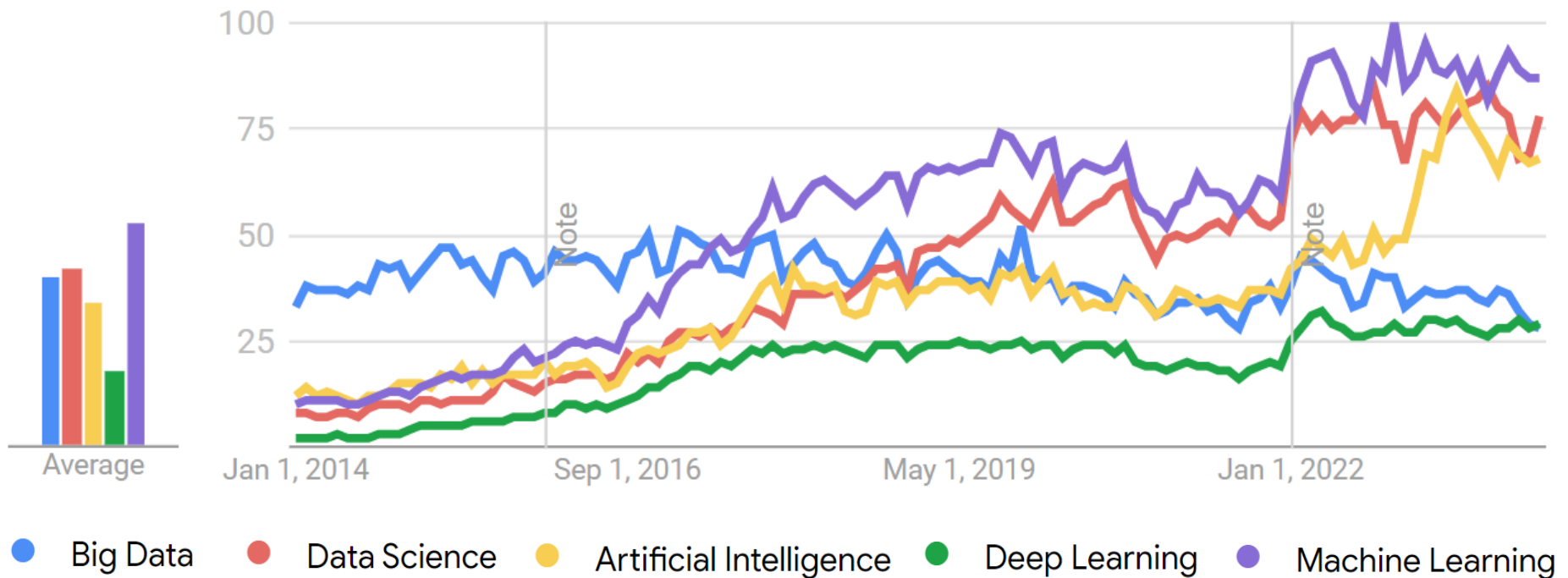# BIG DATA FUNDAMENTALS
## Part I

Le Ngoc Thanh – Nguyen Ngoc Thao

{lnthanh, nnthao}@fit.hcmus.edu.vn

# Big data: The trending term

- Big data is among trending search terms in recent years.



Source: Google Trends, updated 01/2024

# Outline

- What is Big data?
  - The definitions of Big data
  - The V's characteristics of Big data
  - Common Issues in Big data
- Big data case studies
  - The applications of Big data
  - Big data projects in practice
- Motivations and opportunities

# What is Big data?

It is not big. It is just bigger…

# Big data: A definition

- A variety of definitions for Big data are available worldwide.

Big data is a term used to refer to the *study and applications of data sets that are so big and complex* that traditional data-processing application software are inadequate to deal with them. – Wikipedia.
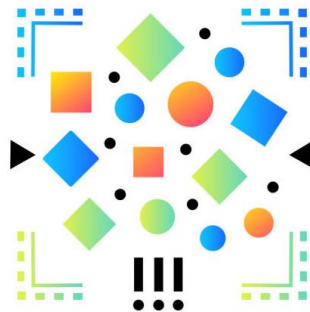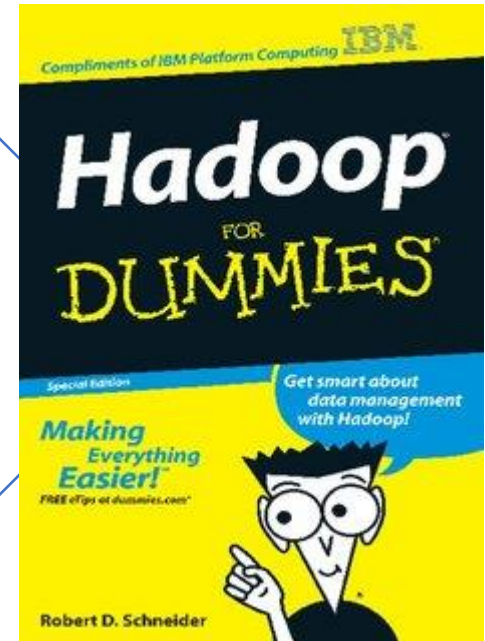
Big data is *high-volume, high-velocity and/or high-variety information assets* that demand cost-effective, innovative forms of information processing that enable enhanced insight, decision making, and process automation. – Garner, 2001.

Big data refers to *the dynamic, large and disparate volumes of data being created by people, tools and machines*; it requires new, innovative and scalable technology to collect, host and analytically process the vast amount of data gathered in order to derive real-time business insights that relate to consumers, risk, profit, performance, productivity management and enhanced shareholder value. – Ernst & Young, 2014.
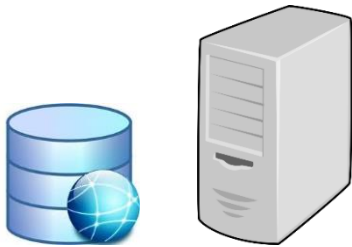
# Big data: A definition

**Big data** is a term that describes at least three separate, but interrelated, trends.

- ✓ Capturing and managing lots of information
- ✓ Working with many new types of data
- ✓ Exploiting these masses of information and new data types with new styles of applications

# Small data vs. Big data

- "Big data" is just the "small data" that grows bigger.
- The new scale of data may require novel approaches for techniques and frameworks.
- It is now able solve new problems or existing problems in a better way.

small data, small computer

bigger data

bigger computer?

or more small computers?

# Technologies in Big data

- Not a single technology but a combination of old and new technologies that helps companies gain actionable insight



There is an urgent need for parallel processing in distributed environment with high scalability.

- Capability to manage a huge volume of disparate data, at the right speed, and within the right time frame → allow for real-time analysis and reaction

# The 5V's characteristics

- The characteristics of Big data are characterized by the 5V's.



VOLUME

VELOCITY

VARIETY

VALUE

VERACITY

Image credit: Dreamtime

# The 5V's characteristics


VOLUME

- Description: The amount of data generated is vast compared to traditional data sources.

- Attributes: Exabyte, zettabyte, yottabytes, etc.

- Drivers: Increase in data sources, higher resolution sensors, scalable infrastructure.
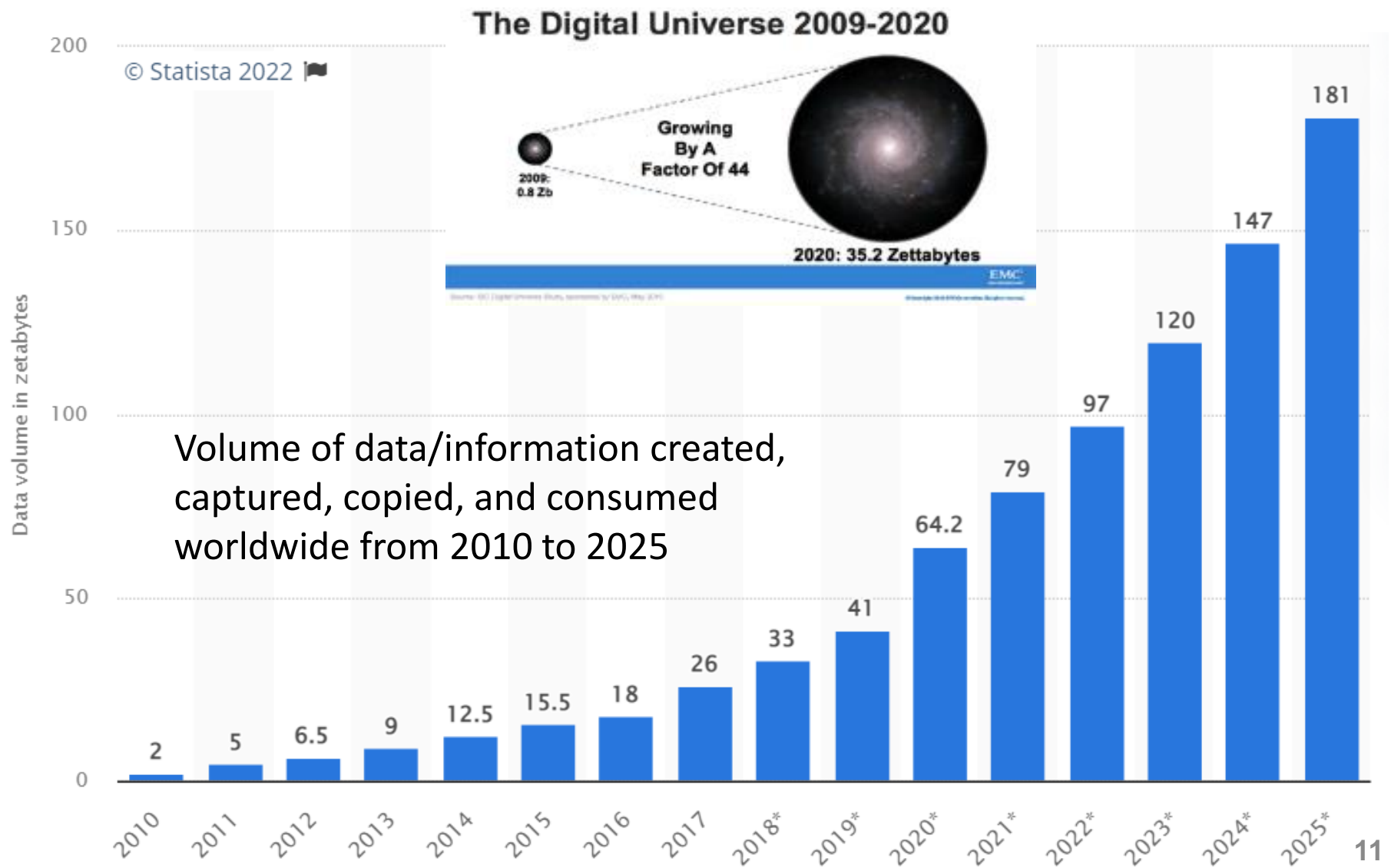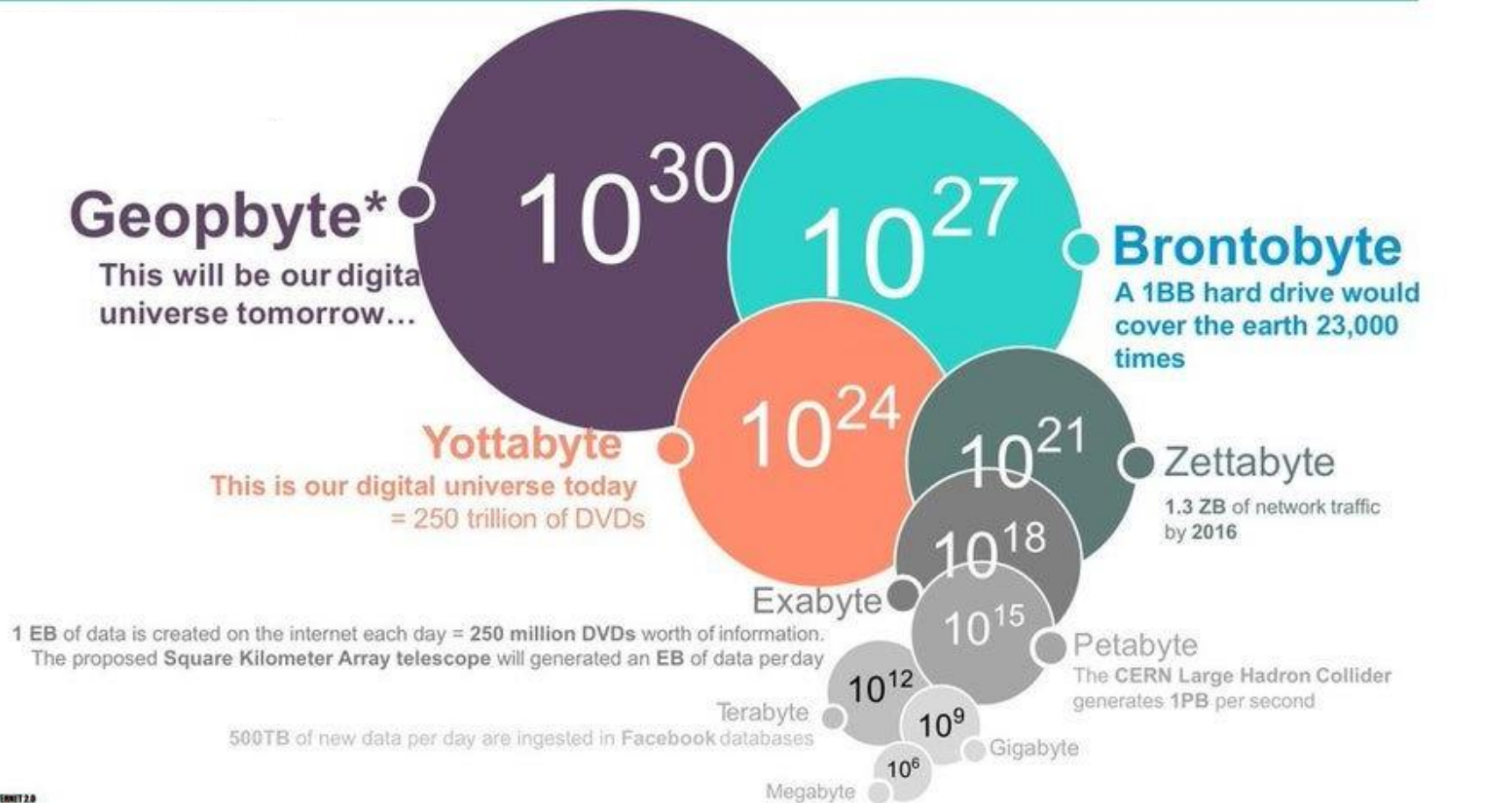
# The explosive growth of data



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

11

# New data units for the Big data era



**Geopbyte\***
This will be our digital universe tomorrow...

$10^{30}$

$10^{27}$

**Brontobyte**
A 1BB hard drive would cover the earth 23,000 times

**Yottabyte**
This is our digital universe today
= 250 trillion of DVDs

$10^{24}$

$10^{21}$

Zettabyte
1.3 ZB of network traffic by 2016

$10^{18}$

Exabyte

1 EB of data is created on the internet each day = 250 million DVDs worth of information.
The proposed Square Kilometer Array telescope will generated an EB of data per day

$10^{15}$

Petabyte
The CERN Large Hadron Collider generates 1PB per second

$10^{12}$

Terabyte
500TB of new data per day are ingested in Facebook databases

$10^{9}$

Gigabyte

$10^{6}$

Megabyte

ZETANET  INTERNET 2.0 SKIPJACK SUPERENCRYP BLOCK

\*The terms Gegobyte and Geobyte are also used in the literature

# 3 Important Statistics About How Much Data Is Created Every Day

**Finances**Online
REVIEWS FOR BUSINESS

## 1 How much data is generated every minute?

**41,666,667**
messages shared
by WhatsApp users

**1,388,889**
video / voice calls made
by people worldwide

**404,444**
hours of video streamed
by Netflix users

**347,222**
stories posted by Instagram users

**150,000**
messages shared by Facebook users

**147,000**
photos shared by Facebook users

## 2 Estimated Data Consumption from 2021 to 2024

| Year | Data |
|------|------|
| 2021 | 74 ZETTABYTES |
| 2022 | 94 ZETTABYTES |
| 2023 | 118 ZETTABYTES |
| 2024 | 149 ZETTABYTES |

## 3 Data Growth in 2021

**2 TRILLION**
searches on Google by the end of 2021

**1.134 TRILLION MB**
volume of data created every day

**3,026,626**
emails sent every second, 67% of which are spam

**278,108 PETABYTES**
global IP data per month by the end of 2021

**230,000**
new malware versions created every day

**82%**
share of video in total global internet traffic at the end of 2021

13

# Why does data become big now?

- Key enablers of appearance and growth of data are

Increase of storage capacity

Availability of data

Increase of processing power

# The 5V's characteristics


VELOCITY

- Description: Data is being generated extremely fast, a process that never stops; and the speed at which data is transformed into insight

- Attributes: Batch; near/real-time; streams

- Drivers: Improved connectivity; competitive advantage; precomputed information

# Real-time and/or fast data

Mobile devices
(tracking all objects all the time)

Scientific instruments
(collecting all sorts of data)

Sensor technology
and networks
(measuring all kinds of data)

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

- Innovations and their progresses are 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.
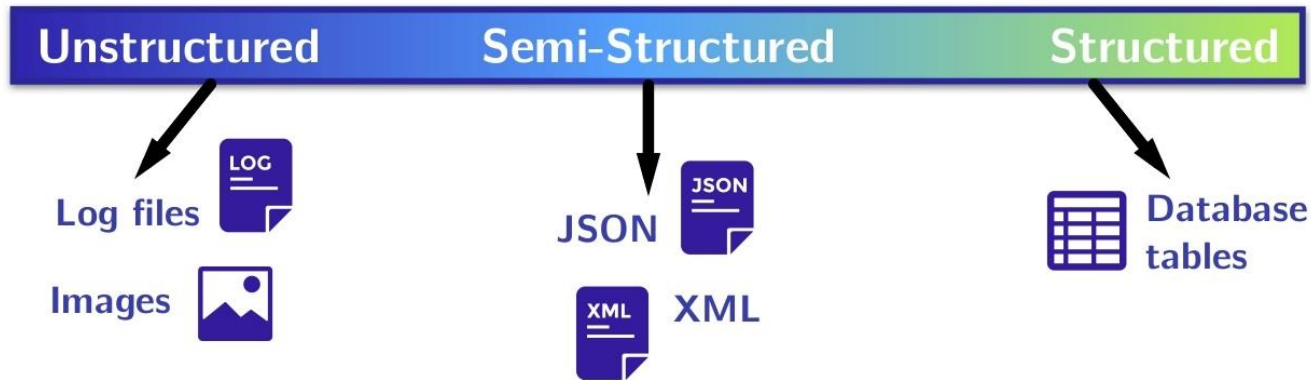
# Real-time analytics/decision requirement

Product Recommendations that are <u>Relevant</u> & <u>Compelling</u>

Influence Behavior

Learning why Customers switch to competitors and their offers; in time to Counter

Customer

Improving the Marketing Effectiveness of a Promotion while it is still in Play

Friend Invitations to join a Game or Activity that expands business

Preventing Fraud as it is O<u>ccurring</u> & Preventing more proactively
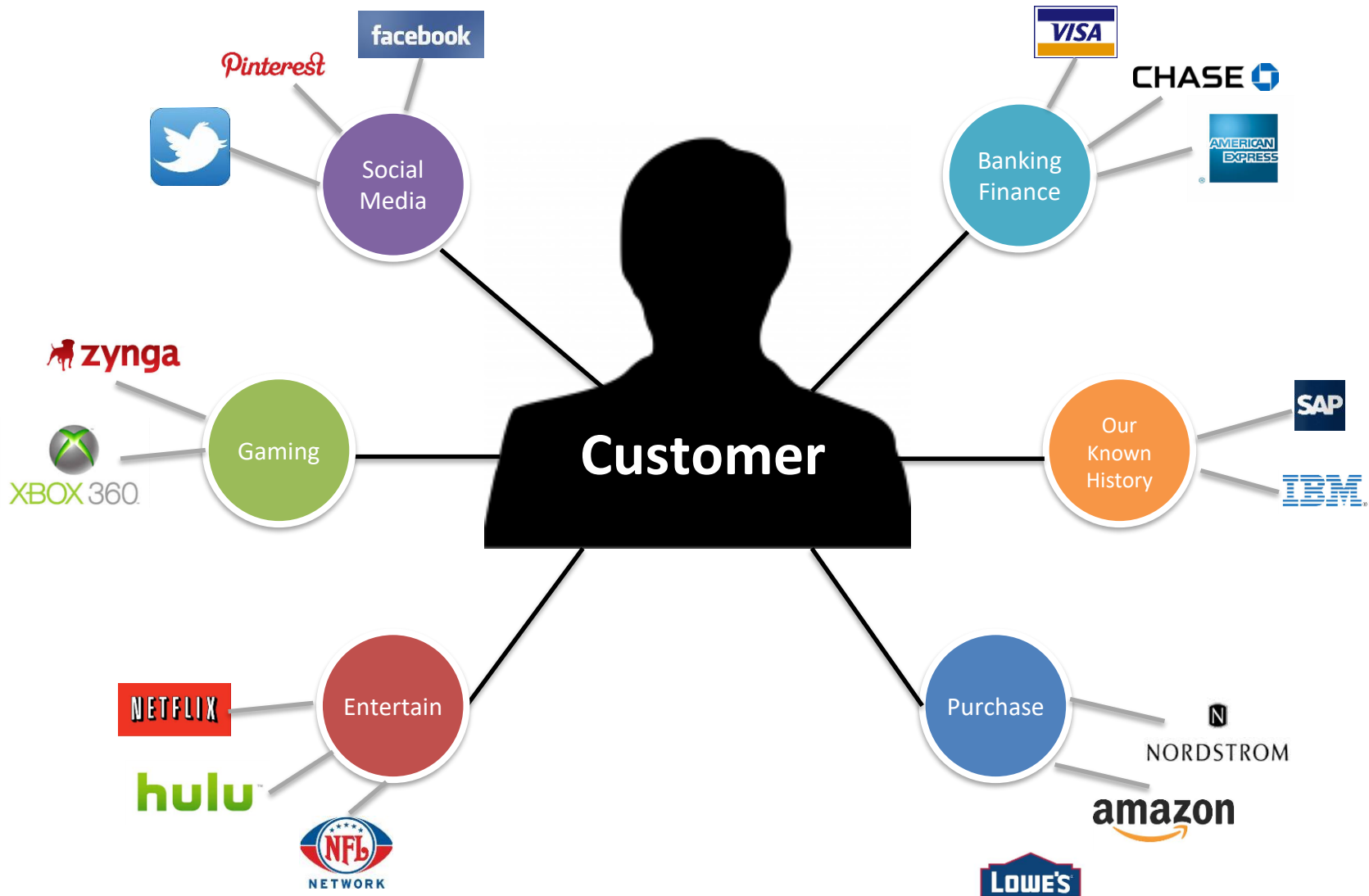
Data processed per day (Updated in January 2022)

# The 5V's characteristics



- Description: Data comes from different sources, machines, people, and processes both from outside and inside the organizations

- Attributes: Degree of structure; complexity

- Drivers: Mobile; social media; video; genomics; IoT

# A single view to the customer

# The 5V's characteristic



**VERACITY**

- Description: Quality and origin of data

- Attributes: Consistency; completeness; integrity; ambiguity

- Drivers: Cost; need of traceability and justification



**VALUE**

- The ability and need to turn data into value

- Value is not only profit but also medical or social benefits, or personal satisfaction (customer, employee, etc.).

# Common issues related to the 5V's

- As the data volume increases, the value of different data records will decrease in proportion to age, type, richness, and quantity among other factors.

- It is hard to handle complex data by existing traditional analytic systems.

  - Big data with relational databases, statistics/visualization packages

  - Massively parallel software running on tens, hundreds, or even thousands of computing units.

  - Data analytics with data that is constantly in motion.

# Issues of personnel

- There is a considerable gap between Business leaders and IT professionals.



- Business leaders concern about adding value to their business and getting more and more profit.
- Meanwhile, IT leaders focus on the technicalities of the storage and processing only.
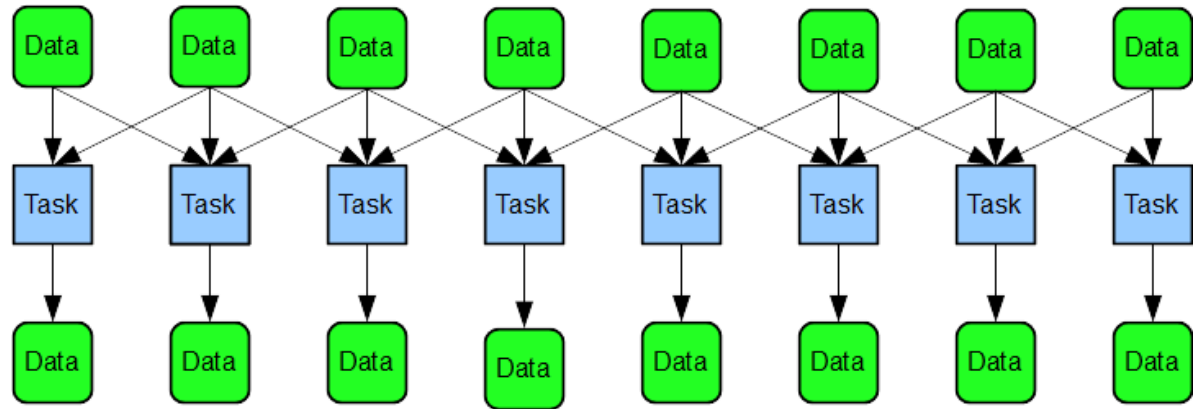
# Issues of storage and transport

- Assume that we have an exabyte ($10^{18}$) of data.

- Modern hard disks can store several TBs ($10^{12}$)

  $\rightarrow$ thousands of disks are required for an exabyte.

- A single computer system would be unable to directly attach the requisite number of disks.

- Accesses to that data also overwhelms the networks.

  - E.g., a 1GB/second network operating at 80% efficiency and a sustainable bandwidth of 100 MB/second would take approximately 14,465 days to transfer an exabyte of data.

# Issues of data management

- Possibly the most difficult problem

- Issues of access, utilization, updating, governance, and reference (in publication) are major stumbling blocks.


- Data sources are varied by size, format, and by method of collection.

  - What, when, where, who, why and how it was collected.

- It is impractical to validate every data item in a huge source.

# Issues of processing power

- Extensive parallel processing and new analytics algorithms are required.



Assume that an exabyte of data need to be processed and it is chunked into blocks of 8 words → 1 exabytes = 1K petabytes.

Assuming a processor expends 100 instructions on one block at 5 gigahertz → 1K petabytes would require a processing time of 635 years.

# Big data case studies

The more data, the better decisions, and then the better outcomes…

# Big Data use case categories

**Big Data Exploration**

Find, visualize, understand all big data to improve decision making

**Enhanced 360° View of the Customer**

Extend existing customer views (MDM, CRM, etc) by incorporating additional internal and external information sources

**Security/Intelligence Extension**

Lower risk, detect fraud and monitor cyber security in real-time
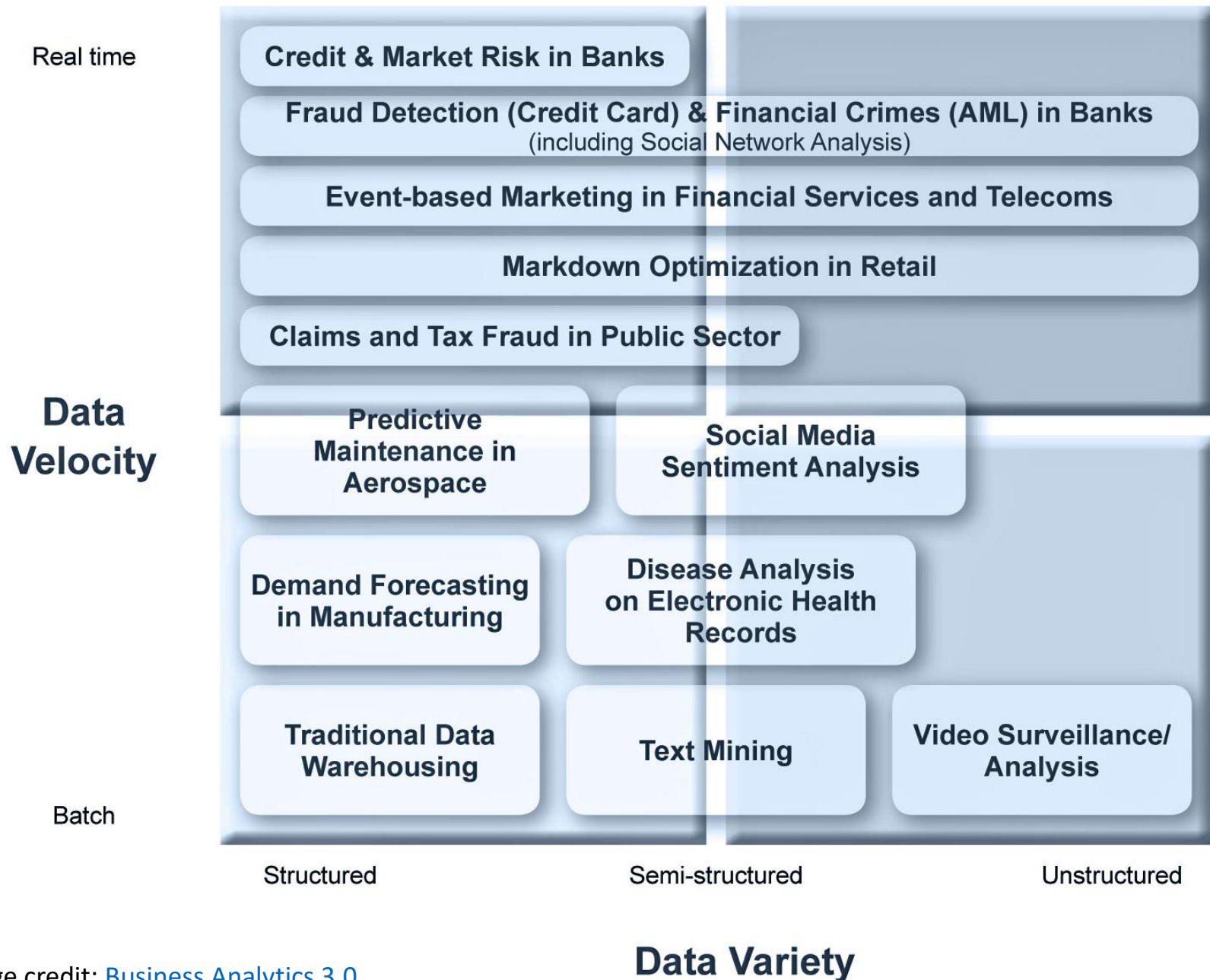
**Operations Analysis**

Analyze a variety of machine data for improved business results
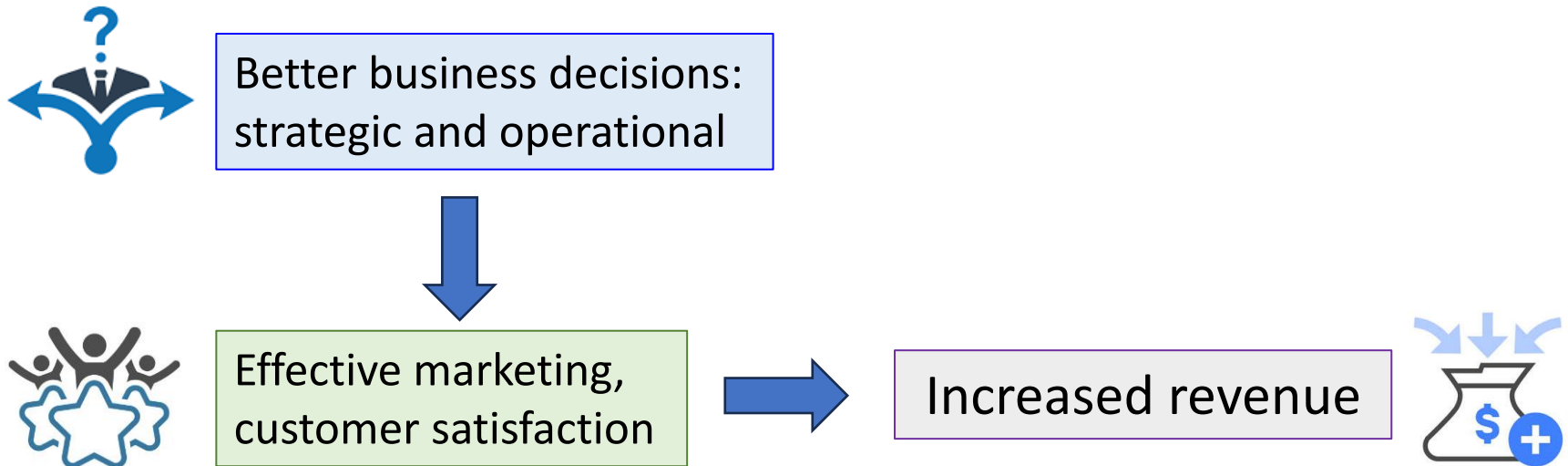
**Data Warehouse Augmentation**

Integrate big data and data warehouse capabilities to increase operational efficiency

# Potential Use Cases for Big Data Analytics



| | Structured | Semi-structured | Unstructured |
|---|---|---|---|
| **Real time** | Credit & Market Risk in Banks | | |
| | Fraud Detection (Credit Card) & Financial Crimes (AML) in Banks (including Social Network Analysis) | | |
| | Event-based Marketing in Financial Services and Telecoms | | |
| | Markdown Optimization in Retail | | |
| | Claims and Tax Fraud in Public Sector | | |
| **Data Velocity** | Predictive Maintenance in Aerospace | Social Media Sentiment Analysis | |
| | Demand Forecasting in Manufacturing | Disease Analysis on Electronic Health Records | |
| **Batch** | Traditional Data Warehousing | Text Mining | Video Surveillance/ Analysis |

**Data Variety**

Image credit: Business Analytics 3.0

# Big data analytics

- This analytics exploits a **large** amount of data for interesting data relationships to gain competitive advantage.

- Appropriate information: hidden patterns or correlations

- Competitive advantage:

Better business decisions: strategic and operational

Effective marketing, customer satisfaction

Increased revenue

# Big data analytics

- Big data is more real-time in nature than traditional data warehouse applications.

  - Conventional architectures are ill-suited for big data apps (e.g., Exadata, Teradata)

- There are many challenges in handling Big data, mainly about the technology bottleneck and the lack of experts.



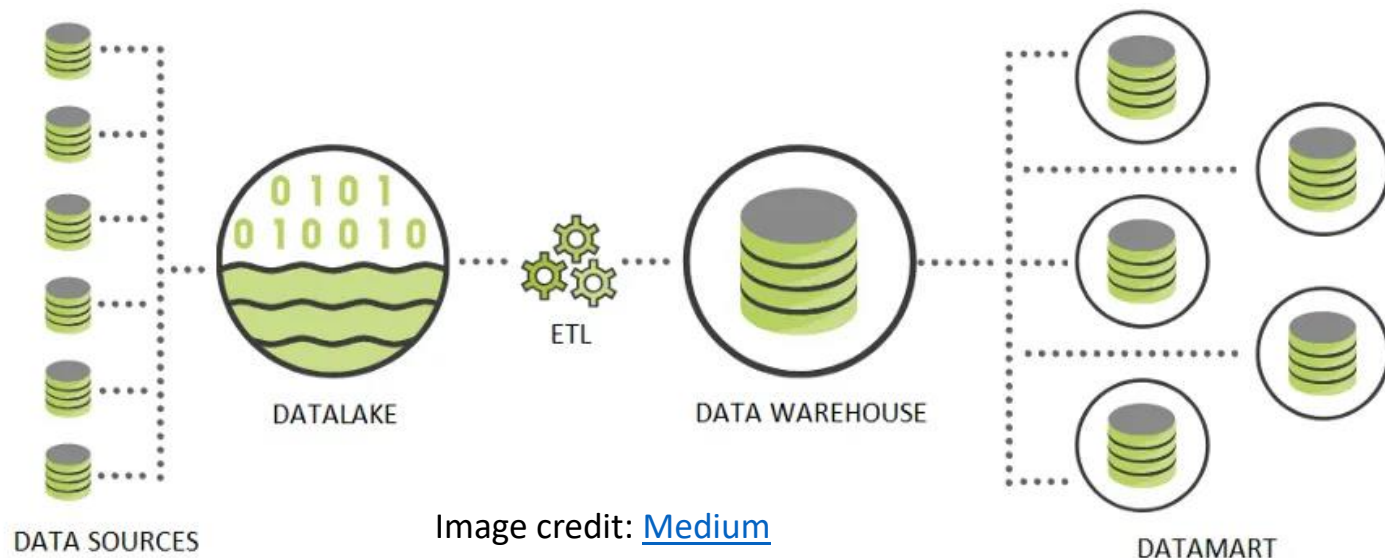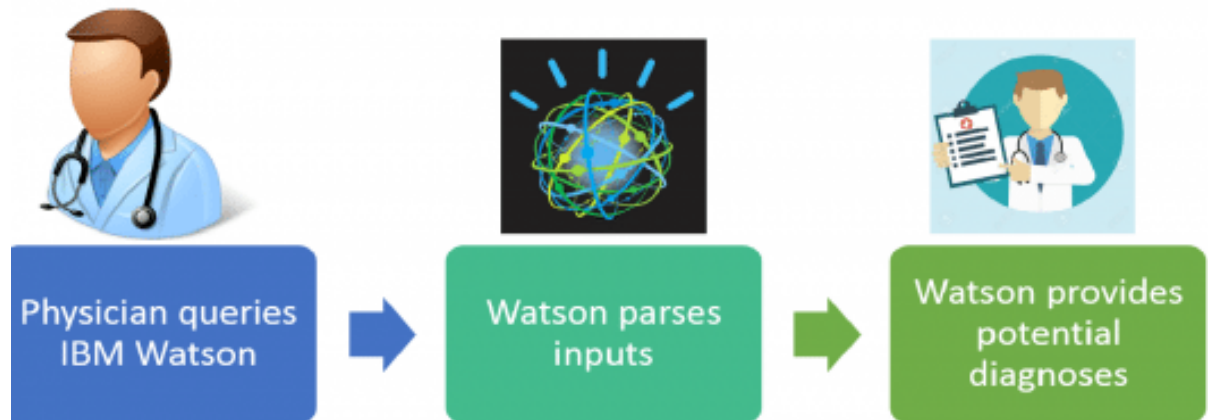DATA SOURCES  DATALAKE  ETL  DATA WAREHOUSE  DATAMART

0 1 0 1
0 1 0 0 1 0

Image credit: Medium

# Big data in Healthcare

- 80% of medical data is unstructured and clinically relevant, residing in multiple places.

    - Individual EMRs, labs and imaging systems, physician notes, medical correspondence, etc.

- Big data may help increase access to healthcare.

    - Build sustainable healthcare systems, collaborate to improve care and outcomes.



Physician queries IBM Watson → Watson parses inputs → Watson provides potential diagnoses

# Big data in Healthcare



VinBigData's Genomics project.
Homepage: https://genome.vinbigdata.org/ (Updated 2021)

# KTH: Reducing traffic congestion



KTH Swedish Royal Institute of Technology Reducing Traffic Congestion

Capabilities Utilized:

**Stream Computing**

- Deployed real-time Smarter Traffic system to predict and improve traffic flow.
- Analyzes streaming real-time data gathered from cameras at entry/exit to city, GPS data from taxis and trucks, and weather information.
- Predicts best time and method to travel such as when to leave to catch a flight at the airport

Significant benefits:

- Enables ability to analyze and predict traffic faster and more accurately than ever before
- Provides new insight into mechanisms that affect a complex traffic system
- Smarter, more efficient, and more environmentally friendly traffic

34

**SONAR:**
**~10-100 KB**
per second

**RADAR:**
**~10-100 KB**
per second

**GPS:**
**~50 KB**
per second

**CAMERAS:**
**~20-40 KB**
per second

**LIDAR:**
**~10-70 KB**
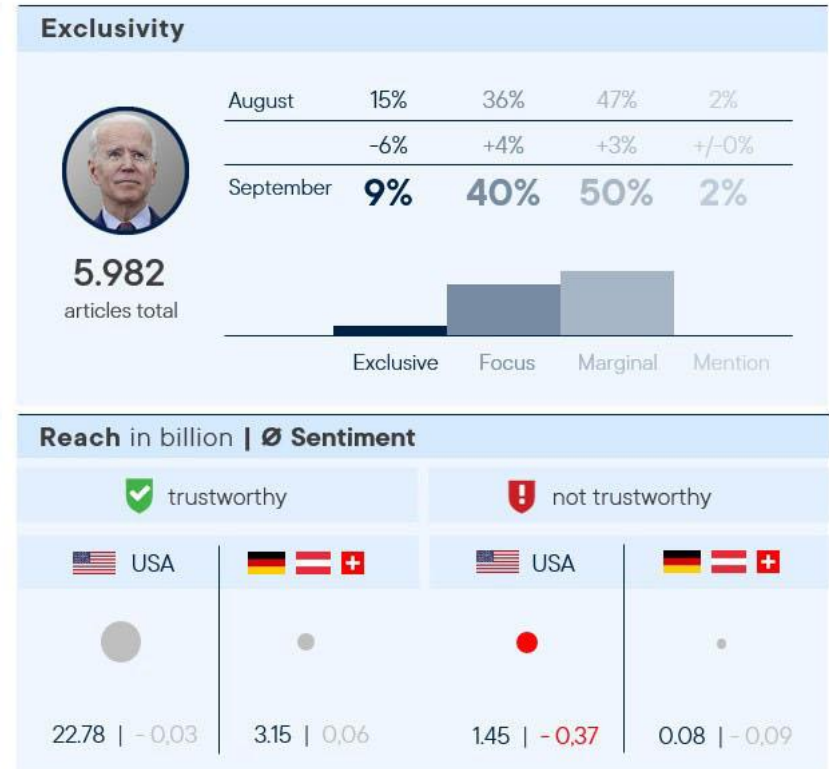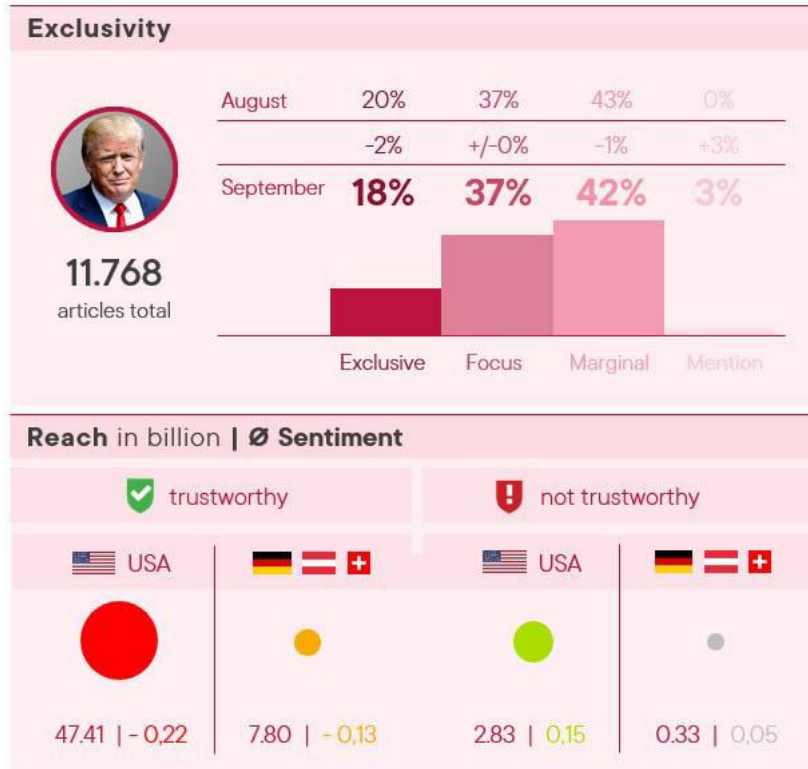per second

**AUTONOMOUS VEHICLES GET**
**4,000 GB**
**DATA EACH DAY**

# Big data in social network analysis



MEDIA RESONANCE

**Trump Dominates the News, Biden has More Positive Coverage in Credible Media**

The trend has even strengthened compared to the previous month. The more trustworthy the media outlet, the more positive is its coverage of Biden. DACH media report neutrally to positively about Biden and predominantly negatively about Trump.

# Big data in social network analysis



Shevtsov, Alexander, et al. "Analysis of Twitter and YouTube during US elections 2020."
*arXiv e-prints* (2020): arXiv-2010.

# Motivations and opportunities

A new horizon that changes our lives…

# New insight into data

- Why deal with more data? ***New insights.***

- The insights are for people throughout the enterprise, not top-level executives only.

  - People involve may include the CEO, marketing staffs, data analyst and programmers.

- They can be transformed into actionable intelligence

  $\rightarrow$ better business, better service to customers

  $\rightarrow$ more profit gained

# Big data analytics: Applications


**Smarter Healthcare**


**Multi-channel sales**


**Finance**


**Log Analysis**


**Homeland Security**


**Traffic Control**


**Telecom**


**Search Quality**


**Manufacturing**


**Trading Analytics**


**Fraud and Risk**


**Retail: Churn, NBO**

# Big data analytics: Applications

**amazon**

310 million active users globally

12 million items across all its categories and services

Over 1 billion monthly active users

**TikTok**

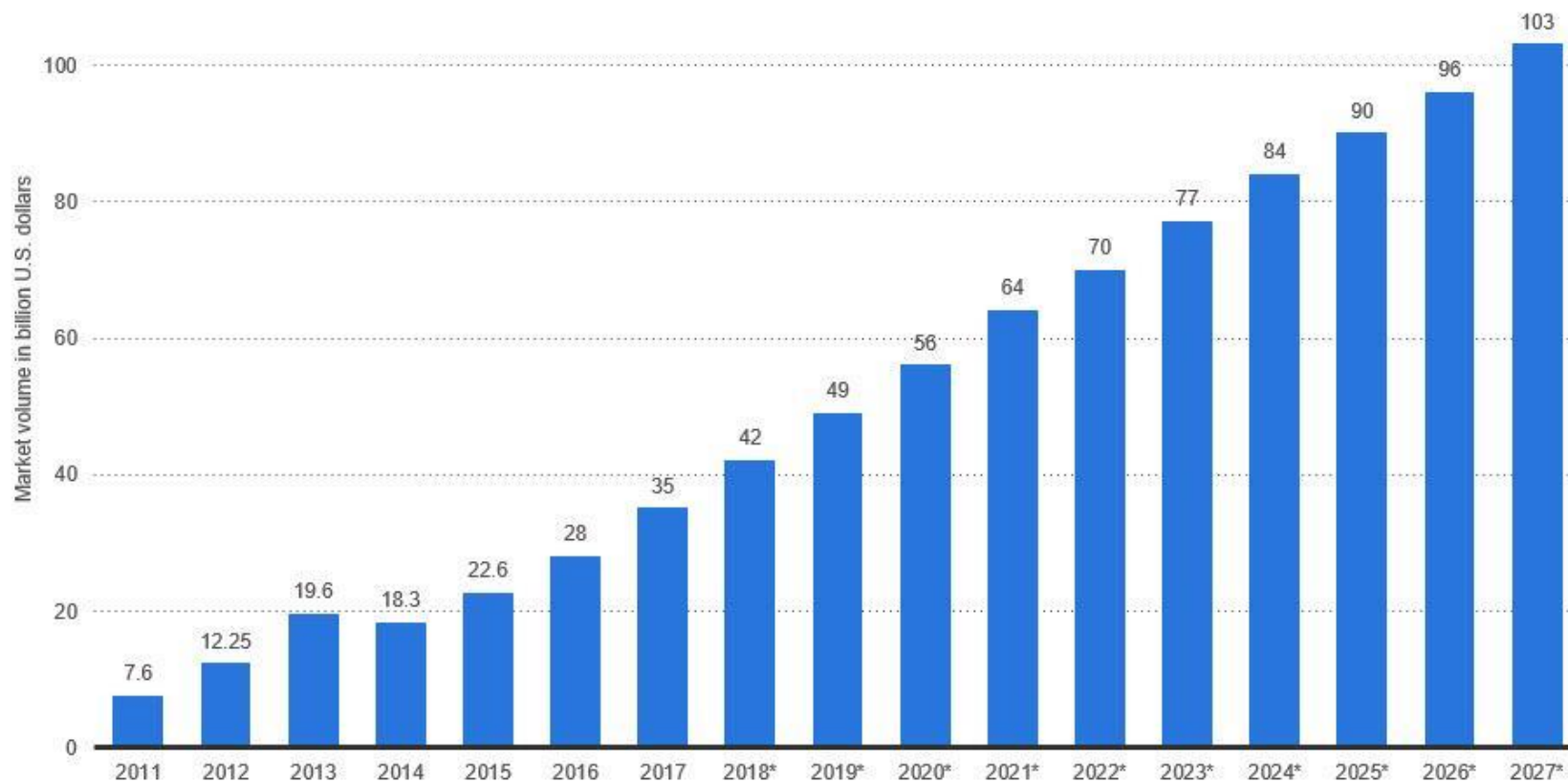There were a total of 8.6 billion videos uploaded in TikTok in 2021. The number in 2023 approximates 14.4 billion videos.

**STEAM®**

20,919,374 players online, 5,163,503 players in-game

105,281 tracked Steam games

Forecast Revenue Big Data Market Worldwide 2011-2027

# Big Data Market Size Revenue Forecast Worldwide From 2011 To 2027 (in billion U.S. dollars)

Market volume in billion U.S. dollars

| Year | Value |
|------|-------|
| 2011 | 7.6 |
| 2012 | 12.25 |
| 2013 | 19.6 |
| 2014 | 18.3 |
| 2015 | 22.6 |
| 2016 | 28 |
| 2017 | 35 |
| 2018* | 42 |
| 2019* | 49 |
| 2020* | 56 |
| 2021* | 64 |
| 2022* | 70 |
| 2023* | 77 |
| 2024* | 84 |
| 2025* | 90 |
| 2026* | 96 |
| 2027* | 103 |

statista

# Big Data Market Forecast Worldwide from 2011 to 2026, by segment (in billion U.S. dollars)



**Legend:** Professional Services ■ Apps & Analytics ■ Storage ■ Compute ■ SQL ■ Data Management ■ Networking ■ NoSQL ■ Hadoop

# Salary by Job roles in India (July 2023)



| Job Role | Salary |
|---|---|
| Data Architect | 28,50,000 |
| Business Analyst | 13,00,000 |
| Machine Learning Engineer | 16,00,000 |
| Statistician | 16,80,000 |
| Quantitative Analyst | 14,00,000 |
| Data Analyst | 14,90,000 |
| Database Administrator | 13,00,000 |
| AI Engineer | 19,90,000 |
| Data Engineer | 13,80,000 |
| Risk Analyst | 16,00,000 |
| Marketing Analyst | 12,00,000 |

ANALYTIXLABS

* The currency unit is INR (Indian Rupee)

# Salary by tools in India (July 2023)



Kuberflow — 27.9
PySpark — 22.4
Kubernetes — 20.7
Qlikview — 18.9
AWS — 18.2
Tableau — 17.8
SAS — 16.6
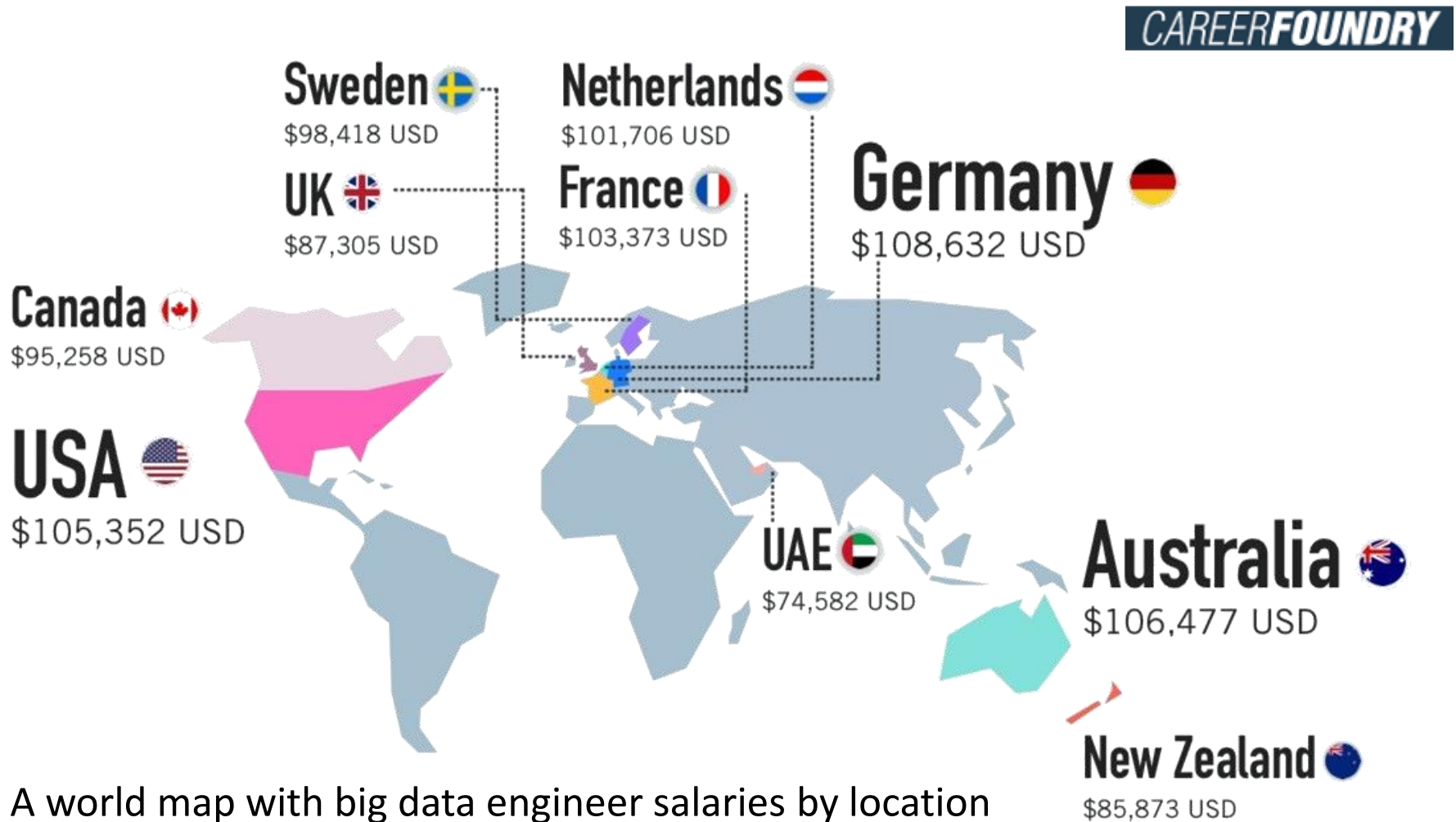SQL — 15.4

ANALYTIXLABS

* The currency unit is Lakhs in INR (Indian Rupee)

# Big data engineer salaries by location



A world map with big data engineer salaries by location (updated in January 2023)

Sweden $98,418 USD
Netherlands $101,706 USD
UK $87,305 USD
France $103,373 USD
Germany $108,632 USD
Canada $95,258 USD
USA $105,352 USD
UAE $74,582 USD
Australia $106,477 USD
New Zealand $85,873 USD

CAREERFOUNDRY

# Average big data engineer salary

CAREER**FOUNDRY**

(updated in January 2023)

| | |
|---|---|
| Meta | $229K |
| Google | $205K |
| Microsoft | $183K |
| amazon | $167K |
| Apple | $170K |