# Data Engineering Platforms For Analytics

## Foundations of Data Systems

# Modules

- Welcome

- Introduction to data

- Data Preparation

- Database Management Systems

- Types of Data bases

# Welcome

Introductions, Syllabus review

# Syllabus

- Session 1 – Foundations of data systems
- Session 2 – Relational Databases
- Session 3 – Structured Query Language ( SQL )
- Session 4 – Analytical Data Platforms
- Session 5 – Business Intelligence
- Session 6 – Cloud Data Pipelines
- Session 7 – NoSQL - Document and Graph Databases
- Session 8 – NoSQL - Columnar and Key Value Databases
- Session 9 – Introduction to Blockchain
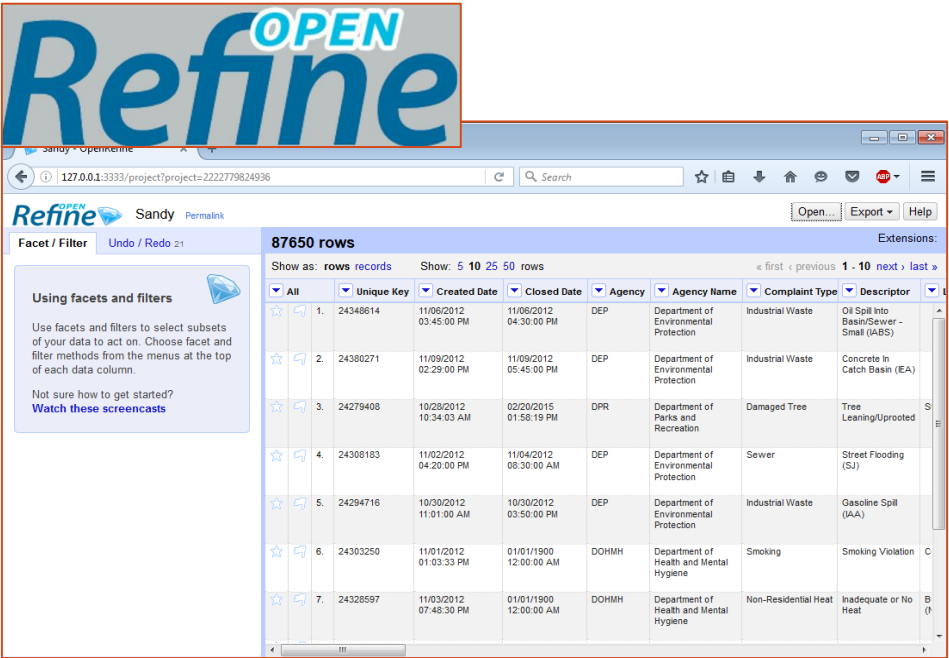- Session 10 –  Team Project Presentations

# Assignments & Final Project

- 4 Assignments ( Individual )
  - 100 points each
  - Bi-weekly submission ( Saturdays 11:59 PM CST )

- Final Project ( Max 4 students )
  - Business Use Case - 10%
  - Data Analysis & Preparation - 25%
  - Data Modelling & Design - 25 %
  - Tools / Database concepts -20%
  - Dashboards and insights - 20%

# Primary Datasets

- Classic Models ( In-Class exercises )
  - Sessions 2, 3, 4

- Sakila ( Assignments )
  - Assignments 2,3,4

- Final Project ( Sample datasets )
  - IRI Dataset
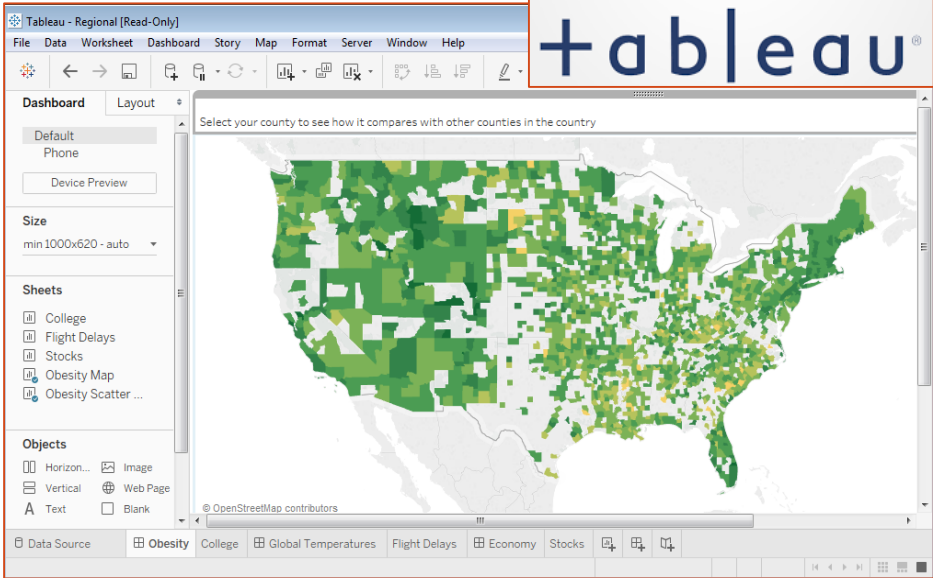  - https://toolbox.google.com/datasetsearch/
  - https://data.cityofchicago.org/
  - https://opendata.cityofnewyork.us/
  - https://data.gov.in/catalogs/
  - https://github.com/awesomedata/awesome-public-datasets/
  - https://www.springboard.com/blog/free-public-data-sets-data-science-project/

# Software Tools

# Software Tools - Continued…

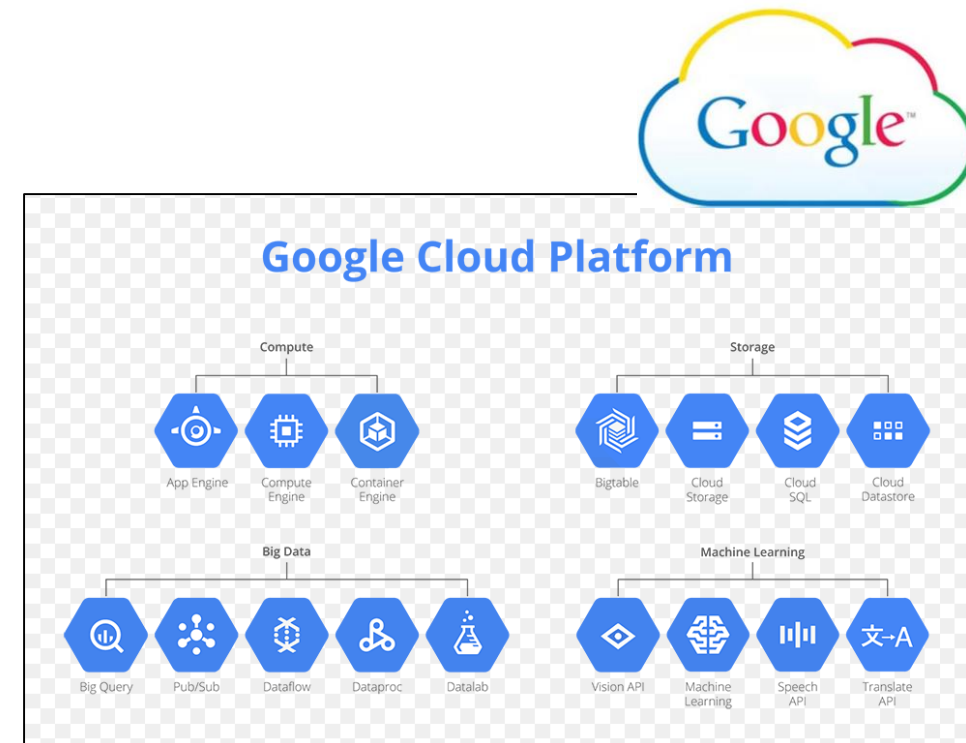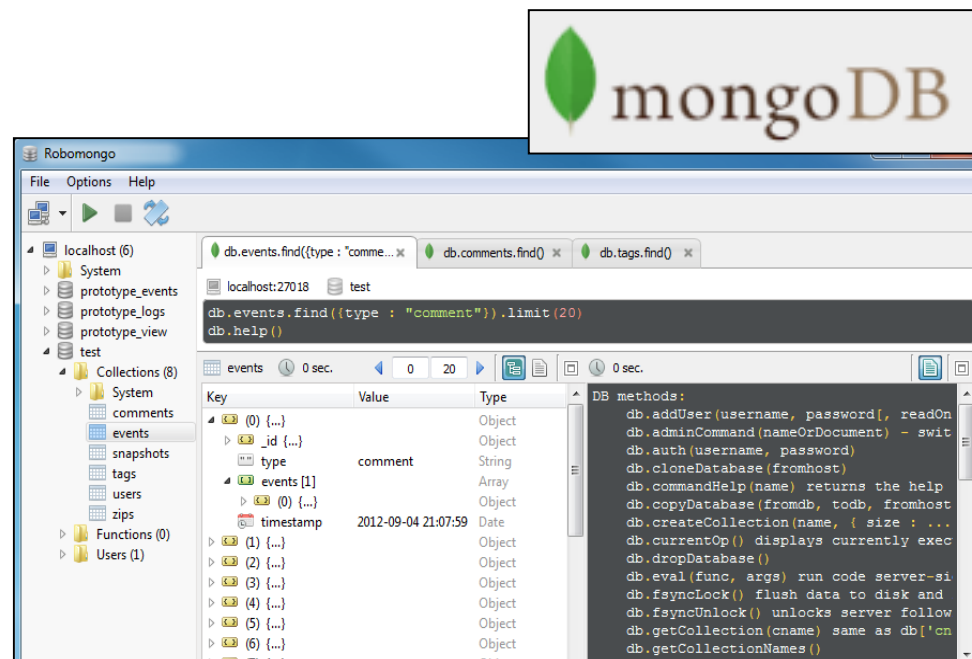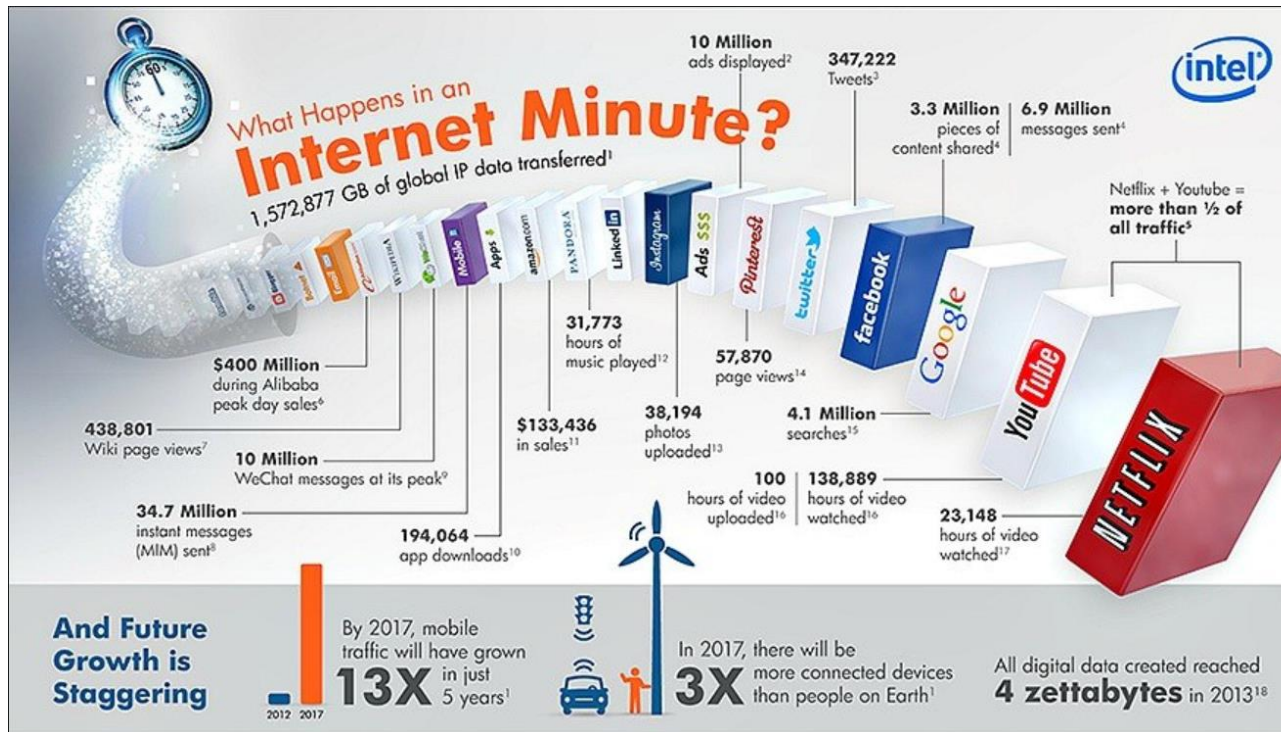# Software Tools - Continued…

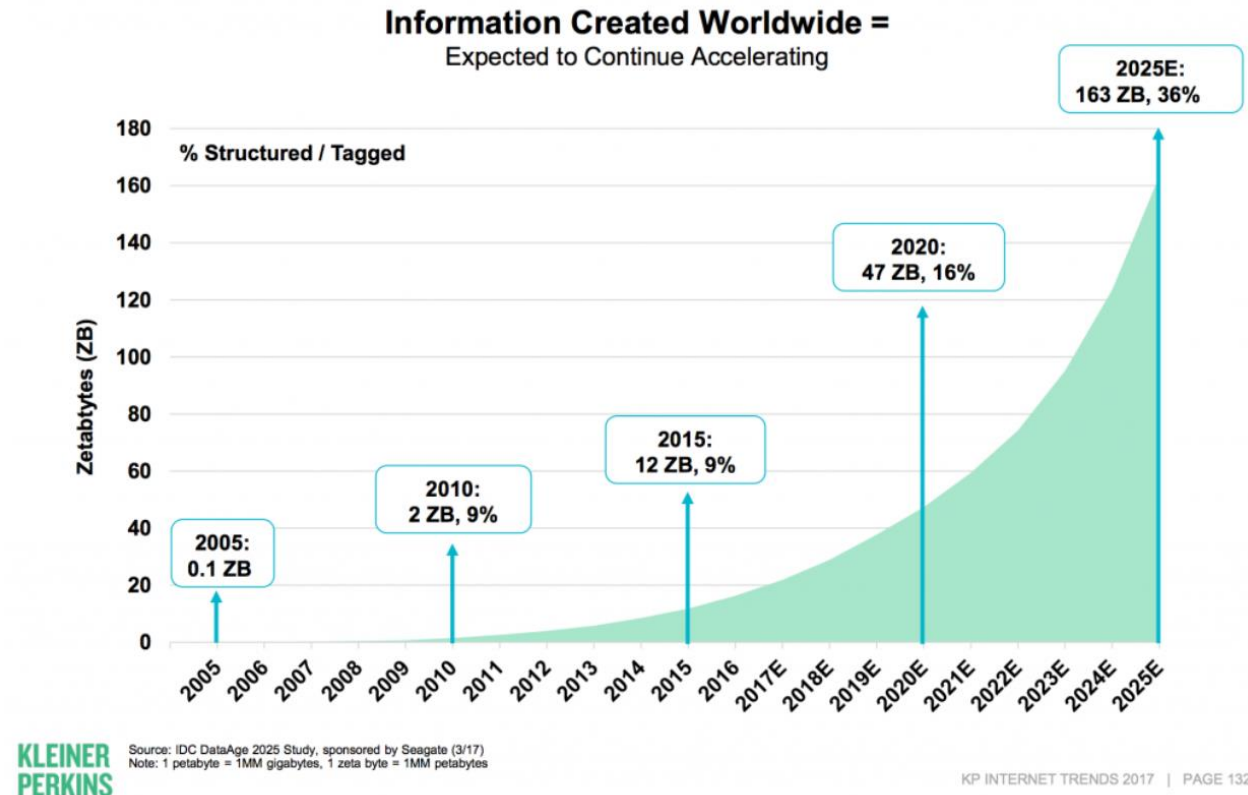# Software Tools - Continued…
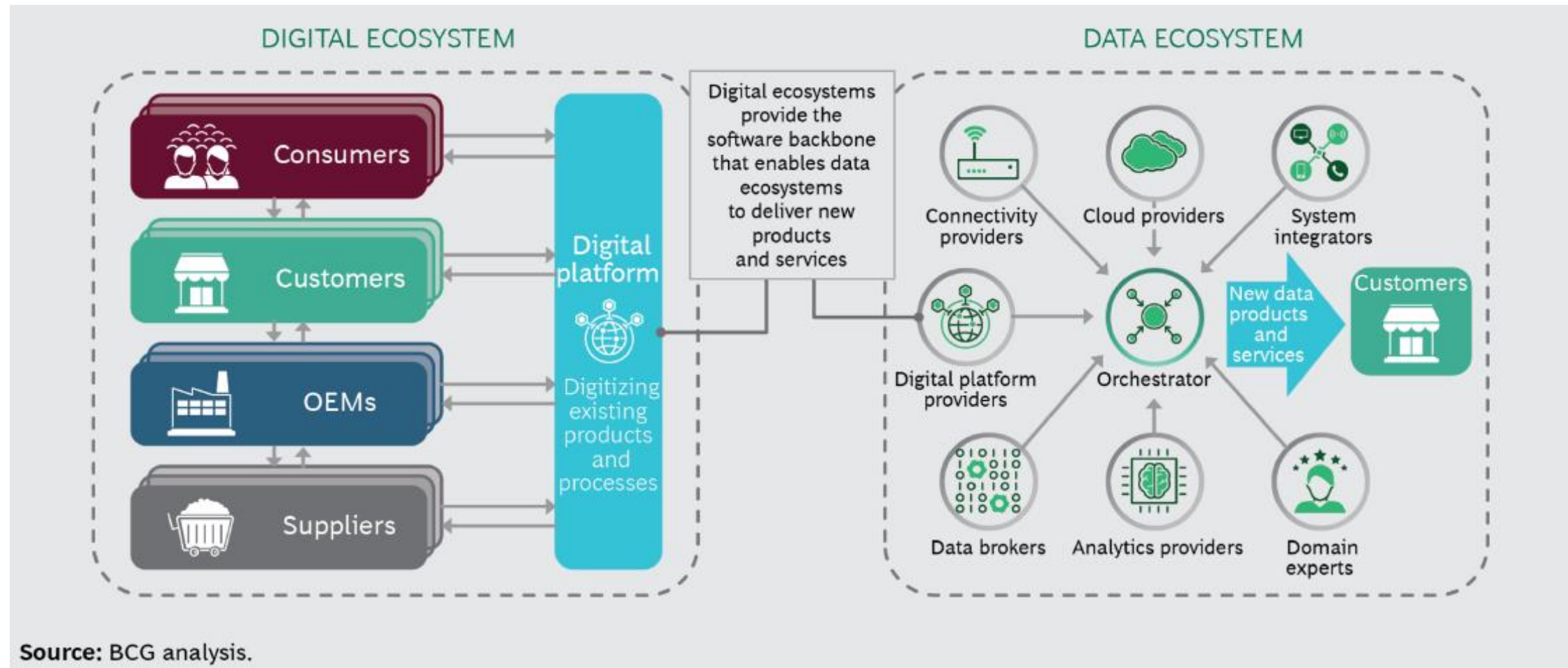
# What Happens in an Internet Minute

# Growth of data

- The New York Stock Exchange generates about 4–5 terabytes of data per day.

- Facebook hosts more than 240 billion photos, growing at 7 petabytes per month.

- Ancestry.com, the genealogy site, stores around 10 petabytes of data.

- The Internet Archive stores around 18.5 petabytes of data.

- The Large Hadron Collider near Geneva, Switzerland, produces about 30 petabytes of data per year.

Source: Hadoop The Definitive Guide

**Information Created Worldwide =**
Expected to Continue Accelerating

% Structured / Tagged

2025E: 163 ZB, 36%

2020: 47 ZB, 16%

2015: 12 ZB, 9%

2010: 2 ZB, 9%

2005: 0.1 ZB

Zetabytes (ZB)

2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017E 2018E 2019E 2020E 2021E 2022E 2023E 2024E 2025E

KLEINER PERKINS

Source: IDC DataAge 2025 Study, sponsored by Seagate (3/17)
Note: 1 petabyte = 1MM gigabytes, 1 zeta byte = 1MM petabytes

KP INTERNET TRENDS 2017 | PAGE 132

# Digital and Data Ecosystems

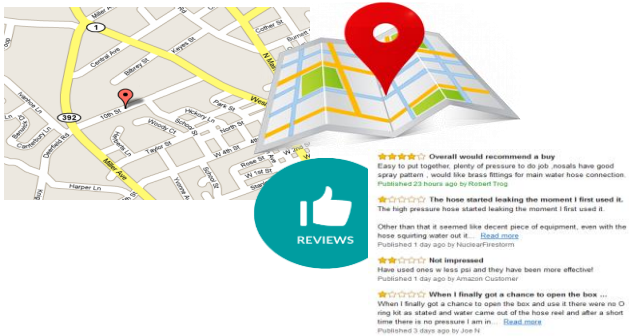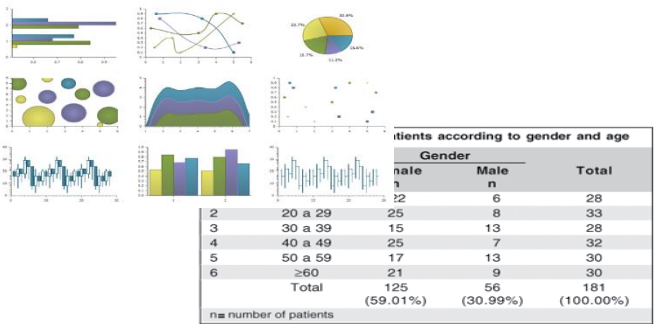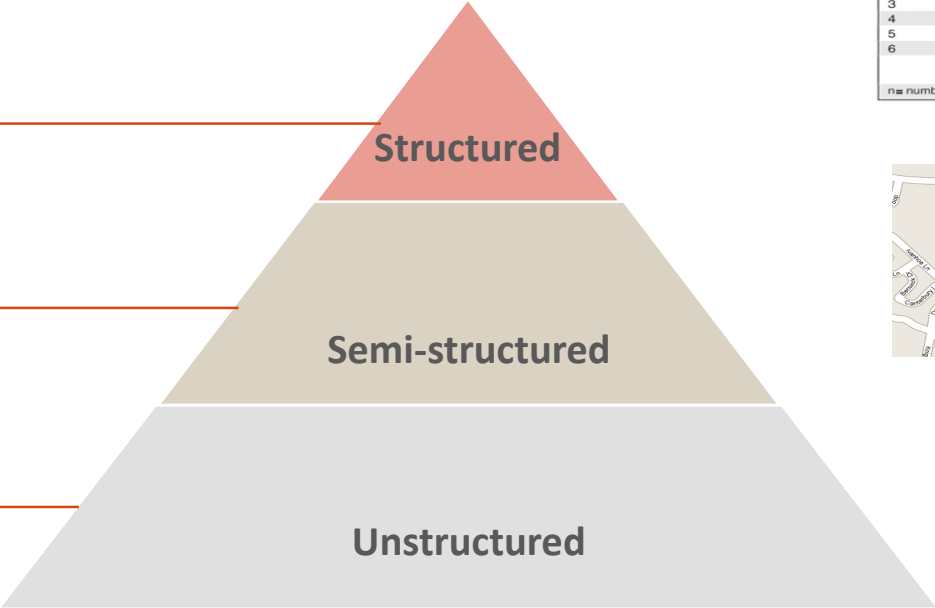

Source: BCG analysis.

# Data Classification

Data Growth, Data Types, Formats

# Data Formats, Classification

Well defined data types, fixed fields e.g. numeric, currency, date, address, text

Textual data with patterns e.g. XML, EDI, logs, Sensors

No inherent structure. File formats e.g. audio, video, pdf

**Structured**

**Semi-structured**

**Unstructured**

# Types of data



**Categorical**
Classified or ordered
Nominal
     Descriptions or labels - No sense of sequence
     Ex: Red/blue/yellow OR M/F, etc
Ordinal
     Meaningful Order
     Ex: 1st (98.3%)/2nd (97.2%)/3rd(91.4%)

**Numerical Data**
Can be measured
Discrete
     Quantitative data with whole numbers
     Ex: number of students in a class
Continuous
     Can take on fractions and decimals
     Ex: height / weight

# Binary Format

- Bit - smallest unit of information

<div align="center">

0/1

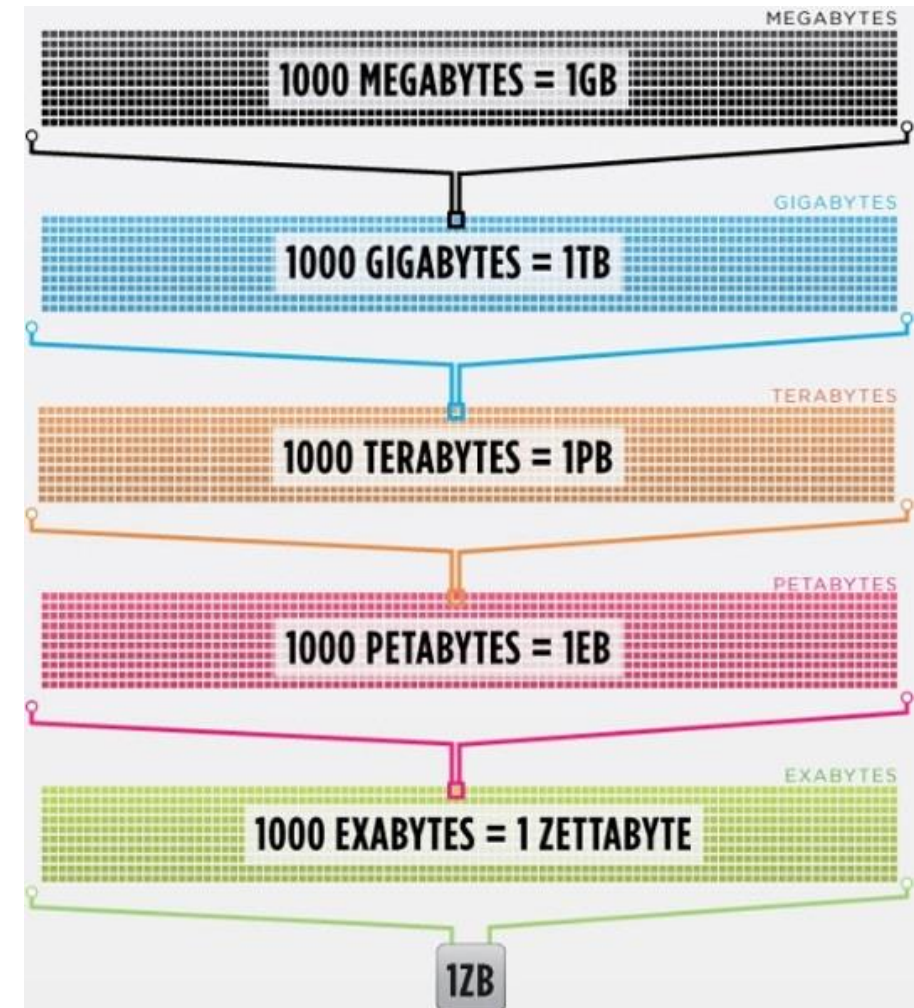</div>

- Byte - 8 bits of information

<div align="center">

0/1   0/1   0/1   0/1   0/1   0/1   0/1   0/1

</div>

- 1 Kilo Byte        $= 2^{10}$ Bytes  = 1024 bytes

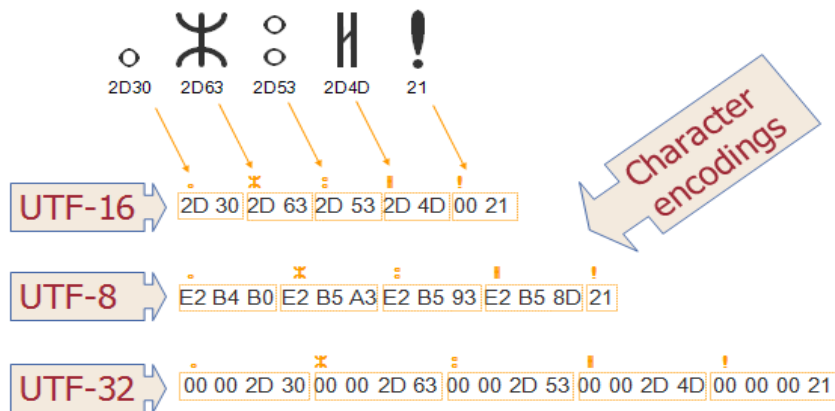- 1 Kilo Byte        $= 10^{3}$ Bytes = 1000 bytes

# Data Sizes

| Data Size | Bytes | Examples |
|-----------|-------|----------|
| Kilobyte (KB) | $10^3$ | Emails, Browser cookies |
| Megabyte (MB) | $10^6$ | Music (mp3), docs, Floppy disk |
| Gigabyte (GB) | $10^9$ | Movies, large docs, Flash drives, DVD |
| Terabyte (TB) | $10^{12}$ | Movies collection, Hard-drive |
| Petabyte (PB) | $10^{15}$ | Networking websites, Financial, Retail |
| Exabyte (EB) | $10^{18}$ | Large search and video websites |
| Zettabyte (ZB) | $10^{21}$ | All data on the internet |

MEGABYTES

1000 MEGABYTES = 1GB

GIGABYTES

1000 GIGABYTES = 1TB

TERABYTES

1000 TERABYTES = 1PB

PETABYTES

1000 PETABYTES = 1EB

EXABYTES

1000 EXABYTES = 1 ZETTABYTE

1ZB

https://www.engadget.com/

# Character Sets and Encodings

- Character Set is a collection of supported characters
- Encoding maps a string of characters to a string of bytes



Source: w3c.org

| Data | Encoding | Bits |
|------|----------|------|
| A | ASCII /UTF-8 | 01000001 |
| A | UTF-16 | 00000000 01000001 |
| a | ASCII /UTF-8 | 01100001 |
| 7 | ASCII /UTF-8 | 00000111 |
| bits | UTF-8 | 01100010 01101001 01110100 01110011 |
| あ | UTF-8 | 11100011 10000001 10000010 00001010 |
| ऋ | UTF-8 | 11100000 10100100 10001011 00001010 |

# Data Format – Text files

- Field delimiters separate data fields. Record delimiters separate groups of fields.

- Examples
  - Comma Separated Values ( , )
  - Pipe Separated Values ( | )
  - Tab Separated Values ( tab )
  - Semicolon separated values ( ; )

```
fname,lname,age,salary
nancy,davolio,33,$30000
erin,borakova,28,$25250
tony,raphael,35,$28700
```

CSV file

https://en.wikipedia.org/wiki/Delimiter/

# Data Format - JavaScript Object Notation (JSON)

- Lightweight data-interchange format

- Platform independent

- Built on two structures :
  - Collection of name/value pairs
  - An ordered list of values

```
{
  "message": {
  "date": "07-04-2016",
  "hour": "08:30",
  "to": "Receiver",
  "from": "Sender",
  "body": "This is an email message"
}}
```
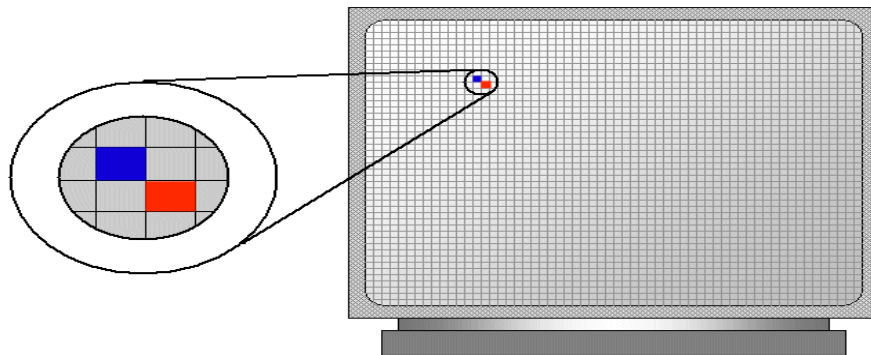
# Data Format - Extensible Markup Language (XML)

- Lightweight data-interchange format

- Platform independent

- Self-descriptive

- Extensible

```xml
<note>
    <date>07-04-2016</date>
    <hour>08:30</hour>
    <to>Receiver</to>
    <from>Sender</from>
    <body>This is an email message</body>
</note>
```
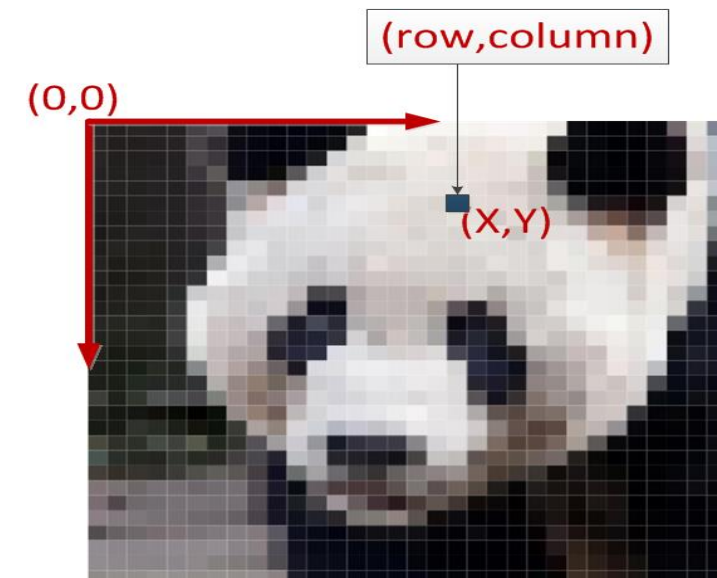
# Digital Images

- Raster Graphics
  - Information stored as pixels (e.g. JPEG)
- Vector Graphics
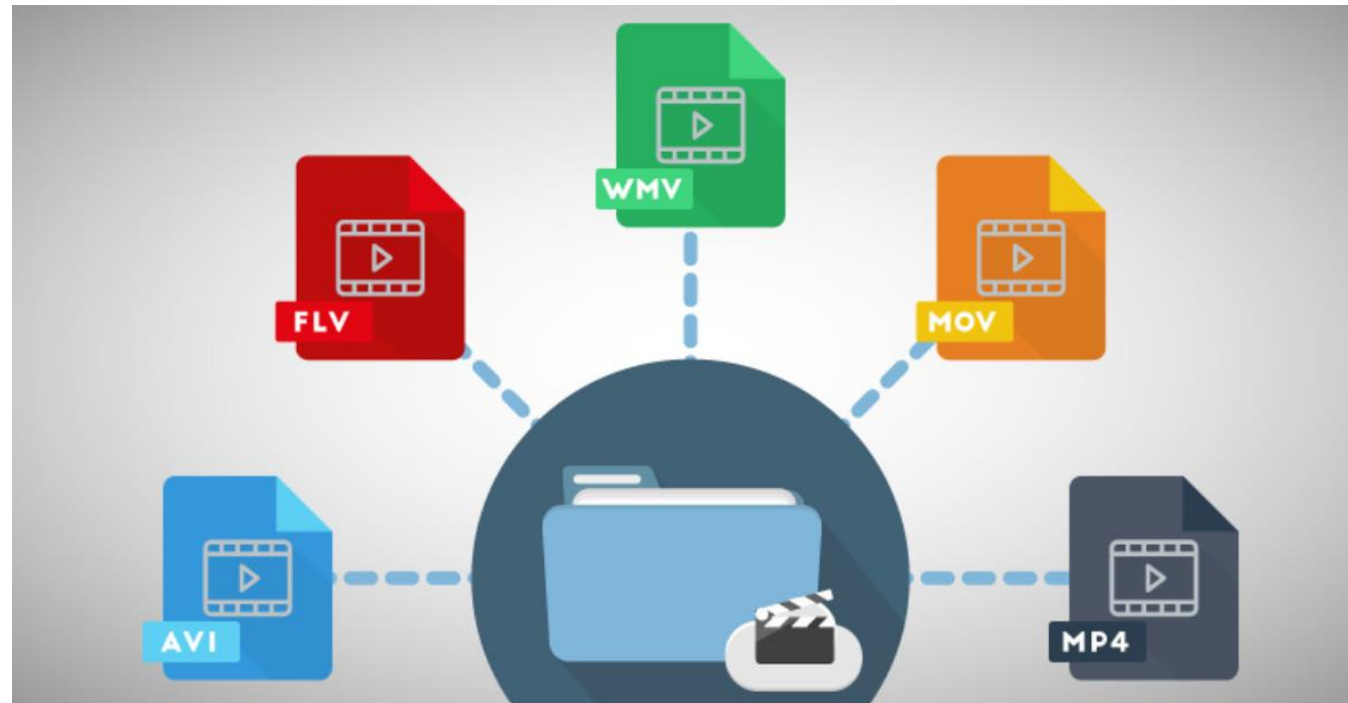  - Information stored as vectors (e.g. SVG)

| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|
| 1 | 1 | 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 | 0 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 |
| 1 | 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 1 |

(row,column)

(0,0)

(X,Y)

# Video Formats

- File format for storing digital video data on a computer system. Video is almost always stored in compressed form to reduce the file size.
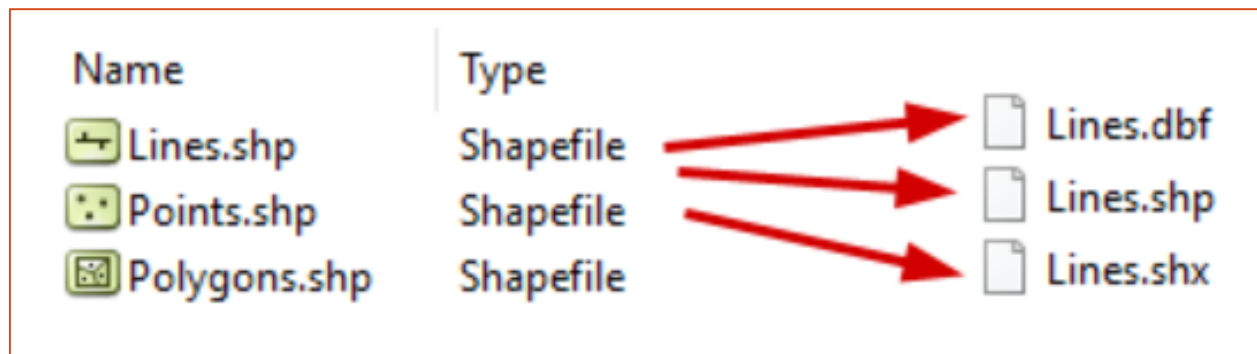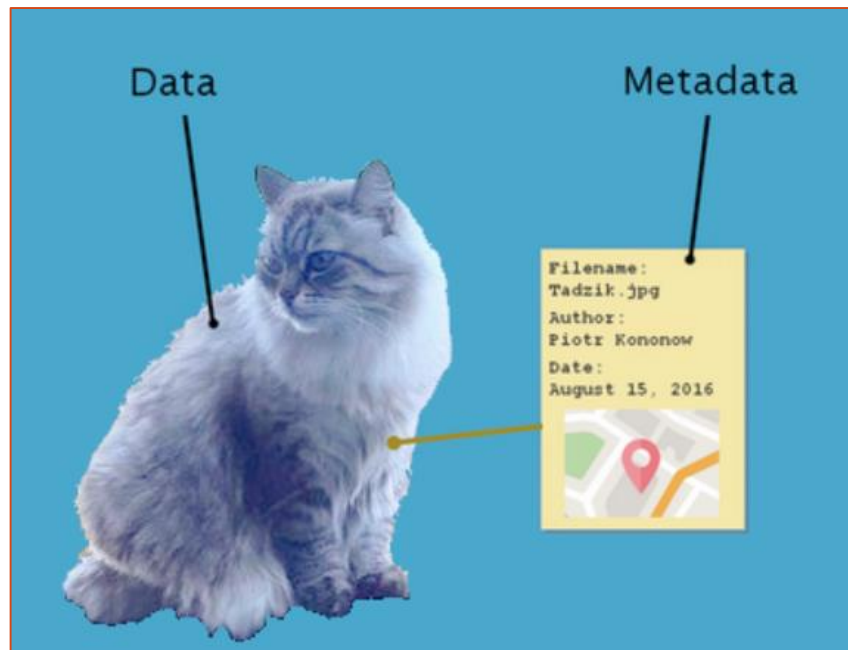
# Spatial data

- Spatial data (geospatial data ) represents the location, size and shape of an object such as a building, lake, mountain and includes attributes that provide more information about the entity being represented.

- The shapefile is the most common geospatial file type. The three required files are SHP is the feature geometry, SHX is the shape index position and DBF is the attribute data.



https://gisgeography.com/gis-formats/

# Metadata

- Data that describes the structure and the properties of the data. It is essential for the proper understanding and use of the data

- Makes finding and working with particular instances of data easier



https://gisgeography.com/gis-formats/

- Data name
- Data create date
- Creator's name
- Data owner
- Data sensitivity
- Group/user permissions
- Source of data
- Construction process of the data

# Data Dictionary

- A "centralized repository of information about data such as meaning, relationships to other data, origin, usage, and format". Also known as metadata repository.

**5.1.1 The actor Table**

The `actor` table lists information for all actors.

The `actor` table is joined to the `film` table by means of the `film_actor` table.

**Columns**

- `actor_id`: A surrogate primary key used to uniquely identify each actor in the table.

- `first_name`: The actor's first name.

- `last_name`: The actor's last name.

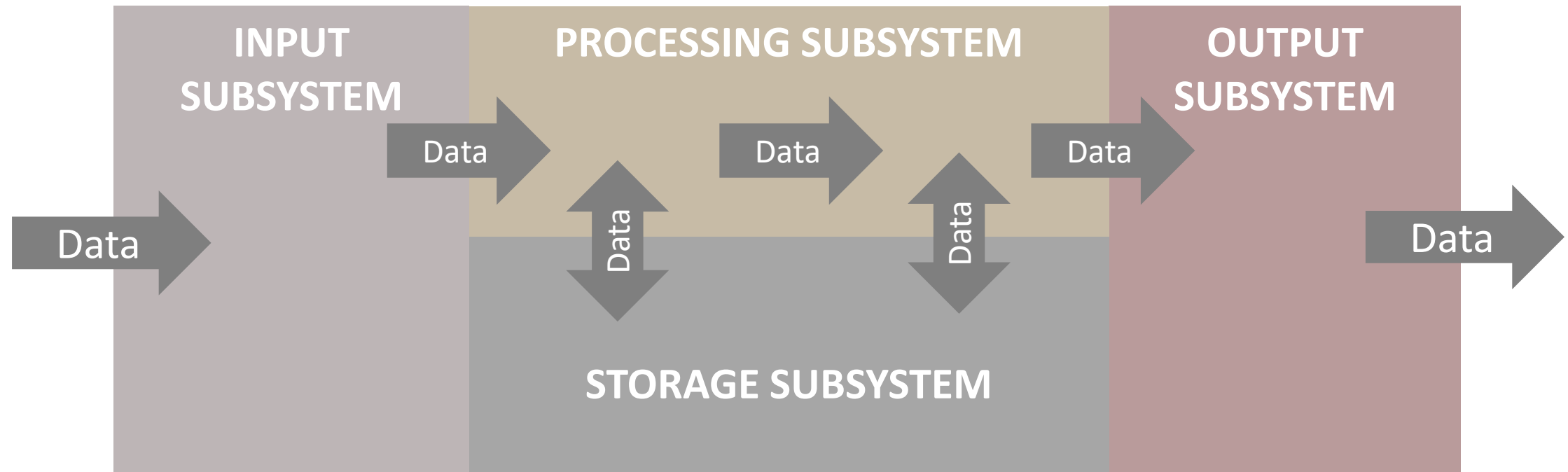- `last_update`: The time that the row was created or most recently updated.

**5.1.2 The address Table**

The `address` table contains address information for customers, staff, and stores.

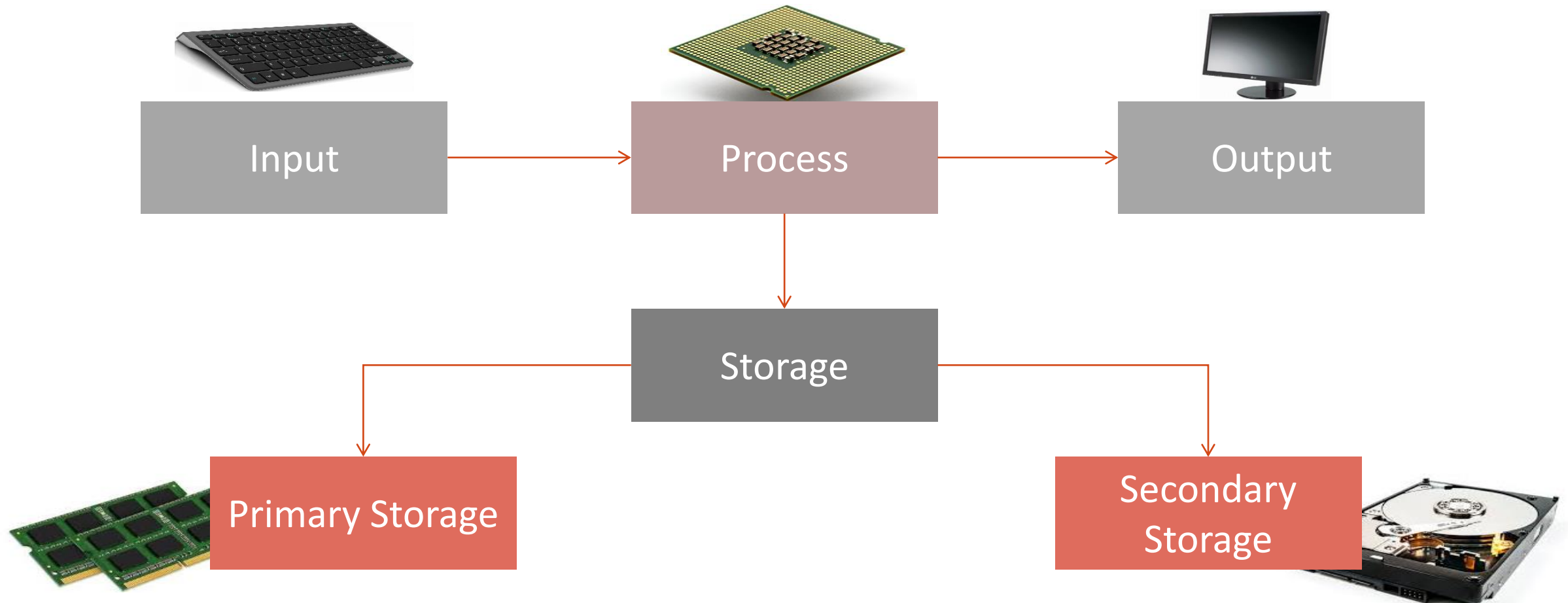The `address` table primary key appears as a foreign key in the `customer`, `staff`, and `store` tables.

# Data Preparation

# Data Processing



Source: wikipedia

# Data Processing



Input → Process → Output

Process → Storage

Storage → Primary Storage

Storage → Secondary Storage

# Application Programming Interfaces (APIs)

- An API is code that allows two software programs to communicate with each other

- APIs are software components or building blocks, which are used by other systems. Usually contain clearly defined methods of communication

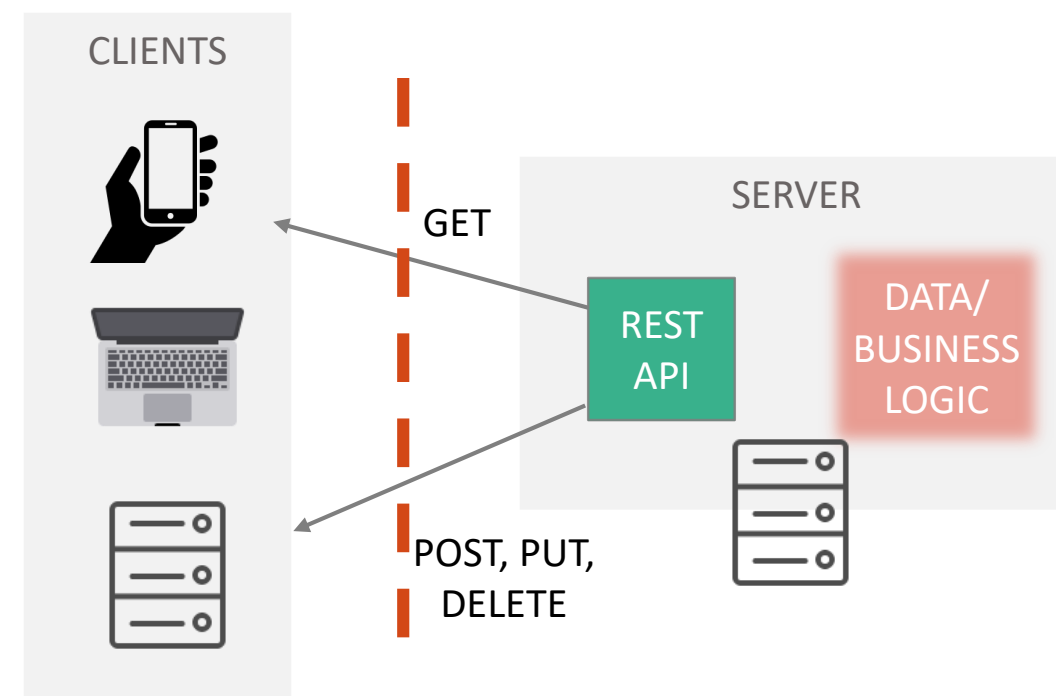- An API may be for a web-based system, operating system, database system, or hardware/software library
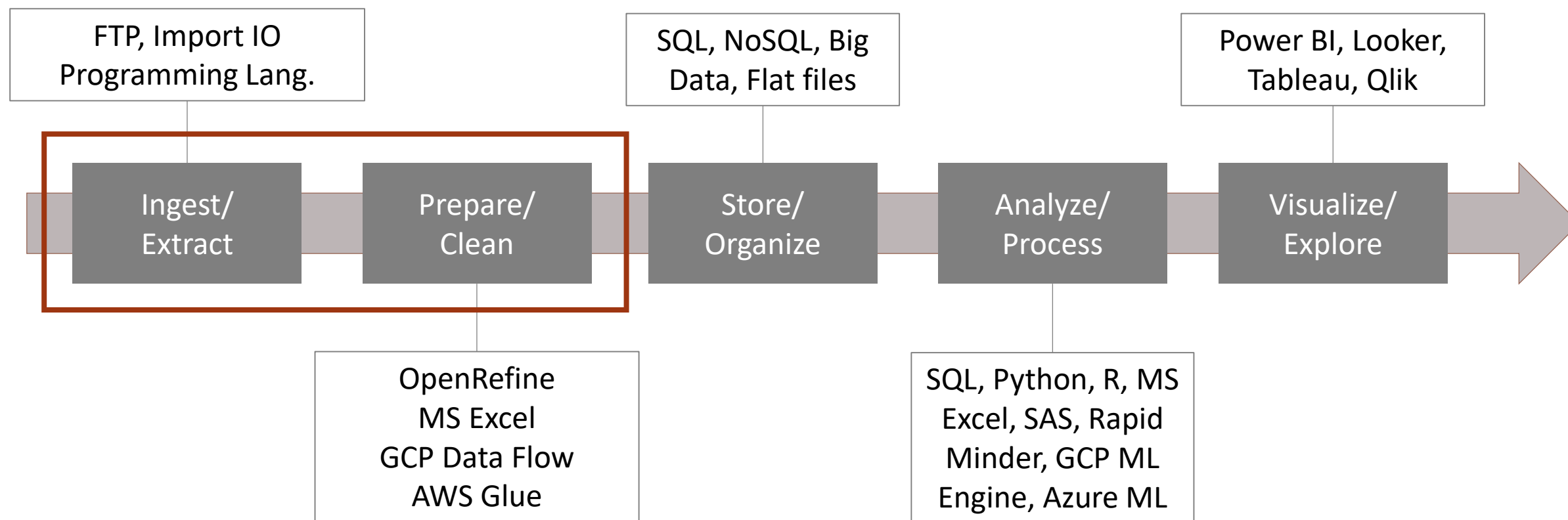
LTX Solutions

# Web Services

- Web Services are services or APIs offered over the internet

- REST/SOAP are widely used web services protocols

- Data exchanged between client and the server is typically in the JSON/XML format

- E.g. Weather Service, Stock Quote Service, etc.
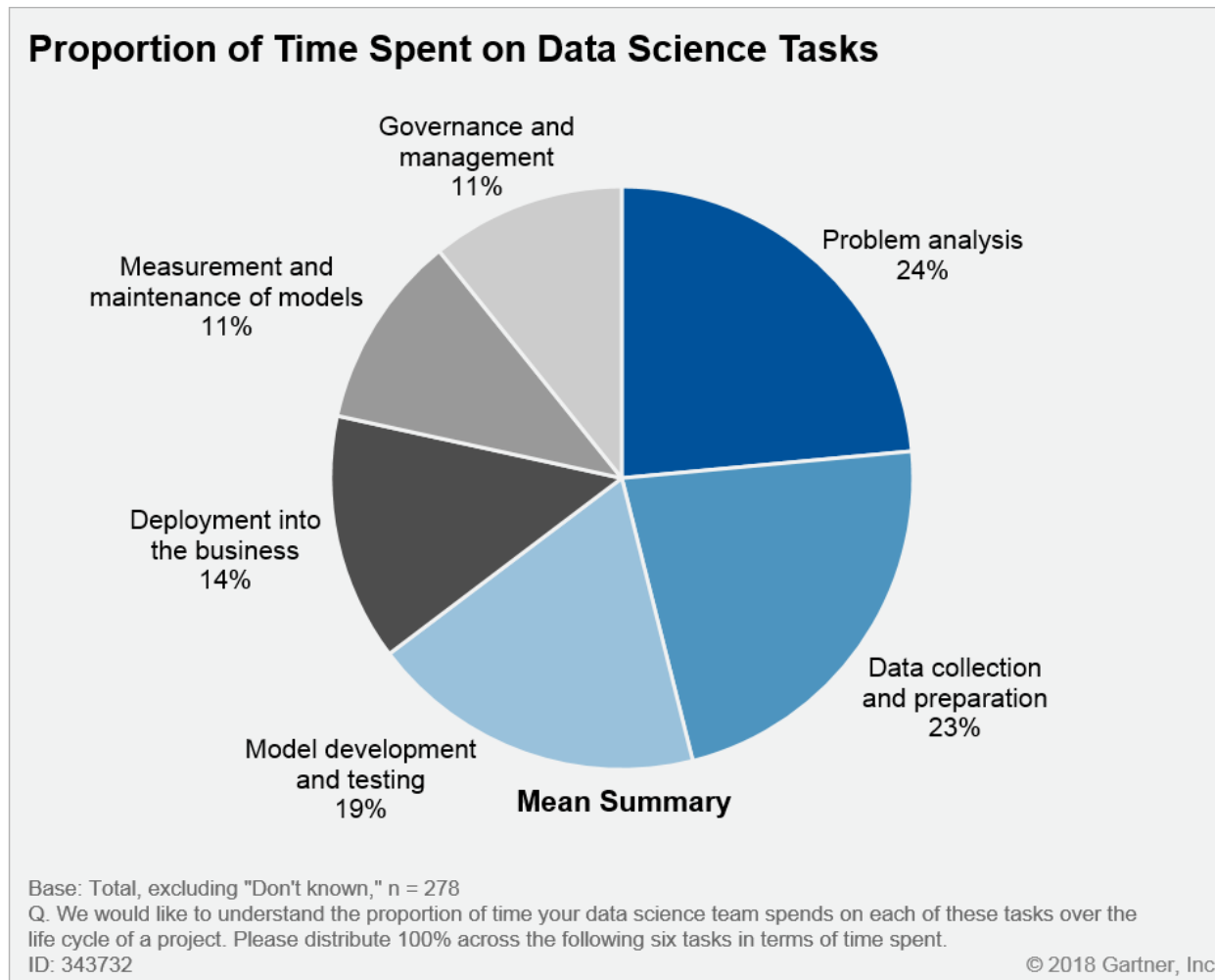
- A curated list of web services https://www.programmableweb.com/category/all/apis

# Data Analysis Pipeline – Ingest, Prepare

| FTP, Import IO Programming Lang. | | SQL, NoSQL, Big Data, Flat files | | Power BI, Looker, Tableau, Qlik |
|---|---|---|---|---|
| **Ingest/ Extract** | **Prepare/ Clean** | **Store/ Organize** | **Analyze/ Process** | **Visualize/ Explore** |
| | OpenRefine MS Excel GCP Data Flow AWS Glue | | SQL, Python, R, MS Excel, SAS, Rapid Minder, GCP ML Engine, Azure ML | |

# Data Science Tasks



**Proportion of Time Spent on Data Science Tasks**

Governance and management — 11%
Problem analysis — 24%
Measurement and maintenance of models — 11%
Deployment into the business — 14%
Data collection and preparation — 23%
Model development and testing — 19%

**Mean Summary**

Base: Total, excluding "Don't known," n = 278
Q. We would like to understand the proportion of time your data science team spends on each of these tasks over the life cycle of a project. Please distribute 100% across the following six tasks in terms of time spent.
ID: 343732

© 2018 Gartner, Inc.

# Data Ingestion

- Process of retrieving data from various sources for further processing
- Various methods of performing data ingestion:
  - File transfers (FTP, Downloads)
  - Web service (REST API)
  - Event based messaging (Publish-Subscribe, Queues, Topics)
  - Extraction Transformation Load (ETL)
  - Web Mining - scraping/crawling

# Exercise – Web Scraping

## Read and process data from a website

# Raw Data

- Can we interpret the data ?

- Does data have errors ?
  - Typos, multiple formats, inconsistencies, missing values

- Does it meet domain standards ?
  - Currency, decimals, censored fields, etc.

> 112, ND748579, 2014-08-20, 33347.88
> 114, GG31455, 2016-05-20, 45864.03
> 114, MA765515, 2014-12-15, 82261.22

# Sources of Errors

- Manual data entry

- Duplicate data entry

- Measurement related errors

- Absence of well defined standards

- Inconsistent data formatting

- Numeric approximations
  - software and hardware constraints



DOW 9,869.62
▼ 998.50 / 9.2%

At 2:45 pm on May 6, 2010, Wall Street essentially had a heart attack. In just minutes, the stock market plunged 1000 points, for reasons traders, analysts, and business media could not explain. The "flash crash" wiped out $1.1 Trillion of investor dollars and even though most of that was quickly regained, it left the market badly shaken.

..It appears that a single keystroke error was to blame. The letter "B" was inserted in a sell order instead of the letter "M". Billion was input where Million should have been and it triggered a ripple effect through the automated financial markets.
https://ungerboeck.com

# Data Quality

- Data Quality
  - Fitness of the data for its intended use in operations, decision making and planning

- Impact of data quality :
  - Consistency - all copies of the data are consistent with each other
  - Completeness – required fields are not missing or incomplete
  - Accuracy - data is correct and has been verified
  - Validity - data adheres to types, formats and business rules
  - Timeliness - data is not out of date
  - Integrity - data is appropriately referenced

# Data Transformation

- Trim white spaces

- Handling missing or null values

- Formatting text, numeric and date fields

- Standardization and consistency

- Correcting domain values
  - Selecting certain attributes
  - Sorting and aggregation
  - Deriving calculated values
  - Statistical adjustments

# OpenRefine

- Open source data cleaning application

- Compatible with a variety of formats
  - CSV, TSV, JSON, XML, Excel, RDF

- More powerful than spreadsheets; better visualization than programming

- Functions:
  - Import/Export, Faceting, Transforming, Clustering, Reconciling

# Exercise - OpenRefine

## Data preparation using OpenRefine

# Database Management Systems

# Data Analysis Pipeline – Store, Analyze

FTP, Import IO
Programming Lang.

SQL, NoSQL, Big
Data, Flat files

Power BI, Looker,
Tableau, Qlik

| Ingest/Extract | Prepare/Clean | Store/Organize | Analyze/Process | Visualize/Explore |
|---|---|---|---|---|

OpenRefine
MS Excel
GCP Data Flow
AWS Glue

SQL, Python, R, MS
Excel, SAS, Rapid
Minder, GCP ML
Engine, Azure ML

# Database

- Shared collection of logically related data and descriptions of this data
- Modeled after real-world systems and entities



Student Database

# History of Databases

| Information Mgmt. Sys(IBM) | Relational Model (E.F.Codd) | DB2, Oracle, INGRES, RDb, etc. | Datawarehouse | No-SQL Databases |

**1960    1970    1980    1990    2000    2010**

| Flat-file databases | Network Model (CODASYL) | Relational Databases - System R (IBM) | Object Oriented Models | Cloud Databases |

# Principle of Data Independence

Logical Database Schema

- A logical schema is a conceptual model of the data. It is primarily concerned with understanding the business entities, their attributes and their relationships
- Logical schema is design-centric in order to meet the business requirements

Physical Database Schema

- A physical model is concrete in the sense that it is implemented in the DBMS
- It contains physical objects such as data files, redo logs, control files etc., that reside on a database system

The ability to modify a schema definition in one level without affecting a scheme definition in a higher level is called **data independence.** The principle of data independence helps keep the logical model separate from the physical schema so that the database model/design can be isolated from the physical constraints

# Logical Architecture (ANSI-SPARC)



ANSI-SPARC Architecture

# Database Management System (DBMS)

- Aggregate of data, software, hardware and users that helps an organization manage its operational data

- Provide efficient and reliable methods of data retrieval

- Involves  monitoring, administration, and maintenance of the databases

# DBMS Benefits

- Program-Data Independence

- Efficiently retrieve, manipulate, store data

- Concurrent Access & Crash Recovery

- Minimize data duplication

- Data Integrity & Security

- Data Administration

# DBMS Components



Administrators
Designers, Users

Query Processor
Memory Mgr.
File Systems

DDL, DCL, DML

Authentication,
Install, Backup

Internal memory
External memory

# DBMS Software Architecture

Obtains results from the database that satisfies queries compiled by the QP

**Memory Manager**

**Query Processor**

User Queries

Converts user queries into instructions the DBMS can process efficiently

**Transaction Manager**

Ensures transactions satisfy ACID
Provides recovery from failures

Storage

# DBMS Users



Administrators

Maintenance
Security
Optimization

Designers

Entity Modeling
Data Formatting
Database Content

**DBMS**

Applications and
Business Users

# Client-Server Architecture

# 2-tier Architecture (Logical)

# 3-tier Architecture (Logical)



**Presentation Tier**

**Application Tier**

**Data Tier**

GUI/Frontend

Business Logic

Databases/Backend

1

2

3

# Database Transactions

- A unit of work performed against a database

- Generally results in a change of database state

- Reliable and independent of other transactions

- Either executes in entirety or is rolled back

- ACID properties

# ACID Properties

- Atomicity
  - When an update occurs to a database, if part of the transaction fails, then the entire transaction fails. Atomicity requires that each transaction be an "all or nothing".

- Consistency
  - Ensures that any transaction beings the database from one valid state to another. A consistency constraint is a predicate on data which serves as a precondition, post-condition, and transformation condition on any transaction

- Isolation
  - Ensures that concurrent execution of transactions results in a system state that would be obtained if the transactions were executed serially, i.e. one after the other

- Durability
  - Ensures that the system is able to recover committed transaction updates if either the system crashes or fails with errors

# Transaction Example



ATM

Withdrawal request

Account Balance

1) Verify account details

$ 1000

2) Accept withdrawal request

$ 1000

3) Check balance

$ 1000

Crash

4) Update balance

$ 900

5) Dispense money

# Transaction Lifecycle

# Exercise

## MySQL DBMS Review

# Types of Databases

Users, Location, Function, Supported Data Types

# Database Classification

- Databases are usually classified based on the following attributes:
  - Number of users
    - Single user, multi-user, large scale
  - Primary function
    - Transactional vs analytical
  - Database location(s)
    - Centralized vs distributed
  - Data type and structure
    - Structured, semi-structured

# Database Types - Users

Databases

- Single-user
- Multi-user
  - Workgroup
  - Enterprise

# Database Types - Function



Databases

Operational
(OLTP)

Analytical
(OLAP)

# Database Types - Location(s)

```
Databases
    ├── Centralized
    └── Distributed
            ├── Partitioned
            ├── Replicated
            └── Partitioned Replicated
```

# Distributed Databases



Partitioned
Non Replicated

Non Partitioned
Replicated

Partitioned
Replicated

# Distributed Databases

Partitioned Non-Replicated

increasing flexibility,
parallelism, availability

increasing cost,
complexity

Non-Partitioned Replicated

Partitioned Replicated

# CAP Theorem

Clients get the same view of data

All requests will receive a response

Consistency     CA     Availability

CP     AP

Partition Tolerance

System continues to work despite message loss or partial failure

*It is impossible for a distributed computer system to simultaneously provide all three of the following guarantees: Consistency, Availability Partition tolerance* - Eric Brewer

# Database Types - Structure

# Relational Databases

- Composed of tables with columns and rows

- Data stored across multiple related tables

- Support SQL like query languages and SQL standards

Columns

Rows

| Student Id | Name | Course | Grade |
|------------|------|--------|-------|
| 1001 | Ash | Algebra | A |
| 1002 | Jeff | Physics | B |
| 1003 | Judy | English | A |
| 1004 | Ram | Spanish | C |

# No-SQL or Non-Relational Databases

- Store semi-structured data
- Flexible schema
- Increased scalability
- Eventual consistency



SQL Database | Non-SQL Database

Relational

Key-Value

Column-Family

Analiticals (OLAP)

Graph

Document

https://www.netsolutions.com/

# Key-Value Store

- Stores a Dictionary (or Hash) data structure
- Each value is associated with unique key

| Key | Value |
|---|---|
| 1001 | Ash , Algebra, A |
| 1002 | Jeff , Physics, B |
| 1003 | Judy , English, A |
| 1004 | Ram , Spanish, C |

Key → Value

# Document Store

- Data stored as documents

Documents

Student_1001
{ id: "1001",
  name: "Ash",
  course: "Algebra",
  grade:  "A"}

Student_1002
{ id: "1002",
  name: "Jeff",
  course: "Physics",
  grade:  "B"}

Student_1003
{ id: "1003",
  name: "Judy",
  course: "English",
  grade:  "A"}

Student_1004
{ id: "1004",
  name: "Ram",
  course: "Spanish",
  grade:  "C"}

# Column Store

- Entities stored in columns rather than rows
- Column oriented databases

| Column1 | Column2 | Column3 | Column4 |
|---------|---------|---------|---------|
| 1001 | 1002 | 1003 | 1004 |
| Ash | Jeff | Judy | Ram |
| Algebra | Physics | English | Spanish |
| A | B | A | C |

# Graph Databases

- Entities stored as graphs - Nodes, Edges, Properties



Source: Wikipedia

# Popular Databases

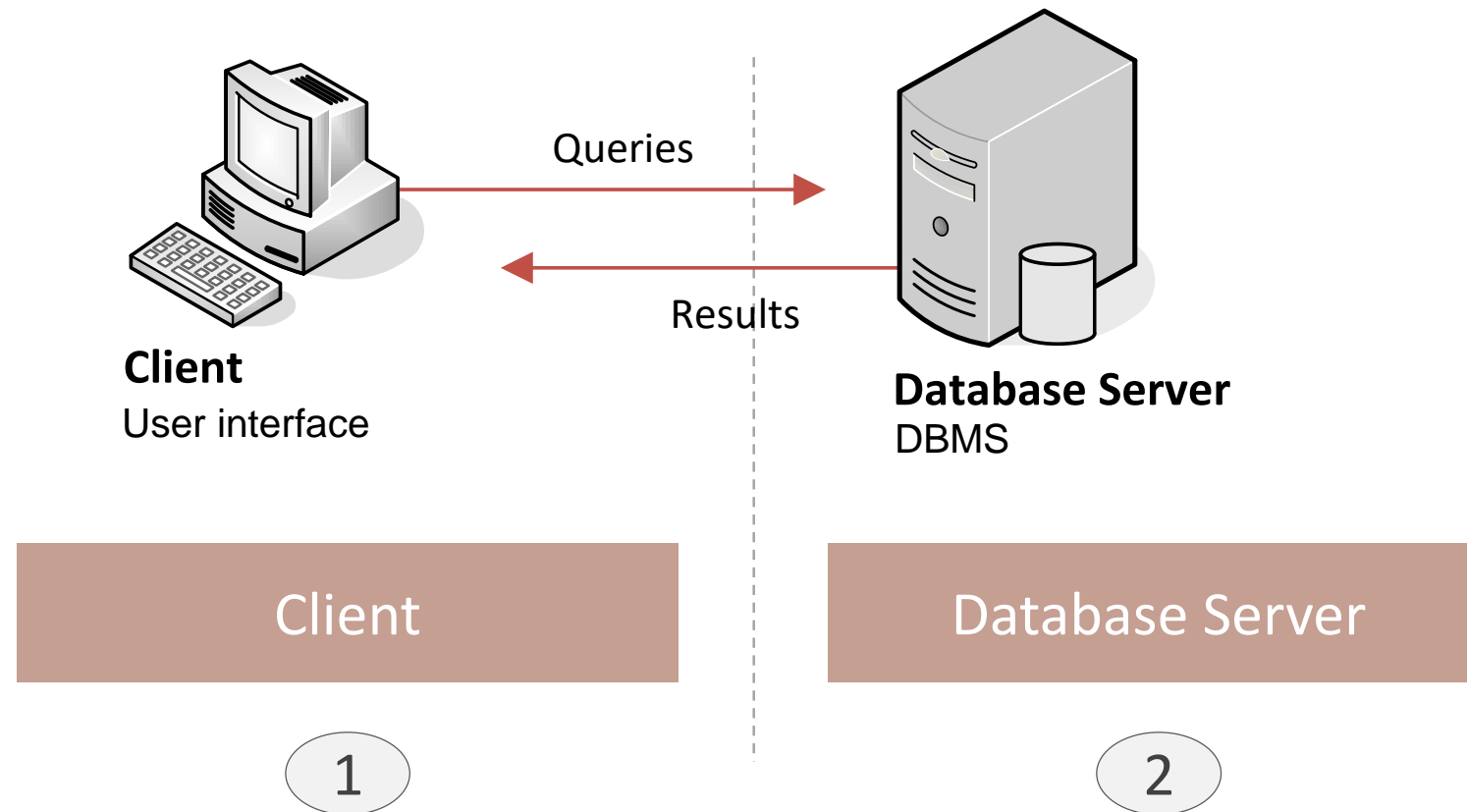| Database Type | Database Names | | | |
|---|---|---|---|---|
| Relational | MySQL | ORACLE | PostgreSQL | IBM DB2 |
| Key-Value | redis | riak | MEMCACHED | |
| Column | APACHE HBASE | Cassandra | Google BigQuery | |
| Document | mongoDB | CouchDB relax | amazon DynamoDB | |
| Graph | neo4j | TITAN | OrientDB | INFINITEGRAPH |

# Appendix

# Evolution of storage



Source: bbc.co.uk

# 2-tier Architecture (Physical)

# 3-tier Architecture (Physical)



User request

Queries

Response

Results

**Client**
User interface

**Application Server**
Business components

**Database Server**
DBMS

Client

Application Server

Database Server

1

2

3