



DS221

# Data Structures, Algorithms & Data Science Platforms

Instructor: Chirag Jain  
(slides from Prof. Simmhan)



# Big Data Concepts

What, Where, Why?

# What is Big Data?



# The term *is* fuzzy ... Handle with care!



Wordle of "Thought Leaders'" definition of Big Data, © Jennifer Dutcher, 2014  
<https://datascience.berkeley.edu/what-is-big-data/>

# Data Generation View

“

*“Big data refers to the approach to data of ‘collect now, sort out later’...The low cost of storage and better methods of analysis mean that you generally don’t need to have a specific purpose for the data in mind before you collect it.”*

***Rohan Deuskar, CEO and Co-Founder, Stylitics***



Wordle of “Thought Leaders” definition of Big Data, © Jennifer Dutcher, 2014  
<https://datascience.berkeley.edu/what-is-big-data/>

# Data Systems View

“



*“Big data is when your business wants to use data to solve a problem, answer a question, produce a product, etc., but the standard, simple methods break down on the size of the data set, causing time, effort, creativity, and money to be spent crafting a solution to the problem that leverages the data without simply sampling or tossing out records.”*

*John Foreman, Chief Data Scientist, MailChimp*

# Data Analysis View



“

*“While the use of the term is quite nebulous ...  
I've understood “big data” to be about  
analysis for data that's really messy or where  
you don't know the right questions or queries  
to make — analysis that can help you find  
patterns, anomalies, or new structures amidst  
otherwise chaotic or complex data points.”*

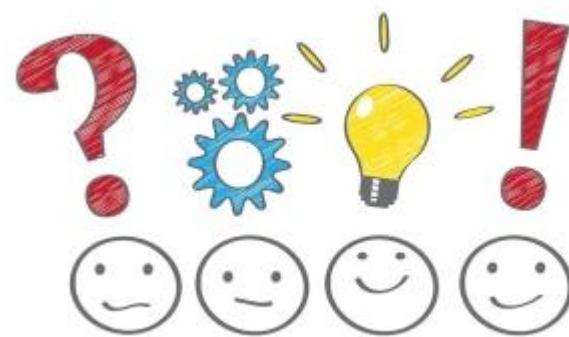
*Philip Ashlock, Chief Architect, Data.gov*

Wordle of “Thought Leaders” definition of Big Data, © Jennifer Dutcher, 2014  
<https://datascience.berkeley.edu/what-is-big-data/>

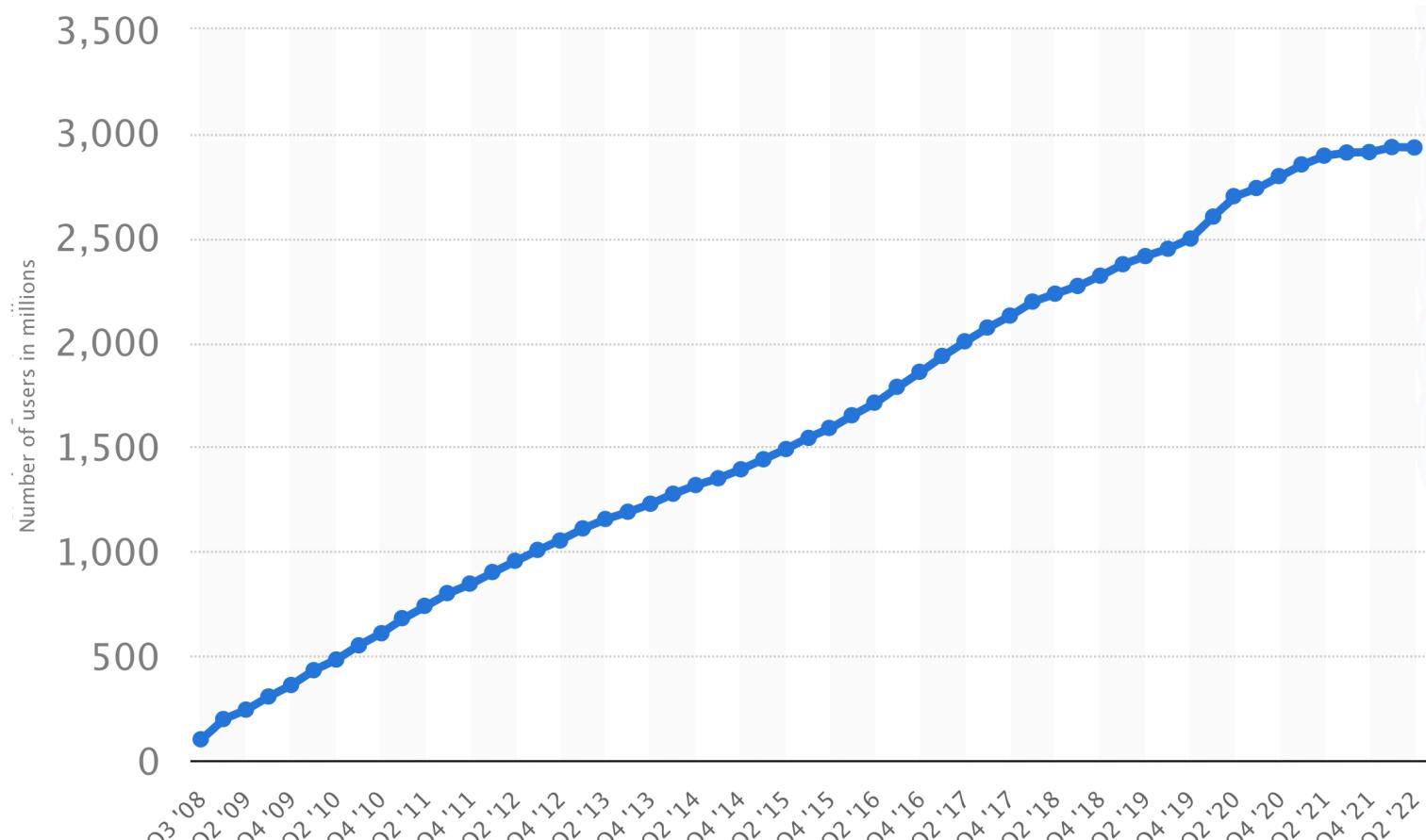
## So...What is Big Data?

Data whose characteristics exceeds the capabilities of conventional *algorithms, systems and techniques* to derive useful **value**.

<https://www.ted.com/ideas/what-is-big-data>



# Facebook Active Users

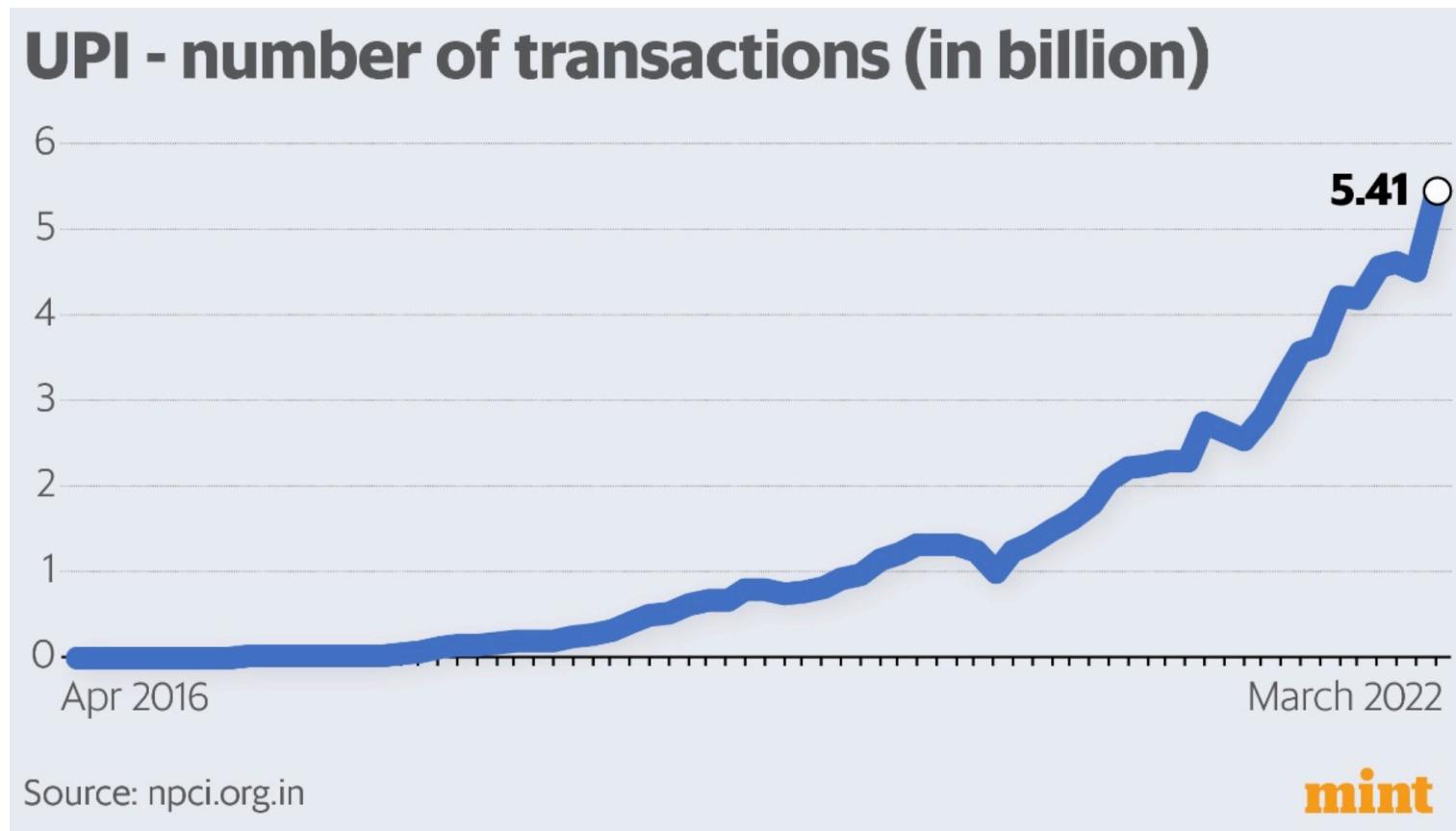


© Statista 2020

Over 2 billion users each month

# FinTech: UPI Transactions

Monthly UPI Transactions



<https://www.npci.org.in/product-statistics/upi-product-statistics>

# Big Data and Science

AAAS [Become a Member](#)

# Science

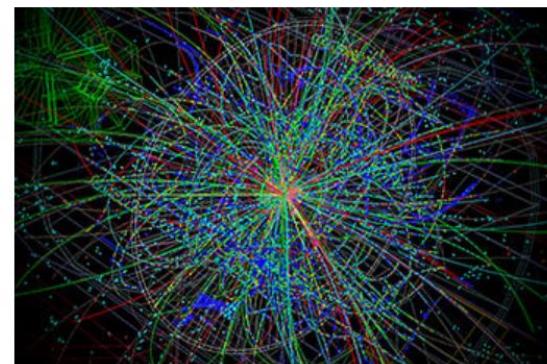
Contents ▾ News ▾ Careers ▾ Journals ▾

## AI's early proving ground: the hunt for new particles

Particle physicists began fiddling with artificial intelligence (AI) in the late 1980s, just as the term “neural network” captured the public’s imagination. Their field lends itself to AI and machine-learning algorithms because nearly every experiment centers on finding subtle spatial patterns in the countless, similar readouts of complex particle detectors—just the sort of thing at which AI excels. “It took us several years to convince people that this is not just some magic, hocus-pocus, black box stuff,” says Boaz Klima, of Fermi National Accelerator Laboratory (Fermilab) in Batavia, Illinois, one of the first physicists to embrace the techniques. Now, AI techniques number among physicists’ standard tools.

Particle physicists strive to understand the inner workings of the universe by smashing subatomic particles together with enormous energies to blast out exotic new bits of matter. In 2012, for example, teams working with the world’s largest proton collider, the Large Hadron Collider (LHC) in Switzerland, discovered the long-predicted Higgs boson, the fleeting particle that is the linchpin to physicists’ explanation of how all other fundamental particles get their mass.

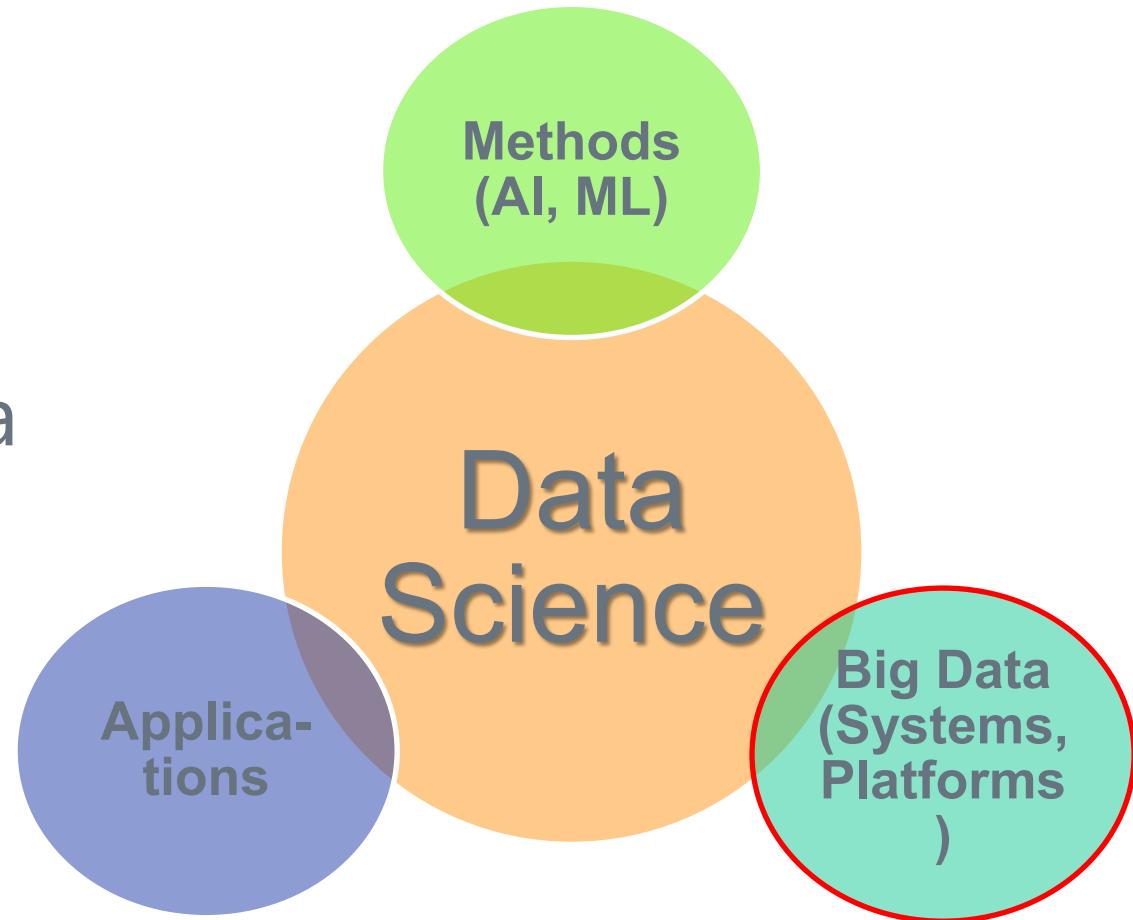
Such exotic particles don’t come with labels,



Neural networks search for fingerprints of new particles in the debris of collisions at the LHC. © 2012 CERN, FOR THE BENEFIT OF THE ALICE COLLABORATION

# Data Science

Inter-disciplinary domain at the intersection of data analysis methods, Big Data Systems and data-driven applications.



**40 ZETTABYTES**

[ 43 TRILLION GIGABYTES ]  
of data will be created by  
2020, an increase of 300  
times from 2005

**6 BILLION PEOPLE**  
have cell phones



## Volume SCALE OF DATA

It's estimated that  
**2.5 QUINTILLION BYTES**  
[ 2.3 TRILLION GIGABYTES ]  
of data are created each day



Most companies in the  
U.S. have at least  
**100 TERABYTES**  
[ 100,000 GIGABYTES ]  
of data stored

Modern cars have close to  
**100 SENSORS**  
that monitor items such as  
fuel level and tire pressure

## Velocity ANALYSIS OF STREAMING DATA

The New York Stock Exchange  
captures

**1 TB OF TRADE  
INFORMATION**  
during each trading session



By 2016, it is projected  
there will be

**18.9 BILLION  
NETWORK  
CONNECTIONS**

- almost 2.5 connections  
per person on earth



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

2005  
2020

# The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume**, **Velocity**, **Variety** and **Veracity**.

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015  
**4.4 MILLION IT JOBS**  
will be created globally to support big data,  
with 1.9 million in the United States



As of 2011, the global size of  
data in healthcare was  
estimated to be

**150 EXABYTES**  
[ 161 BILLION GIGABYTES ]



**30 BILLION  
PIECES OF CONTENT**

are shared on Facebook  
every month



## Variety DIFFERENT FORMS OF DATA

By 2014, it's anticipated  
there will be  
**420 MILLION  
WEARABLE, WIRELESS  
HEALTH MONITORS**

**4 BILLION+  
HOURS OF VIDEO**  
are watched on  
YouTube each month



**400 MILLION TWEETS**  
are sent per day by about 200  
million monthly active users

**1 IN 3 BUSINESS  
LEADERS**

don't trust the information  
they use to make decisions



Poor data quality costs the US  
economy around

**\$3.1 TRILLION A YEAR**



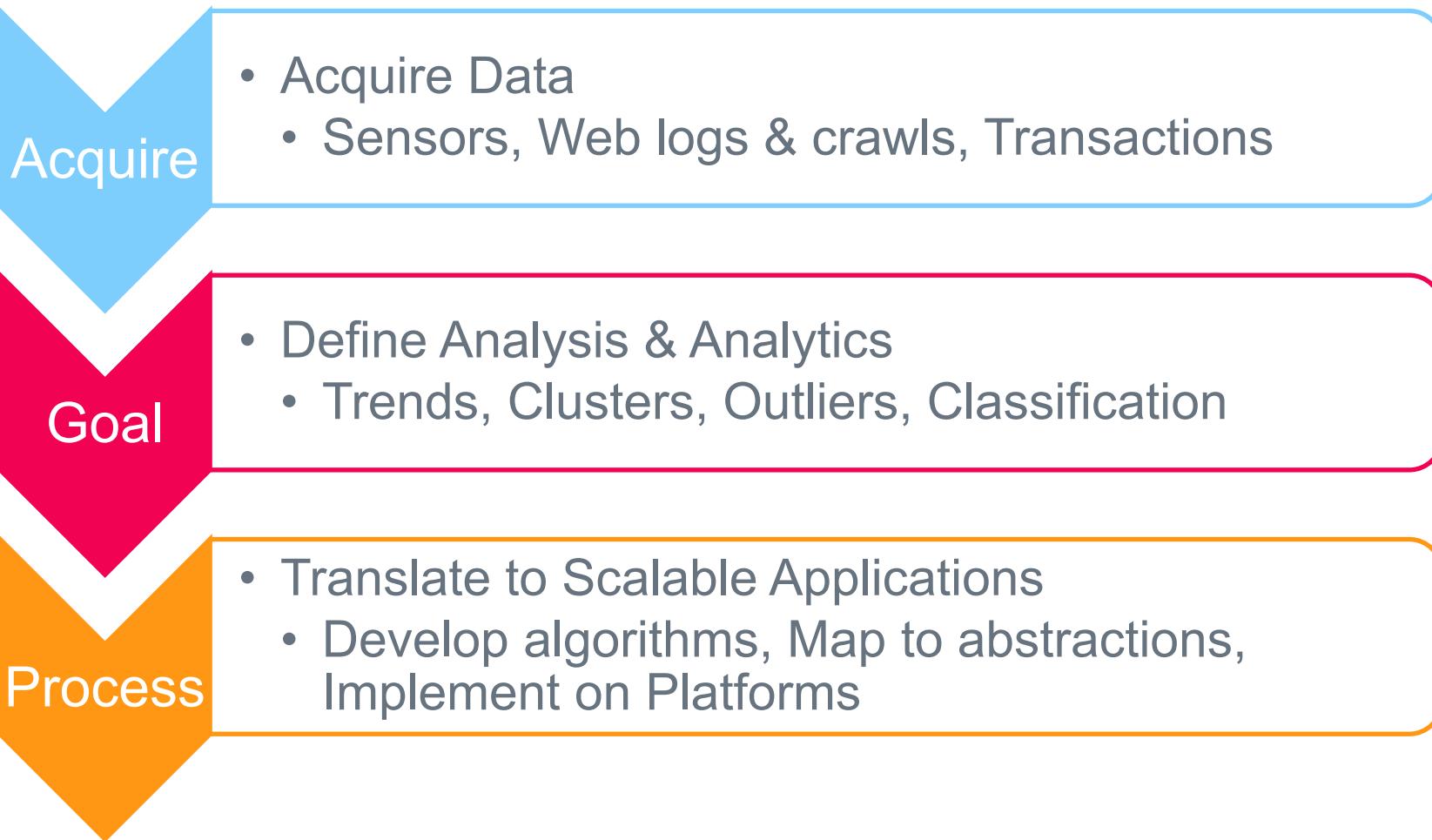
## Veracity UNCERTAINTY OF DATA

**27% OF  
RESPONDENTS**

in one survey were unsure of  
how much of their data was  
inaccurate

IBM

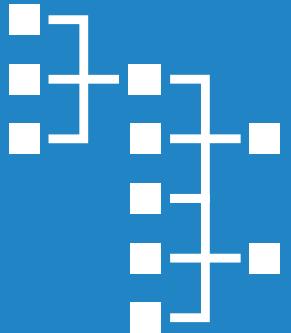
# Data Analysis Lifecycle



# Additional Reading

“

- ▷ A Survey of Big Data Research, H Fang, et al., *IEEE Network*, September/October 2015,  
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4617656/>
- ▷ Beyond the hype: Big data concepts, methods, and analytics, A. Gandomi and M. Haider, *International Journal of Information Management*, Volume 35, Issue 2, 2015, <https://doi.org/10.1016/j.ijinfomgt.2014.10.007>
- ▷ Uncertainty in big data analytics: survey, opportunities, and challenges. R.H. Hariri, et al. *J Big Data* **6**, 44 (2019).  
<https://doi.org/10.1186/s40537-019-0206-3>

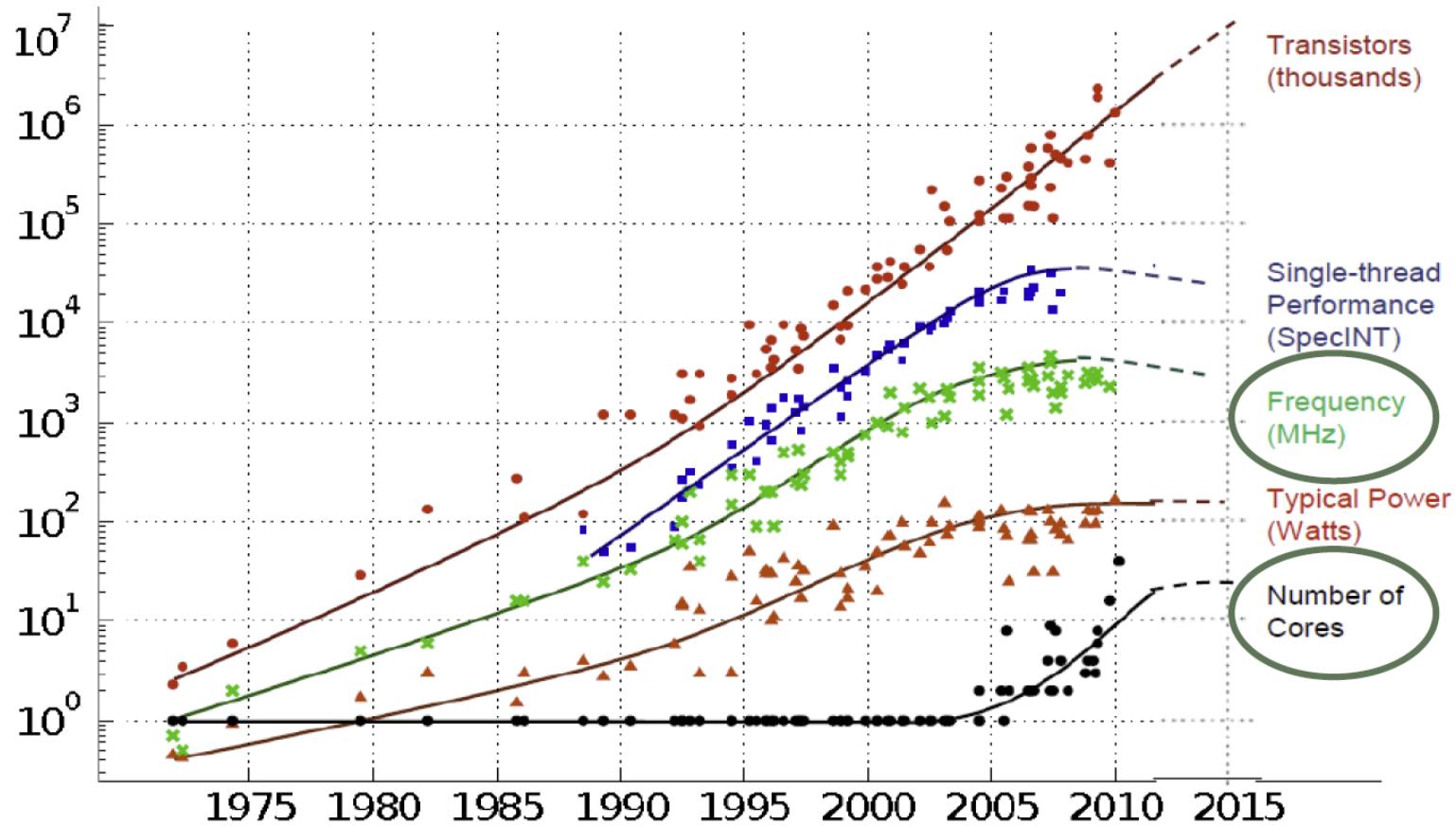


# Distributed Systems

Helping Data Science Scale

# Resources: Scale Up

Power and Heat Problems Led to Multiple Cores and Prevent Further Improvements in Speed



Jeffrey Funk, NUS

Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten  
Dotted line extrapolations by C. Moore

24

Source: Chuck Moore, Data Processing in Exascale-Class Systems, April 27, 2011. Salishan Conference on High-Speed Computing

# Resources: Scale Out

- ▷ More servers, working collaboratively
  - Connected over the network
  - High-speed LAN
  - Ethernet LAN
  - WAN
- ▷ Software is critical to enable this

# Distributed Systems

- ▷ Distributed Computing
  - Clusters of machines
  - Connected over network
- ▷ Distributed Storage
  - Disks attached to clusters of machines
  - Network Attached Storage
- ▷ *How can we make effective use of multiple machines?*
- ▷ **Commodity** clusters vs. **HPC** clusters
  - Commodity: Available off the shelf at large volumes
  - Lower Cost of Acquisition
  - Cost vs. Performance
    - Low disk bandwidth, and high network latency
    - CPU typically comparable (Xeon vs. i3/5/7)
    - Virtualization overhead on Cloud
- ▷ *How can we use many machines of modest capability?*

# Cloud Computing

- ▷ Large Data Centers
- ▷ 1000's of racks, 100k servers
- ▷ Virtualized infrastructure
  - On-demand: Rent VMs by the minute
  - 100 machines for 10mins costs same as 1 machine for 1000 minutes
- ▷ Makes use of economies of scale. Much cheaper than on-premises servers.
  - Reduce operations costs (energy, personnel)
  - Ensures capital costs (servers) are fully used
- ▷ *Have you used the Cloud?*



# Scalability

- ▷ **System Size:** Higher performance when adding more machines
- ▷ **Software:** Can framework and middleware work with larger systems?
- ▷ **Application:** As problem size grows (compute, data), can the system keep up?
- ▷ **Vertical vs Horizontal:** ?

# Scalability Metric

- ▷ If the *problem size* is  $x$  and the number of *processors available* is  $p$ , **Speedup** is:

$$\text{Speedup}(p, x) = \frac{\text{time}(1, x)}{\text{time}(p, x)}$$

- ▷ **Speedup Efficiency:** How good is the speedup relative to perfect linear speedup?

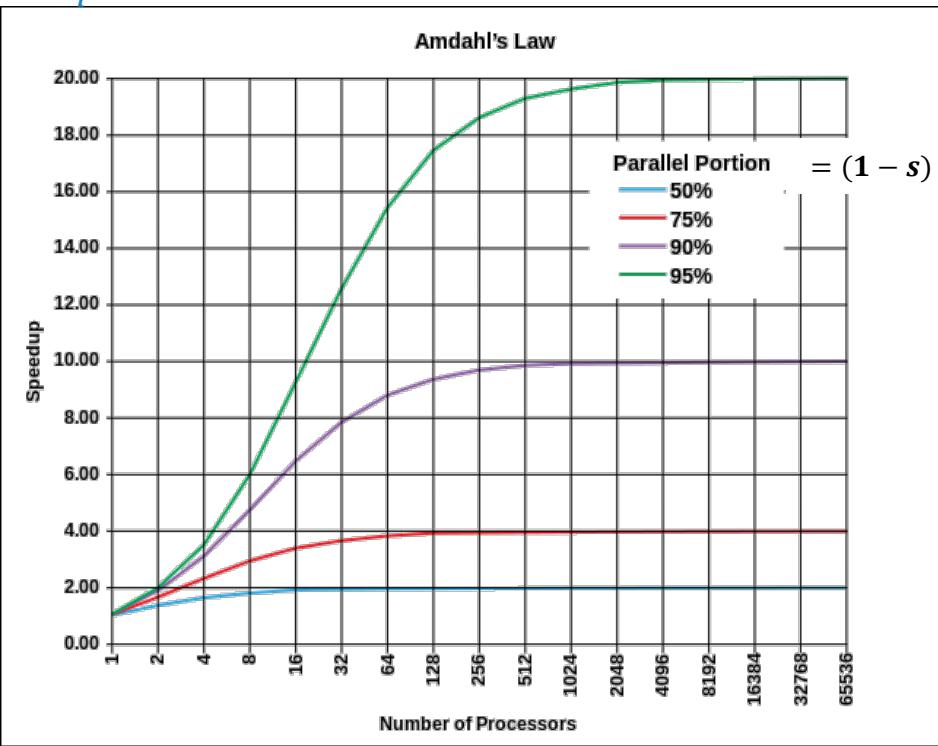
$$\text{Speedup Efficiency} = \frac{\text{Speedup}(p, x)}{p}$$

# Amdahl's Law of Strong Scaling

- ▷ Amdahl's Law for Application Scalability
  - Total problem size is fixed
  - *Speedup limited by sequential bottleneck*
- ▷  $s$  is *serial* fraction of application, cannot be parallelized
- ▷  $(1 - s)$  is fraction of application that can be *parallelized*
- ▷  $p$  is number of processors
- ▷  $t$  is time taken to process input of size  $x$  on 1 processor
- ▷ Then the speedup is limited by the serial fraction

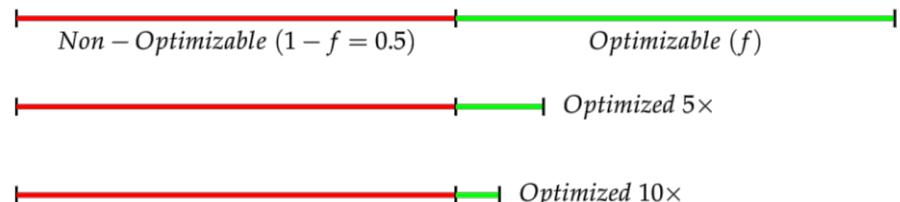
$$\text{Speedup}(p, x) = \frac{\text{time}(1, x)}{\text{time}(p, x)} = \frac{t}{s.t + \frac{(1-s)}{p}.t} = \frac{1}{s + \frac{(1-s)}{p}}$$

$$\frac{1}{s + \frac{s}{p}}$$



© Daniels220 at English Wikipedia

# Amdahl's Law



© Martha A. Kim

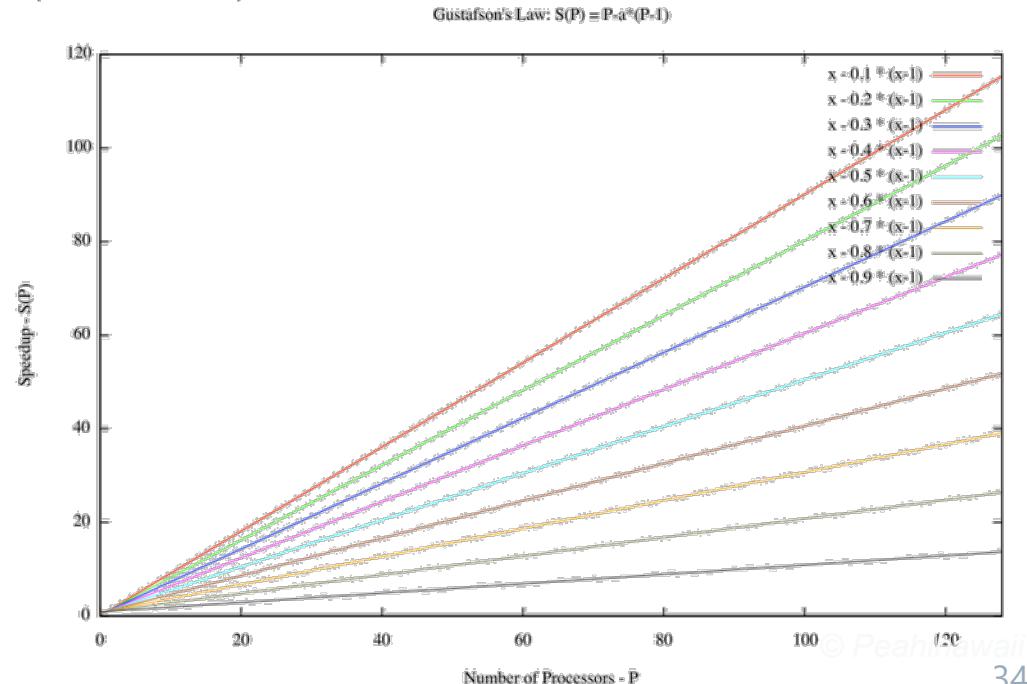
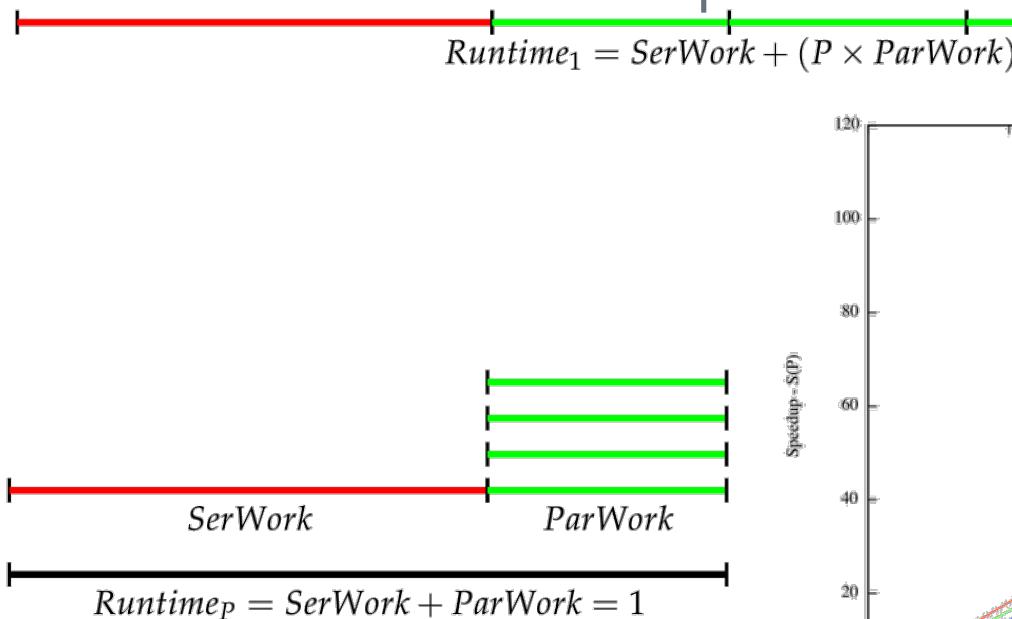
Meaningful only for small values of  $p$ , small fractions of  $s$

# Gustafson's Law of weak scaling

- ▷ Gustafson's Law of Application Scalability
  - Problem size ( $p.x$ ) increases with number of processors ( $p$ )
  - "Scaled speedup"
- ▷  $Speedup(p, x.p) = \frac{time(1,p.x)}{time(p,p.x)} = \frac{s.p.t + (1-s).p.t}{s.p.t + \frac{(1-s).p.t}{p}}$
- ▷ If  $s.p.t \approx s.t$  ...i.e., only parallelizable fraction of input work increases with number of processors
- ▷  $Speedup(p, x.p) = \frac{s.t + (1-s).p.t}{s.t + \frac{(1-s).p.t}{p}} = \frac{s + (1-s).p}{s + (1-s)}$   
 $= s + (1 - s).p$

# Gustafson's Law of weak scaling

- As input work increases, all incremental work can be parallelized. Total time remains constant as we increase number of processors.



# Scalability

- ▷ Strong vs. Weak Scaling
- ▷ **Strong Scaling:** How the performance varies with the # of processors for a *fixed total problem size*
- ▷ **Weak Scaling:** How the performance varies with the # of processors for a *fixed problem size per processor*
  - Big Data platforms are intended for “Weak Scaling”

# Additional Reading

“

- ▷ The evolution of distributed computing systems: from fundamental to new frontiers. Lindsay, D., Gill, S.S., Smirnova, D. et al. *Computing* (2021).  
<https://doi.org/10.1007/s00607-020-00900-y>
- ▷ Big Data computing and clouds: Trends and future directions, M. D. Assunção, et al., *Journal of Parallel and Distributed Computing*, Volumes 79–80, 2015, Pages 3-15, <https://doi.org/10.1016/j.jpdc.2014.08.003>

- . TRY GOOGLE COLAB BEFORE NEXT CLASS
  - . <https://colab.research.google.com/>
- .
- . Will be sharing Spark Notebook