# BIG-DATA & DATA PIPELINES

**CSCI 5408:**
**Data Management, Warehousing, and Analytics**
**Prepared By: Suhaib Qaiser (suhaibqaiser@dal.ca)**

# *Recap from last lecture...*

**Q1. How can we scale horizontally and vertically in Heroku?**

**Q2. What are Salesforce limitations as a platform?**

**Q3. How is Assignment 2 coming along?**

**Q4. What is the difference between Latency and Response delay?**

**Q5. Is latency on Web Servers bigger than Database servers? Compare**
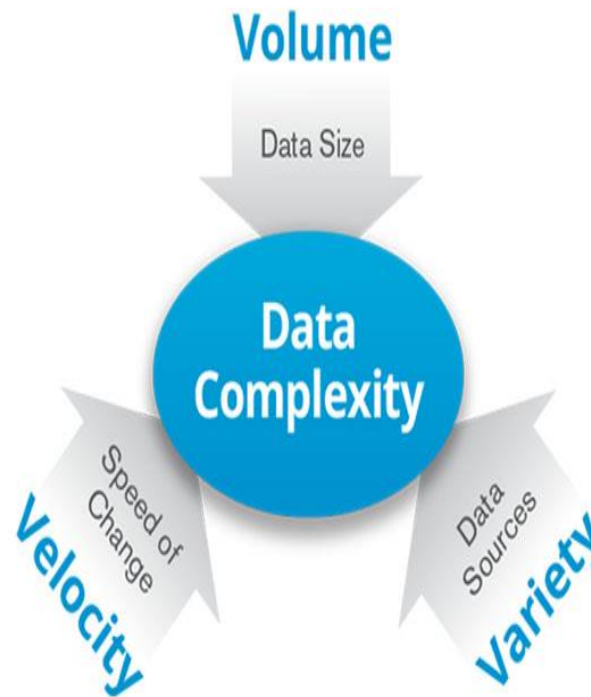
# *BIG DATA*

# *BIG DATA*

## DEFINITION

Big Data refers to data that because of its size, speed or format, that is, its volume, velocity or variety, cannot be easily stored, manipulated or analyzed with traditional methods like spreadsheets, relational databases or common statistical software



**Definition:** Big data is the confluence of the three trends consisting of Big Transaction Data, Big Interaction Data and Big Data Processing

**BIG TRANSACTION DATA**

ONLINE TRANSACTION PROCESSING (OLTP)

ONLINE ANALYTICAL PROCESSING (OLAP) & DW APPLIANCES

Oracle
DB2
Britton-Lee
Ingres
Informix
Sybase
SQLServer

Teradata
Redbrick
EssBase
Sybase IQ
Netezza
Greenplum
DataAllegro
Asterdata
Vertica
Paraccel

**BIG INTERACTION DATA**

SOCIAL MEDIA DATA

OTHER INTERACTION DATA

Clickstream
image/Text
Scientific

• Genomic/pharma
• Medical

Machine/Device

Sensors/meters/
RFID tags
CDR/mobile

**BIG DATA INTEGRATION**

**BIG DATA PROCESSING**

# *BIGDATA*

## What is BIGDATA

One way of describing big data is by looking at the three V's of volume, velocity, and variety. These come from an article written by Doug Laney in 2001, and they're taken as the most common characteristics of big data.
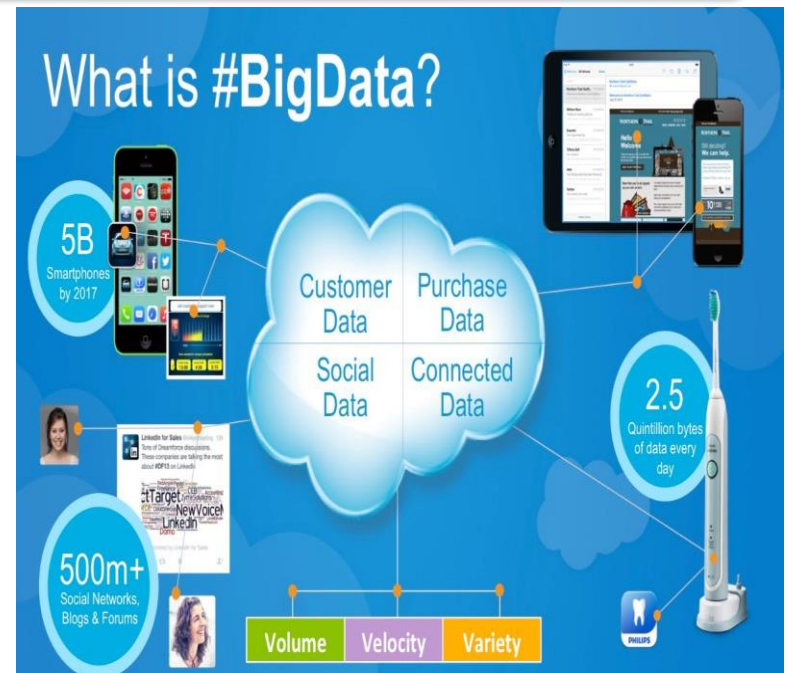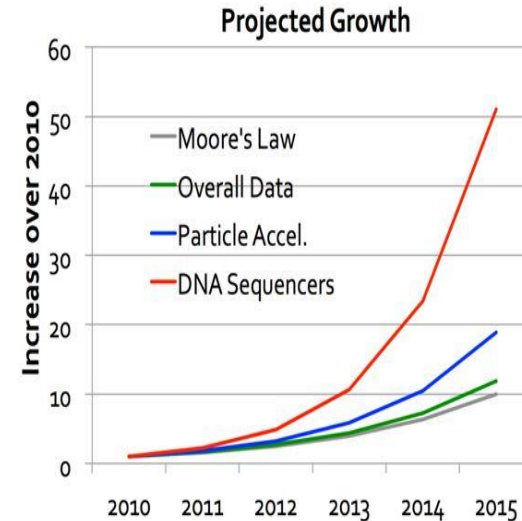
# *BIGDATA*

## What is BIGDATA

one way of describing big data is by looking at the three V's of volume, velocity, and variety. These come from an article written by Doug Laney in 2001, and they're taken as the most common characteristics of big data.

# *BIGDATA*

**VOLUME**

- big data is data that's just too big to work on your computer

## Trends: Big Data

## Moore's Law

- Moore's Law: computers double in speed every 24 months.
  - Applies also to quality-adjusted microprocessor prices, memory capacity, disks, networks, sensors, and the number and size of pixels in cameras.

- How fast would a bicyclist be if Moore's Law applied?
  - 50 years of doubling every **24** months
  - 30 km/h originally
  - ~1 billion km/h now

- Three lessons
  - Growth rates in computing are dramatic and difficult to imagine.
  - Hardware advances are important information technology drivers.
  - Humans don't really improve – they are the constants.

**Projected Growth**

- Moore's Law
- Overall Data
- Particle Accel.
- DNA Sequencers

Increase over 2010

## Data grows faster than Moore's Law!

[Kathy Yelick LBNL, VLDBJ 2012, Dhruba Borthakur]

# *BIGDATA*

## What is BIGDATA

**VOLUME**

One way of describing big data is by looking at the three V's of volume, velocity, and variety. These come from an article written by Doug Laney in 2001, and they're taken as the most common characteristics of big data.

| FEATURES | MAXIMUM LIMITS |
|---|---|
| Open workbooks | Limited by available memory and system resources |
| Worksheet size | 65,536 rows and 256 columns |
| Column Width | 255 Characters |
| Row Height | 409 points |
| Page breaks | 1,000 horizontal and vertical |

**Excel's Limit**

**Version 2003**

| FEATURES | MAXIMUM LIMITS |
|---|---|
| Open workbooks | Limited by available memory and system resources |
| Worksheet size | 1,048,576 rows by 16,384 columns |
| Column Width | 255 Characters |
| Row Height | 409 points |
| Page breaks | 1,026 horizontal and vertical |

**Excel's Limit**

**Version 2007**

# *BIGDATA*

## What is BIGDATA

**VOLUME**

| On the other hand, if you're looking at photos or video and you need to have all of the information in memory at once, you have an entirely different issue. | iPhone takes photos at two or three megabytes per photo and video at about 18 megabytes per minute, or one gigabyte per hour | video camera you could do up to 18 gigabytes per minute. And instantly you have very big data |
|---|---|---|

# *BIGDATA*

## What is BIGDATA

### VELOCITY

- Velocity, this is when data is coming in very fast

**Iris Data**



the most familiar data set for teaching the statistical procedure, cluster analysis, is the Iris data collected by Edgar Anderson and analyzed by Ronald Fisher, both of whom published their papers in 1936

Data from a social media platform, like Twitter, you may have to deal with the so-called "fire horse" That works out to 500,000,000 tweets per day and about 200,000,000,000 tweets per year.

| | |
|---|---|
| | 6000 TWEETS PER SECOND |
| | 5000 MILLION TWEETS PER DAY |
| | 200 BILLION TWEETS PER YEAR |

# *BIGDATA*

## What is BIGDATA

**VELOCITY**

- Velocity, this is when data is coming in very fast

The kind of constant influx of data, better known as streaming data, presents special challenges for analysis, because the data set, itself, is a moving target.

REAL TIME STREAM BIG DATA PROCES

If you're accustomed to working with static data sets, in a program like SPSS or R, the demands and complexities of streaming data can be very daunting, to say the least.
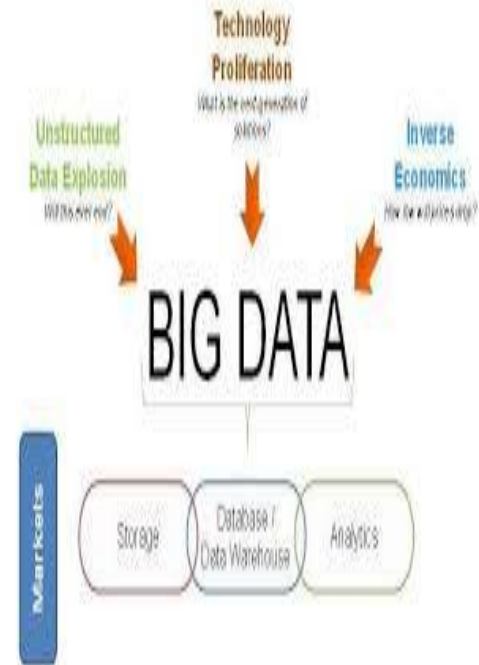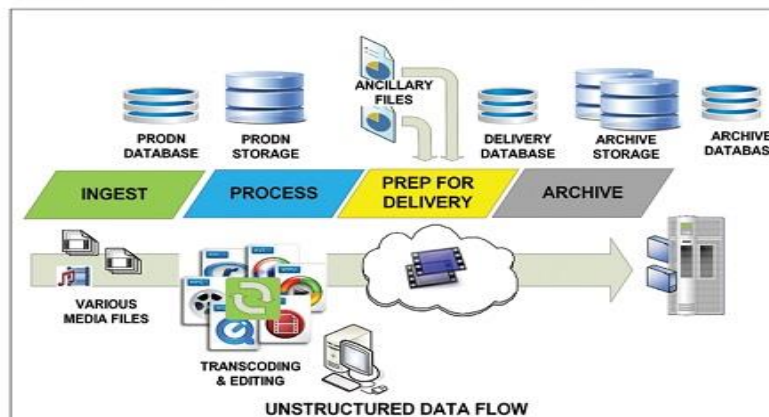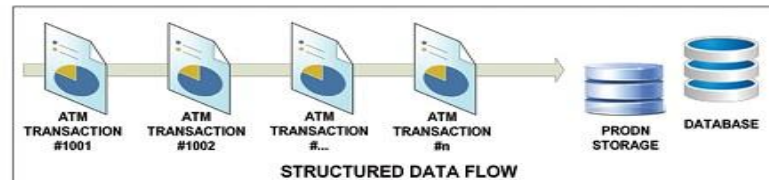
Streaming RSSI Values

# BIGDATA

## What is BIGDATA

# VARIETY

- The Structured Data that are rows and columns of a nicely formatted data set in a spread sheet, and many data sheets in many different formats. Data in unstructured text,like books and blog posts and comments on news articles and tweets. One researcher has estimated that 80 percent of enterprise data may be unstructured.



Structured Data

Unstructured Data

STRUCTURED DATA FLOW

UNSTRUCTURED DATA FLOW

INGEST — PROCESS — PREP FOR DELIVERY — ARCHIVE

Technology Proliferation

Unstructured Data Explosion

Inverse Economics

BIG DATA

Storage | Database / Data Warehouse | Analytics

# *BIGDATA*
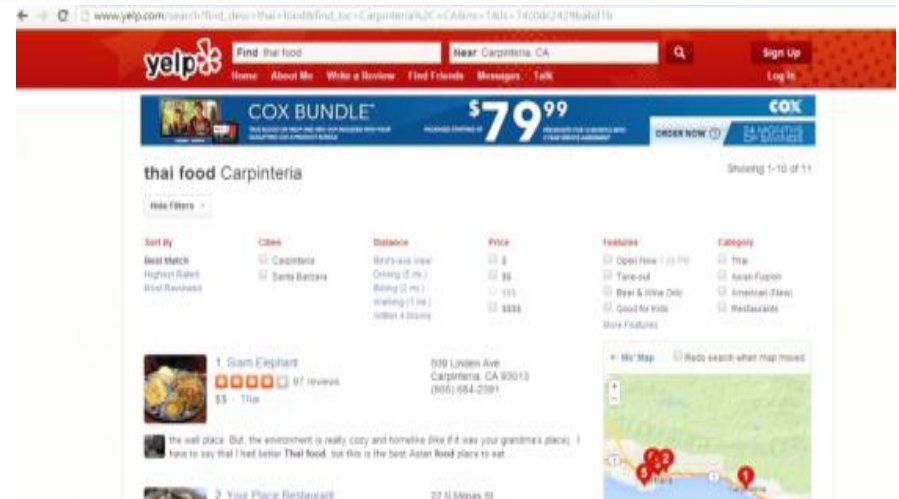
## How Big is the BIGDATA

# BIGDATA How big is Big Data

## Understanding Big Data For Consumers

For Consumers big data plays an enormous role in providing valuable services big data can provide important conveniences and functionality for consumers

some common applications of big data for consumers that you may be using already without being aware of the sophistication of the big data analysis

Apple iPhone or iPad is what Siri can do. So for instance, aside from saying what's the weather like, and Siri actually knows what it is you mean, and where you are, and what time you're talking about, it can do things like look for restaurants of a particular kind of food and see if they have reservations available
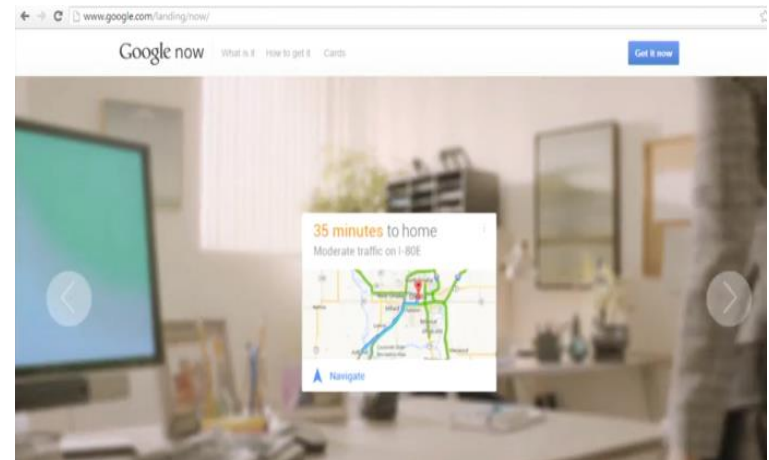**www.yelp.com/events**

# BIGDATA How big is Big Data

## Understanding Big Data For Consumers

For Consumers big data plays an enormous role in providing valuable services big data can provide important conveniences and functionality for consumers

some common applications of big data for consumers that you may be using already without being aware of the sophistication of the big data analysis

Netflix makes specific suggestions for other movies you might like
www.netflix.com/browse

# *BIGDATA*

## Understanding Big Data For Consumers

**For Consumers big data plays an enormous role in providing valuable services big data can provide important conveniences and functionality for consumers**

**some common applications of big data for consumers that you may be using already without being aware of the sophistication of the big data analysis**

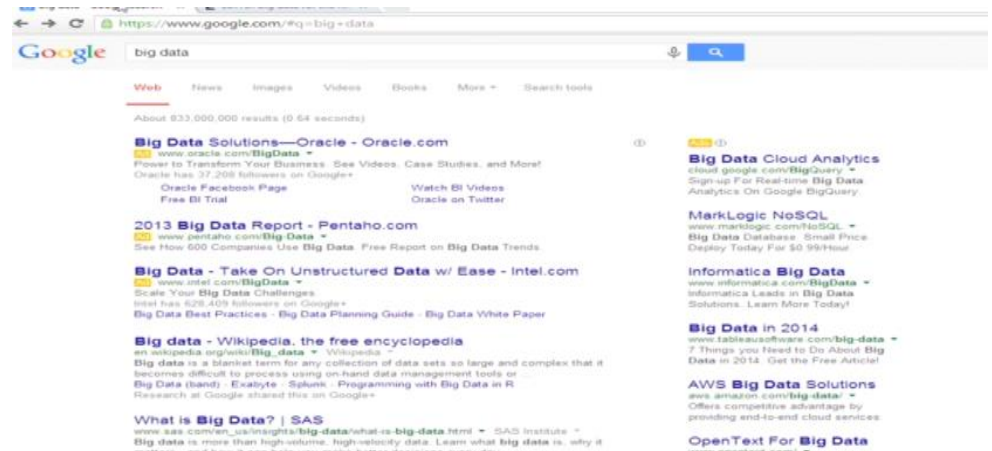**Google Now makes recommendations before you ask for them:**

**https://play.google.com/store/apps/details?id=com.google.android.launcher&hl=en**

# BIGDATA How big is Big Data

## Understanding Big Data For Business

For the business world, big data is revolutionizing the way people do commerce

most people have encountered big data in commerce, and that's in the results for Google ad searches

**Google Searches**
**www.google.com**

# *BIGDATA*

## *Understanding Big Data for Business*

Another interesting place is what's called predictive marketing. This is when big data is used to help decide who the audience would be for something before they actually get there.

**Predictive Marketing**

Predicts major life events
Looks at consumer behavior
Uses demographic info
Can purchase more data

Linked in

This is trying to predict, for example, major life events, like for instance graduating, or getting married, or getting a new job, or having a child, or any number of events that are often associated with a whole series of commercial transactions.

# *BIGDATA* How big is Big Data

## *Understanding Big Data for Business*

For the business world, big data is revolutionizing the way people do commerce

The use of big data in commerce is for fraud detection.

specifically means, how are you making the purchase? Are you online, and what website are you using? They can use geo location. Where are you physically located in the world? They can look at the IP address. What computer are you using to access the website. They can look at the log in time.

**Fraud Detection**

Point of sale
Geolocation and IP address
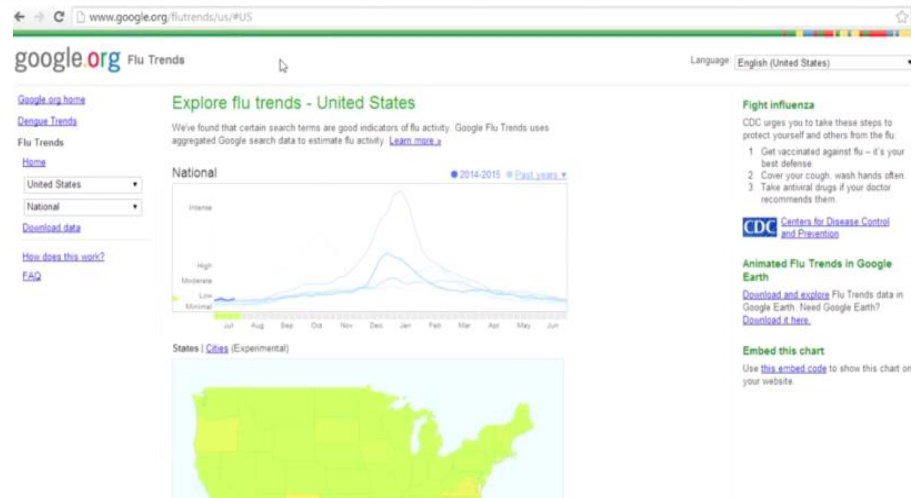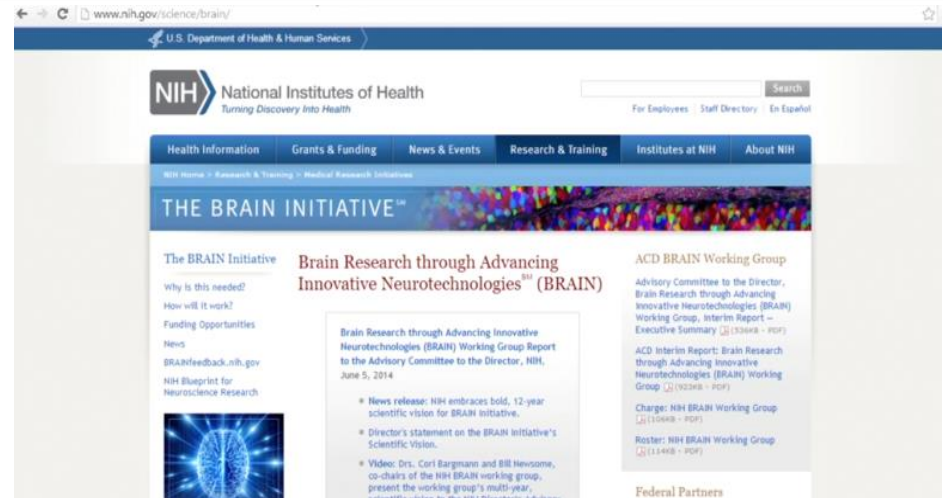Login time
Biometrics

# BIGDATA How big is Big Data

## Understanding Big Data for Research

Big data has been revolutionizing aspects of scholarship and research

examples of where big data has influenced scientific progress..

The first one we want to look at is Google flu trends where they were able to find that search patterns for flu related words were actually able to identify outbreaks of the flu in the United States much faster than there search that the Center for Disease Control could do.

**https://www.google.com/publicdata/explore?ds=z3bsqef7ki44ac_**

# *BIGDATA How big is Big Data*

## *Understanding Big Data for Research*

**Big data has been revolutionizing aspects of scholarship and research**

**examples of where big data has influenced scientific progress..**

**The National Institutes of Health created the Brain Initiative as a way of taking enormous numbers of brain scans to create a full map of brain functioning**.

[https://www.nih.gov/research-training](https://www.nih.gov/research-training)

# *BIGDATA How big is Big Data*

## *Understanding Big Data for Research*

**Big data has been revolutionizing aspects of scholarship and research**

**examples of where big data has influenced scientific progress..**

A GROUP of researchers created an application on Facebook that used scientifically valid measure of personality
https://applymagicsauce.com/demo.html

# BIG DATA

## BIG DATA AND DATA SCIENCES
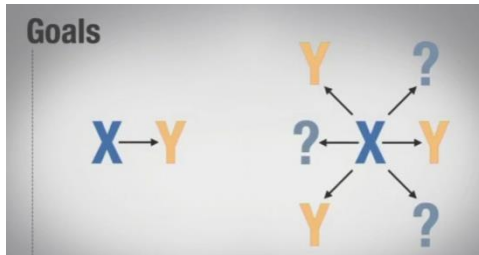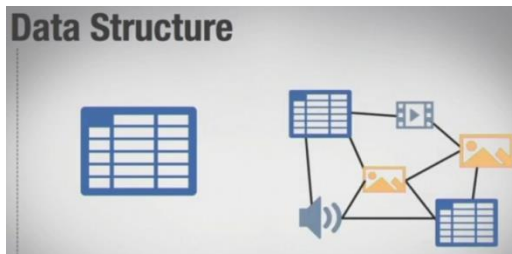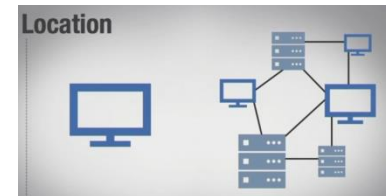
# Big Data and Data Science

## Ten Ways Big data is different from Small Data

**Goals**

The first is goals. Small data is usually gathered for a specific goal. Big data on the other hand may have a goal in mind when it's first started, but things can evolve or take unexpected directions.

The second is location. Small data is usually in one place, and often in a single computer file . Big data on the other hand can be in multiple files in multiple servers on computers in different geographic locations.
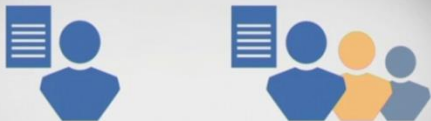
**Location**

**Data Structure**

Third, the data structure and content.

Small data is usually highly structured like an Excel spreadsheet, and it's got rows and columns of data. Big data on the other hand can be unstructured, it can have many formats in files involved across disciplines, and may link to other resources.

# Big Data and Data Science
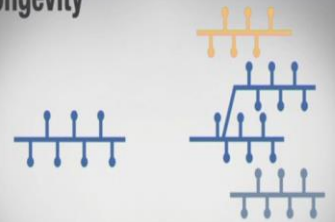
## Ten Ways Big data is different from Small Data

**Data Preparation**

Fourth, data preparation. Small data is usually prepared by the end user for their own purposes, but with big data the data is often prepared by one group of people, analyzed by a second group of people, and then used by a third group of people, and they may have different purposes, and they may have different disciplines. .

Fifth, longevity. Small data is usually kept for a specific amount of time after the project is over because there's a clear ending point. but with big data each data project, because it often comes at a great cost, gets continued into others, and things are going to stay there for a very long time. They may be added on to in terms of new data at the front, or contextual data of things or additional variables, or linking up with different files. So it has a much longer compared to a small data set.

**Longevity**

**Measurements**

Small data is typically measured with a single protocol using set units and it's usually done at the same time. With big data on the other hand, because you can have people in very different places, in very different times, different organizations, and countries, you may be measuring things using different protocols, and you may have to do a fair amount of conversion to get things consistent

# Big Data and Data Science

## Ten Ways Big data is different from Small Data


Reproducibility

Number seven is reproducibility. Small data sets can usually be reproduced in their entirety if something goes wrong in the process. Big data sets on the other hand, because they come in so many forms and from different directions, it may not be possible to start over again if something's gone wrong. Usually the best you can hope to do is to at least identify which parts of the data project are problematic and keep those in mind as you work around them. .

Number eight is stakes. On small data, if things go wrong the costs are limited, it's not an enormous problem, but with big data, projects can cost hundreds of millions of dollars, and losing the data or corrupting the data can doom the project, possibly even the researcher's career or the organization's existence.
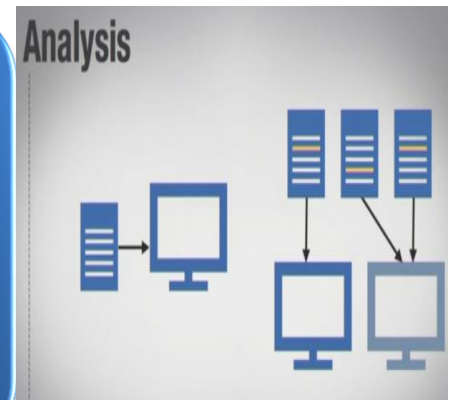

Stakes

# Big Data and Data Science

## Ten Ways Big data is different from Small Data



The ninth is what's called introspection, and what this means is that the data describes itself in an important way.

The final characteristic is analysis. With small data it's usually possible to analyze all of the data at once in a single procedure from a single computer file. With big data however, because things are so enormous and they're spread across lots of different files and servers, you may have to go through extraction, reviewing, reduction, normalization, transformation, and other steps and deal with one part of the data at a time to make it more manageable, and then eventually aggregate your results.

# *QUIZ*

Q1. Was dataset given to you in Assignment 2 an example of Big data?

Q2. Can Casandra hold Big Data?

Q3. What conventional tools fail to work with Big data applications?

Q4. What is the role of Hadoop in Big Data?

Q5. How can we use Casandra for Big Data?