

BIG DATA PROCESSING

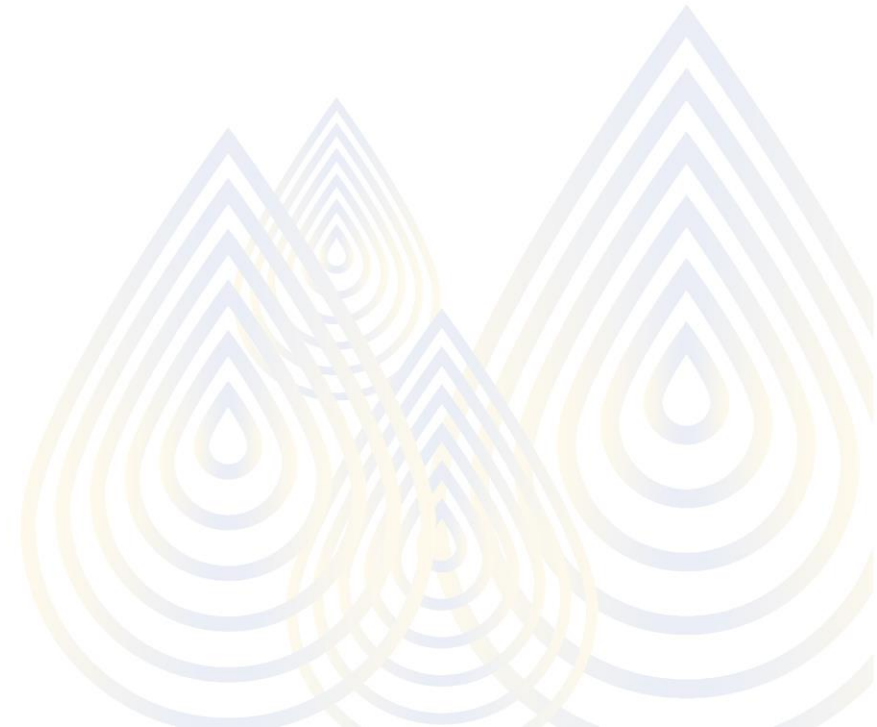
EGCI 466

Week 1



Google classroom

- rdhi44x5



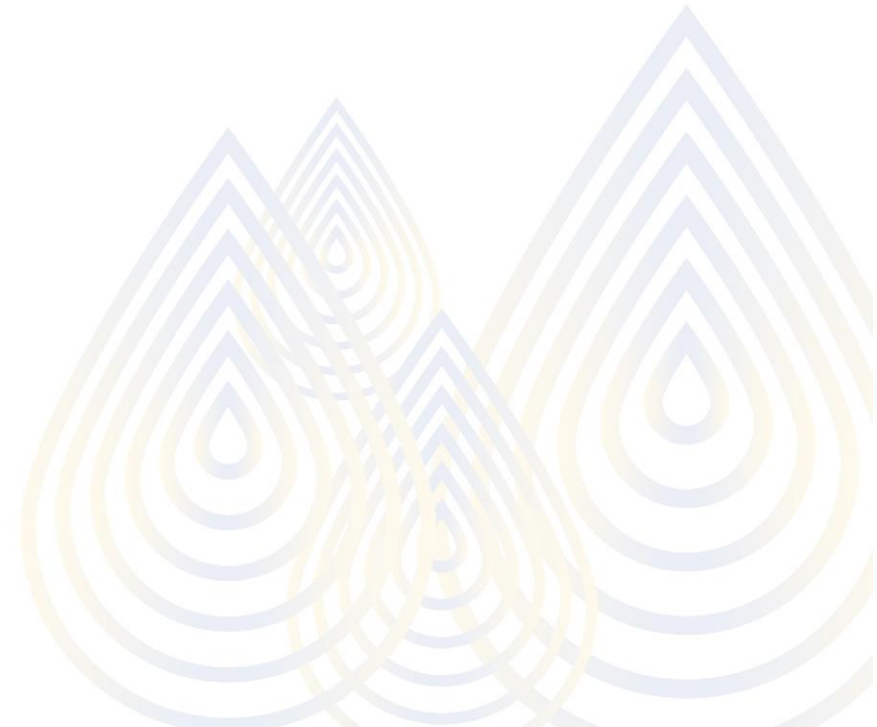
Course Outline

- [Schedule](#)

Evaluation	Percentage		
Weekly Lab Report/ Cloud Skill	10%		
Participation	10%		
Homework/Assignment	20%		
Exam/ Quizzes	30%		TBC
Midterm Project	15%		TBC
Final Project/ Presentation	25%		
	110%		

Learning Objectives (Today)

- Understand what "Big Data" means
- Explore real-world Big Data challenges
- Overview of key Big Data tools and architecture
- Set expectations for labs and projects



Why big data

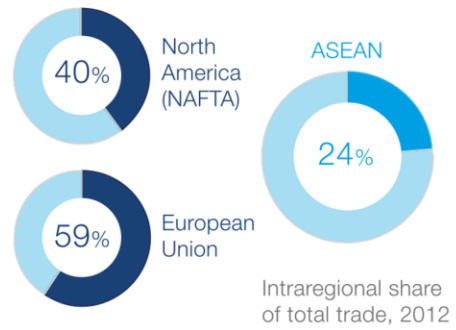
Three global trends create opportunities to transform Southeast Asia by 2030.

1

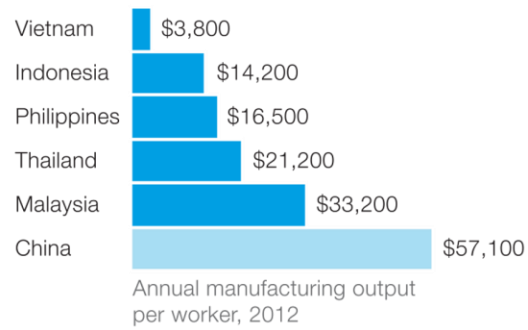
Capturing a greater share of global flows

Up to \$615 billion in annual economic value

The ASEAN¹ Economic Community (AEC) sets the stage for **greater intraregional trade**



To attract more global production, Southeast Asia must **raise labor productivity**



2

Riding the urbanization wave

Up to \$930 billion in annual economic value

An **expanding** consumer class

2013  81 million households

2030  163 million households

\$7 trillion in investment needed for infrastructure, housing, and commercial space



3

Deploying disruptive technologies

Up to \$625 billion in annual economic value



Mobile Internet



Big data



Internet of Things



Automation of knowledge work

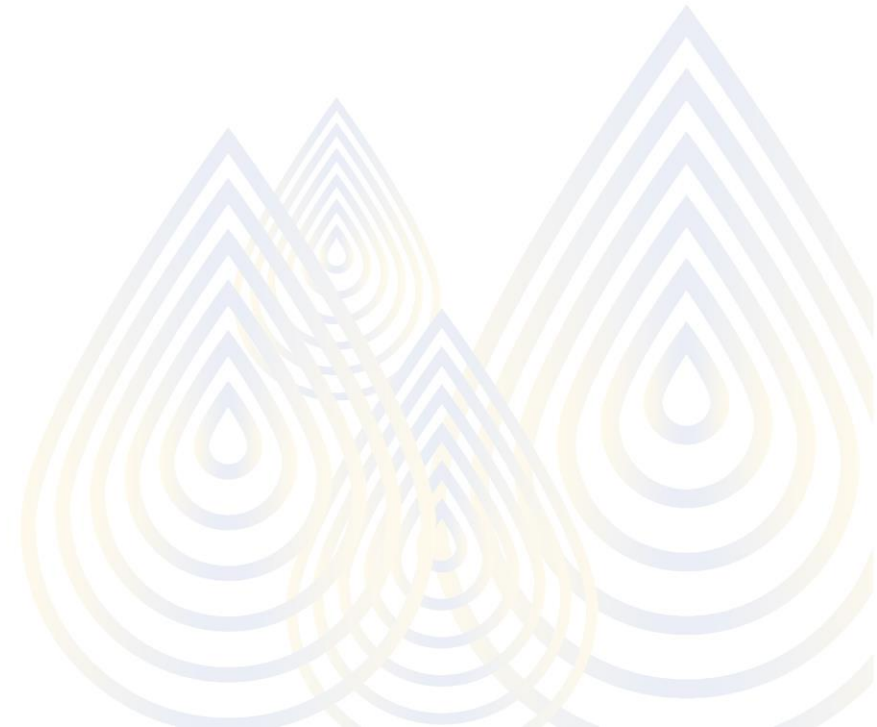


Cloud

¹Association of Southeast Asian Nations.

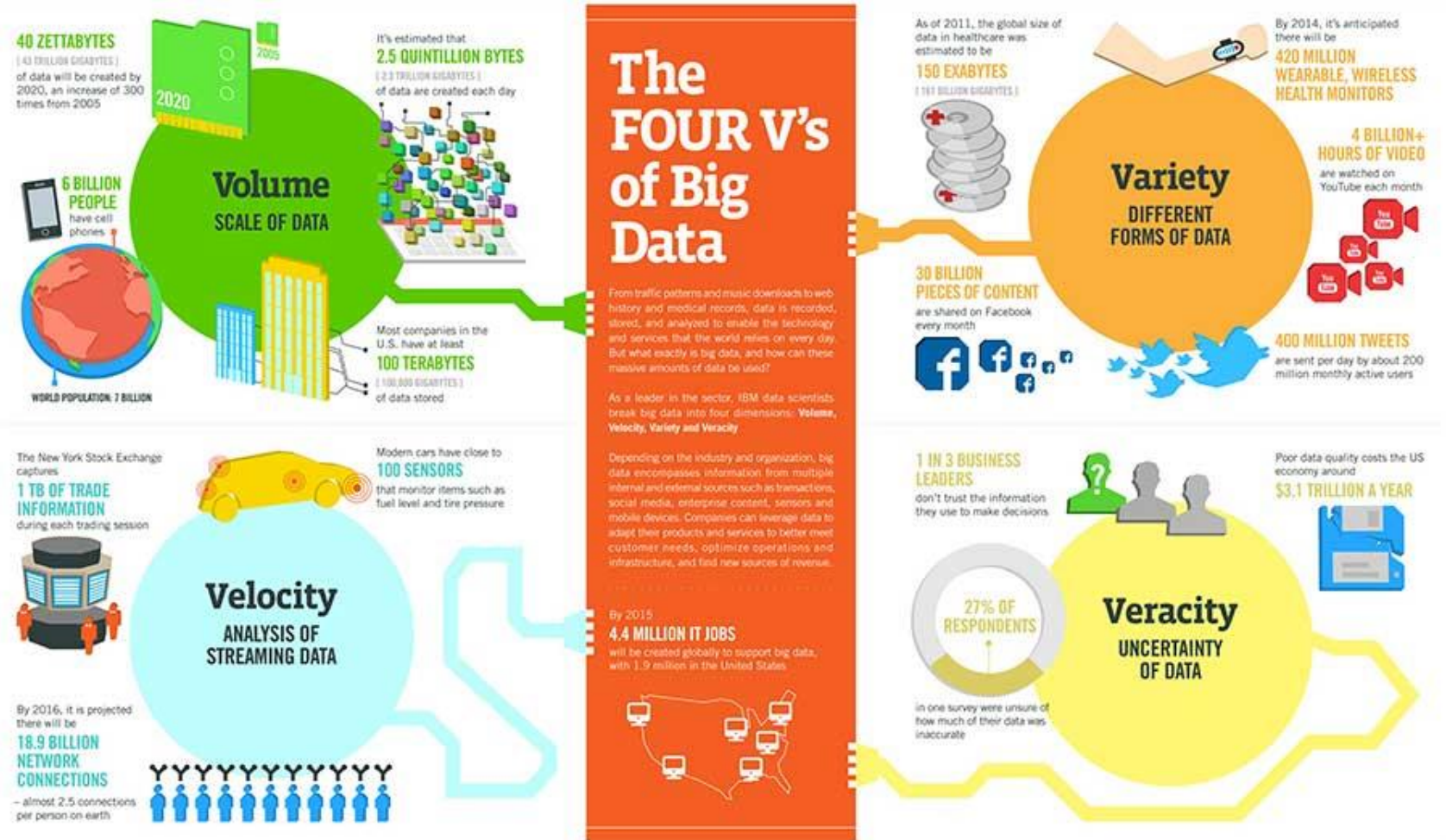
What is big data?

- 4 Vs. → 5 Vs.
- Volume
- Variety
- Velocity
- Veracity
 - Value
 - Variability
 - Valence



What is big data?

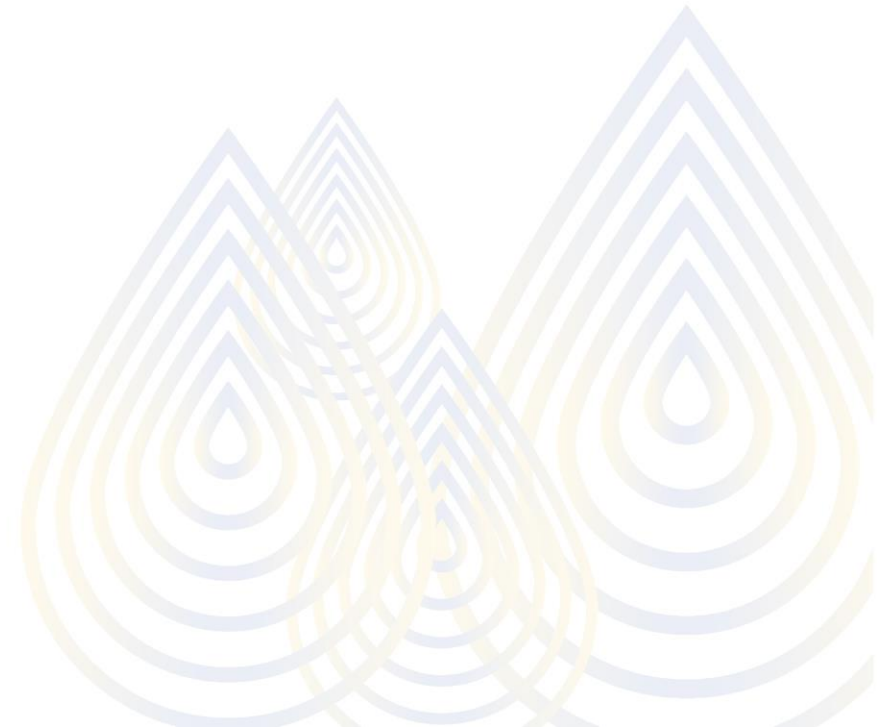
- 4 Vs. → 5 Vs.
- Volume
- Variety
- Velocity
- Veracity
 - Value
 - Variability



Applications

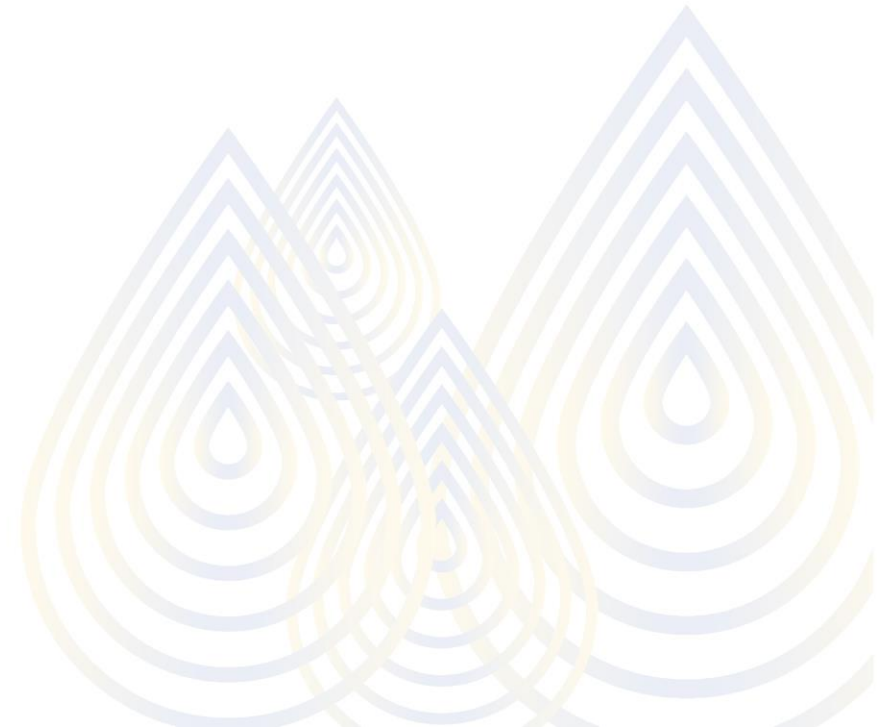
Big Data → Better Model → Better Precision

Personalized marketing



Personalized marketing

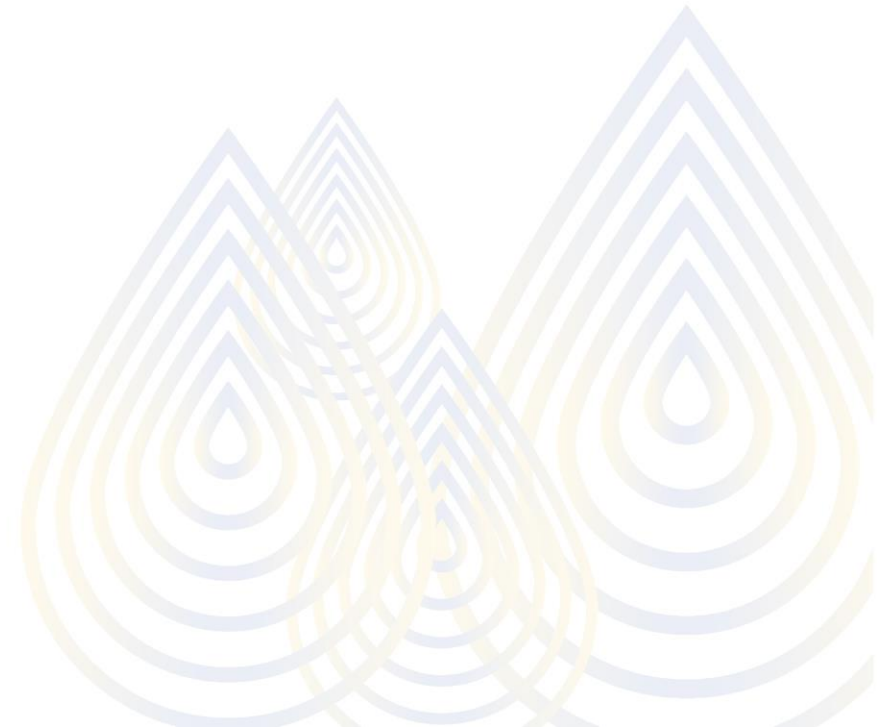
- Recommendations Engines
 - Netflix, [Amazon](#), Shopee
- Sentiment Analysis
- Real-time Fraud Detection
- Smart Cities with IoT sensors



Mobile Advertising

- Location based advertising
- Geolocation data

Is this a big data? If so why?





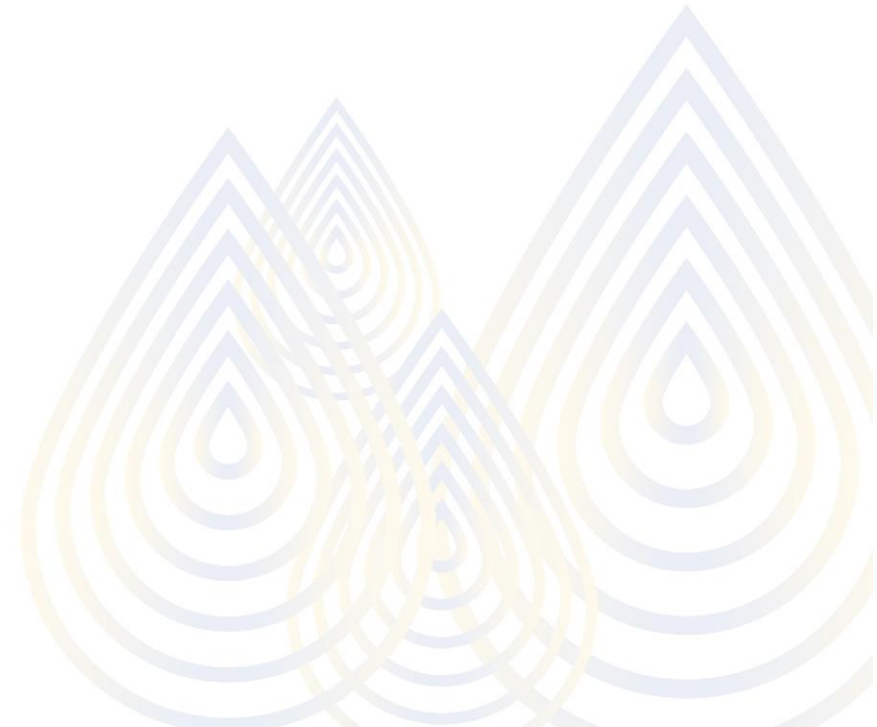
Mahidol University
Wisdom of the Land

Smart city



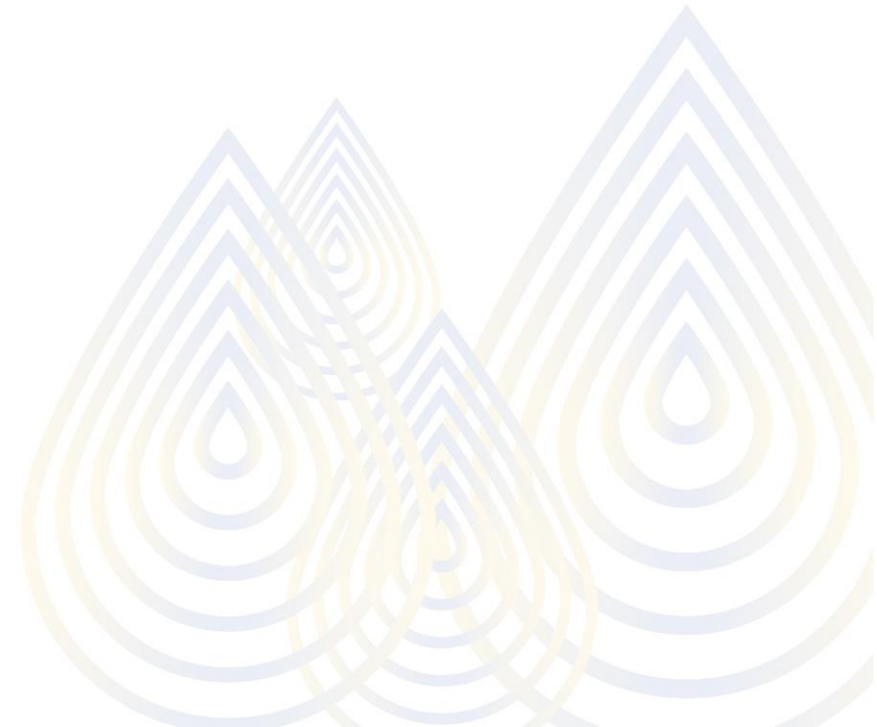
Fun Fact about big data

https://www.slideshare.net/BernardMarr/big-data-25-facts/2-Every_2_dayswe_create_as



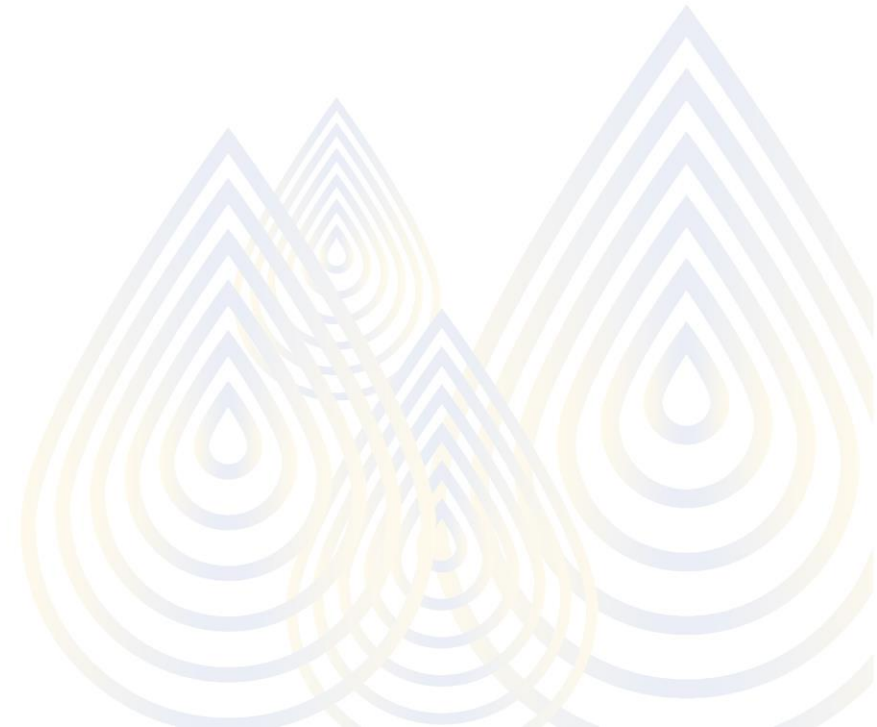
Where do data are from??

- Machine generated
- Human Generated
- Organization generated



Machine generated data

- Big plane → big data
 - 0.5 terabytes / flight
- Smart phones/ Smart Watch
 - GPS, HR, O2,

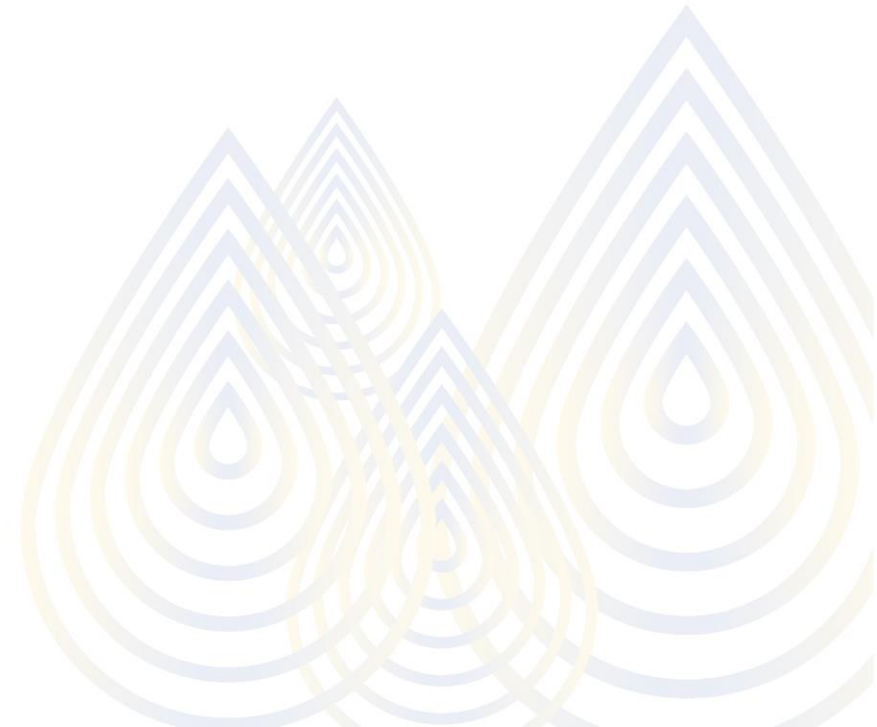


Why it is useful?

- Detect mal-function
- Notify of unhealthy behaviour

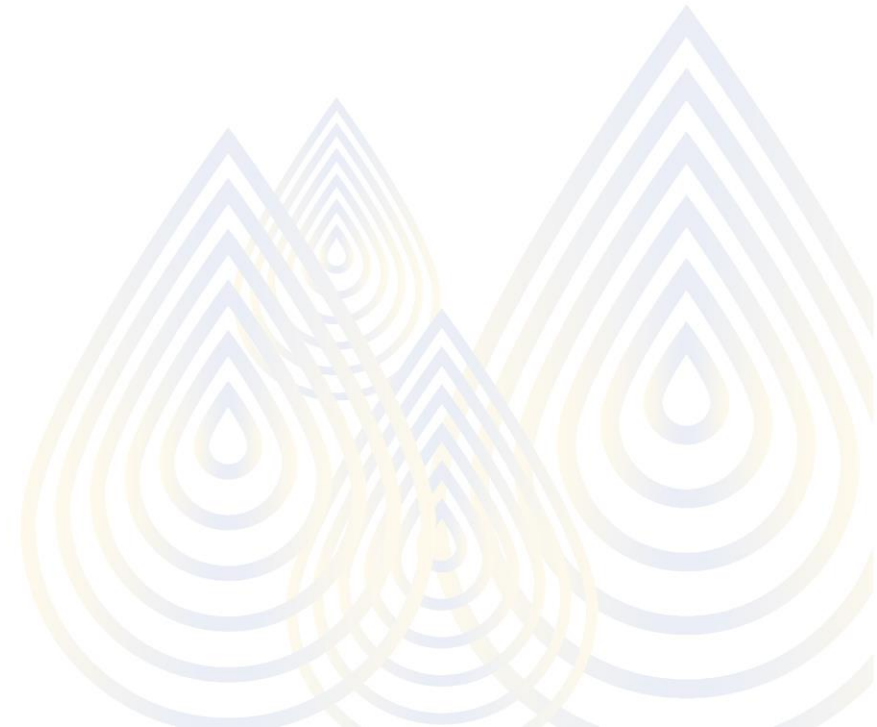
Business

- Fraud detection
- Business Planing
- System moniting/control
- Customer statisfication

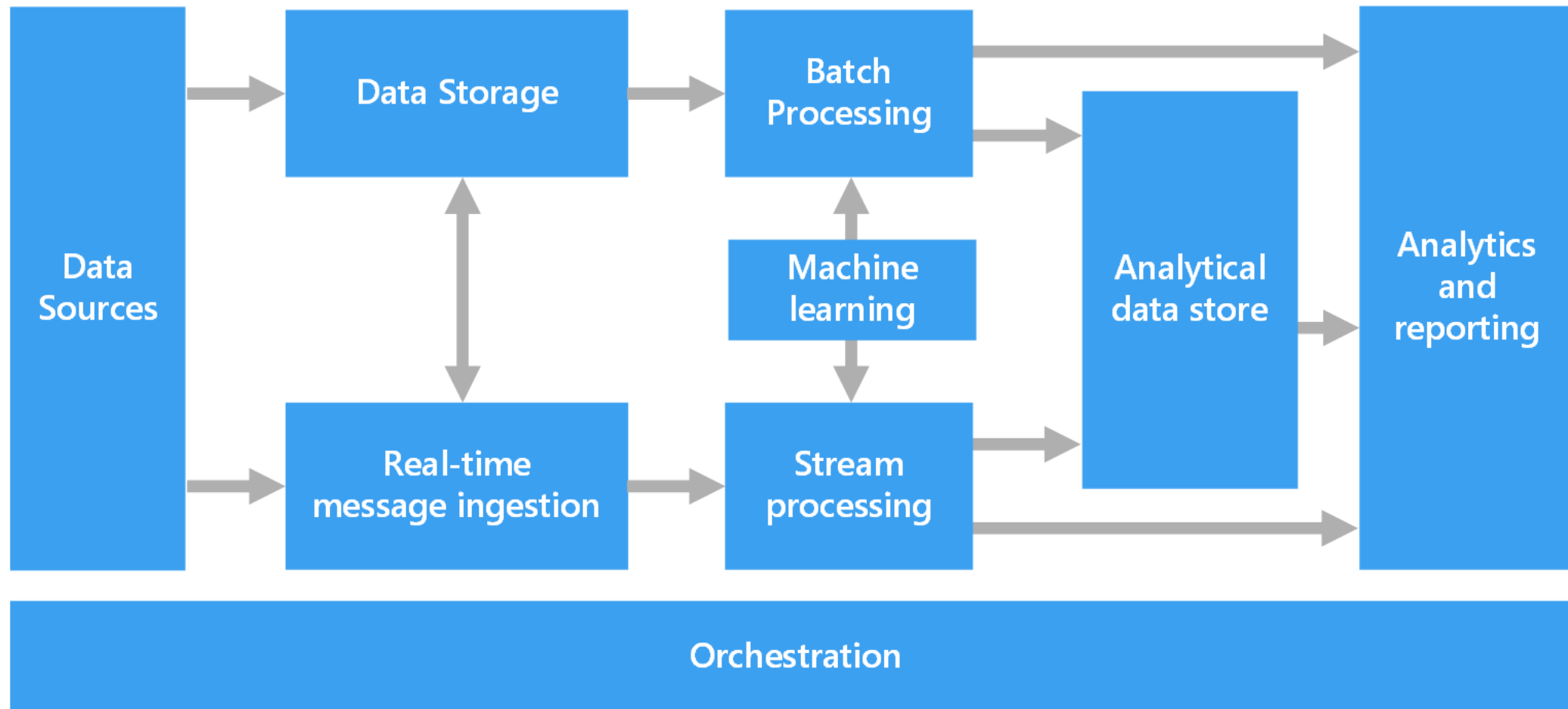


Why traditional System fail?

- Single-node databases can't handle petabytes
- Slow processing times
- Limited scalability and fault tolerance
- Need for distributed computing



Topic covered



Course covered

- Hadoop Basics: HDFS, MapReduce
- Cloud Simulation
- Spark Processing: RDDs, DataFrames, SQL
- NoSQL Databases: MongoDB, Cassandra, Neo4J
- Real-Time Streaming: Spark Streaming
- Midterm Project:
 - Build a mini Big Data pipeline with no data processing and literature review
- Final Project:
 - Build a mini Big Data pipeline with data processing

Setting Up for Success

- Installations: Spark local, Python, Docker (optional)
- Join Discord group (optional)

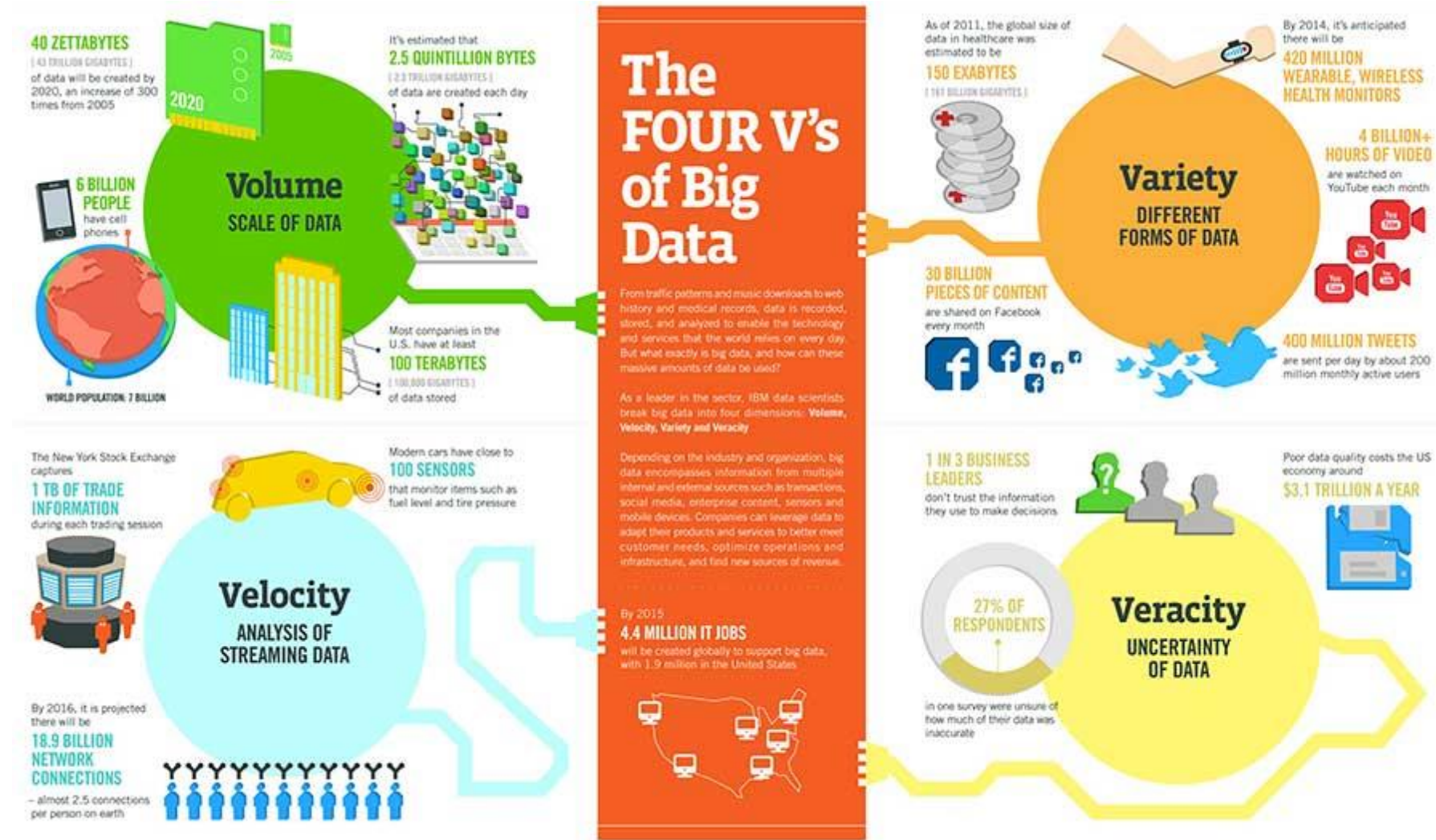
<https://discord.gg/VyWp8pSw>

- Join Cloudskill (I will add your email)
- Apply for Google cloud credit (I will provide education credit later)





Group Discussion



Discuss 4 Vs

- What is it?
- Where do you think the data come from?
- What is the big challenge?
- What solutions have been introduced?
- What tools have been used?

