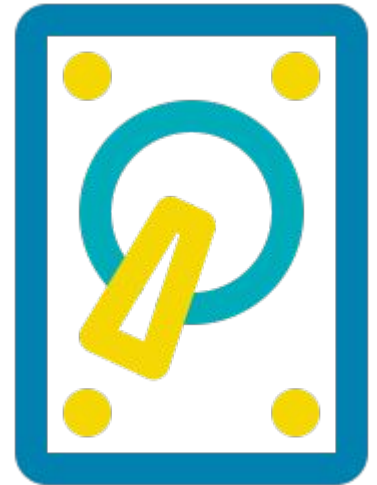
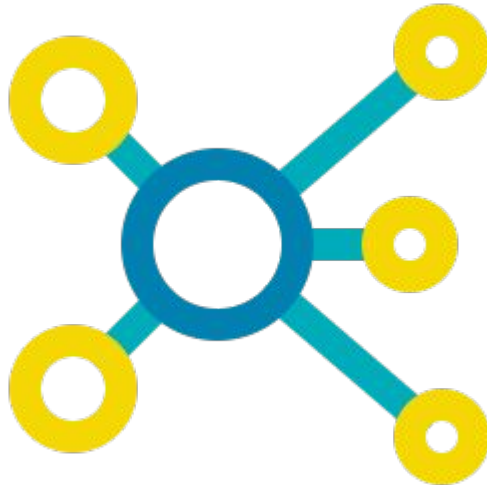


# Overview of Big Data (3Vs)

(or is it 4Vs ...)

Suzanne Little

You may have heard the term but just  
how **big** is “big data”?



# Big data

Characterised by 3 'V's

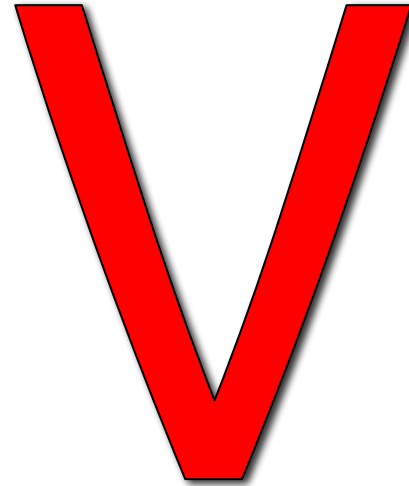
Volume

Variety

Velocity

a 4th V is sometimes added ...

April 1st post [42 Vs of Big data!](#)



# Big data: Volume

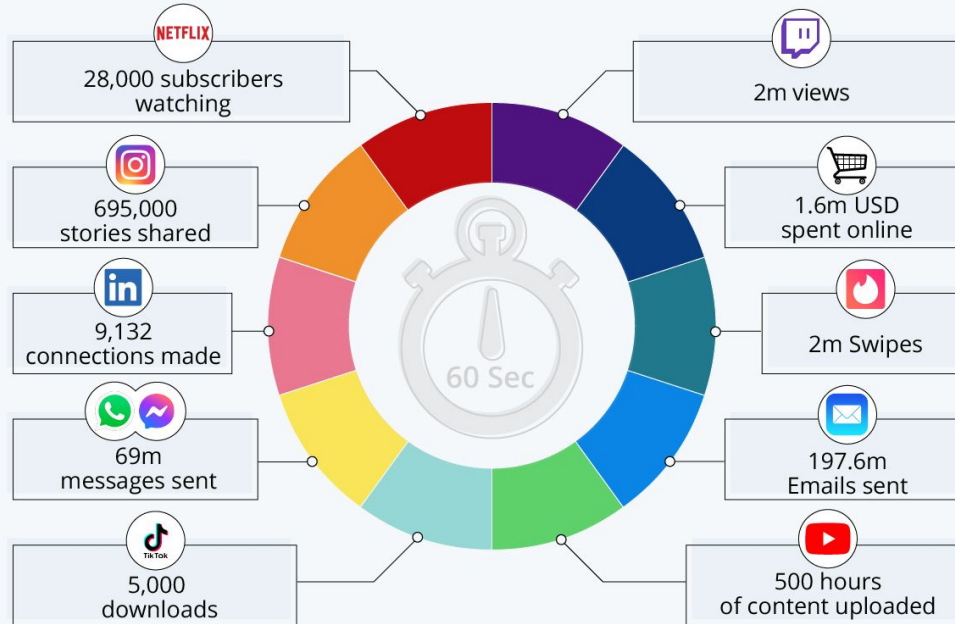
- Refers to the **amount** of data
- For big data, varies from terabytes to petabytes to zettabytes
- Can you open the whole dataset in your PC? Probably not big.
- In 2008, Google was already processing 20,000 terabytes of data (20 petabytes) a day
- In 2018, Google processes 40,000 searches per second!
- Social media produces vast quantities of data per minute!

<https://www.forbes.com/sites/bernardmarr/2018/05/21/how-much-data-do-we-create-every-day-the-mind-blowing-stats-everyone-should-read/#672fd38760ba>



# A Minute on the Internet in 2021

Estimated amount of data created  
on the internet in one minute



Source: Lori Lewis via AllAccess



statista

[https://www.statista.com/chart/25443/  
estimated-amount-of-data-created-on-  
the-internet-in-one-minute/](https://www.statista.com/chart/25443/estimated-amount-of-data-created-on-the-internet-in-one-minute/)



# 2020 *This Is What Happens In An Internet Minute*



# 2019 *This Is What Happens In An Internet Minute*



# Big data: Variety

- Refers to **differing** types and data sources
- Structured, semi-structured and unstructured data
- Organisations may need to combine data from many different sources
- With the proliferation of analytics and sensor data, the variety of data is expanding rapidly
- Q: How many digital cameras do you own?



# Big data: Velocity

- Data in motion -- **dynamic**, temporal
- We may need to process data as it arrives
  - ...because we cannot store such volumes
  - ...because we need timely processing
- Related notion of latency (lag-time)
  - The time between data being generated and processed
  - Some applications, such as fraud detection are highly time-sensitive.

# Big data: The 4th V - Veracity

- Introducing the notion of data **uncertainty**
  - Veracity refers to how reliable the data is
- Is the data correct?
- Is it out-of-date?
- Is it complete?
- Data cleansing helps greatly
  - Using techniques such as data fusion, stochastic models
  - But with the huge volumes of data being generated – errors will slip in

## 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



WORLD POPULATION: 7 BILLION

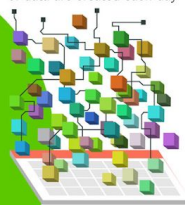
## Volume SCALE OF DATA

2005

2020

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

The New York Stock Exchange captures

**1 TB OF TRADE INFORMATION**

during each trading session



## Velocity ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be

**18.9 BILLION NETWORK CONNECTIONS**

— almost 2.5 connections per person on earth



Modern cars have close to

**100 SENSORS**

that monitor items such as fuel level and tire pressure



# 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**



**27% OF RESPONDENTS**

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

## Veracity UNCERTAINTY OF DATA

# Resources

Video: Introduction to Big Data, [O'Reilly](#)

Caesar Wu and Rajkumar Buyya and Kotagiri Ramamohanarao (2016) “Big Data Analytics = Machine Learning + Cloud Computing”, <https://arxiv.org/abs/1601.03115> or version at <https://learning.oreilly.com/library/view/big-data/9780128093467/B9780128053942000015.xhtml>