

Introduction to Emerging Technologies

Chapter Two

Data Science

Outline

- **An Overview of Data Science**
- Data Types and Data Representation
- Data Value Chain
- Basic Concepts of Big Data

Introduction

- In the previous chapter, the concept of the role of data for emerging technologies was discussed.
- In this chapter, you are going to learn more about data science, data vs. information, data types and data representation, data value chain, and basic concepts of big data.

What is Data Science?

No single definition

Components:

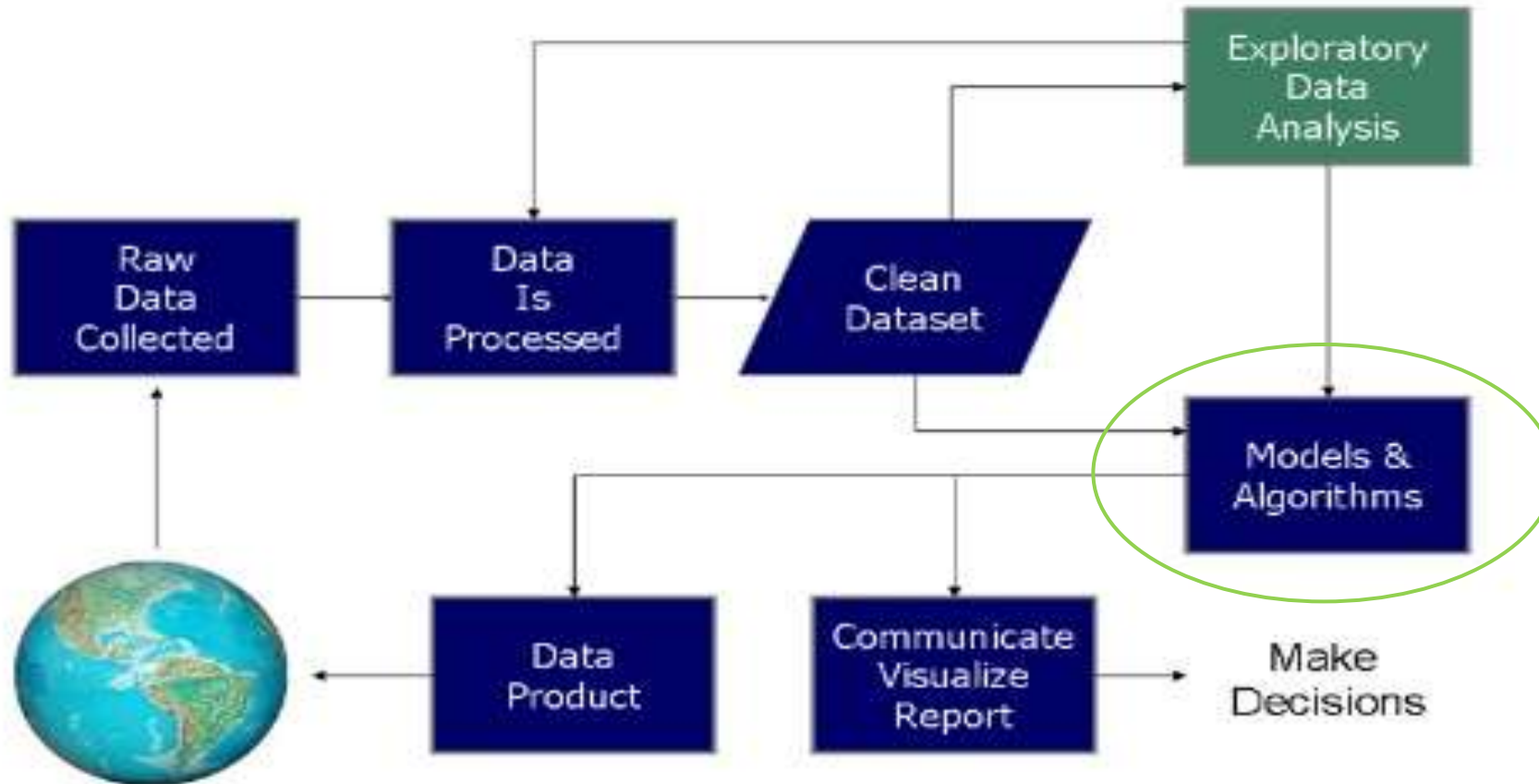
- Data-driven (the more the better)
- Interdisciplinary (math, stat, CS, ...)
- Extract knowledge from observed data

“Data Science refers to an emerging area of work concerned with the collection, preparation, analysis, visualization, management and preservation of large collections of information.”

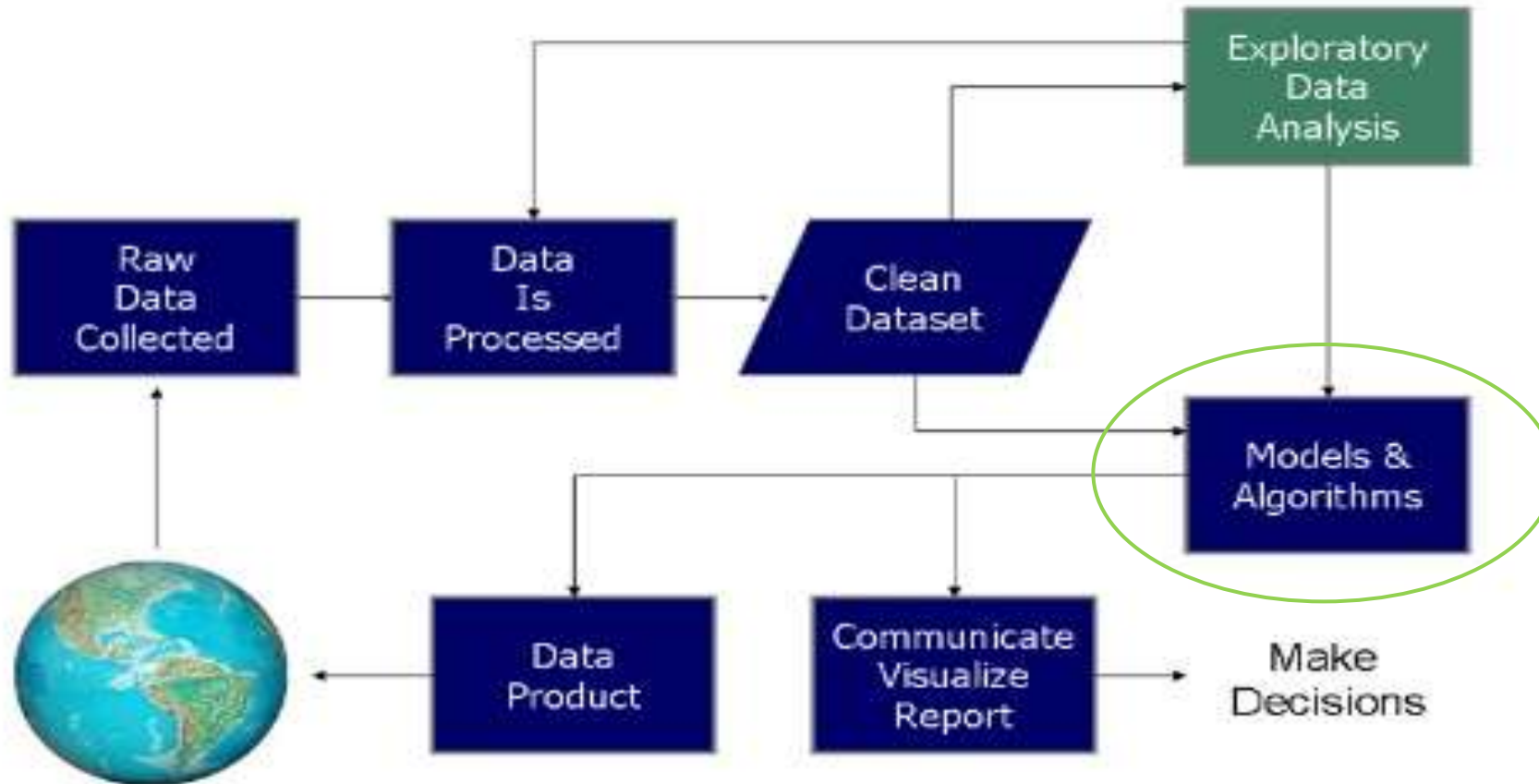
*Data Science is about the whole **processing pipeline** to **extract (and utilize)** information out of data.*

***Data Scientists** understand and care about the whole **data pipeline**.*

Data Science: process



Data Science: process



The field of algorithms has traditionally assumed that the input data to a problem is presented in **random access memory**, which the algorithm can repeatedly access. This is not feasible for problems involving **enormous amounts of data**. The **streaming model** and other models have been formulated to reflect this. In this setting, **sampling** plays a crucial role and, indeed, we have to **sample on the fly**.

DATA SCIENCE IS ABOUT *DATA PRODUCTS*

- **Data-driven apps**
 - Spellcheckers
 - Language Translators
 - Automatic Image Captioning apps
- **Interactive visualizations**
 - Google flu application
 - Global Burden of Disease
- **Online Databases**
 - Enterprise data warehouse

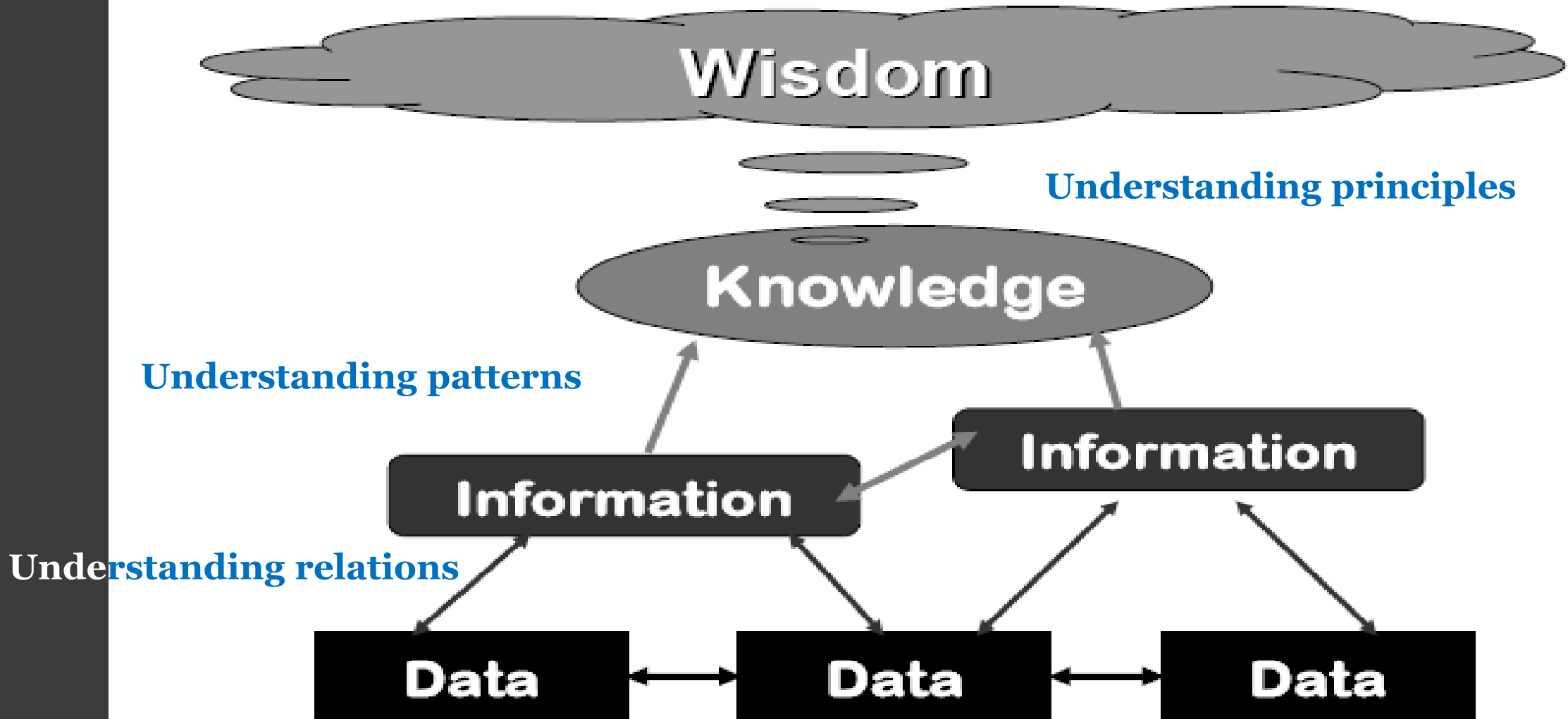
What is data?

A representation of **facts**, **concepts**, or **instructions** in a formalized manner, which should be **suitable for communication, interpretation, or processing** by human or electronic machine.

Data can be described as **unprocessed facts and figures**.

It can also be defined as groups of **non-random symbols** in the form of **text, images, and voice** representing quantities, actions and objects

Information Hierarchy



Definition of Information

- ◎ **Information = organized data**
 - **Formatted, filtered, organized, structured, interpreted, summarized data**
 - **data + relations (context) = information**
 - **Relates to a description, definition or perspective (what, who, when , where)**

Definition of Knowledge

- ◎ **Knowledge = information that has been organized, internalized and integrated with experience, study, or intuition**
 - Case, rule, process, model, ideas
 - Rules and procedures that guide decisions and actions
 - Information + application = knowledge
 - Comprises of strategy, practice, method, or approach (how)

Reading assignment: Read about Data Processing Cycle

Data types and data representations

In computer science and computer programming, a **data type** or simply type is a **structure** of data which tells the compiler or interpreter **how the programmer intends to use the data.**

Common data types include:

Integers, Boolean, Characters, Floating-Point Numbers, Strings.

These data types **define the operations that can be done on the data, the meaning of each data attribute, and the way values of those types can be stored.**

Data Types Continued...

- ⦿ A **data type** defines a collection of data values and a set of predefined operations on those values.
- ⦿ Data types that are not defined in terms of other types are called **primitive data types**. Nearly all programming languages provide a set of primitive data types. Some of the primitive types are merely reflections of the hardware—for example, most integer types.

Data Types Continued...

- ⦿ Programming languages provide high-level **data types** such as truth values, integers, characters, records, and arrays, together with operations over these types. Target machines provide only machine 'types' such as bits, bytes, words, and double-words, together with **low-level arithmetic and logical operations**. To bridge the semantic gap between the source language and the target machine, the implementer must decide how to *represent* the source language's types and operations in terms of the target machine's types and operations.
- ⦿ The **primitive types** of a programming language are those whose values are primitive, i.e., cannot be decomposed into simpler values.

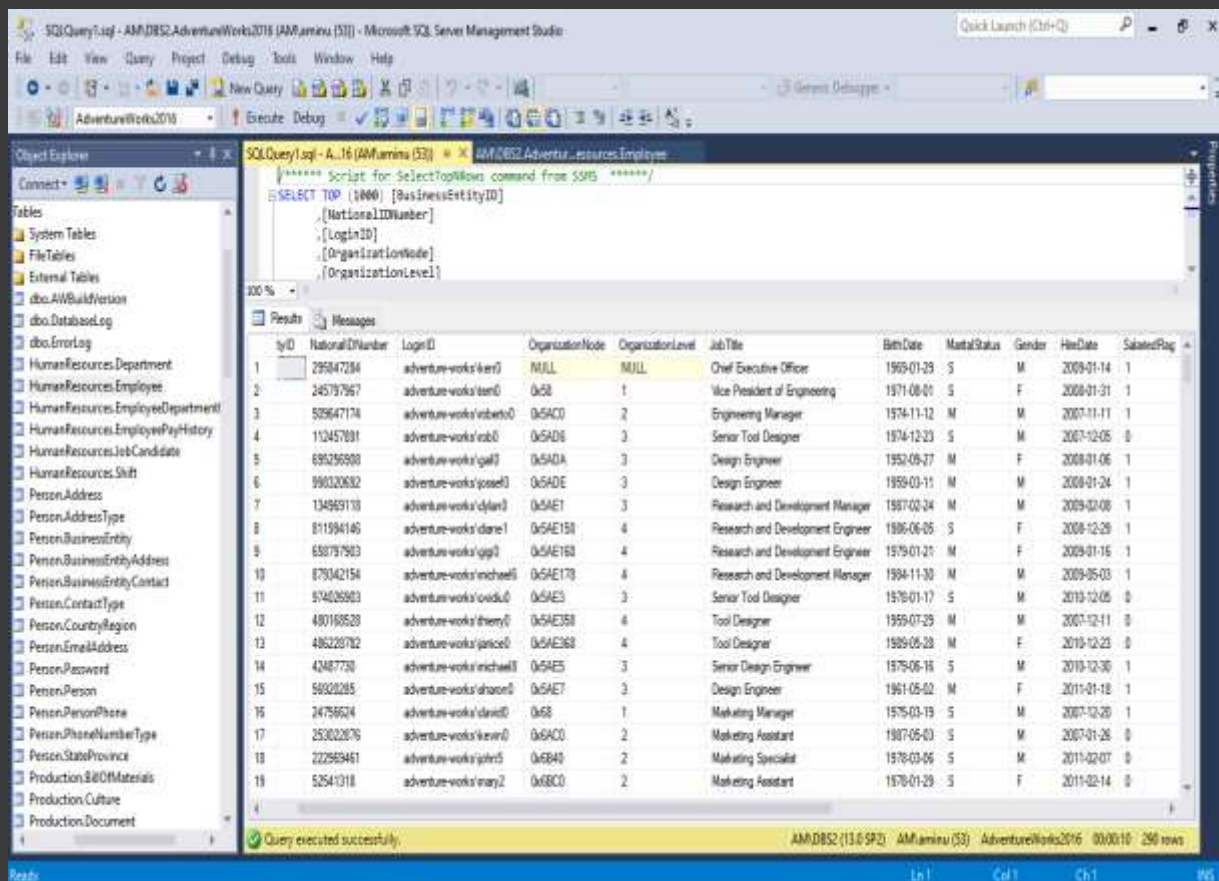
Types of Data from Data Analytics perspective

Structured, Unstructured, and Semi-structured data types



Structured Data -- examples

SQL Data



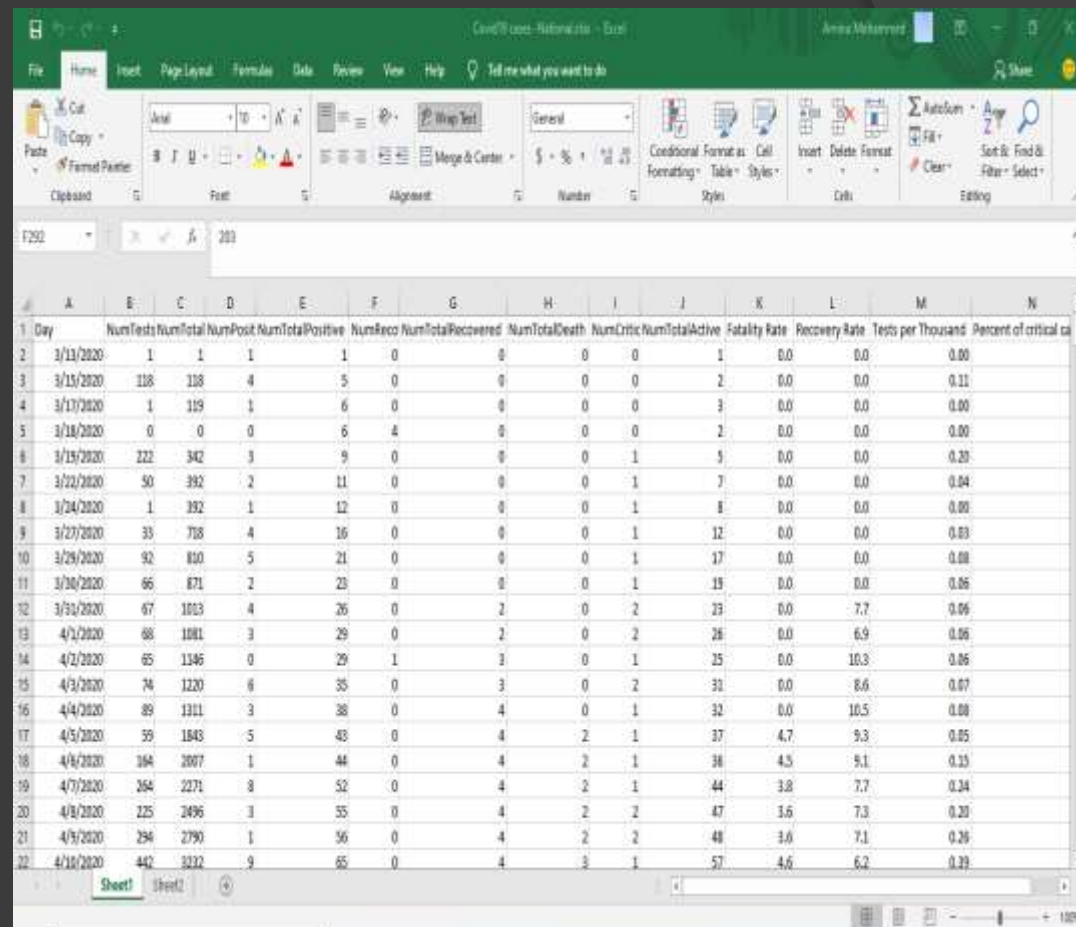
SQLQuery1.sql - AM/DBS2:AdventureWorks2016 (AM/aminu (53)) - Microsoft SQL Server Enterprise Manager

Object Explorer: AdventureWorks2016

Query: `SELECT TOP (1000) [BusinessEntityID], [NationalIDNumber], [LoginID], [OrganizationNode], [OrganizationLevel], [JobTitle], [BirthDate], [MaritalStatus], [Gender], [HireDate], [SalariedFlag]`

EmployeeID	NationalIDNumber	LoginID	OrganizationNode	OrganizationLevel	JobTitle	BirthDate	MaritalStatus	Gender	HireDate	SalariedFlag
1	295847234	adventure-works\ken0	NULL	NULL	Chief Executive Officer	1959-01-29	S	M	2009-01-14	1
2	245797967	adventure-works\ken0	0x58	1	Vice President of Engineering	1971-08-01	S	F	2008-01-31	1
3	529647174	adventure-works\vinet0	0x5AC0	2	Engineering Manager	1974-11-12	M	M	2007-11-11	1
4	112457891	adventure-works\vinet0	0x5AD6	3	Senior Tool Designer	1974-12-23	S	M	2007-12-05	0
5	695256908	adventure-works\gal0	0x5ADA	3	Design Engineer	1952-09-27	M	F	2008-01-06	1
6	990220832	adventure-works\josee0	0x5ADE	3	Design Engineer	1959-03-15	M	M	2008-01-24	1
7	134959113	adventure-works\dylan0	0x5AE1	3	Research and Development Manager	1987-02-24	M	M	2009-02-08	1
8	811994146	adventure-works\diane1	0x5AE158	4	Research and Development Engineer	1986-06-05	S	F	2008-12-29	1
9	658797963	adventure-works\gigi0	0x5AE160	4	Research and Development Engineer	1979-01-21	M	F	2009-01-16	1
10	879342154	adventure-works\michael0	0x5AE170	4	Research and Development Manager	1994-11-30	M	M	2009-05-03	1
11	974026983	adventure-works\ovidiu0	0x5AE3	3	Senior Tool Designer	1976-01-17	S	M	2010-12-05	0
12	480169528	adventure-works\thom0	0x5AE358	4	Tool Designer	1959-07-25	M	M	2007-12-11	0
13	486228702	adventure-works\janice0	0x5AE368	4	Tool Designer	1989-06-28	M	F	2010-12-23	0
14	42487730	adventure-works\michael0	0x5AE5	3	Senior Design Engineer	1975-06-16	S	M	2010-12-30	1
15	58328285	adventure-works\aharon0	0x5AE7	3	Design Engineer	1961-05-02	M	F	2011-01-18	1
16	24756624	adventure-works\clark0	0x58	1	Marketing Manager	1975-03-19	S	M	2007-12-26	1
17	253022676	adventure-works\kevin0	0x5AC0	2	Marketing Assistant	1987-05-03	S	M	2007-01-26	0
18	222959461	adventure-works\johnd0	0x5B40	2	Marketing Specialist	1978-03-06	S	M	2011-02-07	0
19	525413118	adventure-works\mary2	0x5BCC	2	Marketing Assistant	1976-01-29	S	F	2011-02-14	0

Excel File

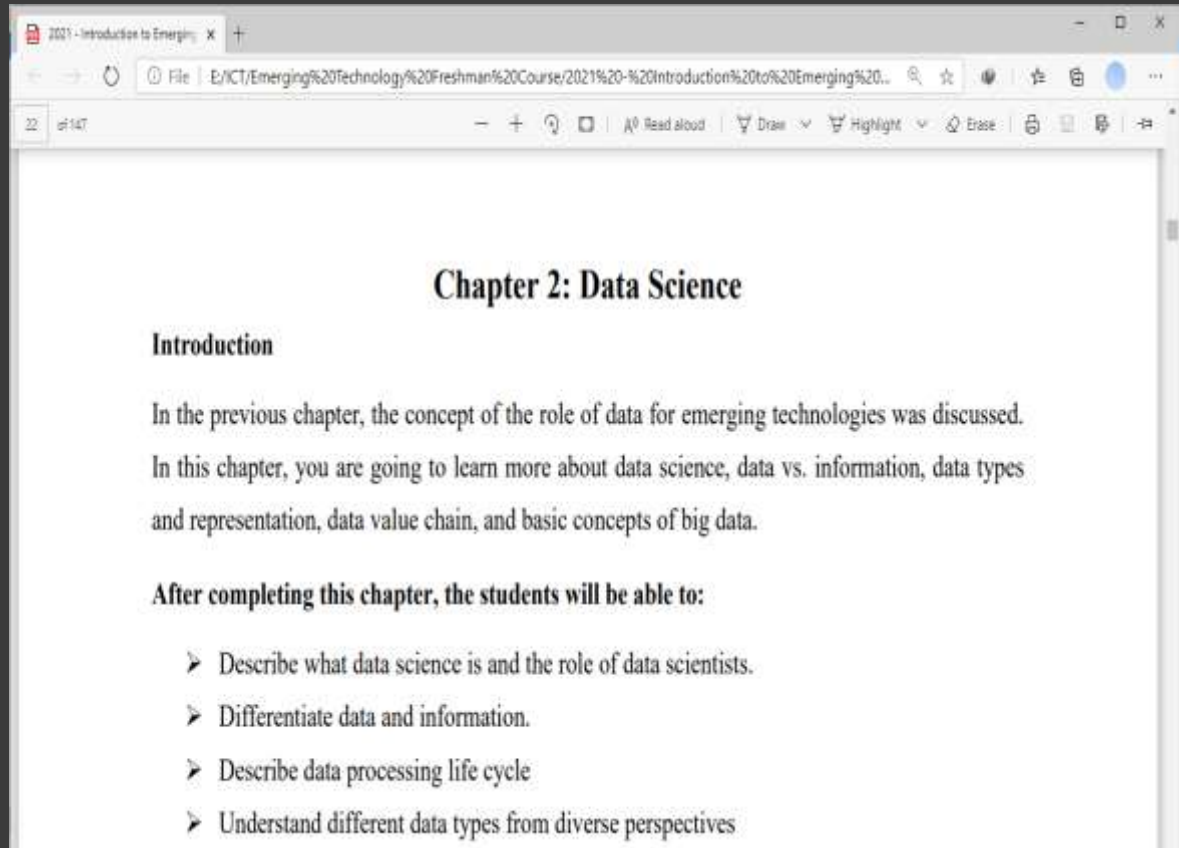


COVID-19 Stats - National.xlsx - Excel

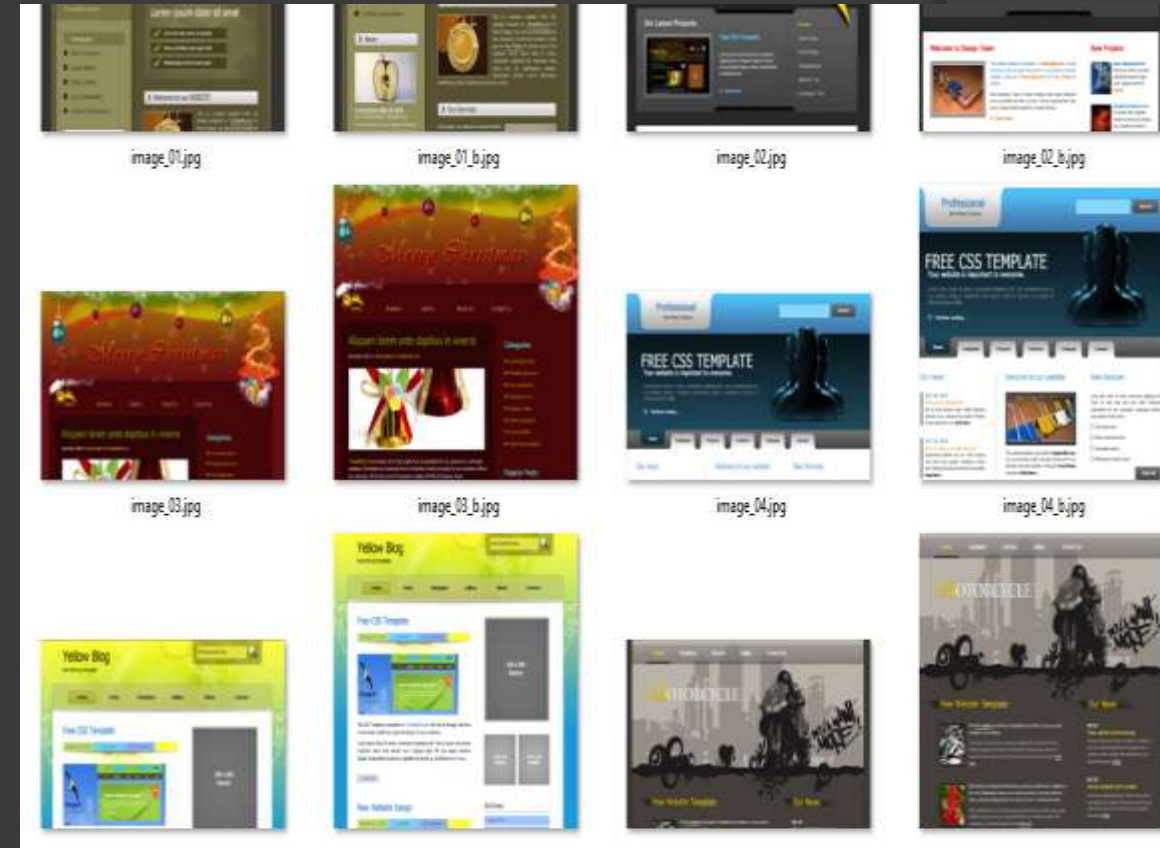
Day	NumTests	NumTotal	NumPost	NumTotalPositive	NumReco	NumTotalRecovered	NumTotalDeath	NumCritic	NumTotalActive	Fatality Rate	Recovery Rate	Tests per Thousand	Percent of critical ca
3/13/2020	1	1	1	1	0	0	0	0	1	0.0	0.0	0.00	
3/15/2020	118	118	4	5	0	0	0	0	2	0.0	0.0	0.11	
3/17/2020	1	119	1	6	0	0	0	0	3	0.0	0.0	0.00	
3/18/2020	0	0	0	6	4	0	0	0	2	0.0	0.0	0.00	
3/19/2020	222	342	3	9	0	0	0	1	5	0.0	0.0	0.20	
3/22/2020	50	392	2	11	0	0	0	1	7	0.0	0.0	0.04	
3/24/2020	1	392	1	12	0	0	0	1	8	0.0	0.0	0.00	
3/27/2020	33	718	4	16	0	0	0	1	12	0.0	0.0	0.03	
3/29/2020	92	810	5	21	0	0	0	1	17	0.0	0.0	0.08	
3/30/2020	66	871	2	23	0	0	0	1	19	0.0	0.0	0.06	
3/31/2020	67	1013	4	26	0	2	0	2	23	0.0	7.7	0.06	
4/1/2020	68	1081	3	29	0	2	0	2	26	0.0	6.9	0.06	
4/2/2020	65	1146	0	29	1	3	0	1	25	0.0	10.3	0.06	
4/3/2020	74	1220	6	35	0	3	0	2	31	0.0	8.6	0.07	
4/4/2020	89	1311	3	38	0	4	0	1	32	0.0	10.5	0.08	
4/5/2020	59	1843	5	43	0	4	2	1	37	4.7	9.3	0.05	
4/6/2020	184	2007	1	44	0	4	2	1	38	4.3	9.1	0.13	
4/7/2020	264	2271	8	52	0	4	2	1	44	3.8	7.7	0.14	
4/8/2020	225	2496	3	55	0	4	2	2	47	1.6	7.3	0.20	
4/9/2020	294	2790	1	56	0	4	2	2	48	3.6	7.1	0.26	
4/10/2020	442	3232	9	65	0	4	3	1	57	4.6	6.2	0.39	

Unstructured Data -- examples

□ Pdf files



□ Images



Semi-structured Data -- examples

Examples of semi-structured data
JSON and XML

```
{ "employees": [  
  { "firstName": "John", "lastName": "Doe" },  
  { "firstName": "Anna", "lastName": "Smith" },  
  { "firstName": "Peter", "lastName": "Jones" }  
]}
```

```
<employees>  
  <employee>  
    <firstName>John</firstName> <lastName>Doe</lastName>  
  </employee>  
  <employee>  
    <firstName>Anna</firstName> <lastName>Smith</lastName>  
  </employee>  
  <employee>  
    <firstName>Peter</firstName> <lastName>Jones</lastName>  
  </employee>  
</employees>
```

Metadata -- example

■ Metadata about an image

Property	Value
Image	
Image ID	
Dimensions	800 x 600
Width	800 pixels
Height	600 pixels
Horizontal resolution	96 dpi
Vertical resolution	96 dpi
Bit depth	24
Compression	
Resolution unit	
Color representation	
Compressed bits/pixel	
Camera	
Camera maker	
Camera model	
F-stop	
Exposure time	
ISO speed	

Data Science Continued

- ⦿ **Data Value Chain**
- ⦿ **Basic Concepts of Big Data**

Data Value Chain

Describes the information flow within a big data system as a **series of steps** needed to **generate value** and **useful insights** from data.

The **Big Data Value Chain** identifies the following key high-level activities:

Data Acquisition, Data Analysis, Data Curation, Data Storage, Data Usage

Data Acquisition

The infrastructure required for data acquisition must

- deliver **low, predictable latency** in both capturing data and in executing queries
- be able to **handle very high transaction volumes**, often in a distributed environment
- support **flexible** and **dynamic data structures**

Data Analysis

Involves **exploring, transforming**, and **modelling** data with the goal of **highlighting relevant data, synthesizing** and **extracting** useful hidden information with high potential from a business point of view.

Related areas include **data mining, business intelligence**, and **machine learning**.

What is synthesizing?

Data Curation

Data curation processes can be categorized into different activities

- **content creation, selection, classification, transformation, validation, and preservation**
- ensuring that data are **trustworthy, discoverable, accessible, reusable, and fit** for purpose

A key trend for the curation of big data **utilizes community and crowd sourcing** approaches (**attribute of IR4**).

Data Storage

Relational databases that guarantee database transactions, **lack flexibility** with regard to **schema changes**, **performance** and **fault tolerance** when data volumes and complexity grow, making them **unsuitable for big data scenarios**

NoSQL technologies have been designed with the **scalability** goal in mind and present a wide range of solutions based on **alternative data models**

Data Usage

Covers the **data-driven business activities** that need **access** to the curated data, its **in-use analytics**, and the tools needed to **integrate** the data analyses within the business activity.

In business decision-making , it can **enhance competitiveness** through **reduction of costs, increased added value**, or any other parameter that can be measured against existing performance criteria

Big Data Value Chain Main Activities

Data Acquisition

- Structured data
- Unstructured data
- Event Processing
- Sensor networks
- Protocols
- Real-time data
- Data streams
- Multimodality

Data Analysis

- Stream mining
- Semantic analysis
- Machine Learning
- Information Extraction
- Linked Data
- Data Discovery
- Whole World semantics
- Ecosystems
- Cross-sectional data analysis

Data Curation

- Data Quality
- Trust/Provenance
- Annotation
- Data Validation
- Human Data Interaction
- Top-Down/Bottom-up
- Human Computation
- Curation at scale
- Incentivization
- Automation
- Interoperability

Data Storage

- In memory DBS
- No SQL DBS
- Cloud storage
- Query Interfaces
- Scalability and Performance
- Data Models
- Consistency
- Availability
- Partition-tolerance
- Security and Privacy
- Standardization

Data Usage

- Decision support
- Prediction
- In-use analytics
- Simulation
- Exploration
- Visualization
- Modelling
- Control
- Domain-specific usage

Basic concepts of big data

Big data is a blanket term for the **non-traditional** strategies and technologies needed to gather, organize, process, and gain **insights** from large **datasets**. While the problem of working with data that exceeds the computing power or storage of a single computer is not new, the **pervasiveness**, **scale**, and **value** of this type of data have greatly expanded in recent years.

What Is Big Data?

A collection of datasets **so large and complex** that it becomes difficult to process using on-hand database management tools or traditional data processing applications and individual machines. *Wikipedia*

“Big data is **high volume, high velocity, and/or high variety** information assets that require **new forms of processing** to enable enhanced **decision making, insight discovery and process optimization**” *Laney 2001*

“When the **size of the data itself becomes part of the problem** and traditional techniques for working with data run out of steam”

The 4 Vs Characterizing Big Data.

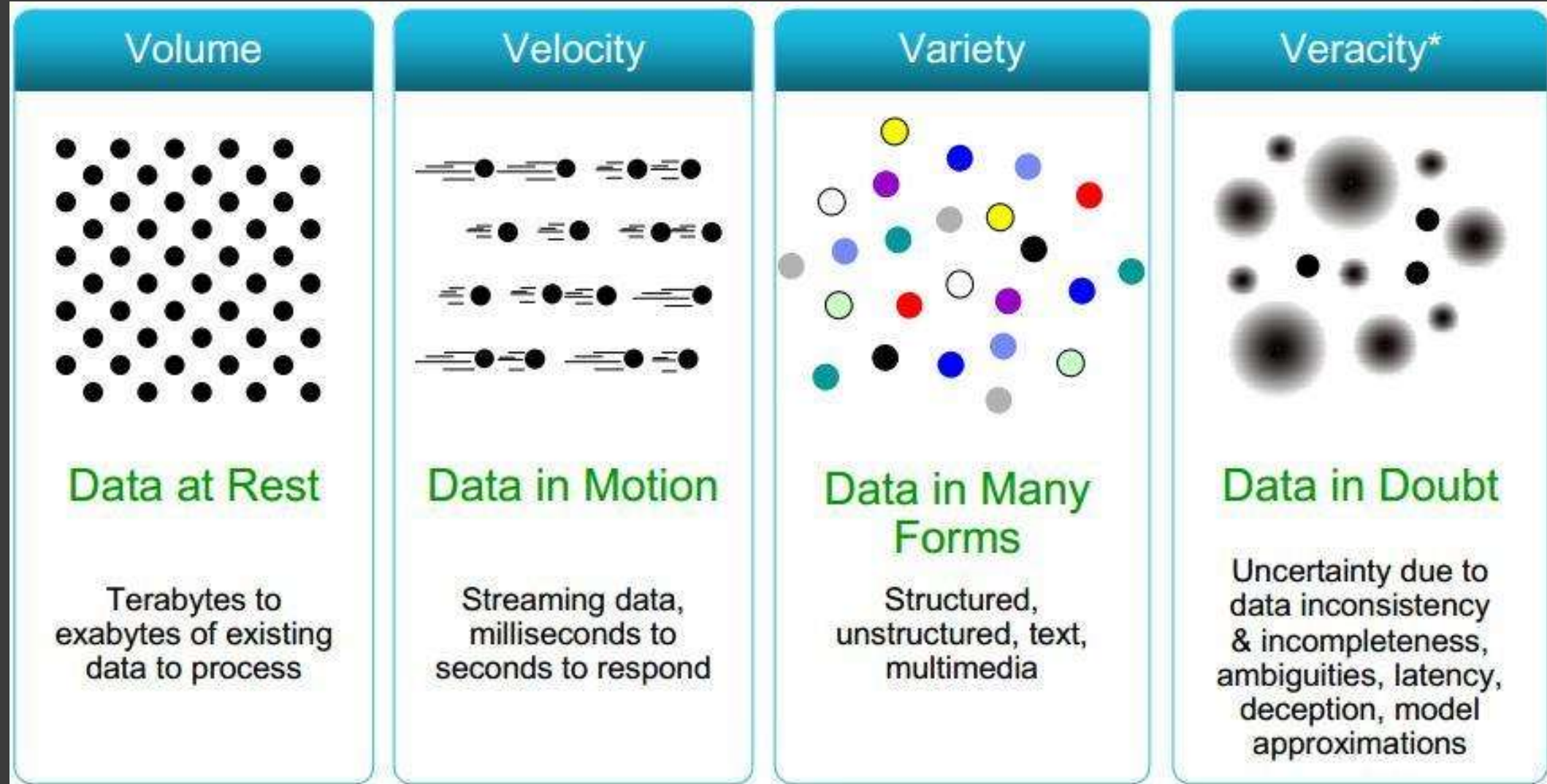
Volume: large amounts of data in yottabytes or zetabytes, (massive datasets)

Velocity: Data is live streaming or in motion

Variety: data comes in many different forms / from diverse sources and formats

Veracity: can we trust the data? How accurate is it? Doubt in data, etc.

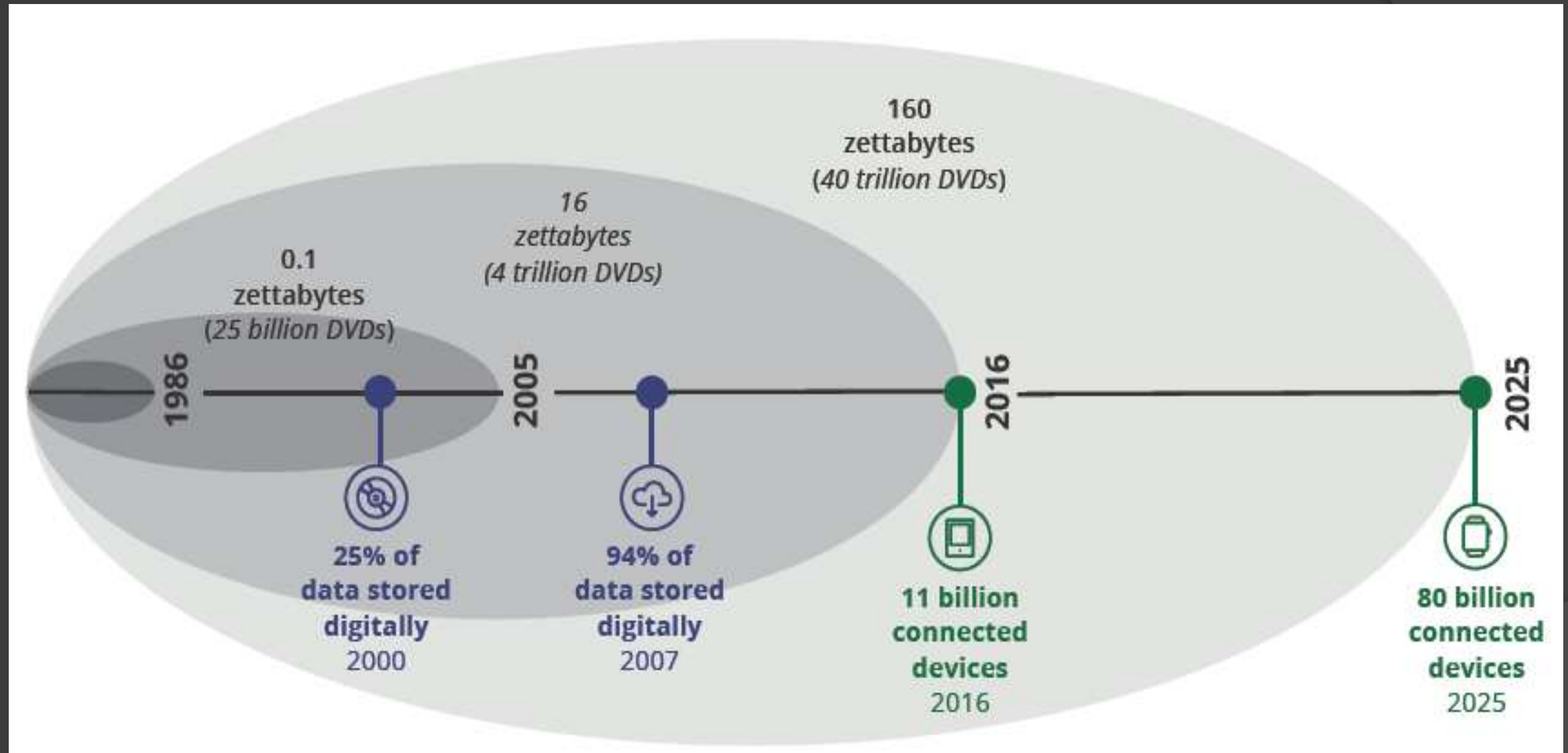
The 4 Vs





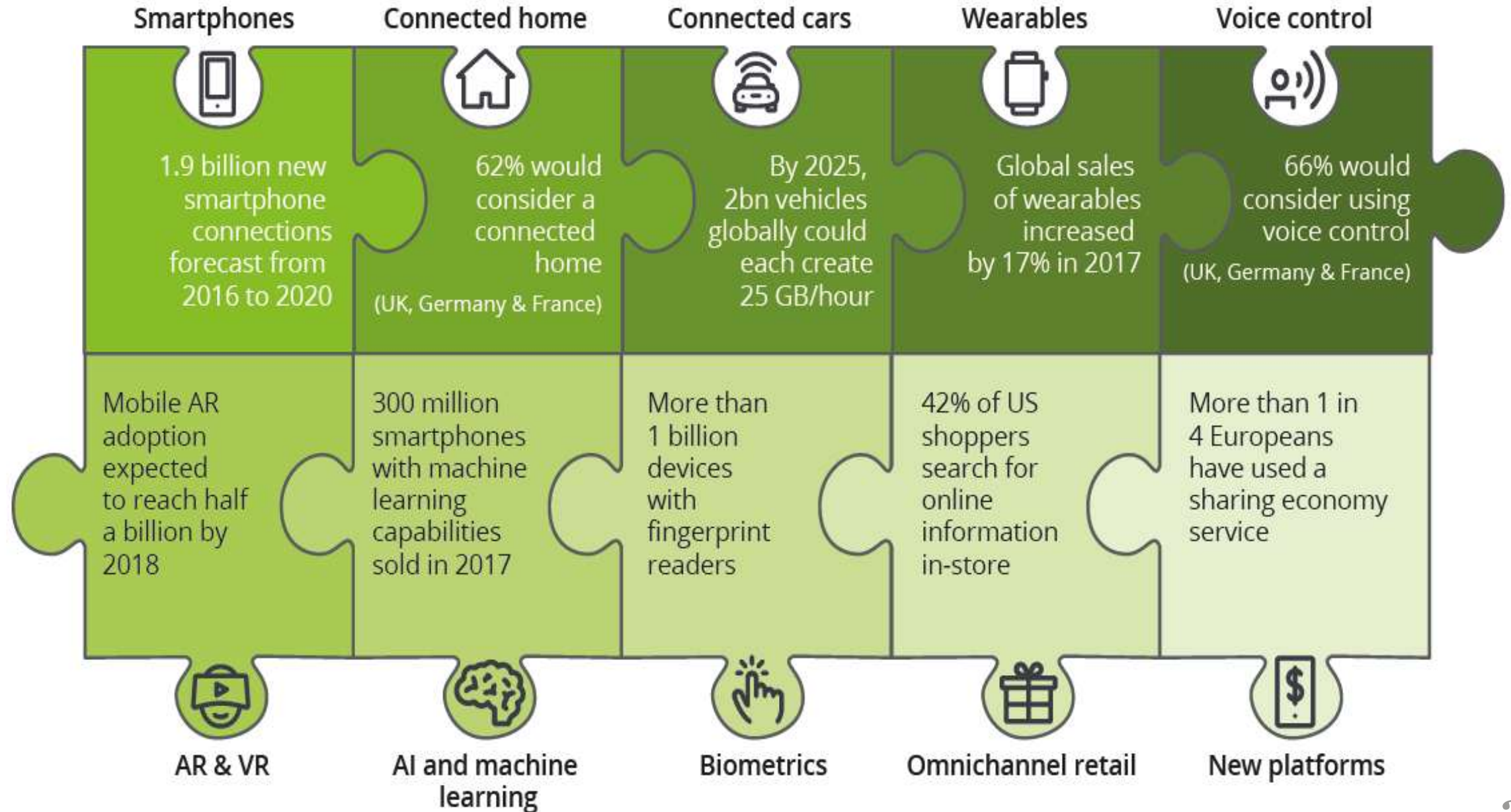
CERN's Large Hadron Collider (LHC) generates 15 PB a year

Illustration of global data growth trends over time



1 zettabyte = 1,000 exabytes = 1,000,000 petabytes = 1,000,000,000 terabytes = 1,000,000,000,000 gigabytes

Drivers of data growth (Big data drivers)



Variety (Complexity)

- ⦿ Relational Data (Tables/Transaction/Legacy Data) in Databases
- ⦿ Text Data (Web)
- ⦿ Semi-structured Data (XML)
- ⦿ Graph Data
 - Social Networks, Semantic Web (RDF), ...
- ⦿ Streaming Data
 - You can only scan the data once
- ⦿ A single application can be generating/collecting many types of data
- ⦿ Big Public Data (online, weather, finance, etc.)

To extract knowledge → all these types of data need to be linked together

Velocity (Speed)

- ⦿ Data is being generated fast and need to be processed fast
- ⦿ Online Data Analytics
- ⦿ Late decisions → missing opportunities
- ⦿ **Example**
 - **Healthcare monitoring:** sensors monitoring your activities and body → any abnormal measurements require immediate reaction



Real-time/Fast Data



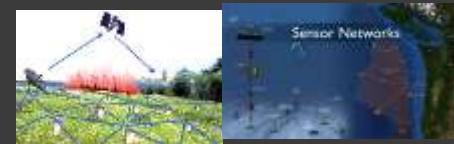
Social media and computer networks
(all of us are generating data)



Scientific instruments
(collecting all sorts of data)



Mobile devices
(tracking all objects all the time)



Sensor technology and networks
(measuring all kinds of data)

- The progress and innovation is no longer hindered by the ability to collect data
- But, by the ability to manage, analyze, summarize, visualize, and discover knowledge from the collected data in a timely manner and in a scalable fashion.

Assignment 1

- ④ Write an article explaining the difference between relational database systems (RDBMS) and HDFS. Include other Hadoop equivalent systems as much as possible in your discussion. Read about other types of file systems (reading assignment).
 - Individual , in 10 days (on or before November 10).
 - Follow Standard formatting rules.
 - Similar works are not allowed; shall be rejected.
 - Lamination not required; Maximum of 5 pages (you decide).

Clustered Computing and Hadoop Ecosystem

- ⦿ Because of the **qualities** of **big data**, individual computers are often inadequate for handling big data at most stages.
- ⦿ To address the **high storage** and **computational needs** of big data, computer **clusters** are a better fit.

Clustered Computing Continued...

- ◎ Big data clustering software **combines** the **resources** of many **smaller machines**, seeking to provide the following benefits:
 - **Resource pooling**: Storage, CPU, Working memory
 - **High availability**: Fault tolerance, Availability guaranties
 - **Easy scalability**: Easy to scale horizontally by adding more machines

Clustered Computing Continued...

- ④ Using clusters requires a **solution** for **managing cluster membership**, coordinating **resource sharing**, and **scheduling actual work** on individual nodes. Cluster membership and resource allocation can be handled by software like **Hadoop's YARN** (Yet Another Resource Negotiator).
- ④ The machines involved in the computing cluster are also typically involved with the management of a distributed storage system (**distributed computing**).

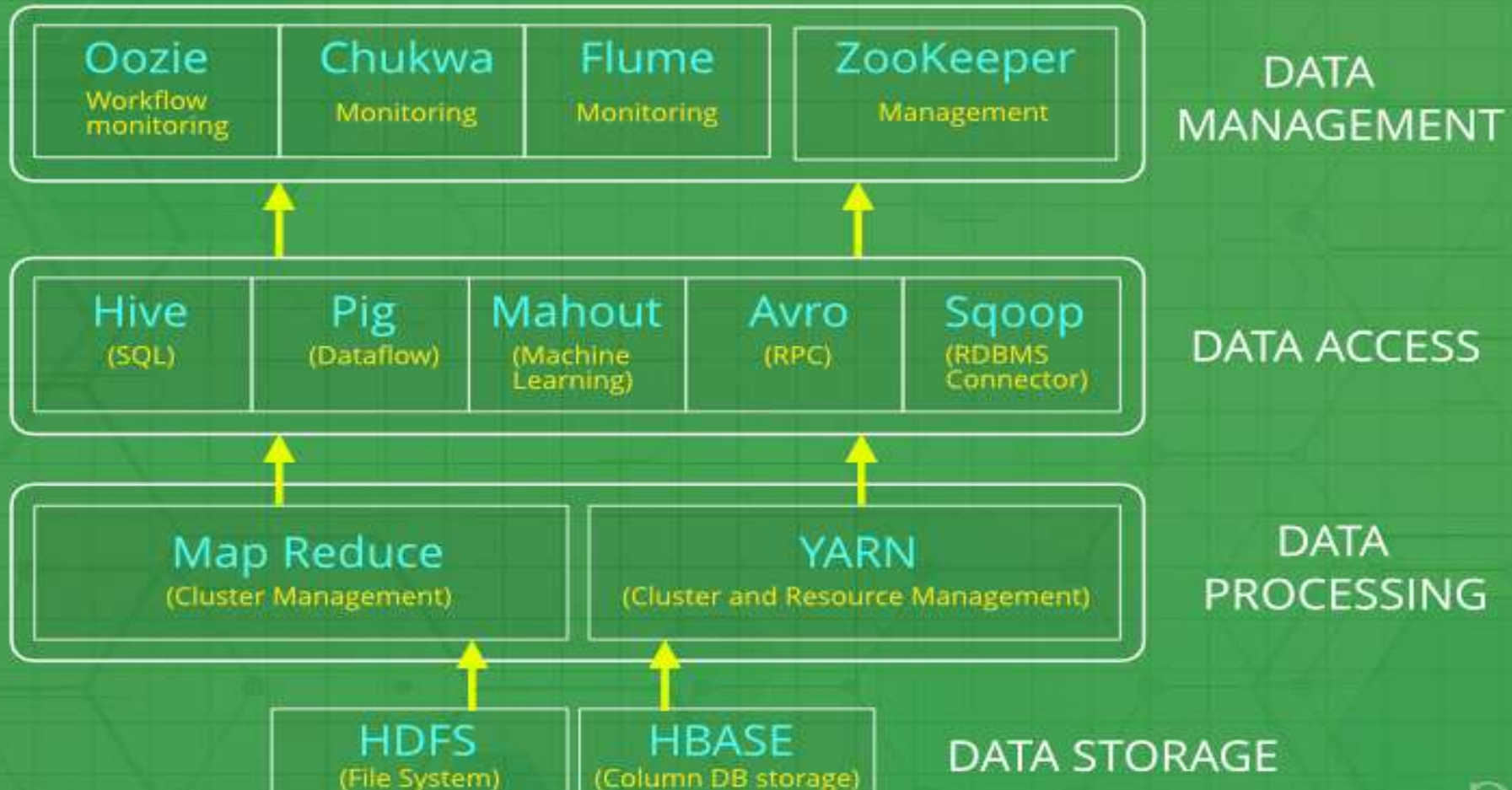
Hadoop and its Ecosystem

- Hadoop is an open-source **framework** intended **to make interaction with big data easier**. It is a framework that allows distributed processing of large datasets across **clusters** of computers using simple programming models. It is inspired by a technical document published by Google.

Hadoop and its Ecosystem Cont...

- ◎ Four key **characteristics** of Hadoop are:
 - **Economical:** highly economical as ordinary computers can be used for data processing.
 - **Reliable**
 - **Scalable**
 - **Flexible**

Hadoop Ecosystem



Hadoop Ecosystem

Big Data Life Cycle with Hadoop

- ④ **Ingesting data into the system:** data is transferred to **Hadoop** from various sources such as relational databases, systems, or local files. **Sqoop** transfers data from RDBMS to **HDFS**, whereas **Flume** transfers event data.
- ④ **Processing data for storage and (out from storage)**
- ④ **Computing and analyzing data**
- ④ **Visualizing the results**