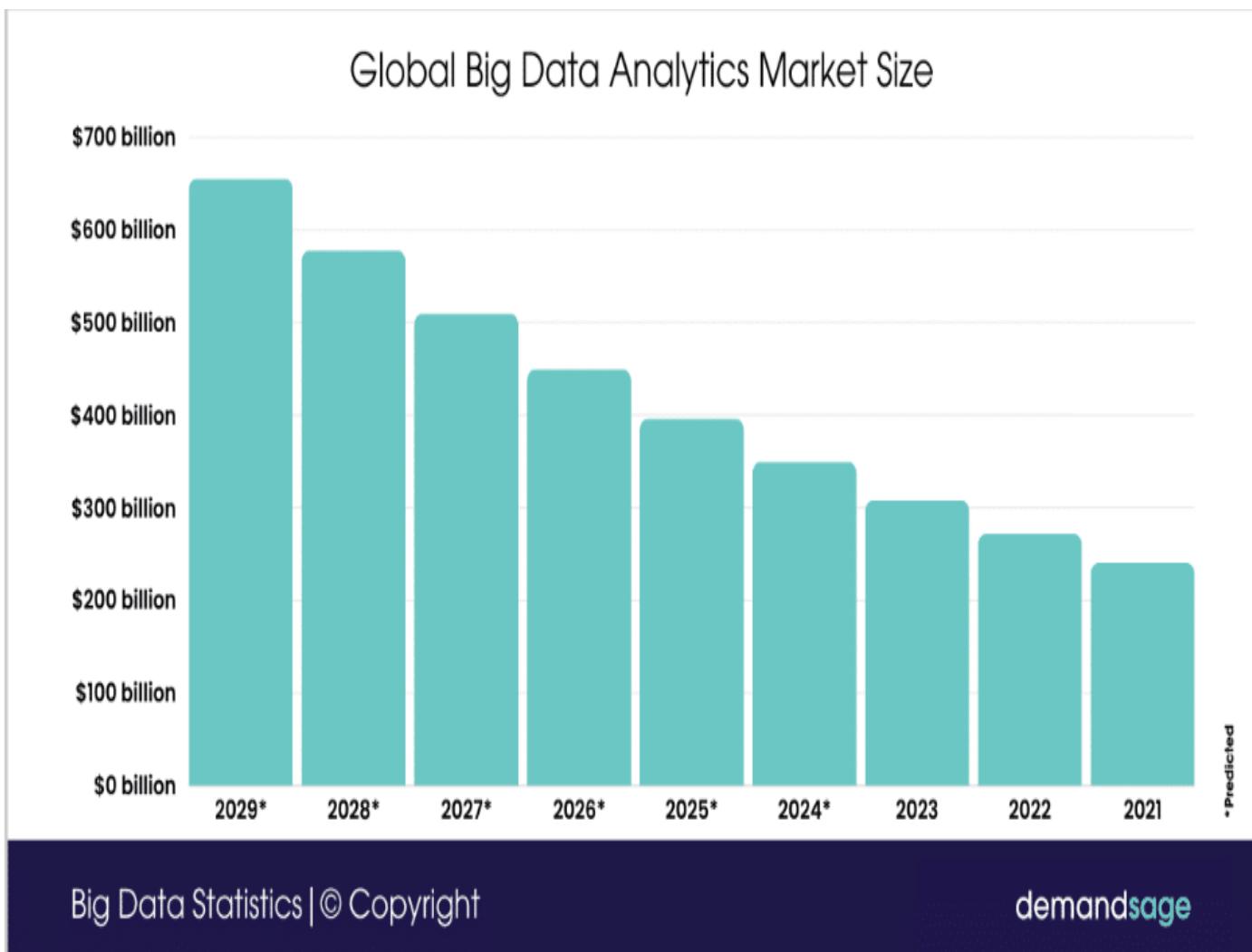
FIGURE 3.4

## Skills on the rise, 2025-2030

Share of employers that consider skills to be increasing, decreasing, or remaining stable in importance. Skills are ranked based on net increase, which is the difference between the share of employers that consider a skill category to be increasing in use and those that consider it to be decreasing in use.



# Big-data Market Size



# OPEN Government Data Act

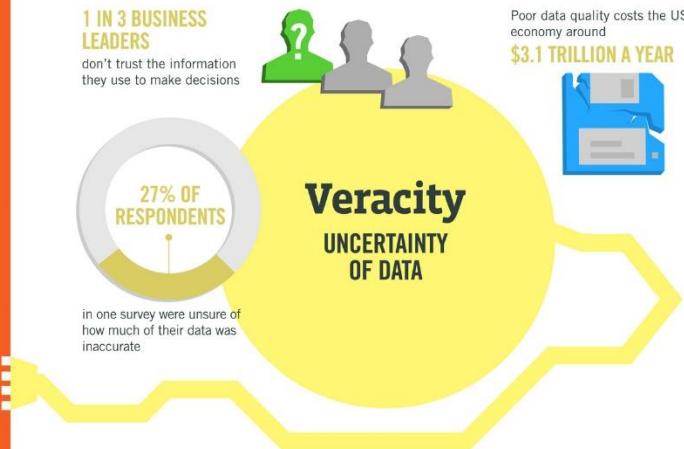
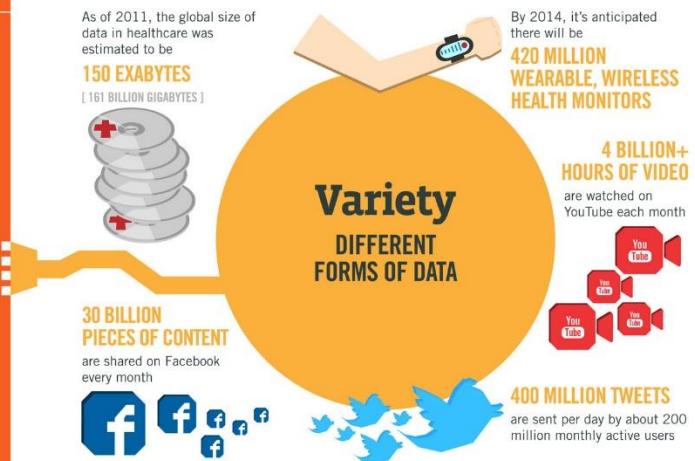
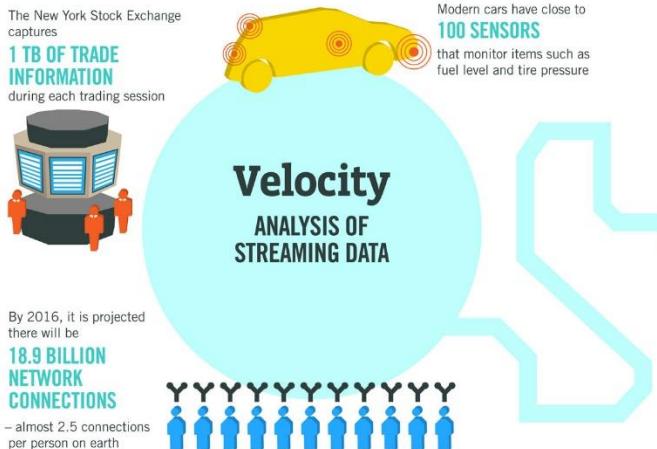
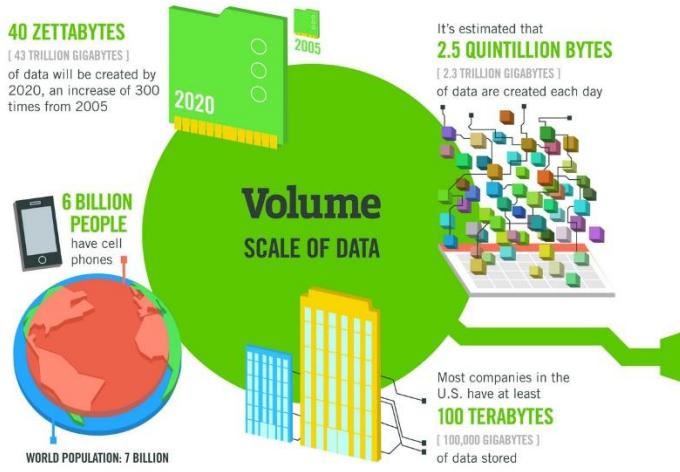
## White House finalizes OPEN Government Data Act guidance, restarts CDO Council

The long-awaited guidance for the data and evidence statute comes six years after it became law.

BY MADISON ALDER • JANUARY 15, 2025



# V's of Big Data



Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS

IBM

# Big Data Vs Big Computation

- Full scans (e.g., log processing)
- Range scans
- Point lookups
- Iterations
- Joins (self, binary, or multiway)
- Proximity queries
- Closures and graph traversals

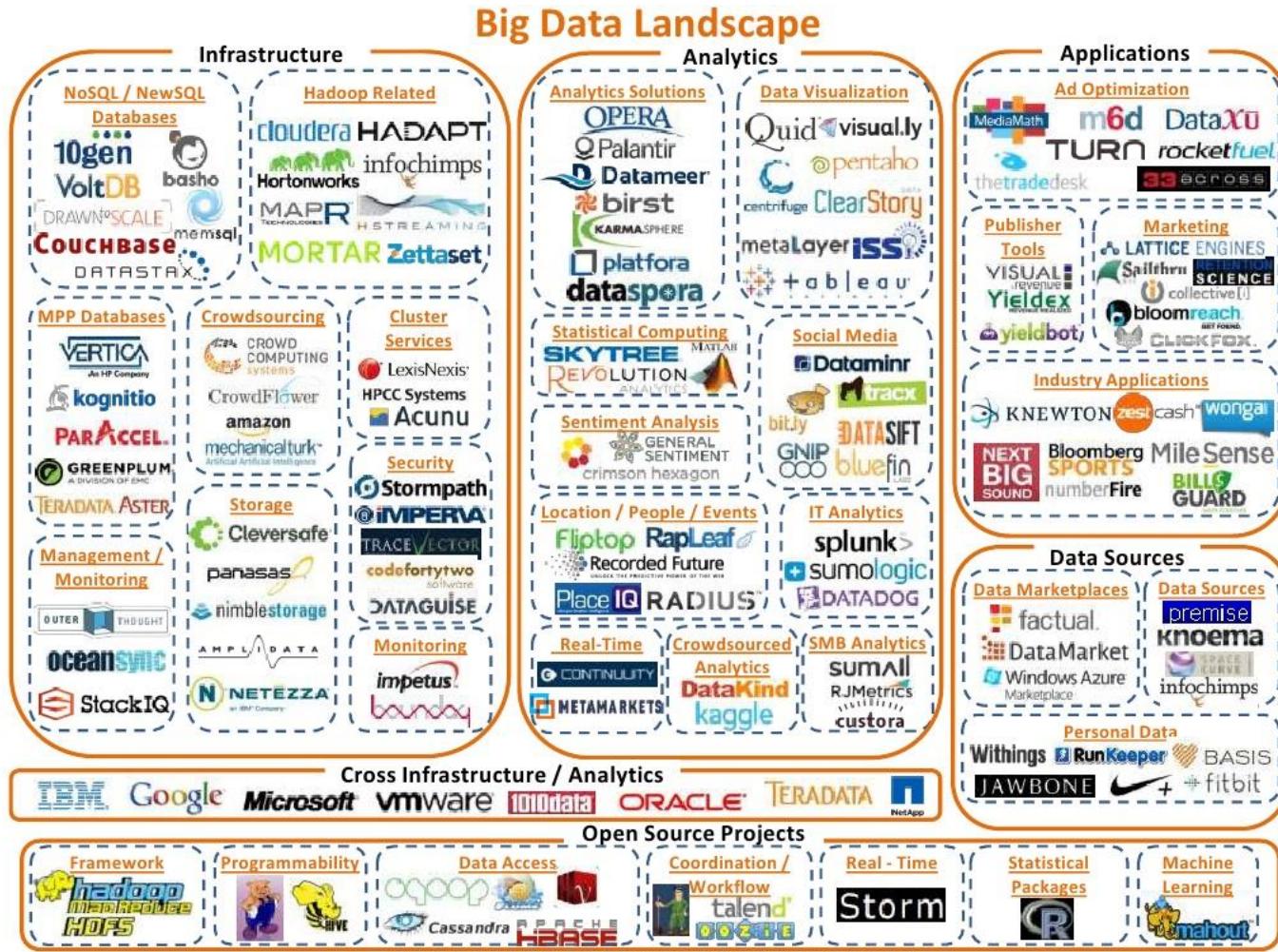
# Big Data Applications

- Web search
- Marketing and advertising
- Data cleaning
- Knowledge base
- Information retrieval
- Internet of Things (IoT)
- Visualization
- Behavioral studies

# Publicly Available Datasets

- Data.gov
- UCR Star [<https://star.cs.ucr.edu>]
- Flickr API
- Kaggle.com
- Industry: Yelp, Yahoo! Webscope, Zillow, ...

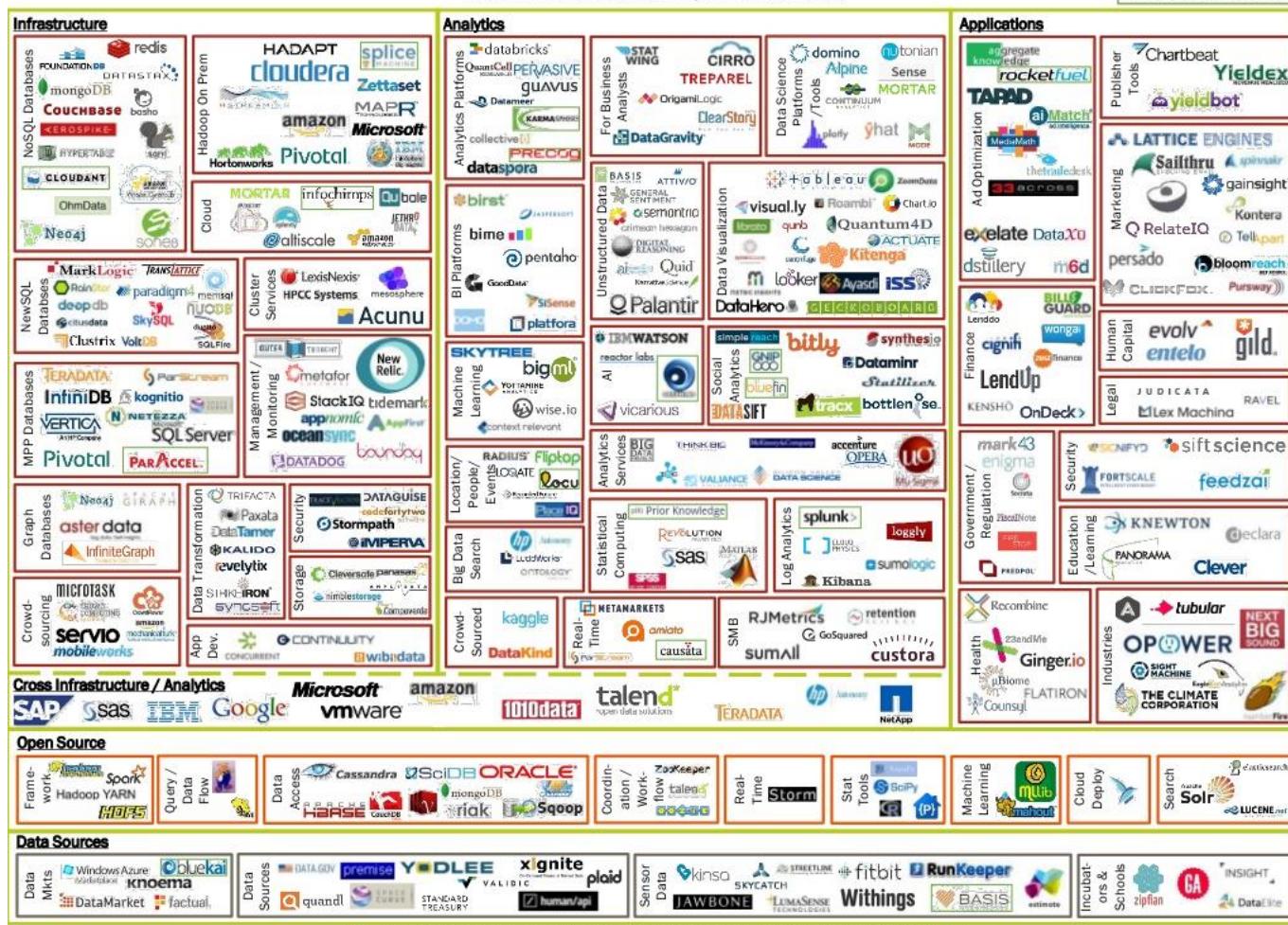
# Big Data Landscape 2012



# Big Data Landscape 2014

## BIG DATA LANDSCAPE, VERSION 3.0

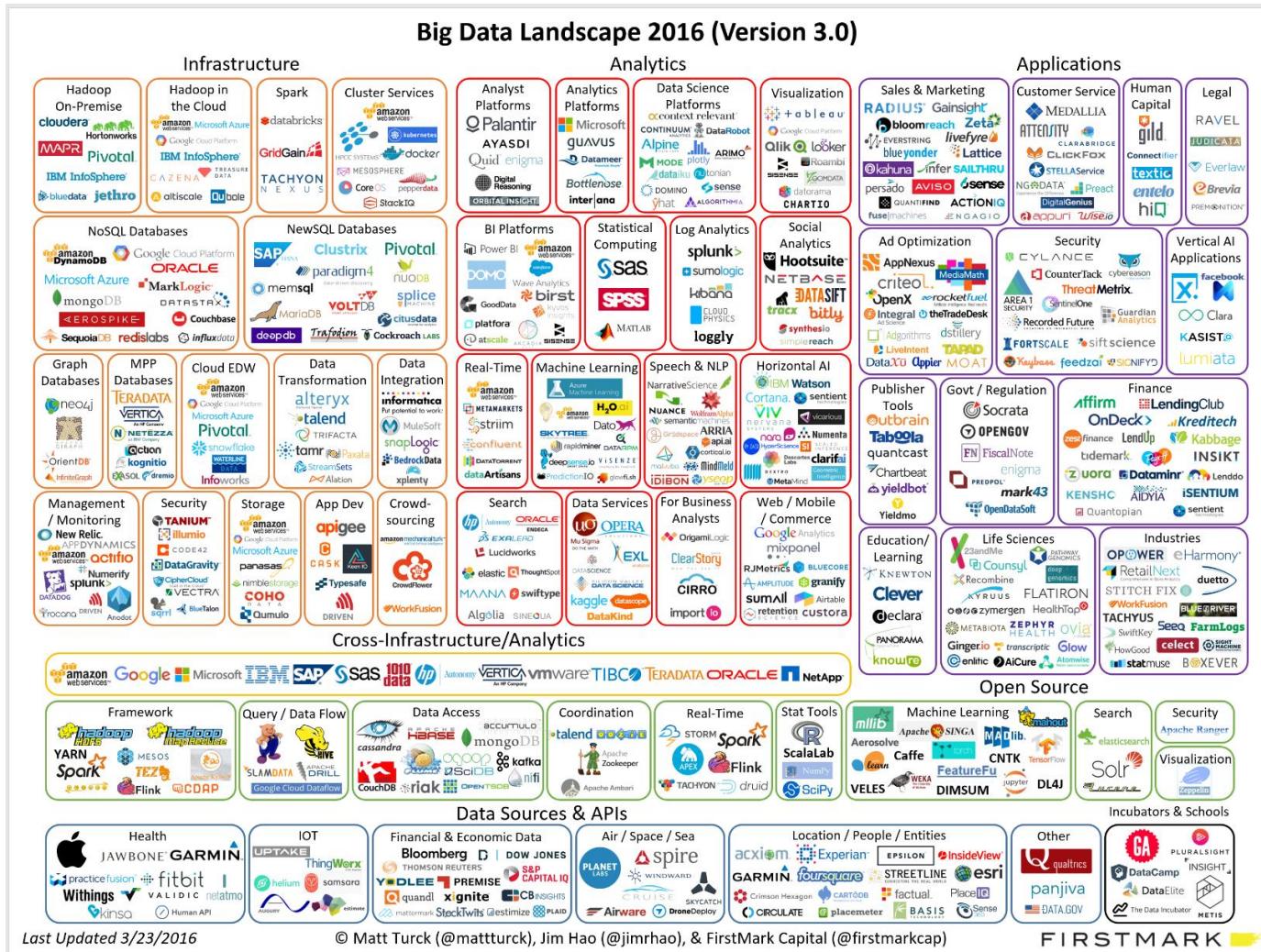
Exited: Acquisition or IPO



© Matt Turck (@mattturck), Sutian Dong (@sutiandong) & FirstMark Capital (@firstmarkcap)

<http://mattturck.com/2014/05/11/the-state-of-big-data-in-2014-a-chart/>

# Big Data Landscape 2016



<http://mattturck.com/2016/02/01/big-data-landscape/>

# Big Data Landscape 2018



V1 – Last updated 6/19/2018

© Matt Turck (@mattturck), Demilade Obavomi (@demi\_obavomi), & FirstMark (@firstmarkcap) mattturck.com/bigdata2018

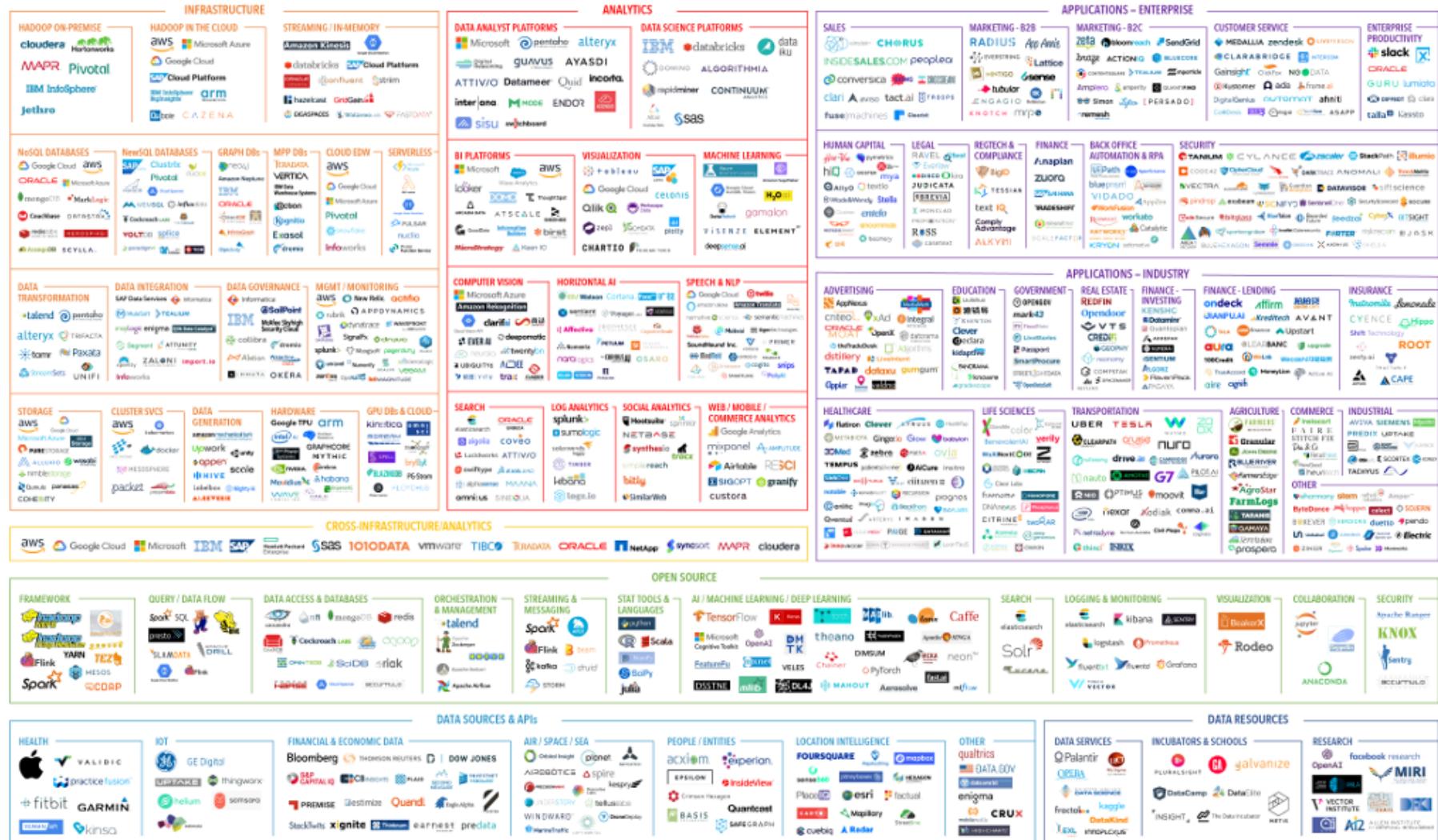
FIRSTMARK  
EARLY STAGE VENTURE CAPITAL



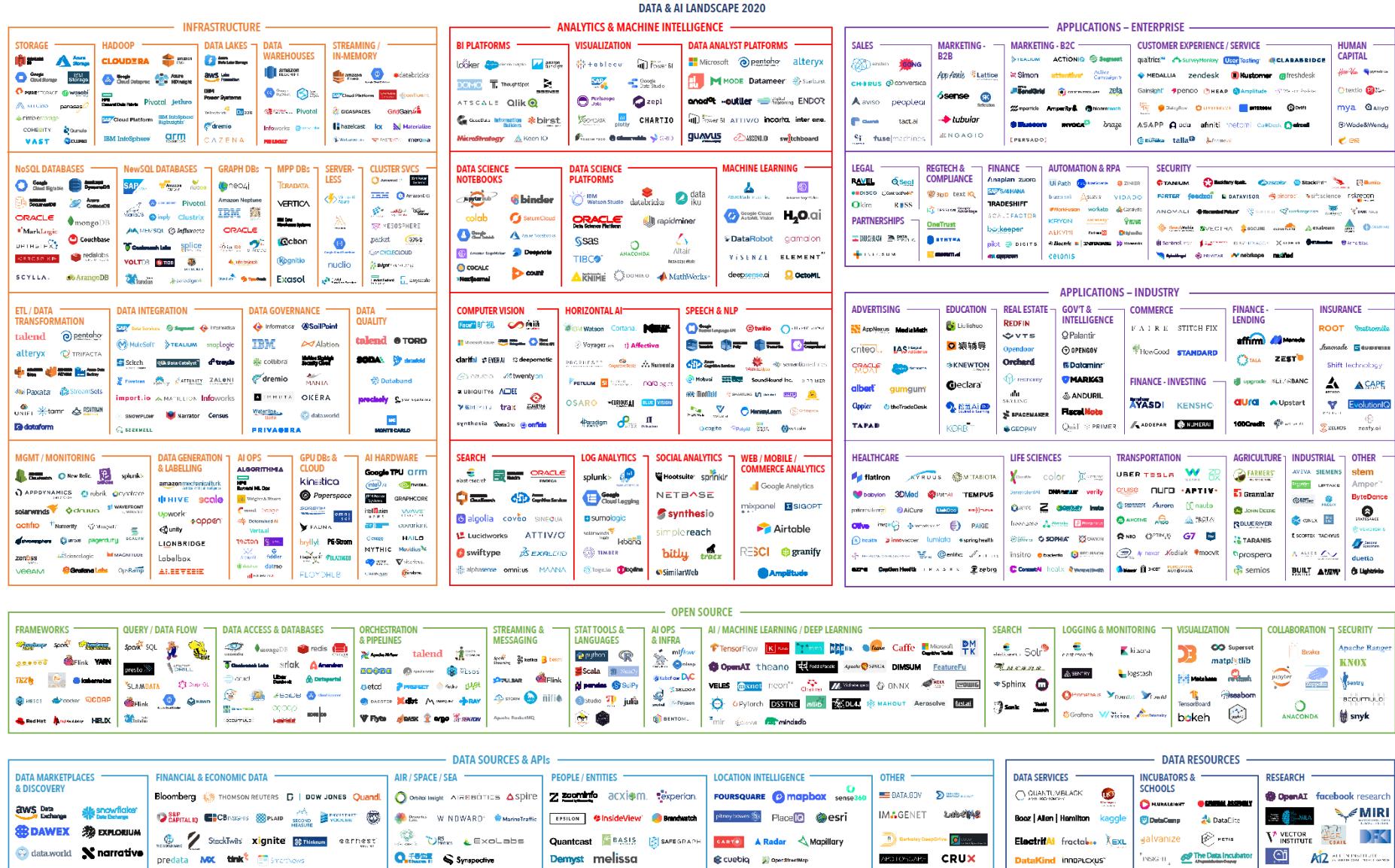
**CELEBRATING 30 YEARS**  
Marlan and Rosemary Bourns  
College of Engineering

# Data & AI Landscape 2019

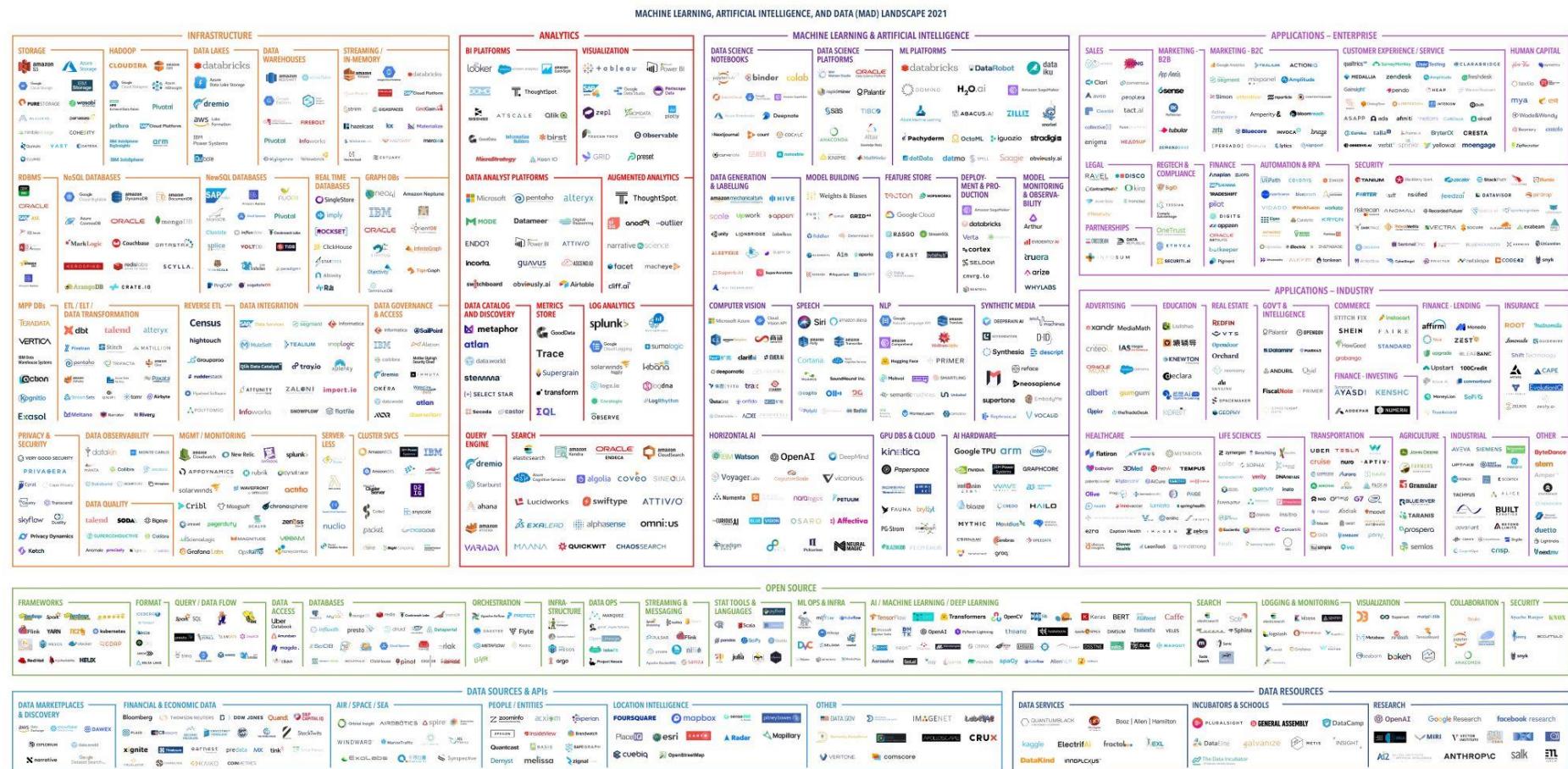
DATA & AI LANDSCAPE 2019



# Data & AI Landscape 2020



# Machine Learning, AI, and Data (MAD) Landscape 2021

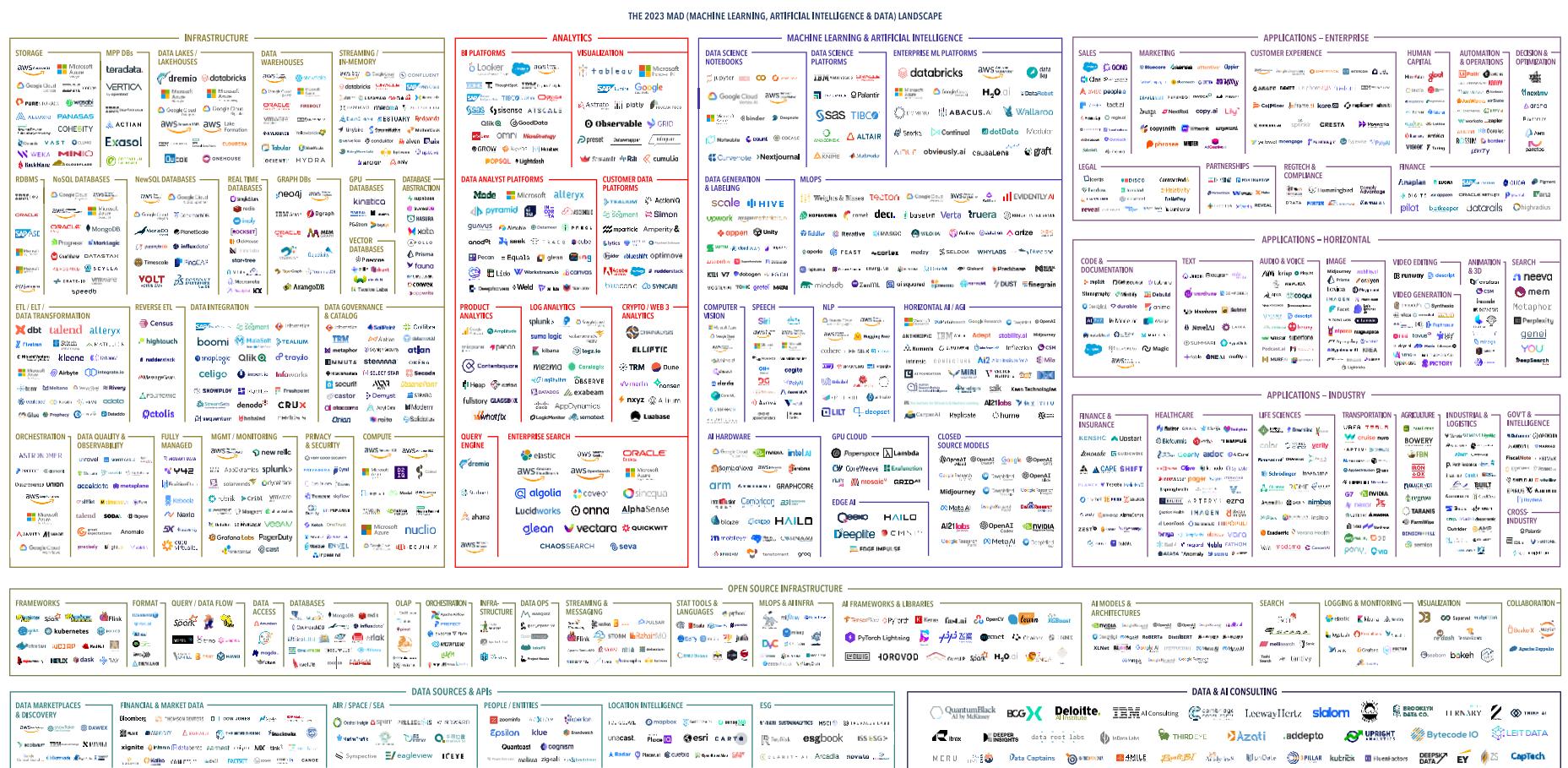


Mountain 3.0 November 2001

© Matt Turck (@mattturck), John Wu (@john\_d\_wu) & FirstMark (@firstmarkcap)

[www.ncbi.nlm.nih.gov/pmc/articles/323031/](https://www.ncbi.nlm.nih.gov/pmc/articles/323031/)

# Machine Learning, AI, and Data (MAD) Landscape 2023



Version 1.0 - Feb 2023

© Matt Turck (@mattturck), Kevin Zhang (@ykevinzhang) & FirstMark (@firstmarkcap)

Blog post: mattturck.com/MAD2023

Interactive version: [MAD.firstmarkcap.com](http://MAD.firstmarkcap.com)

Comments? Email MAD2023@firstmarkcap.com



# Components of Big Data

# Components of Big Data

## Big-data Libraries

MLlib (Machine Learning), GraphX, Visualization

## Structured/ Semi-structured Data Processing

SparkSQL, Pig, SQL++, HiveQL

## Distributed Computing

MapReduce (Hadoop and Google), Resilient Distributed Dataset (Spark), Hyracks (AsterixDB)

## Big Data Distributed Storage

Hadoop Distributed File System, Cloud storage systems (Amazon S3 and Google File System), Key-value stores

## Cloud Services

Amazon Web Services, Microsoft Azure, and Google Cloud Platform

## Coordination/Cluster Management

Oozie, Yarn, Kubernetes

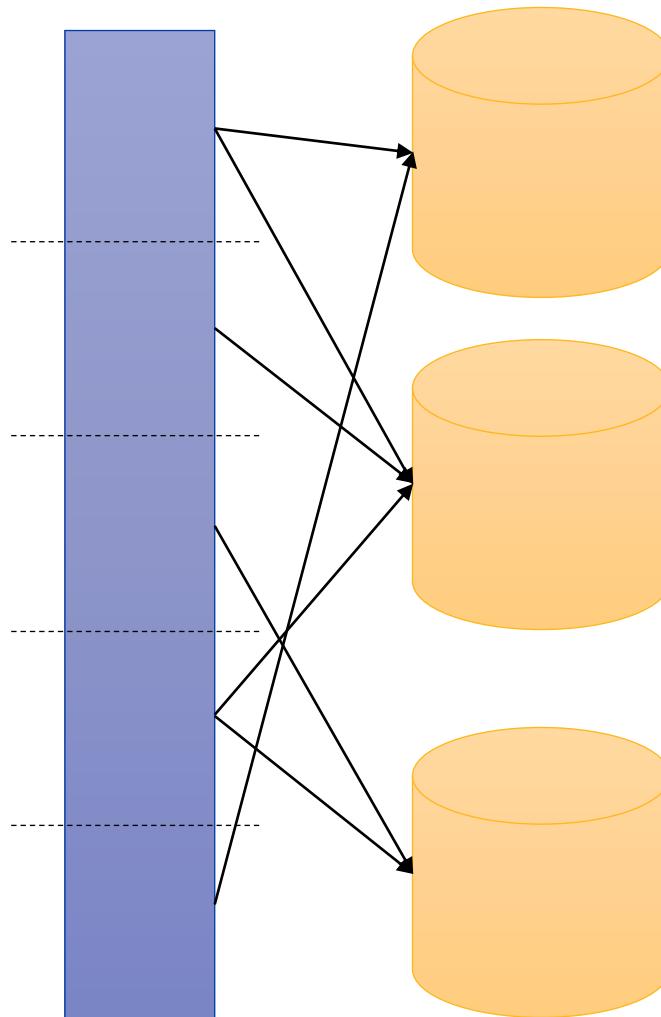


Marlan and Rosemary Bourns  
College of Engineering

# Big Data Storage

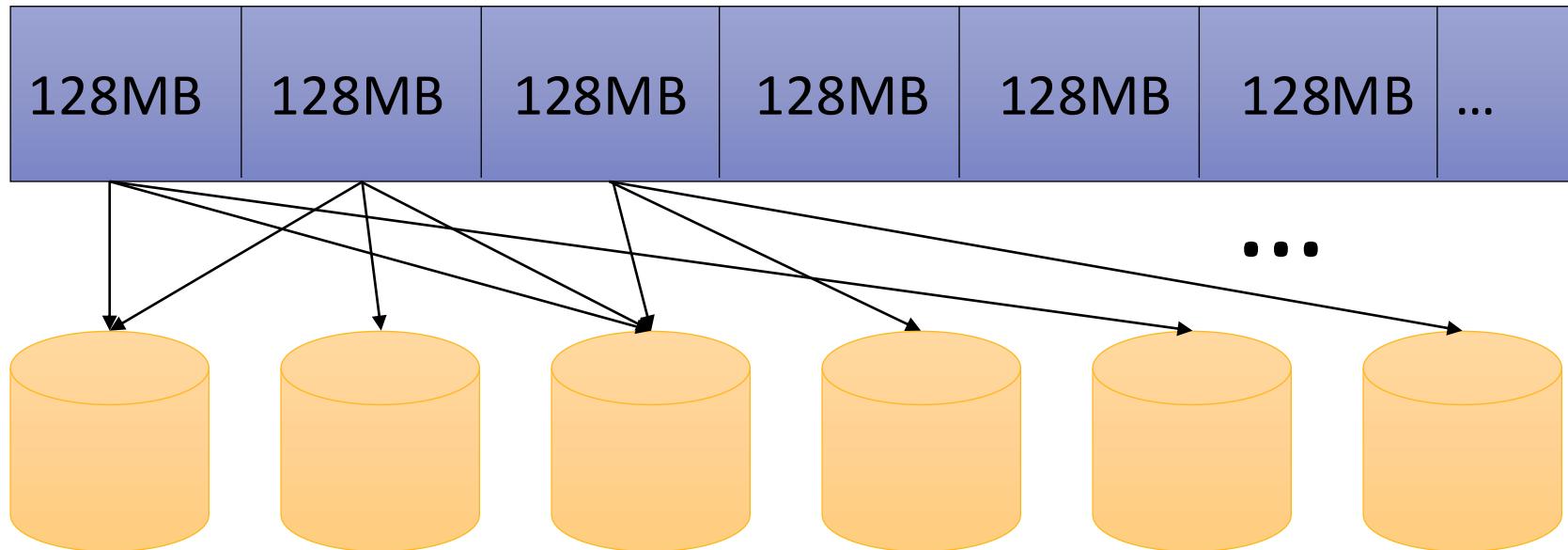
# Storage of Big Data

- Data is growing faster than Moore's Law
- Too much data to fit on a single machine
- Partitioning
- Replication
- Fault-tolerance



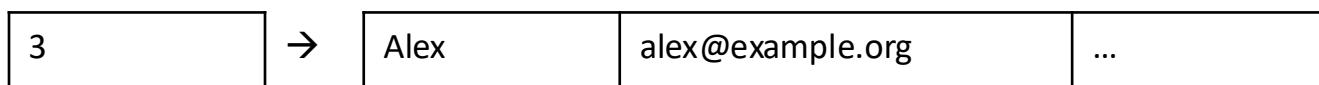
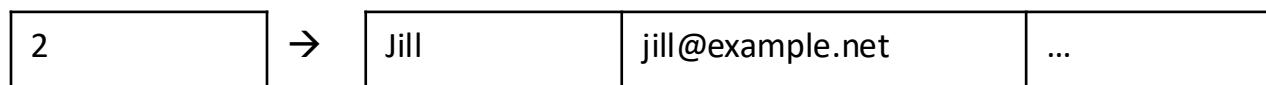
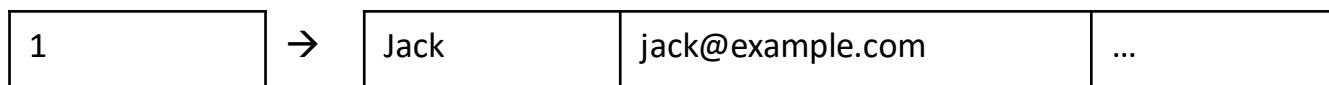
# Hadoop Distributed File System (HDFS)

- The most widely used distributed file system
- Fixed-sized partitioning
- 3-way replication
- Write-once read-many
- See also: GFS, Amazon S3, Azure Blob Store

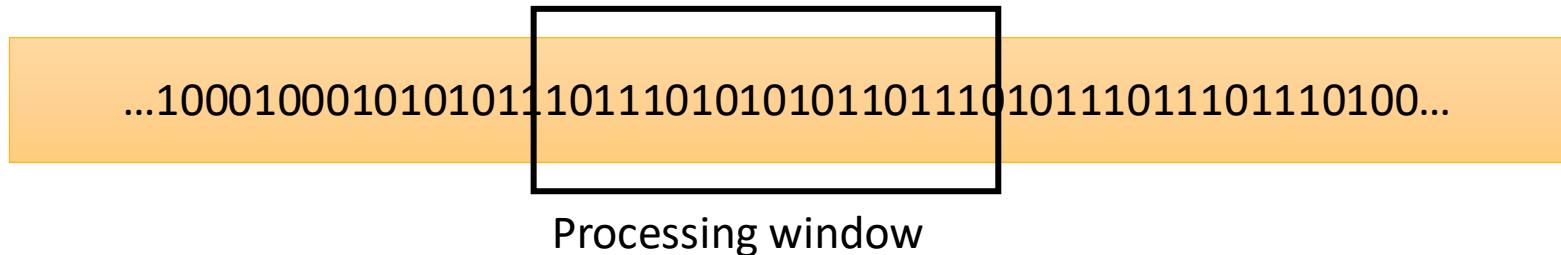


# Key-value Stores

| ID | Name | Email            | ... |
|----|------|------------------|-----|
| 1  | Jack | jack@example.com |     |
| 2  | Jill | jill@example.net |     |
| 3  | Alex | alex@example.org |     |



# Streaming



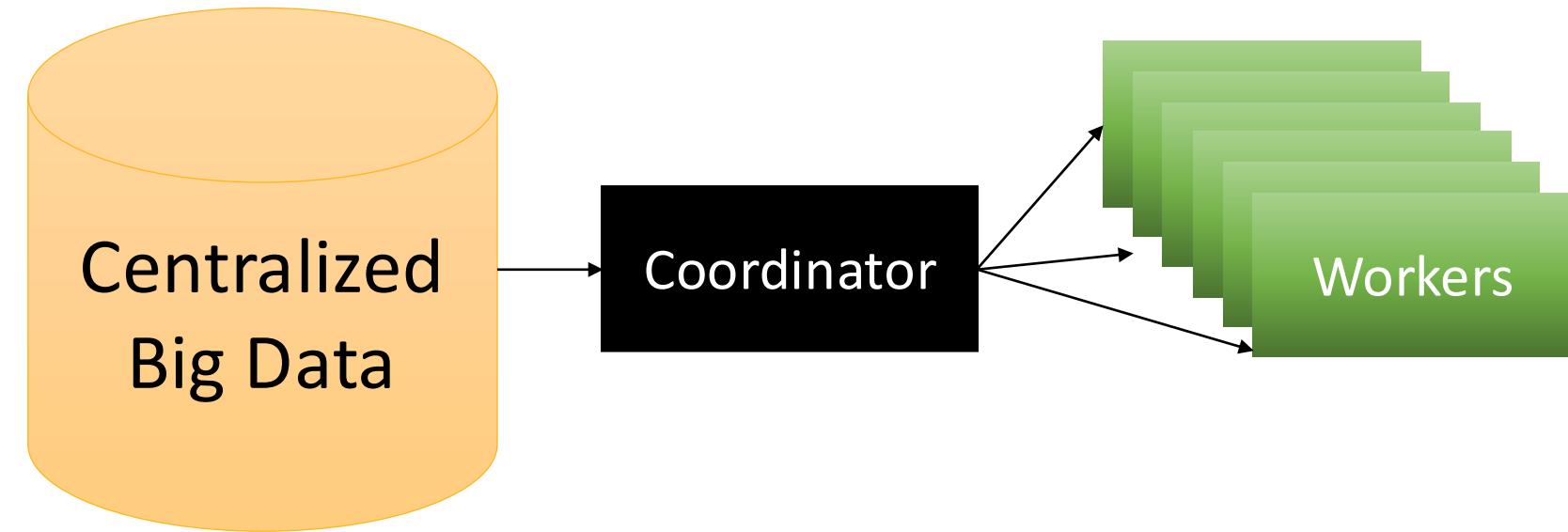
- Sub-second latency for queries
- One scan over the data
- (Partial) preprocessing
- Continuous queries
- Eviction strategies
- In-memory indexes



Marlan and Rosemary Bourns  
College of Engineering

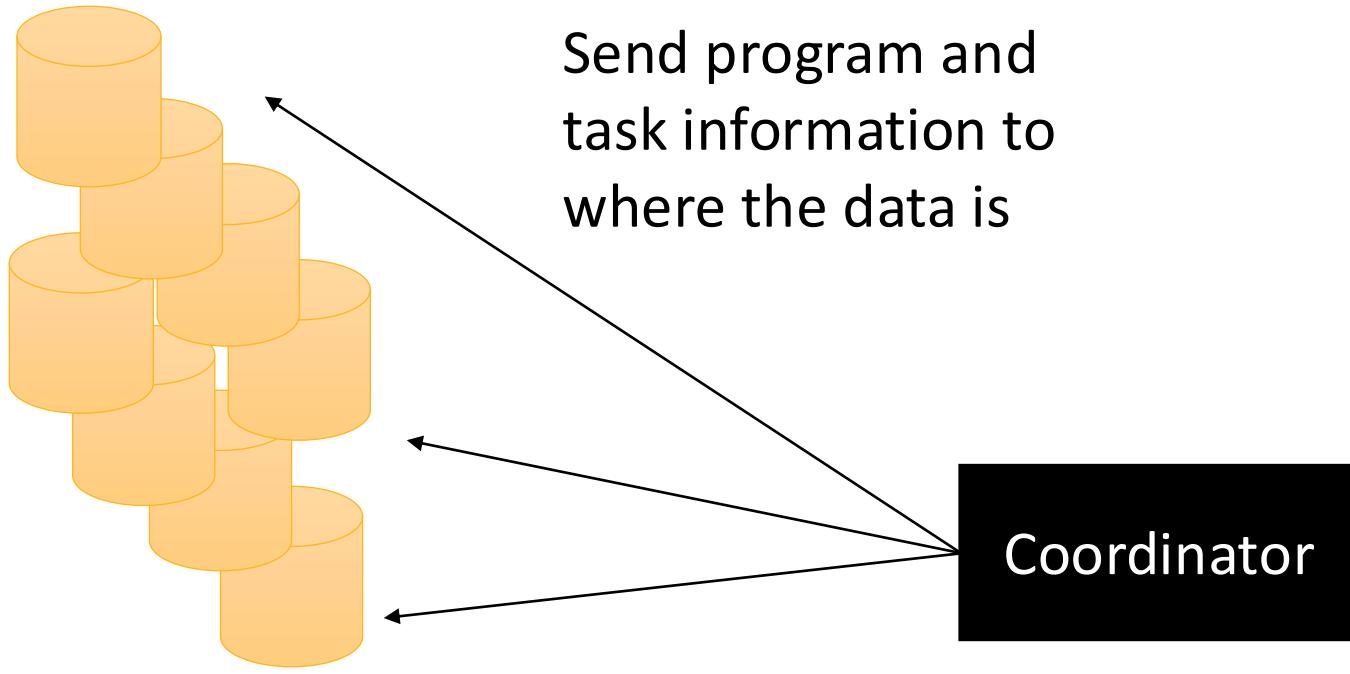
# Big Data Computation Models

# Traditional Distributed Computing



Ship data to computation paradigm  
e.g., High performance computing (HPC)

# Big-data Computing



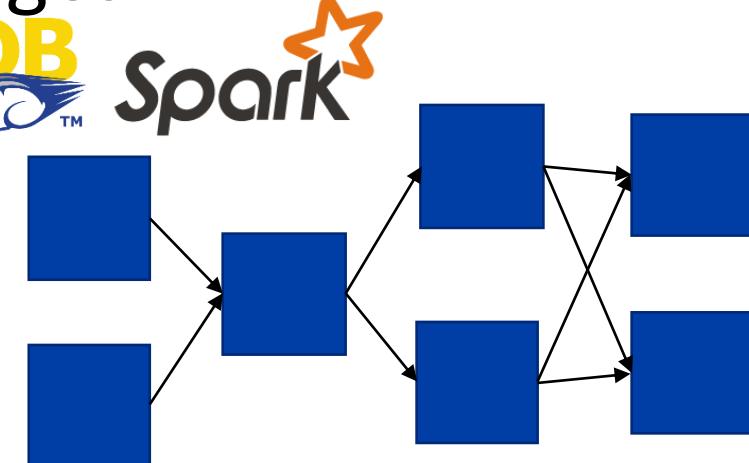
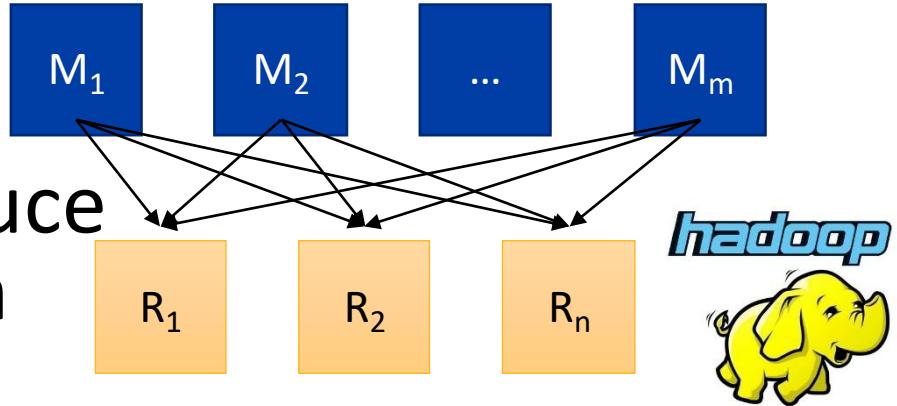
Storage/Compute  
Nodes

Coordinator

Ship compute to data paradigm

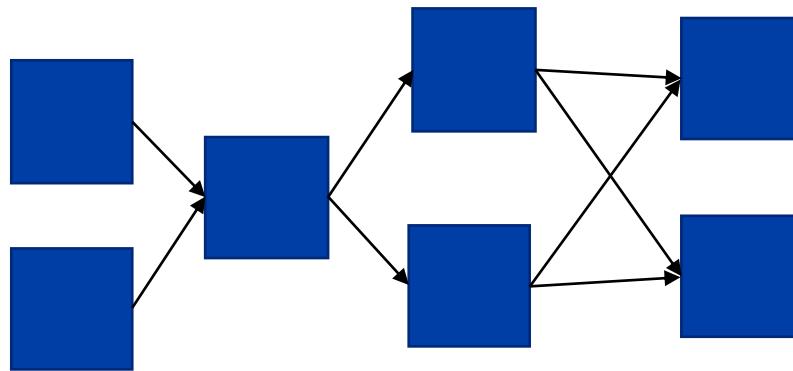
# Task Execution

- MapReduce
  - Map-Shuffle- Reduce
  - Resiliency through materialization
- Resilient Distributed Datasets (RDD)
  - Directed-Acyclic-Graph (DAG)
  - In-memory processing
  - Resiliency through lineages
- Hyracks
- Stragglers
- Load balance



# Provenance

- Debugging in distributed systems is painful



- We need to keep track of transformations on each record



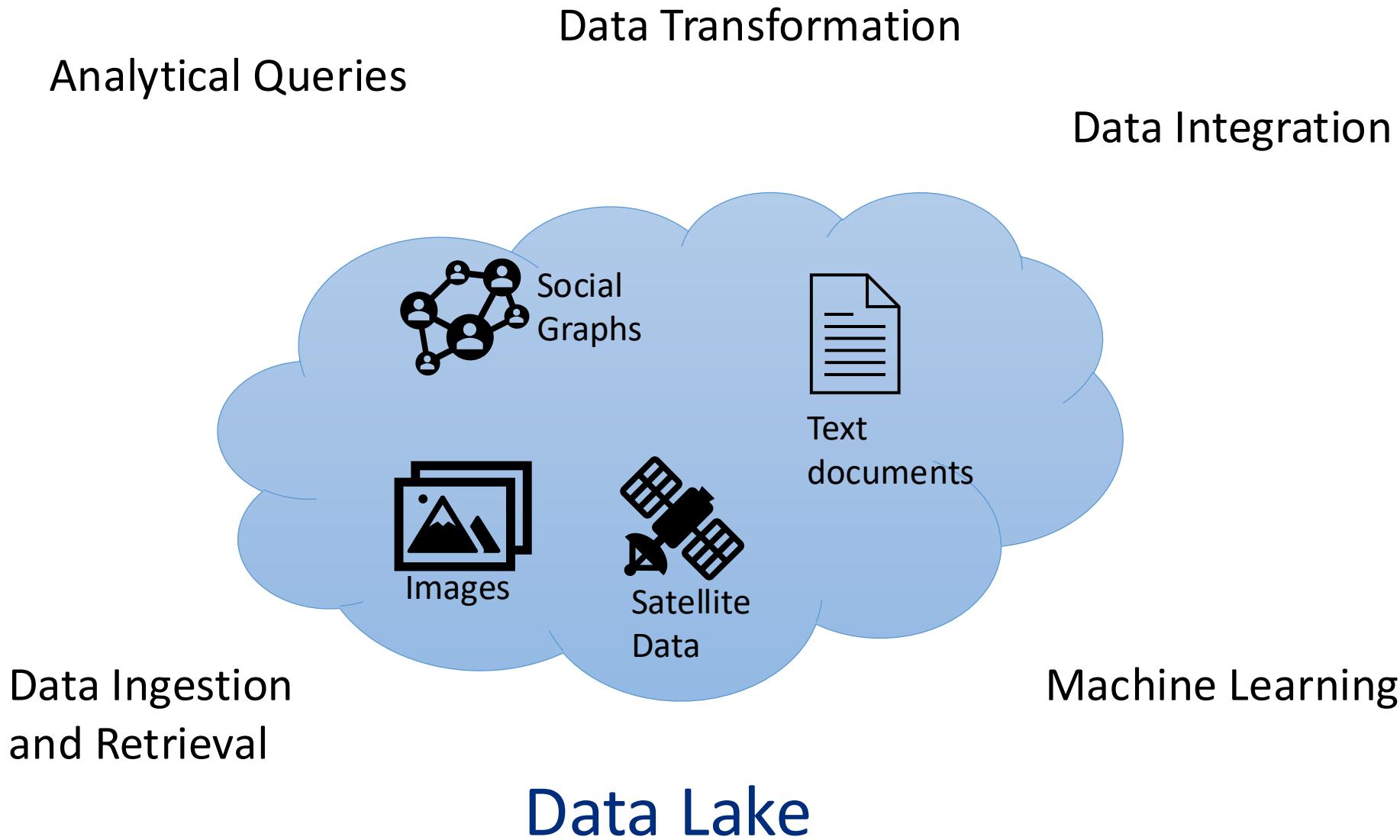
Marlan and Rosemary Bourns  
College of Engineering

# (Semi-)Structured Big-data Processing

# Structured Data Processing

- There is a need for processing structured and semi-structured data
- The relational model and SQL are still popular
- Let the big-data system know about the structure of the data and processing
- Allow the system to optimize query processing
- Examples: Algebricks, SparkSQL, and Pig

# Data Lakes



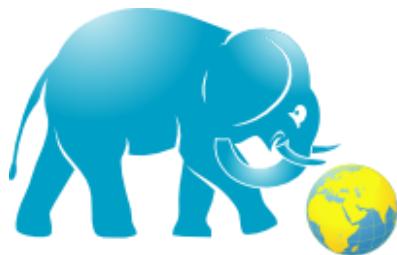


Marlan and Rosemary Bourns  
College of Engineering

# Big-data Applications

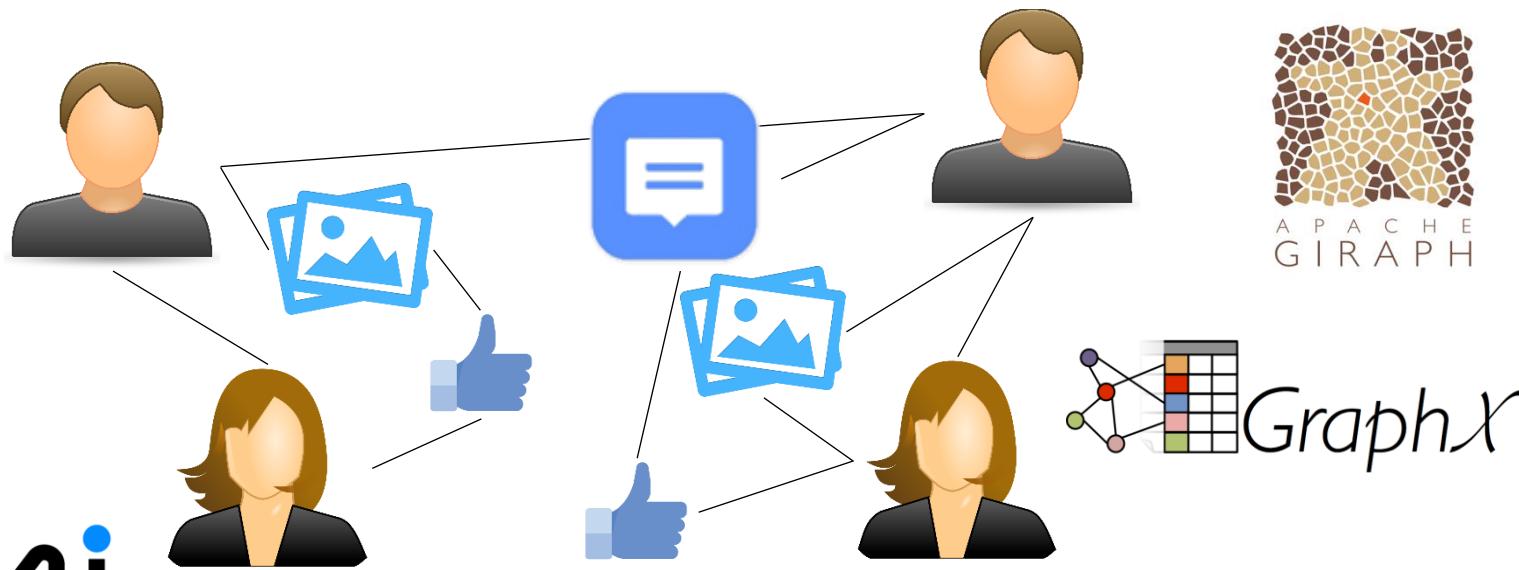
# Big Spatial Data

- GPS-enabled devices
- Satellite images
- Autonomous vehicles
- Scientific applications



# Big Graphs

- Motivated by social networks
- Billions of nodes and trillions of edges
- Tens of thousands of insertions per second
- Complex queries with graph traversals



# Machine Learning

- The rise of machine learning and data-driven models
- Modern generative AI models rely on big-data for training
- Data preparation and model training are costly for big data
- Use big-data processing to scale up the process



# Streaming

- Internet-of-Things (IoT)
- GPS-enabled devices
- Surge of social media

 **Spark**  
*Streaming*

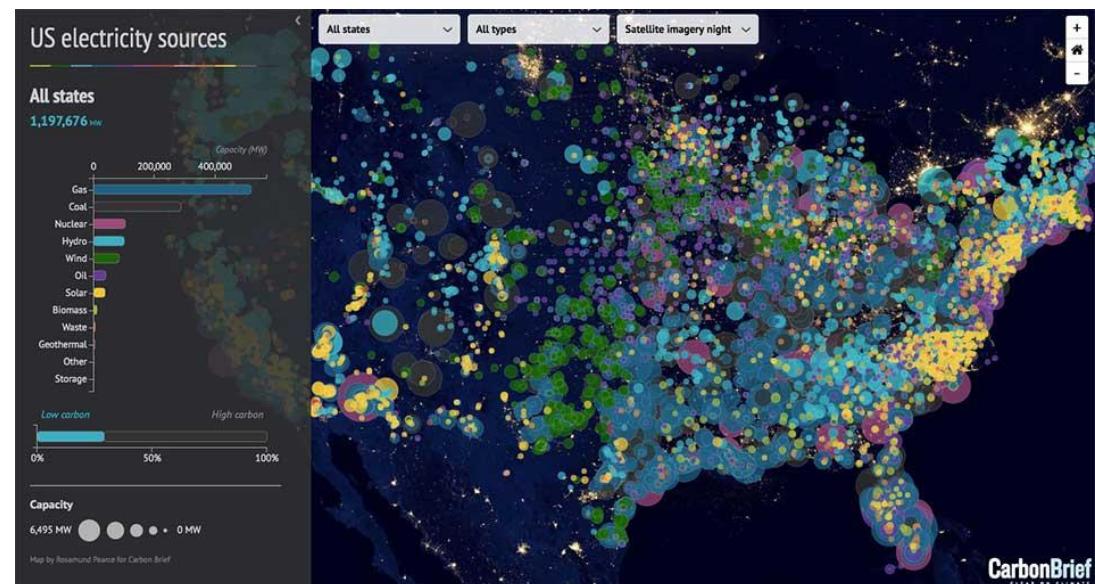
 APACHE **kafka**®

 **Flink**

 APACHE **STORM™**

# Big-data Visualization

- Turn big-data into insights
- Effective way to communicate with non-technical users



# Summary

- Data is getting big, so as computation
- Characteristics of big-data
- Big-data is a new technology for easy handling of large-scale data
- Big-data is a collection of systems, each focusing on a specific aspect of big-data, e.g., storage, processing, analysis, and visualization

# Next Steps

- Think about project ideas and try to find interesting or relevant datasets
- (Required) Use the available discussion space on Canvas to initiate and discuss project ideas
  - This is your opportunity to find teammates
- (Required) Check your knowledge about the prerequisites using the short quiz on Canvas
- (Required) Check the syllabus and finish the course overview assignment