# 02- Introduction to Big Data

**Dr Shafaq Khan**
**School of Computer Science**
**University of Windsor**

# Submission Deadline

**Project groups**: Please add project lead details by highlighting it in the share document. Also, add an extra line below group title, specifying the project title.


**Project proposal submission deadline**: Sec 2: Jan 30; Sec 3: Jan 31; Sec 4: Feb 1

University of Windsor

# Agenda

- Evolution of Data
- Data Classification
- Introduction to Big data
- Big Data Analytics (BDA)
- The Hadoop Ecosystem
- Apache Spark
- The Architecture of NoSQL Databases
- **Workshop on MongoDB**

3

# Introductory Questions

?? Why is the amount of data increasing tremendously?

?? How is data classified?

?? What is Big Data?

?? Why RDBMS is not suitable for Big Data?

?? Why Big Data Analysis is helpful?

?? How do we store and process big data?

University of Windsor

# Evolution of Data

Evolution of technology



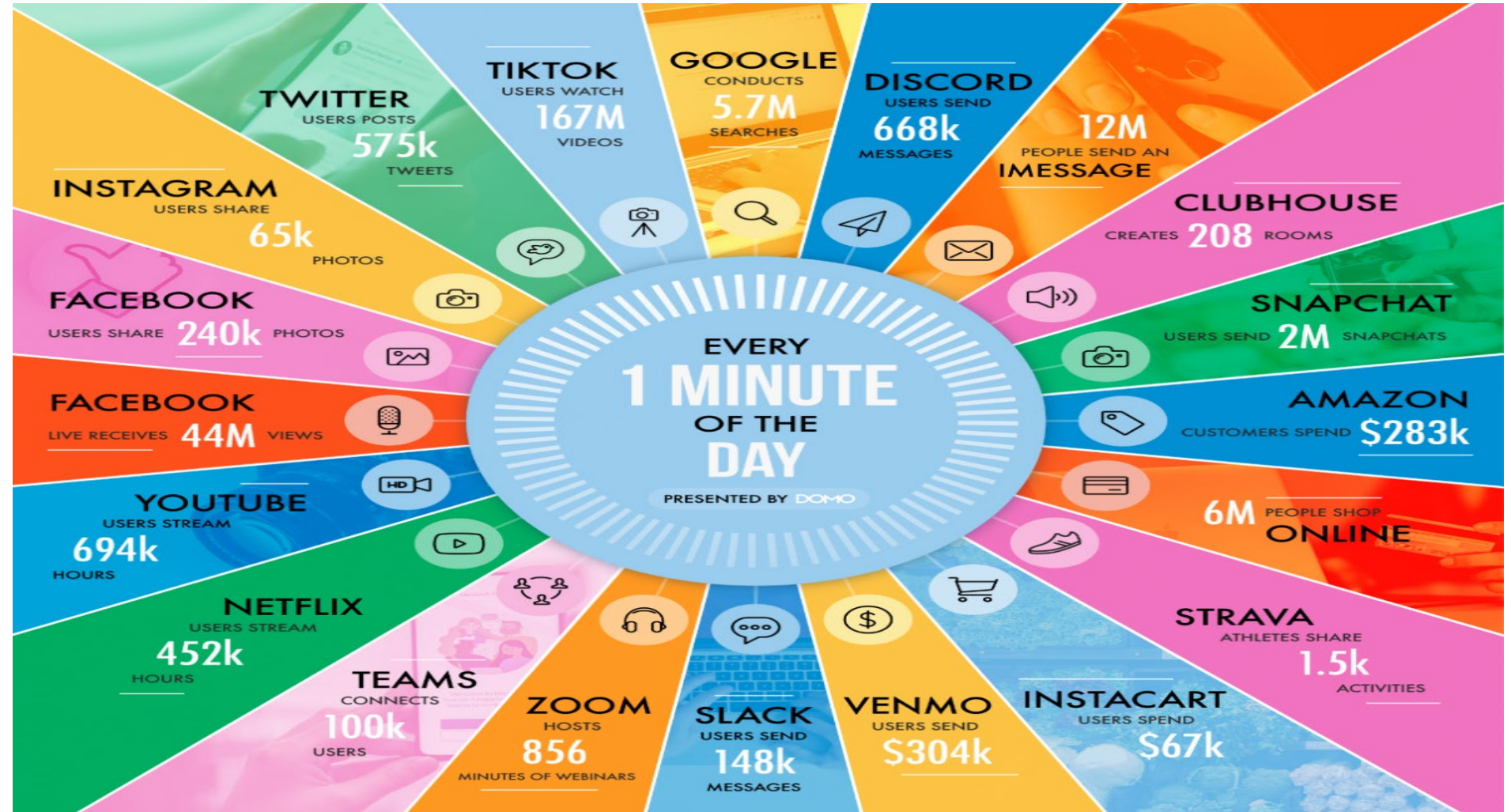As of November 2022, 59.5% of all website traffic comes from people using mobile devices.
Source: https://explodingtopics.com/blog/mobile-internet-traffic

University of Windsor

# Evolution of Data

Evolution of technology

Social Media

University of Windsor

# Evolution of Data

## IOT

Evolution of technology



Social Media





| Year | Connected devices in billions |
|------|------|
| 2019 | 8.6 |
| 2020 | 9.76 |
| 2021 | 11.28 |
| 2022* | 13.14 |
| 2023* | 15.14 |
| 2024* | 17.08 |
| 2025* | 19.08 |
| 2026* | 21.09 |
| 2027* | 23.14 |
| 2028* | 25.21 |
| 2029* | 27.31 |
| 2030* | 29.42 |

Connected devices in billions
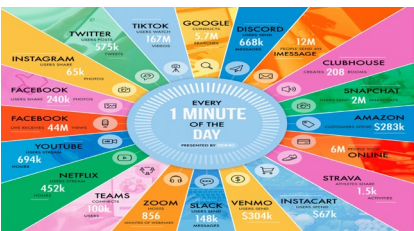
© Statista 2023



Number of Internet of Things (IoT) connected devices worldwide from 2019 to 2021, with forecasts from 2022 to 2030

University of Windsor

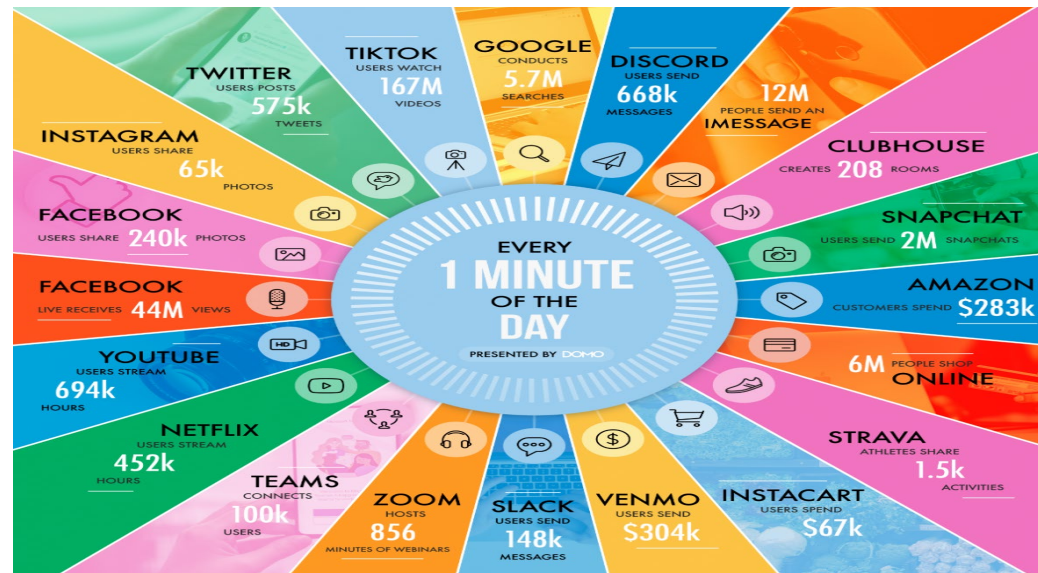# Evolution of Data

Evolution of technology

IOT



Social Media

HEALTHCARE
EDUCATION
RETAIL
MEDIA
GOVERNMENT
.
.
.

https://www.iottechtrends.com/history-of-iot/



EVERY 1 MINUTE OF THE DAY
PRESENTED BY DOMO

TWITTER USERS POSTS 575k TWEETS
TIKTOK USERS WATCH 167M VIDEOS
GOOGLE CONDUCTS 5.7M SEARCHES
DISCORD USERS SEND 668k MESSAGES
12M PEOPLE SEND AN IMESSAGE
INSTAGRAM USERS SHARE 65k PHOTOS
FACEBOOK USERS SHARE 240k PHOTOS
CLUBHOUSE CREATES 208 ROOMS
SNAPCHAT USERS SEND 2M SNAPCHATS
FACEBOOK LIVE RECEIVES 44M VIEWS
AMAZON CUSTOMERS SPEND $283k
YOUTUBE USERS STREAM 694k HOURS
6M PEOPLE SHOP ONLINE
NETFLIX USERS STREAM 452k HOURS
STRAVA ATHLETES SHARE 1.5k ACTIVITIES
TEAMS CONNECTS 100k USERS
ZOOM HOSTS 856 MINUTES OF WEBINARS
SLACK USERS SEND 148k MESSAGES
VENMO USERS SEND $304k
INSTACART USERS SPEND $67k

https://dailyinfographic.com/how-much-data-is-generated-every-minute
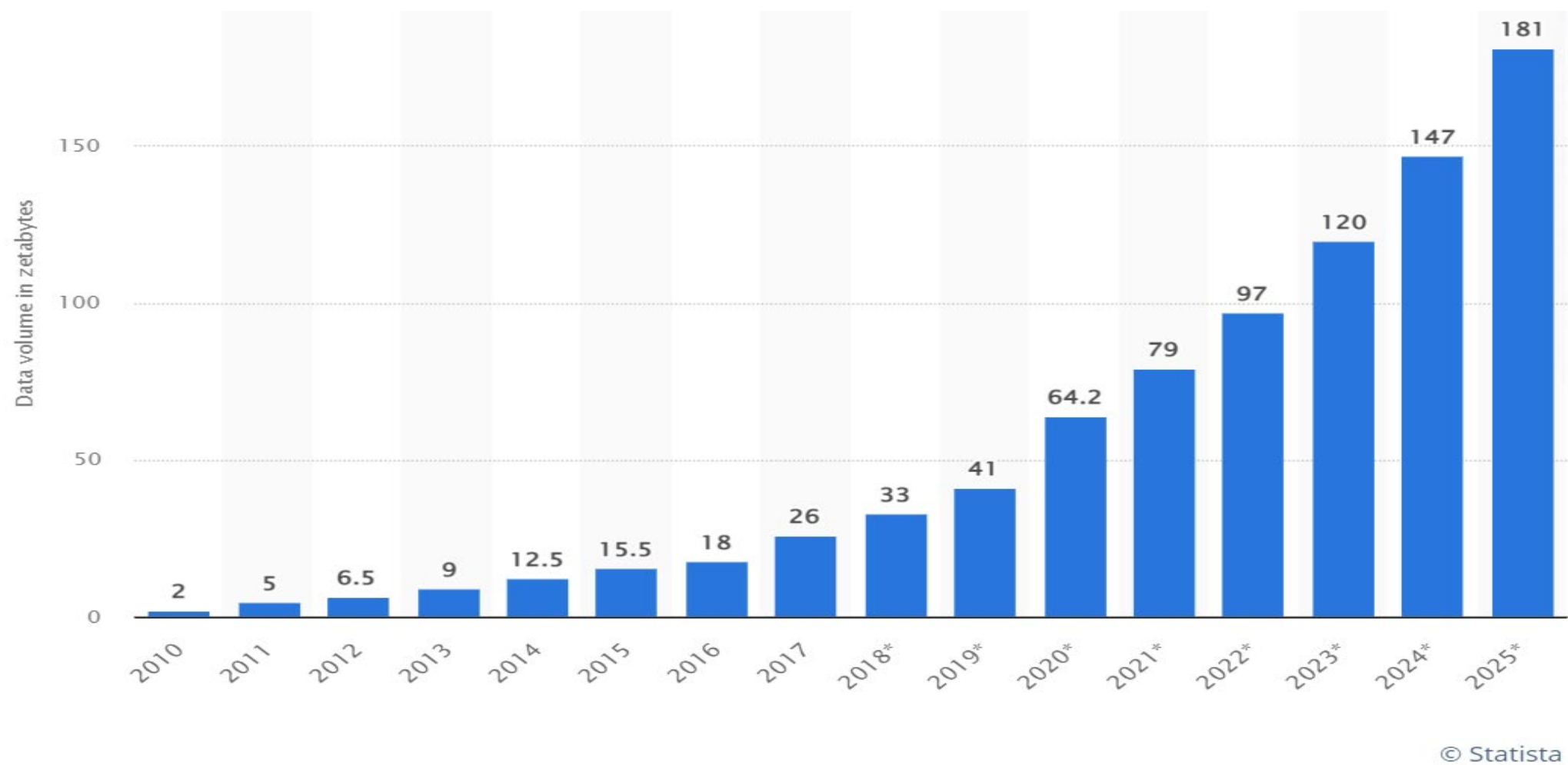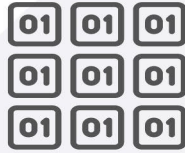
University of Windsor

# Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2025

https://www.statista.com/statistics/871513/worldwide-data-created/

University of Windsor

# Classification of Data

## Structured data

**Characteristics**

Predefined data models
Easy to search
Text-based
Shows what's happening

**Resides in**

Relational databases
Data warehouses

**Stored in**

Rows and columns

**Examples**

Dates, phone numbers, social security numbers, customer names, transaction info

## Unstructured data

**Characteristics**

No predefined data models
Difficult to search
Text, pdf, images, video
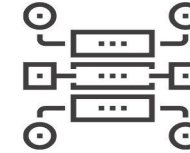Shows the why

**Resides in**

Applications
Data warehouses and lakes

**Stored in**

Various forms

**Examples**

Documents, emails and messages, conversation transcripts, image files, open-ended survey answers

## Semi-structured data

**Characteristics**

Loosely organized
Meta-level structure that can contain unstructured data
HTML, XML, JSON

**Resides in**

Relational databases
Tagged-text format

**Stored in**

Abstracts & figures

**Examples**

Server logs, tweets organized by hashtags, emails sorting by folders (inbox; sent; draft)

LEVITY

University of Windsor

# Historical Interpretation of Big Data
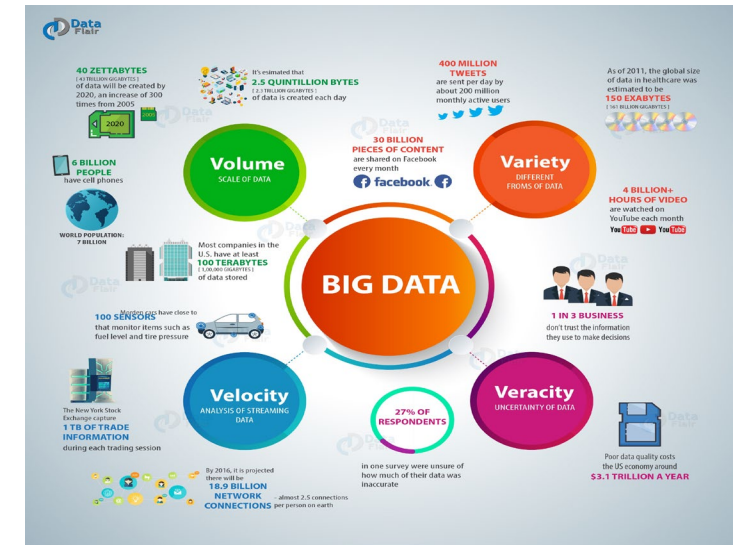
**Douglas Laney's 3Vs:**

He noticed that due to surging of e-commerce activities, data has grown along three dimensions, namely:

1. **Volume**: means the incoming data stream and cumulative volume of data
2. **Velocity**: represents the pace of data used to support interaction and generated by interactions
3. **Variety**: signifies the variety of incompatible and inconsistent data formats and data structures



https://dzone.com/articles/why-is-big-data-in-buzz

**IBM —4Vs:**

IBM added another attribute or "V" for "Veracity" on the top of Douglas Laney's 3Vs notation:

1. Volume stands for the scale of data
2. Velocity denotes the analysis of streaming data
3. Variety indicates different forms of data
4. **Veracity** implies the uncertainty of data. Accuracy and trustworthiness of data is termed as veracity

University of Windsor

# Historical Interpretation Of Big Data
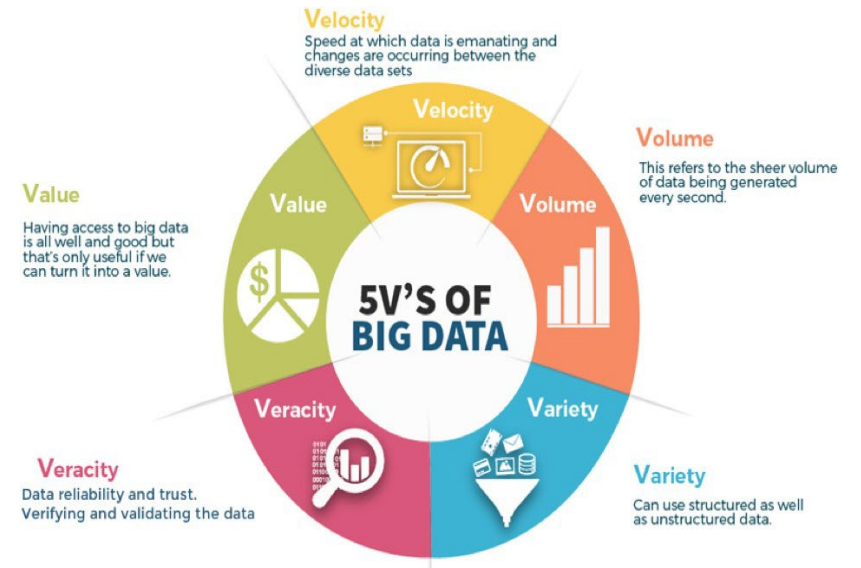
**Yuri Demchenko's 5Vs:**

Yuri added the value dimension along with the IBM 4Vs' definition in 2013:

1. Volume stands for the scale of data
2. Velocity denotes the analysis of streaming data
3. Variety indicates different forms of data
4. Veracity implies the uncertainty of data
5. **Value** refers to how useful the data is in decision making. we must be able to analyze the data to generate valuable knowledge that helps in decision making.

**Microsoft —6Vs:**

For the sake of maximizing the business value, Microsoft extended Douglas Laney's 3Vs attributes to 6 Vs:

1. Volume stands for scale of data
2. Velocity denotes the analysis of streaming data
3. Variety indicates different forms of data
4. Veracity focuses on trustworthiness of data sources
5. **Variability** refers to the complexity of data set. In comparison with "Variety" (or different data format), it means the number of variables in data sets.
6. **Visibility** emphasizes that you need to have a full picture of data in order to make informative decision



https://www.techentice.com/the-data-veracity-big-data/

University of Windsor

# Big Data Analytics (BDA)

"Big data analytics is the often-complex process of examining big data to uncover information -- such as hidden patterns, correlations, market trends and customer preferences -- that can help organizations make informed business decisions"

- **Explanatory analytics** focuses on discovering and explaining data characteristics based on existing data

- **Predictive analytics** focuses on predicting future data outcomes with a high degree of accuracy

- **Big Data Analytics-Benefits**

  - Effective marketing,

  - New revenue opportunities,

  - Customer personalization

  - Improved operational efficiency

  - Enhanced competitive edge

Over 97 percent of organizations are investing in big data and artificial intelligence.
Source: https://www.simplilearn.com/big-data-analytics-and-ai-to-manage-pandemics-article

On 16th Jan 2022,  858 Data Scientist jobs available.
Source: https://ca.indeed.com/Data-Science-Analyst-jobs?vjk=5d891b57d22a8a7b

University of Windsor

# Examples Of Big Data Analytics in Industries







**Analytics in retail** enables companies to create customer recommendations based on their purchase history, resulting in personalized shopping experiences and improved customer service. It also helps with forecasting trends and making strategic decisions based on market analysis.

**Healthcare big data analytics** drive quicker responses to emerging diseases and improve direct patient care, the customer experience, and administrative, insurance and payment processing.
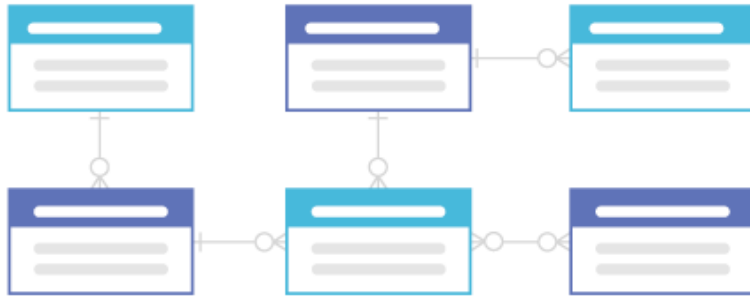
**Financial analytics** improve customer targeting using customer analytics. Businesses can make better informed underwriting decisions and provide better claims management while mitigating risk and fraud.

University of Windsor
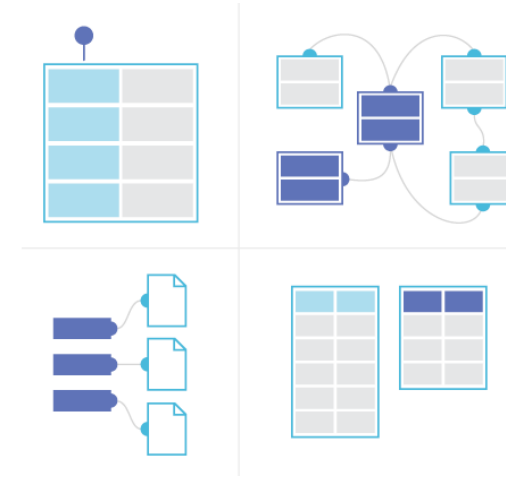
# Questions Reflect the Bottom Line Of BI

1. How to store massive data (such as in PB or EB scale currently) or information in the available resources

2. How to access these massive data or information quickly

3. How to work with datasets in variety formats: structured, semi-structured, and unstructured

4. How to process these datasets in a full scalable, fault tolerant, and flexible manner

5. How to extract BI interactively and cost-effectively

University of Windsor

# Relational Data to Big Data

RDBMS

**Drawback** with RDBMS for Large dataset
- Processing large data may fail
- Can't store unstructured data
- High cost.

Non-Relational Database

University of Windsor

# Hadoop Ecosystem

- De facto standard for most Big Data stora

- Java-based framework for distributing and processing very large data sets across clusters of computers

- Most important components:

  - **Hadoop Distributed File System (HDFS)**: Low-level distributed file processing system that can be used directly for data storage

  - **MapReduce:** Programming model that supports processing large data sets
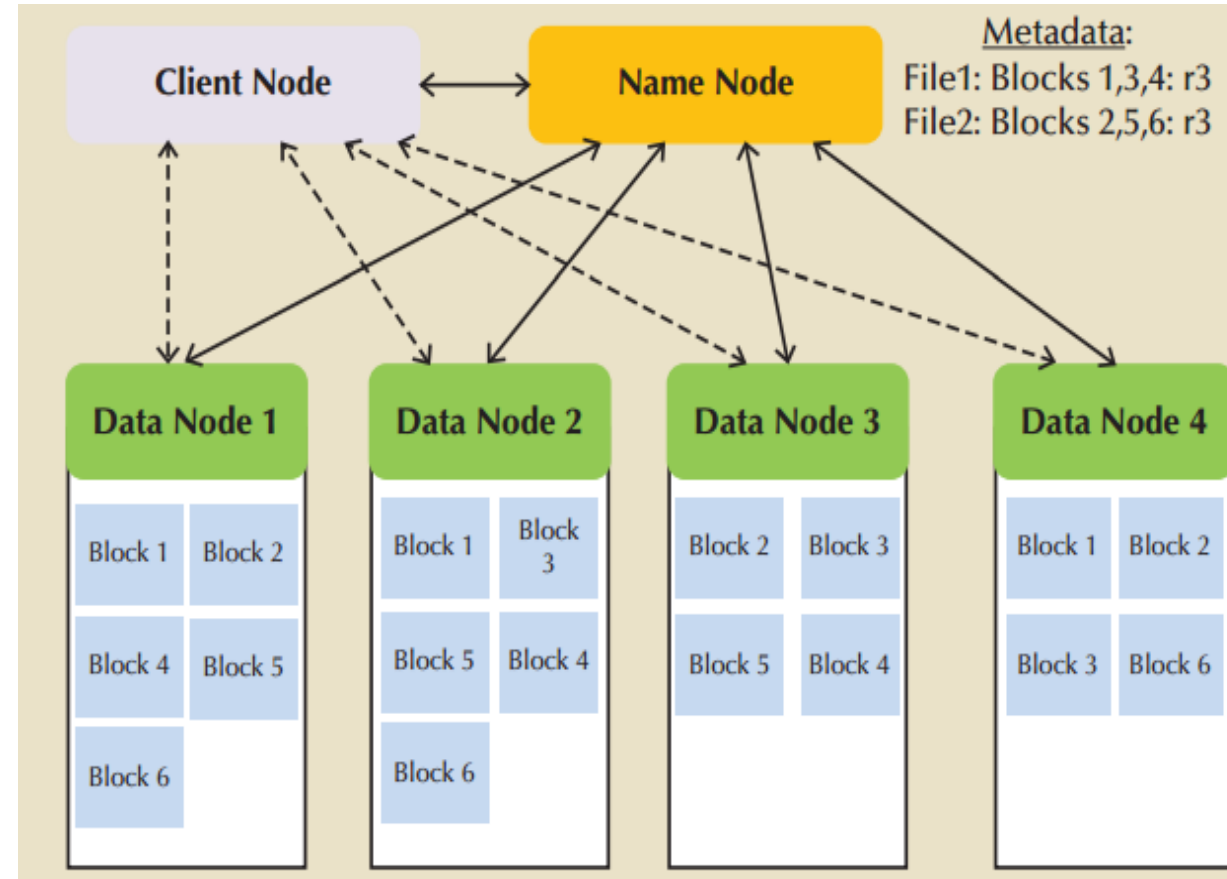
# Hadoop Distributed File System (HDFS)

- Approach based on several key assumptions:
  - *High volume* - Default block sizes is 64 MB and can be configured to even larger values
  - *Write-once, read-many* - Hadoop we can store all kinds of data once which can be accessed any number of times. Model simplifies concurrent issues and improves data throughput
  - *Streaming access* - Hadoop is optimized for batch processing of entire files as a continuous stream of data
  - *Fault tolerance* – HDFS is designed to replicate data across many different devices so that when one fails, data is still available from another device
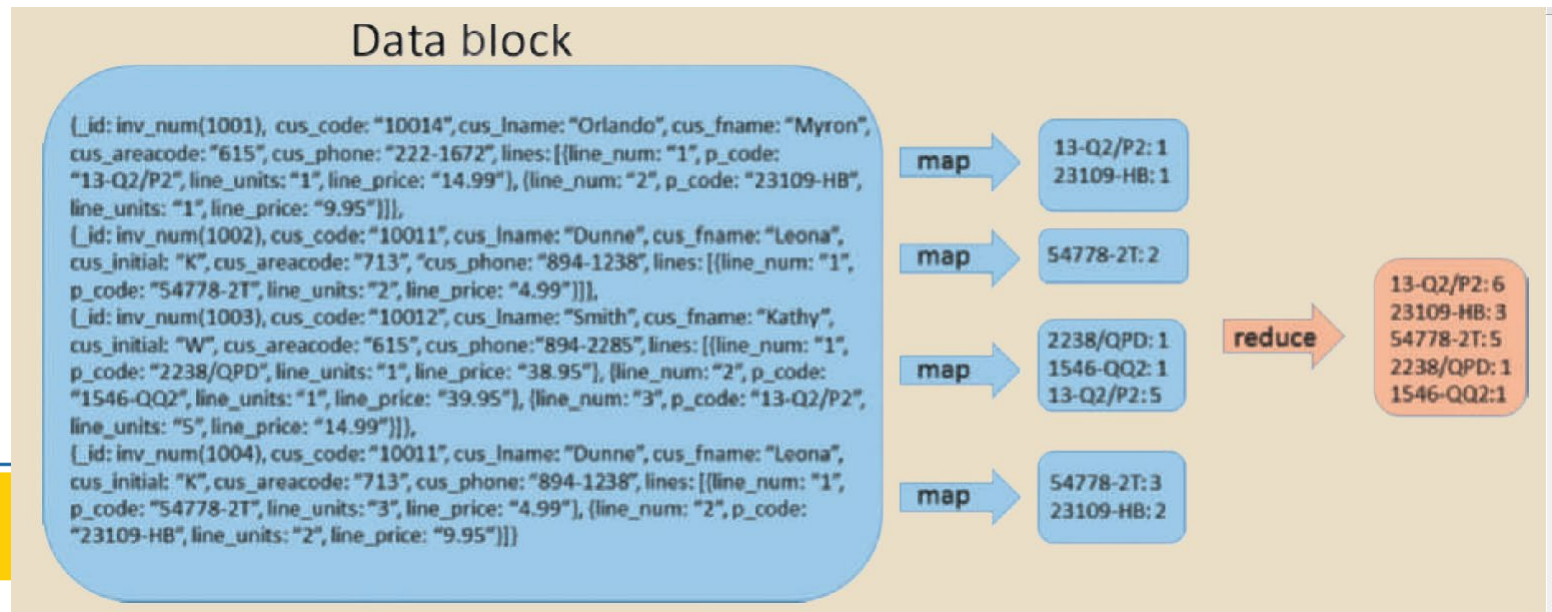
University of Windsor

# Hadoop Distributed File System (HDFS)

- Uses several types of nodes (computers):
  - Data node store the actual file data
  - Name node contains file system metadata
  - Client node makes requests to the file system as needed to support user applications

# MapReduce

- Framework used to process large data sets across clusters
  - Breaks down complex tasks into smaller subtasks, performing the subtasks and producing a final result
  - **Map** function takes a collection of data and sorts and filters it into a set of key-value pairs
    - **Mapper** program performs the map function
  - **Reduce** summaries results of map function to produce a single result
    - **Reducer** program performs the reduce function

# MapReduce

# A Sample of the Hadoop Ecosystem

**Map Reduce Simplification Applications:**

*Hive* is a data warehousing system that sites on top of HDFS and supports its own SQL-like language

*Pig* compiles a high-level scripting language (Pig Latin) into MapReduce jobs for executing in Hadoop



**Data Ingestion Applications:**

*Flume* is a component for ingesting data in Hadoop

*Sqoop* is a tool for converting data back and forth between a relational database and the HDFS

**Direct Query Applications:**

*HBase* is a column-oriented NoSQL database designed to sit on top of the HDFS that quickly processes sparse datasets

*Impala* was the first SQL-on-Hadoop application

University of Windsor

# SPARK

Spark is a unified analytics engine for large-scale data processing. Spark is a fast- and general-purpose computation platform based on large clusters.

It was developed by the UC Berkeley RAD Lab (now called as AMP Lab).

**Speed**: Uses In-Memory Processing- 100x faster than disk access

**Ease of Use:** Applications can be written in Java, Scala,Python, R etc.

**Generality**: Can combine SQL, Streaming and Complex Analytics

**Access Diverse Data Stores**: Hadoop, Apache Mesos(Cluster Management) , Kubernetes (Application Deployment) etc.

Open source:  Apache Spark

# NoSQL

- Name given to non-relational database technologies developed to address Big data challenges.
- There are literally hundreds of products that can be considered as being under the broadly defined term NoSQL.
- Most of these fit roughly into one of four categories:

| NoSQL DATABASES | | |
| --- | --- | --- |
| **NoSQL CATEGORY** | **EXAMPLE DATABASES** | **DEVELOPER** |
| Key-value database | Dynamo<br>Riak<br>Redis<br>Voldemort | Amazon<br>Basho<br>Redis Labs<br>LinkedIn |
| Document databases | MongoDB<br>CouchDB<br>OrientDB<br>RavenDB | MongoDB, Inc.<br>Apache<br>OrientDB Ltd.<br>Hibernating Rhinos |
| Column-oriented databases | HBase<br>Cassandra<br>Hypertable | Apache<br>Apache (originally Facebook)<br>Hypertable, Inc. |
| Graph databases | Neo4J<br>ArangoDB<br>GraphBase | Neo4j<br>ArangoDB, LLC<br>FactNexus |

University of Windsor

# NoSQL

**Key-value (KV) databases** store data as a collection of key-value pairs organized as **buckets** which are the equivalent of tables

| Bucket = Customer | |
| --- | --- |
| **Key** | **Value** |
| 10010 | "LName Ramas FName Alfred Initial A Areacode 615 Phone 844-2573 Balance 0" |
| 10011 | "LName Dunne FName Leona Initial K Areacode 713 Phone 894-1238 Balance 0" |
| 10014 | "LName Orlando FName Myron Areacode 615 Phone 222-1672 Balance 0" |

Ex: Redis, Oracle NoSQL

University of Windsor

# NoSQL

**Document databases** store data in key-value pairs in which the value components are tag-encoded documents grouped into logical groups called **collections**

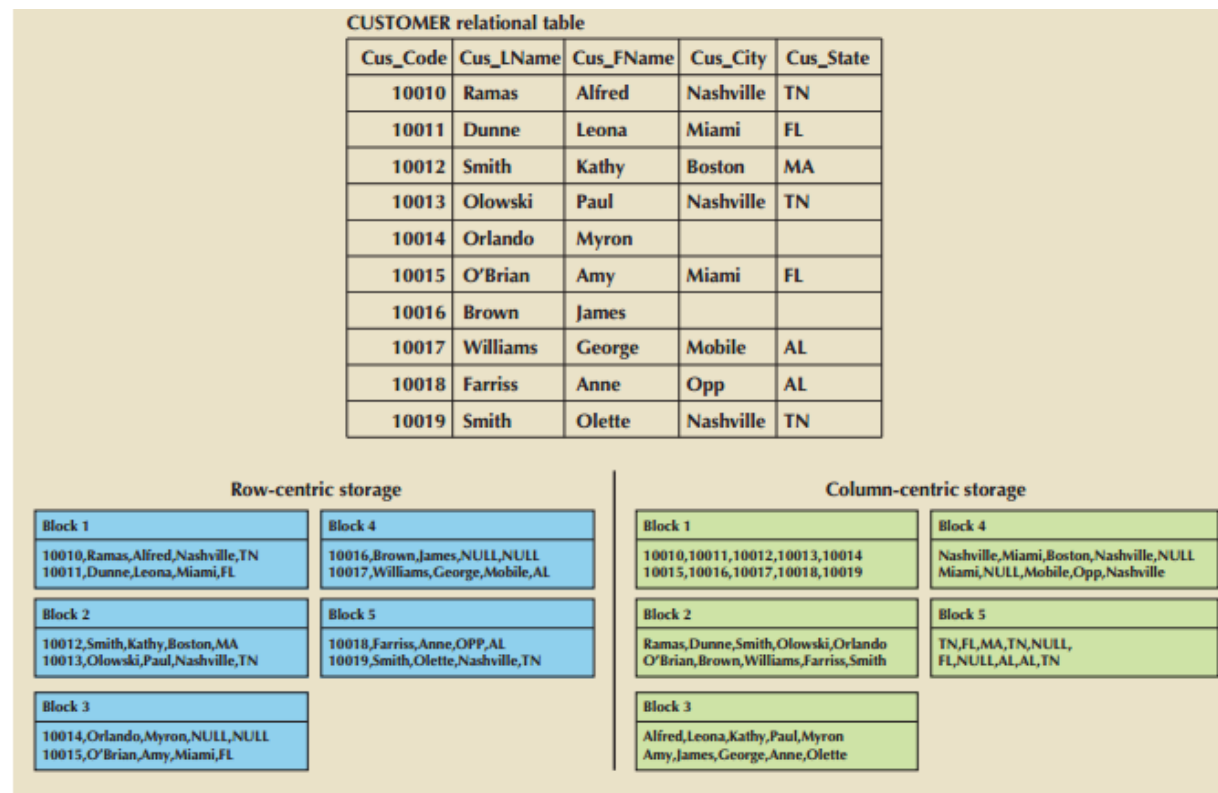| Collection = Customer | |
|---|---|
| **Key** | **Document** |
| 10010 | {LName: "Ramas", FName: "Alfred", Initial: "A", Areacode: "615", Phone: "844-2573", Balance: "0"} |
| 10011 | {LName: "Dunne", FName: "Leona", Initial: "K", Areacode: "713", Phone: "894-1238", Balance: "0"} |
| 10014 | {LName: "Orlando", FName: "Myron", Areacode: "615", Phone: "222-1672", Balance: "0"} |

Ex: MongoDB

University of Windsor

# NoSQL

**Column-oriented databases** refers to two technologies:

- **Column-centric storage**: Data stored in blocks which hold data from a single column across many rows
- **Row-centric storage:** Data stored in block which hold data from all columns of a given set of rows

**CUSTOMER relational table**

| Cus_Code | Cus_LName | Cus_FName | Cus_City | Cus_State |
|----------|-----------|-----------|----------|-----------|
| 10010 | Ramas | Alfred | Nashville | TN |
| 10011 | Dunne | Leona | Miami | FL |
| 10012 | Smith | Kathy | Boston | MA |
| 10013 | Olowski | Paul | Nashville | TN |
| 10014 | Orlando | Myron | | |
| 10015 | O'Brian | Amy | Miami | FL |
| 10016 | Brown | James | | |
| 10017 | Williams | George | Mobile | AL |
| 10018 | Farriss | Anne | Opp | AL |
| 10019 | Smith | Olette | Nashville | TN |

Ex: Google's BigTable, HBase, and Cassandra.

**Row-centric storage**

| Block 1 | Block 4 |
|---------|---------|
| 10010,Ramas,Alfred,Nashville,TN<br>10011,Dunne,Leona,Miami,FL | 10016,Brown,James,NULL,NULL<br>10017,Williams,George,Mobile,AL |

| Block 2 | Block 5 |
|---------|---------|
| 10012,Smith,Kathy,Boston,MA<br>10013,Olowski,Paul,Nashville,TN | 10018,Farriss,Anne,OPP,AL<br>10019,Smith,Olette,Nashville,TN |

| Block 3 |
|---------|
| 10014,Orlando,Myron,NULL,NULL<br>10015,O'Brian,Amy,Miami,FL |

**Column-centric storage**

| Block 1 | Block 4 |
|---------|---------|
| 10010,10011,10012,10013,10014<br>10015,10016,10017,10018,10019 | Nashville,Miami,Boston,Nashville,NULL<br>Miami,NULL,Mobile,Opp,Nashville |

| Block 2 | Block 5 |
|---------|---------|
| Ramas,Dunne,Smith,Olowski,Orlando<br>O'Brian,Brown,Williams,Farriss,Smith | TN,FL,MA,TN,NULL,<br>FL,NULL,AL,AL,TN |

| Block 3 |
|---------|
| Alfred,Leona,Kathy,Paul,Myron<br>Amy,James,George,Anne,Olette |

University of Windsor

28

## CUSTOMER relational table

| Cus_Code | Cus_LName | Cus_FName | Cus_City | Cus_State |
|----------|-----------|-----------|----------|-----------|
| 10010 | Ramas | Alfred | Nashville | TN |
| 10011 | Dunne | Leona | Miami | FL |
| 10012 | Smith | Kathy | Boston | MA |
| 10013 | Olowski | Paul | Nashville | TN |
| 10014 | Orlando | Myron | | |
| 10015 | O'Brian | Amy | Miami | FL |
| 10016 | Brown | James | | |
| 10017 | Williams | George | Mobile | AL |
| 10018 | Farriss | Anne | Opp | AL |
| 10019 | Smith | Olette | Nashville | TN |

## Row-centric storage

**Block 1**

10010,Ramas,Alfred,Nashville,TN
10011,Dunne,Leona,Miami,FL

**Block 2**

10012,Smith,Kathy,Boston,MA
10013,Olowski,Paul,Nashville,TN

**Block 3**

10014,Orlando,Myron,NULL,NULL
10015,O'Brian,Amy,Miami,FL

**Block 4**

10016,Brown,James,NULL,NULL
10017,Williams,George,Mobile,AL

**Block 5**

10018,Farriss,Anne,OPP,AL
10019,Smith,Olette,Nashville,TN

## Column-centric storage

**Block 1**

10010,10011,10012,10013,10014
10015,10016,10017,10018,10019

**Block 2**

Ramas,Dunne,Smith,Olowski,Orlando
O'Brian,Brown,Williams,Farriss,Smith

**Block 3**

Alfred,Leona,Kathy,Paul,Myron
Amy,James,George,Anne,Olette

**Block 4**

Nashville,Miami,Boston,Nashville,NULL
Miami,NULL,Mobile,Opp,Nashville

**Block 5**

TN,FL,MA,TN,NULL,
FL,NULL,AL,AL,TN

# NoSQL

**Graph databases** store data on relationship-rich data as a collection of **nodes** and **edges** (relationships)

**Properties** are the attributes of a node or edge of interest to a user

**Traversal** is a query in a graph database



Ex: Oracle RDF (Resources Description Framework)

30

# Any Questions

University of Windsor

# Recap and Conclusion

- Evolution of Data

- Relational Data to Big Data

- Debates of Big Data Implication

- Historical Interpretation Of Big Data

- Big Data Analytics (BDA)

- Hadoop (HDFS & Map Reduce) help to handle Big data efficiently.

- Spark a powerful open-source unified analytics engine

University of Windsor

# APPENDIX

# JavaScript Object Notation (JSON)

- JSON Syntax Rules
  - Data is in name/value pairs
  - Data is separated by commas
  - Curly braces hold objects
  - Square brackets hold arrays

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