

Introduction to Data Science, Data Analytic and Big Data

Tissana Tanaklang

Software and Solution Development Trainer
Iverson Training Center Co., Ltd.

tissana_t@hotmail.com

- Master of Science Program in Software Engineering King Mongkut's University of Technology Thonburi
- Bachelor of Science Program in Computer Science Naresuan University
- Microsoft Certified Trainer (MCT)
- Microsoft Certified Solutions Associate (MCSA) - Web Application Development
- Microsoft Certified Azure Fundamentals
- Microsoft Certified Azure Data Fundamentals



01 – Introduction to Data



02 – Data Science for Business



03 – Big Data Fundamentals



04 – Big Data Sources



05 – Introduction to Data Analytics



06 – Understanding Machine Learning

01 - Introduction to Data

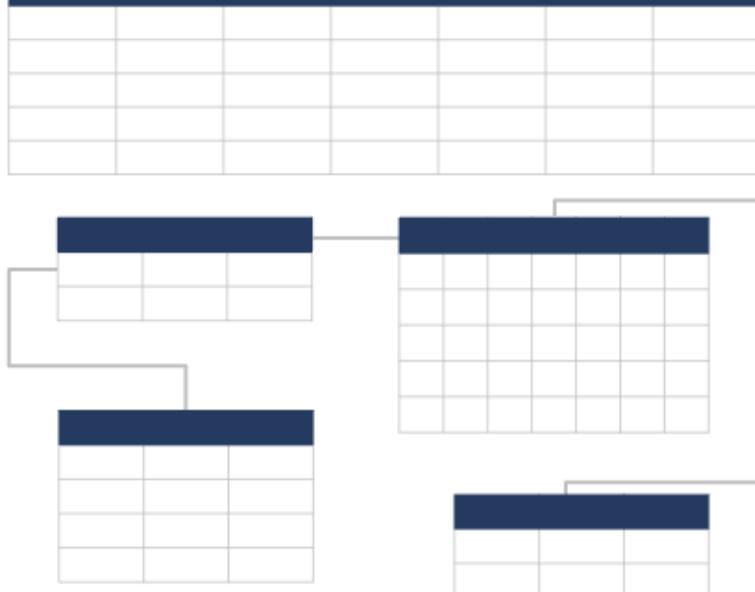


Data is a collection of facts such as numbers, descriptions, and observations used in decision making.

What is data?

Structured

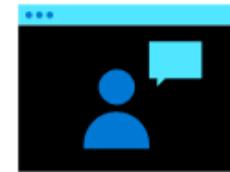
Table



Semi-structured

```
## Document 1 ## {
"customerID": "103248",
"name": { "first": "AAA",
"last": "BBB" }, "address": {
"street": "Main Street",
"number": "101", "city": "Acity",
"state": "NY" },
"ccOnFile": "yes",
"firstOrder": "02/28/2003" }
## Document 2 ## {
"customerID": "103249",
"name": { "title": "Mr",
"forename": "AAA",
"lastname": "BBB" },
"address": { "street": "Another Street",
"number": "202", "city": "Bcity",
"county": "Gloucestershire",
"country-region": "UK" },
"ccOnFile": "yes" }
```

Unstructured



Data Categories

- A transactional system is often what most people consider the primary function of business computing.
- **A transactional system records transactions.**
- A transaction could be financial, such as the movement of money between accounts in a banking system, or it might be part of a retail system, tracking payments for goods and services from customers.

What is a Transactional System?

Customer		
CustomerID	CustomerName	CustomerPhone

Account	
CustomerID	Balance
5558	1000
6023	1500

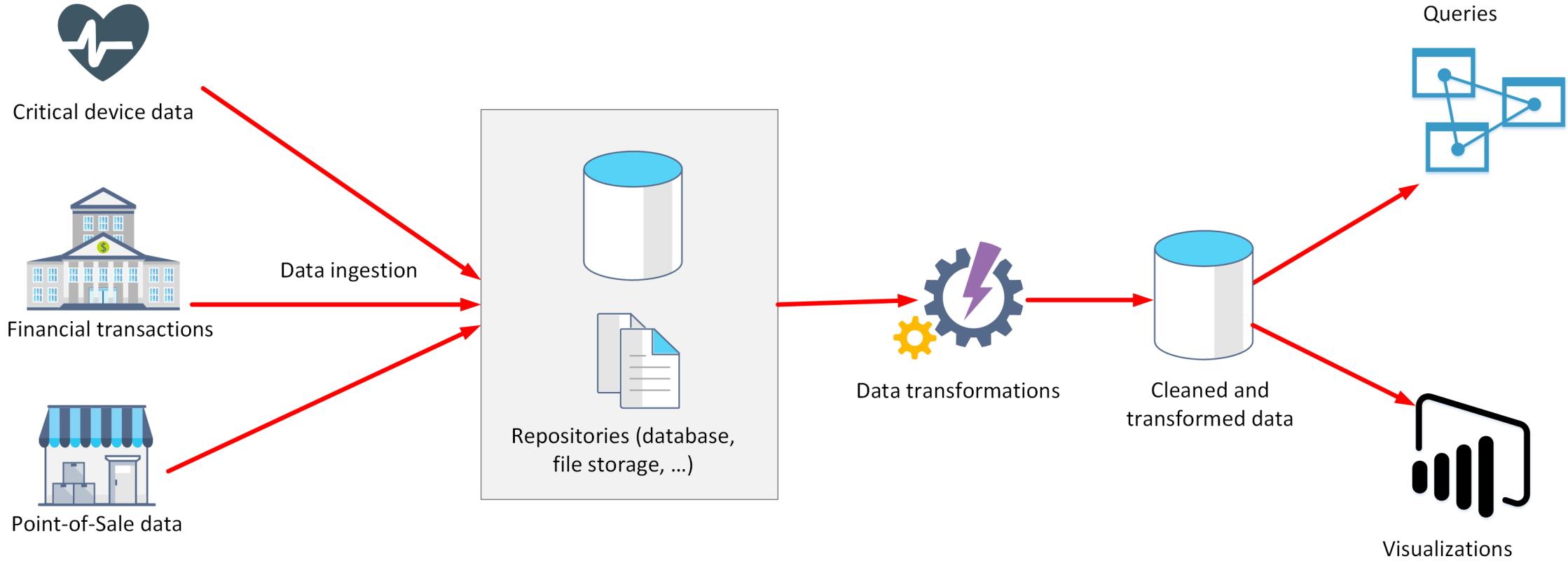
Orders		
OrderID	CustomerID	OrderDate

Transfers					
TransactionID	FromAccount	ToAccount	Transaction Amount	OrderDate	TransactionDescription
982801	6023	5558	500	DD/MM/YY	Transfer 500 from account 6023 to account

```

BEGIN TRANSACTION
UPDATE Account
SET Balance = Balance -500
WHERE CustomerID=6023;
UPDATE Account
SET Balance = Balance +500
WHERE CustomerID=5558;
INSERT INTO Transfers (Fromaccount, ToAccount, TransactionAmount, TransactionDescription)
VALUES (6023,5558,500,'Transfer 500 from account 6023 to account 5558')
COMMIT TRANSACTION
  
```

Transactional workloads



What is an Analytical System?



On-premises data

SQL Server, Oracle,
fileshares, SAP



Cloud data

Azure, AWS, GCP



SaaS data

Salesforce, Dynamics

Data ingestion



Data storage



Data processing



Data visualization



What is an Analytical System?

Batch

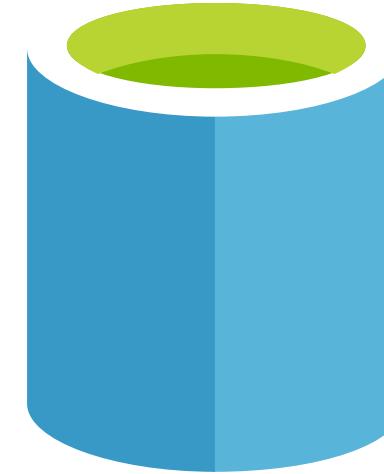


Streaming

101010 101010 101010
010101 010101 010101
101010 101010 101010

101010 101010 101010
010101 010101 010101
101010 101010 101010

101010 101010 101010
010101 010101 010101
101010 101010 101010



Batch Data / Streaming Data

CustomerID	FirstName	LastName	CustomerID	AddressID	AddressID	LineNumber	Text
1	Jay	Adams	1	A	A	1	12
2	Donna	Carreras	2	B	A	2	Park Street
3	Linda	Burnett	3	C	A	3	Some City
4	Frances	Adams	4	A	B	1	The Big House
					B	2	High Road
					B	3	Another City
					B	4	90210
					C	1	Freepost
					C	2	AAA 123

Relational Data

- All data is tabular. Entities are modeled as tables, each instance of an entity is a row in the table, and each property is defined as a column.
- All rows in the same table have the same set of columns.
- A table can contain any number of rows.

Characteristics of a Relational Database

- A primary key uniquely identifies each row in a table. No two rows can share the same primary key.
- A foreign key references rows in another, related table. For each value in the foreign key column, there should be a row with the same value in the corresponding primary key column in the other table.

Characteristics of a Relational Database

- SQL is a standard language for use with relational databases
- SQL standards are maintained by ANSI and ISO
- Proprietary RDBMS systems have their own extensions of SQL such as T-SQL, PL/SQL, pgSQL

Query - SQL

DML	DDL	DCL
<p>Data Manipulation Language</p> <p>Used to query and manipulate data</p> <p>SELECT, INSERT, UPDATE, DELETE</p>	<p>Data Definition Language</p> <p>Used to define database objects</p> <p>CREATE, ALTER, DROP</p>	<p>Data Control Language</p> <p>Used to manage security permissions</p> <p>GRANT, REVOKE, DENY</p>

Query - SQL

Statement	Description
SELECT	Select/read from a table
INSERT	Insert new rows in a table
UPDATE	Edit/Update existing rows in a table
DELETE	Delete existing rows in a table

Query - SQL

```
SELECT EmployeeId, YEAR(OrderDate) AS OrderYear  
FROM Sales.Orders  
WHERE CustomerId = 71  
GROUP BY EmployeeId, YEAR(OrderDate)  
HAVING COUNT(*) > 1  
ORDER BY EmployeeId, OrderYear;
```

Query - SQL

Statement	Description
CREATE	Create a new object in the database, such as a table or a view
INSERT	Modify the structure of an object. For instance, altering a table to add a new column
UPDATE	Remove an object from the database
DELETE	Rename an existing object

Query - SQL

```
CREATE TABLE Mytable  
(Mycolumn1 int NOT NULL PRIMARY KEY, Mycolumn2 VARCHAR(50) NOT  
NULL , Mycolumn2 VARCHAR(10) NOT NULL)
```

Query - SQL

Customers

CustomerID	CustomerName	CustomerPhone
100	Muisto Linna	XXX-XXX-XXXX
101	Noam Maoz	XXX-XXX-XXXX
102	Vanja Matkovic	XXX-XXX-XXXX
103	Qamar Mounir	XXX-XXX-XXXX
104	Zhenis Omar	XXX-XXX-XXXX
105	Claude Paulet	XXX-XXX-XXXX
106	Alex Pettersen	XXX-XXX-XXXX
107	Francis Ribeiro	XXX-XXX-XXXX

Tables

Customers

CustomerID	CustomerName	CustomerPhone
100	Muisto Linna	XXX-XXX-XXXX
101	Noam Maoz	XXX-XXX-XXXX
102	Vanja Matkovic	XXX-XXX-XXXX
103	Qamar Mounir	XXX-XXX-XXXX
104	Zhenis Omar	XXX-XXX-XXXX
105	Claude Paulet	XXX-XXX-XXXX
106	Alex Pettersen	XXX-XXX-XXXX

Orders

OrderID	CustomerName	CustomerPhone
AD100	Noam Maoz	XXX-XXX-XXXX
AD101	Noam Maoz	XXX-XXX-XXXX
AD102	Noam Maoz	XXX-XXX-XXXX
AX103	Qamar Mounir	XXX-XXX-XXXX
AS104	Qamar Mounir	XXX-XXX-XXXX
AR105	Claude Paulet	XXX-XXX-XXXX
MK106	Muisto Linna	XXX-XXX-XXXX

Data is normalized to:

Reduce storage

Avoid data duplication

Improve data quality

Normalization

Customers

CustomerID	CustomerName	CustomerPhone
100	Muisto Linna	XXX-XXX-XXXX
101	Noam Maoz	XXX-XXX-XXXX
102	Vanja Matkovic	XXX-XXX-XXXX
103	Qamar Mounir	XXX-XXX-XXXX
104	Zhenis Omar	XXX-XXX-XXXX
105	Claude Paulet	XXX-XXX-XXXX
106	Alex Pettersen	XXX-XXX-XXXX

Orders

OrderID	CustomerID	SalesPersonID
AD100	101	200
AD101	101	200
AD102	101	200
AX103	103	201
AS104	103	201
AR105	105	200
MK106	105	201

In a normalized database schema:

Primary Keys and Foreign keys are used to define relationships

No data duplication exists (other than key values in 3rd Normal Form (3NF)

Data is retrieved by joining tables together in a query

Relations

Customers

CustomerID	CustomerName	CustomerPhone
100	Muisto Linna	XXX-XXX-XXXX
101	Noam Maoz	XXX-XXX-XXXX
102	Vanja Matkovic	XXX-XXX-XXXX
103	Qamar Mounir	XXX-XXX-XXXX
104	Zhenis Omar	XXX-XXX-XXXX
105	Claude Paulet	XXX-XXX-XXXX
106	Alex Pettersen	XXX-XXX-XXXX

IDX-CustomerRegion

CustomerID	Region
100	France
101	Brazil
102	Croatia
103	Jordan
104	Spain
105	France
106	USA

An index:

Optimizes search queries for faster data retrieval

Reduces the amount of data pages that need to be read to retrieve the data in a SQL Statement

Data is retrieved by joining tables together in a query

Indexes

Customers

CustomerID	CustomerName	CustomerPhone
100	Muisto Linna	XXX-XXX-XXXX
101	Noam Maoz	XXX-XXX-XXXX
102	Vanja Matkovic	XXX-XXX-XXXX
103	Qamar Mounir	XXX-XXX-XXXX
104	Zhenis Omar	XXX-XXX-XXXX
105	Claude Paulet	XXX-XXX-XXXX
106	Alex Pettersen	XXX-XXX-XXXX

Orders

OrderID	CustomerID	SalesPersonID
AD100	101	200
AD101	101	200
AD102	101	200
AX103	103	201
AS104	103	201
AR105	105	200
MK106	105	201
DB205	100	205

Create the definition of a view:

```
CREATE VIEW  
vw_customerorders AS  
SELECT Customers.CustomerID,  
Customers.CustomerName,  
Orders.OrderID FROM  
Customers JOIN Orders on  
Customers.CustomerID =  
Orders.CustomerID
```

Retrieve the orders placed by customer 102 using the view:

```
SELECT CustomerName, OrderID  
from vw_customerorders WHERE  
CustomerID=102
```

A view is a virtual table based on the result set of query:

Views are created to simplify the query

Combine relational data into a single pane view

View



ORACLE



Relational Database Management System - RDBMS

```
## Document for Jay Adams ##
{
  "customerID": "1",
  "name":
  {
    "firstname": "Jay",
    "lastname": "Adams"
  },
  "address":
  {
    "number": "12",
    "street": "Park Street",
    "city": "Some City",
  }
}
```

```
## Document for Frances Adams ##
{
  "customerID": "4",
  "name":
  {
    "firstname": "Francis",
    "lastname": "Adams"
  },
  "address":
  {
    "number": "12",
    "street": "Park Street",
    "city": "Some City",
  }
}
```

Non-Relational Data

```
## Customer 1 ID: 1
Name: Mark Hanson
Telephone: [ Home: 1-999-9999999, Business: 1-888-8888888, Cell: 1-777- 7777777 ]
Address: [ Home: 121 Main Street, Some City, NY, 10110,
           Business: 87 Big Building, Some City, NY, 10111 ]
## Customer 2 ID: 2
Title: Mr
Name: Jeff Hay
Telephone: [ Home: 0044-1999-333333, Mobile: 0044-17545-444444 ]
Address: [ UK: 86 High Street, Some Town, A County, GL8888, UK,
           US: 777 7th Street, Another City, CA, 90111 ]
```

Non-relational collections can have:

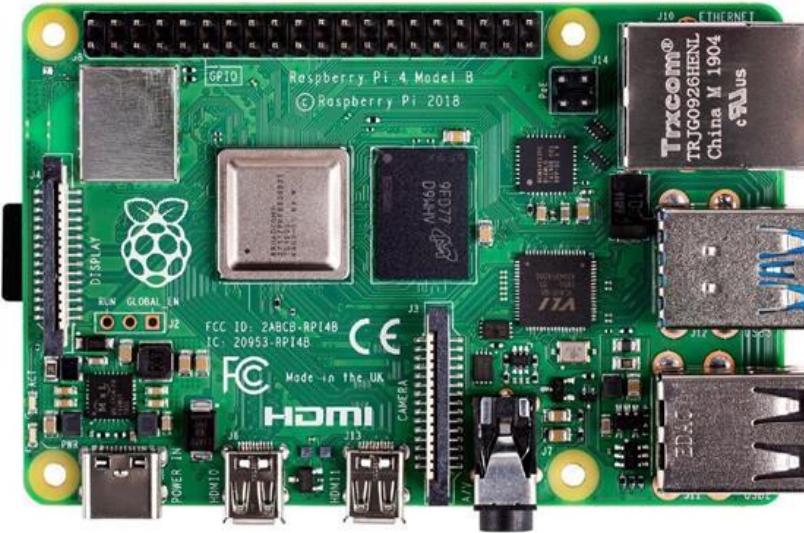
Multiple entities in the same collection or container
with different fields

Have a different,
non-tabular schema

Are often defined by labeling each field with the
name it represents

Non-Relational Data

- IoT and Telematics
- Retail and Marketing
- Gaming
- Web and Mobile



Identify non-relational database use cases

```
{"latitude":37.8267,"longitude":-122.4233,"timezone":"America/Los_Angeles","currently":{"time":1598191217,"summary":"Partly Cloudy","icon":"partly-cloudy-day","nearestStormDistance":5,"nearestStormBearing":58,"precipIntensity":0,"precipProbability":0,"temperature":58.63,"apparentTemperature":58.63,"dewPoint":52.42,"humidity":0.8,"pressure":1011.8,"windSpeed":5.08,"windGust":7.73,"windBearing":210,"cloudCover":0.54,"uvIndex":0,"visibility":9.933,"ozone":291.2},"minutely":{"summary":"Partly cloudy for the hour. ","icon":"partly-cloudy-day","data":[{"time":1598191200,"precipIntensity":0,"precipProbability":0},{"time":1598191260,"precipIntensity":0,"precipProbability":0}, {"time":1598191320,"precipIntensity":0,"precipProbability":0}, {"time":1598191380,"precipIntensity":0,"precipProbability":0}, {"time":1598191440,"precipIntensity":0,"precipProbability":0}, {"time":1598191500,"precipIntensity":0,"precipProbability":0}, {"time":1598191560,"precipIntensity":0,"precipProbability":0}, {"time":1598191620,"precipIntensity":0,"precipProbability":0}, {"time":1598191680,"precipIntensity":0,"precipProbability":0}, {"time":1598191740,"precipIntensity":0,"precipProbability":0}, {"time":1598191800,"precipIntensity":0,"precipProbability":0}, {"time":1598191860,"precipIntensity":0,"precipProbability":0}, {"time":1598191920,"precipIntensity":0,"precipProbability":0}, {"time":1598191980,"precipIntensity":0,"precipProbability":0}, {"time":1598192040,"precipIntensity":0,"precipProbability":0}, {"time":1598192100,"precipIntensity":0,"precipProbability":0}, {"time":1598192160,"precipIntensity":0,"precipProbability":0}, {"time":1598192220,"precipIntensity":0,"precipProbability":0}, {"time":1598192280,"precipIntensity":0,"precipProbability":0}, {"time":1598192340,"precipIntensity":0,"precipProbability":0}, {"time":1598192400,"precipIntensity":0.0026,"precipIntensityError":0.0004,"precipProbability":0.01,"precipType":"rain"}, {"time":1598192460,"precipIntensity":0.0026,"precipIntensityError":0.0005,"precipProbability":0.01,"precipType":"rain"}, {"time":1598192520,"precipIntensity":0,"precipProbability":0}, {"time":1598192580,"precipIntensity":0.0026,"precipIntensityError":0.0005,"precipProbability":0.01,"precipType":"rain"}, {"time":1598192640,"precipIntensity":0.0026,"precipIntensityError":0.0005,"precipProbability":0.01,"precipType":"rain"}, {"time":1598192700,"precipIntensity":0.0027,"precipIntensityError":0.0005,"precipProbability":0.01,"precipType":"rain"}, {"time":1598192760,"precipIntensity":0.0027,"precipIntensityError":0.0005,"precipProbability":0.02,"precipType":"rain"}, {"time":1598192820,"precipIntensity":0.0026,"precipIntensityError":0.0005,"precipProbability":0.02,"precipType":"rain"}, {"time":1598192880,"precipIntensity":0,"precipProbability":0}], "hourly":{},"daily":{},"alerts":[]}}
```

Non-Relational Database use case

<https://api.darksky.net/>

Open API :] สำหรับนักพัฒนา

แสดงค่าประจำวัน :

[//covid19.th-stat.com/api/open/today](https://covid19.th-stat.com/api/open/today)

ข้อมูลสรุปตามช่วงเวลา [เริ่มตั้งแต่วันที่ 01/01/20] :

[//covid19.th-stat.com/api/open/timeline](https://covid19.th-stat.com/api/open/timeline)

ข้อมูลแต่ละเคส :

[//covid19.th-stat.com/api/open/cases](https://covid19.th-stat.com/api/open/cases)

ข้อมูลสรุปจากเคส :

[//covid19.th-stat.com/api/open/cases/sum](https://covid19.th-stat.com/api/open/cases/sum)

แจ้งเตือนพื้นที่ตามคำประกาศ :

[//covid19.th-stat.com/api/open/area](https://covid19.th-stat.com/api/open/area)



กรมควบคุมโรค
DEPARTMENT OF DISEASE CONTROL

Non-Relational Database use case

	Schema	Data relationships	Examples
Structured data	Adheres to a schema, with the same data fields or properties.	Storable in relational database tables, with rows and columns.	Sensor data and financial data.
Semi-structured data	Has an ad hoc schema with less organized fields and properties.	Non-relational or NoSQL data, not storable in tables, rows and column.	Books, blogs, JSON, HTML documents.
Unstructured data	Has no designated schema or data structure.	Non-relational or blob data, with no restrictions on the kinds of data blobs contain.	PDFs, JPGs, videos.

- You might see the term *NoSQL* when reading about non-relational databases.
- NoSQL is a rather loose term that simply means non-relational.
- NoSQL (non-relational) databases generally fall into four categories:
 - key-value stores
 - document databases
 - column family databases
 - graph databases.

What is NoSQL?

A diagram illustrating a key-value store. On the left is a table with four rows. The first row has a light blue header with 'Key' and 'Value' columns. The subsequent three rows have white headers. Each row contains a key in the first column and a binary value in the second column. An arrow points from a callout box labeled 'Opaque to data store' towards the binary value of the first row.

Key	Value
AAAAAA	110100111101010011010111...
AABAB	100110000101100110101110...
DFA766	000000000101010110101010...
FABCC4	1110110110101010100101101...

Opaque to data store

A key-value store is the simplest (and often quickest) type of NoSQL database for inserting and querying data.

Key-Value Stores

Key	Document
1001	{ "CustomerID": 99, "OrderItems": [{ "ProductID": 2010, "Quantity": 2, "Cost": 520 }, { "ProductID": 4365, "Quantity": 1, "Cost": 18 }], "OrderDate": "04/01/2017" }
1002	{ "CustomerID": 220, "OrderItems": [{ "ProductID": 1285, "Quantity": 1, "Cost": 120 }], "OrderDate": "05/08/2017" }

A document database represents the opposite end of the NoSQL spectrum from a key-value store. In a document database, each document has a unique ID, but the fields in the documents are transparent to the database management system. Document databases typically store data in JSON format,

Document Databases

RDBMS	MongoDB
Database	Database
Table	Collection
Tuple/Row	Document
column	Field
Table Join	Embedded Documents
Primary Key	Primary Key (Default key <code>_id</code> provided by mongodb itself)



Document Database

Customer	
PK	CustomerID
FK1	Title FirstName LastName AddressID



Address	
PK	AddressID
	StreetAddress City State ZipCode

Customer Table

CustomerID	Title	FirstName	LastName	AddressID
1	Mr	Mark	Hanson	500
2	Ms	Lisa	Andrews	501
3	Mr	Walter	Harp	500

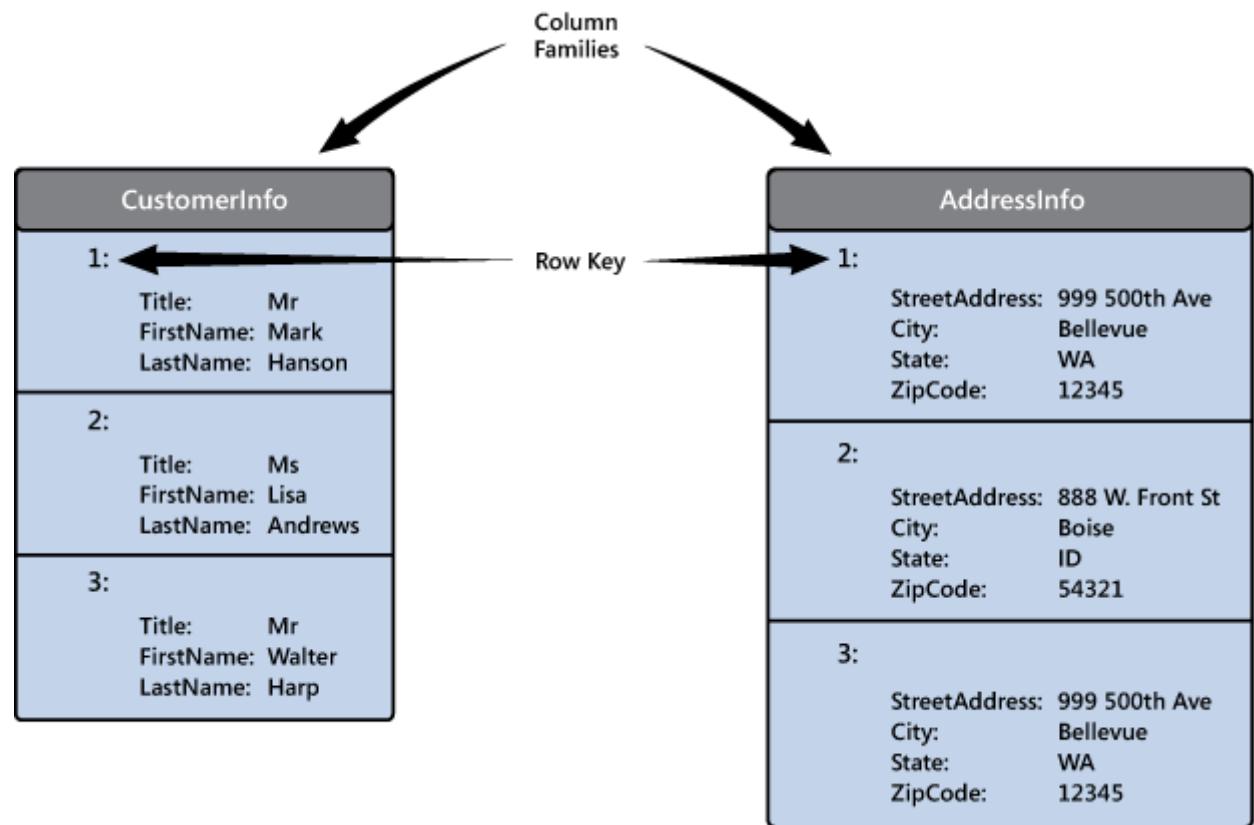
Address Table

AddressID	StreetAddress	City	State	ZipCode
500	999 500th Ave	Bellevue	WA	12345
501	888 W. Front St	Boise	ID	54321

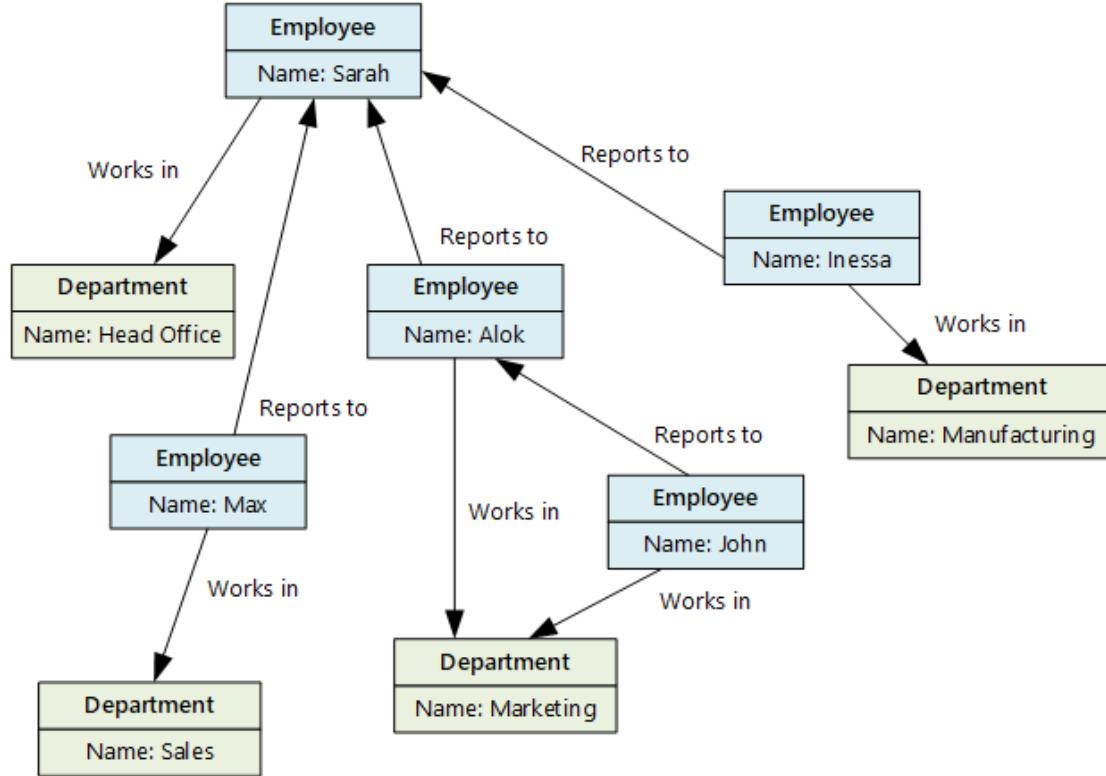
RDBMS is Row-based oriented

Column Family Databases

Row Key	Column Families		
	CustomerInfo		AddressInfo
CustomerID	CustomerInfo:Title	CustomerInfo:FirstName	CustomerInfo:LastName
1	CustomerInfo:Title Mr	CustomerInfo:FirstName Mark	CustomerInfo:LastName Hanson
	AddressInfo:StreetAddress 999 500th Ave	AddressInfo:City Bellevue	AddressInfo:State WA
	AddressInfo:ZipCode 12345		
2	CustomerInfo:Title Ms	CustomerInfo:FirstName Lisa	CustomerInfo:LastName Andrews
	AddressInfo:StreetAddress 888 W. Front St	AddressInfo:City Boise	AddressInfo:State ID
	AddressInfo:ZipCode 54321		
3	CustomerInfo:Title Mr	CustomerInfo:FirstName Walter	CustomerInfo:LastName Harp
	AddressInfo:StreetAddress 999 500th Ave	AddressInfo:City Bellevue	AddressInfo:State WA
	AddressInfo:ZipCode 12345		



Column Family Databases

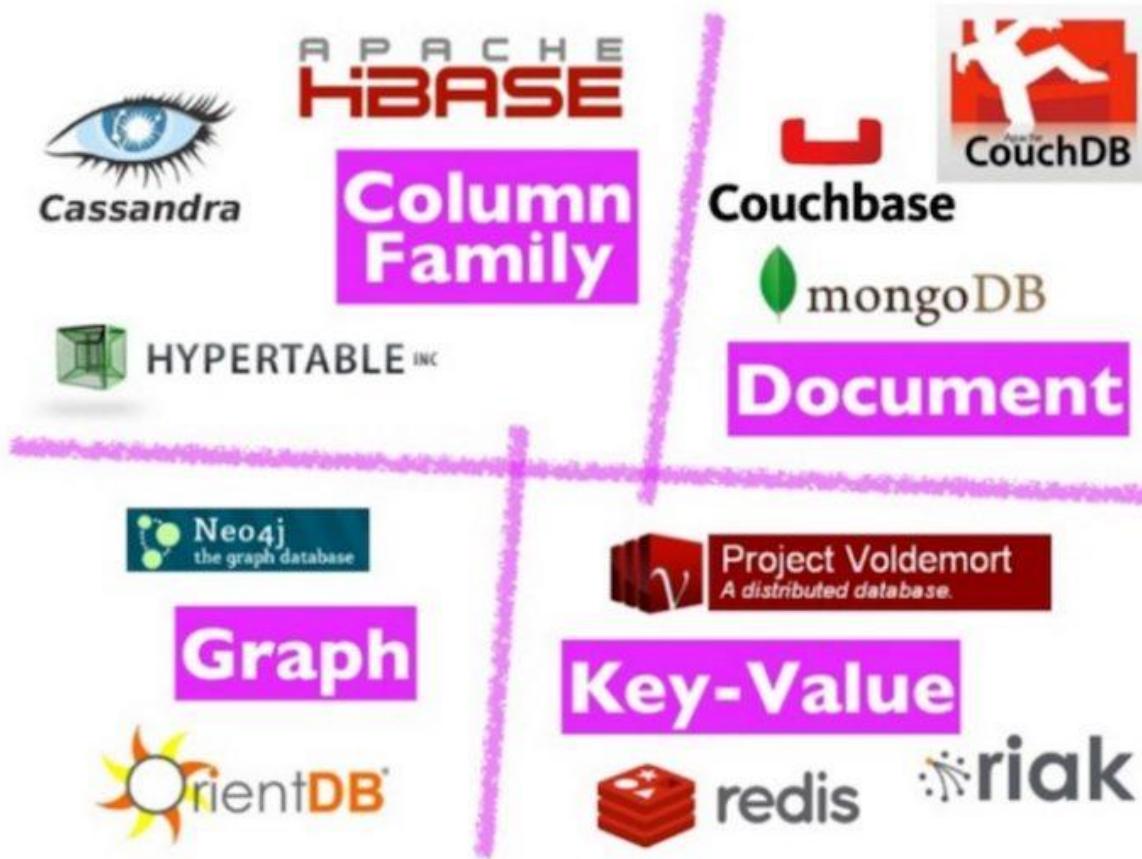


Graph databases enable you to store entities, but the main focus is on the relationships that these entities have with each other.

A graph database stores two types of information: nodes that you can think of as instances of entities, and edges, which specify the relationships between nodes.

Graph Databases

No-SQL Tools



Database Administrator

Database Management

Implements Data Security

Backups

User Access

Monitors performance



Data Engineer

Data Pipelines and processes

Data Ingestion storage

Prepare data for Analytics

Prepare data for analytical processing



Data Analyst

Provides insights into the data

Visual Reporting

Modeling Data for Analysis

Combines data for visualization and analysis



Roles in Data

DATA SCIENTISTS

A.K.A. STATISTICIANS,
DATA MANAGERS

SKILLS

Mathematics,
Programming
Communication

TOOLS

SQL, Python, R

DATA ENGINEERS

A.K.A. DATA ARCHITECTS,
DATABASE ADMINISTRATORS

SKILLS

Programming
Mathematics
Big Data

TOOLS

Hadoop, NoSQL, Python

DATA ANALYSTS

A.K.A. BUSINESS ANALYSTS

SKILLS

Statistics
Communication
Business Knowledge

TOOLS

Excel, Tableau, SQL

Roles in Data

- Installing and upgrading the database server and application tools.
- Allocating system storage and planning storage requirements for the database system.
- Modifying the database structure, as necessary, from information given by application developers.
- Enrolling users and maintaining system security.

Database Administrator tasks and responsibilities

- Ensuring compliance with database vendor license agreement.
- Controlling and monitoring user access to the database.
- Monitoring and optimizing the performance of the database.
- Planning for backup and recovery of database information.
- Maintaining archived data.

Database Administrator tasks and responsibilities

- Backing up and restoring databases.
- Contacting database vendor for technical support.
- Generating various reports by querying from database as per need.
- Managing and monitoring data replication.
- Acting as liaison with users.

Database Administrator tasks and responsibilities

- Developing, constructing, testing, and maintaining databases and data structures.
- Aligning the data architecture with business requirements.
- Data acquisition.
- Developing processes for creating and retrieving information from data sets.
- Using programming languages and tools to examine the data.

Data Engineer tasks and responsibilities

- Identifying ways to improve data reliability, efficiency, and quality.
- Conducting research for industry and business questions.
- Deploying sophisticated analytics programs, machine learning, and statistical methods.
- Preparing data for predictive and prescriptive modeling.
- Using data to discover tasks that can be automated.

Data Engineer tasks and responsibilities

- Making large or complex data more accessible, understandable, and usable.
- Creating charts and graphs, histograms, geographical maps, and other visual models that help to explain the meaning of large volumes of data, and isolate areas of interest.
- Transforming, improving, and integrating data from many sources, depending on the business requirements.

Data Analyst tasks and responsibilities

- Combining the data result sets across multiple sources. For example, combining sales data and weather data provides a useful insight into how weather influenced sales of certain products such as ice creams.
- Finding hidden patterns using data.
- Delivering information in a useful and appealing way to users by creating rich graphical dashboards and reports.

Data Analyst tasks and responsibilities

Ask a question about your data

Total Stores NEW & EXISTING STORES
104

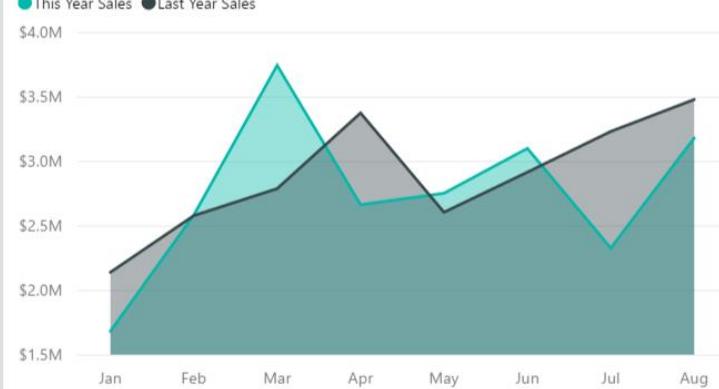
This Year's Sales NEW & EXISTING STORES
\$22M

This Year's Sales BY CHAIN
Lindseys Fashions Direct

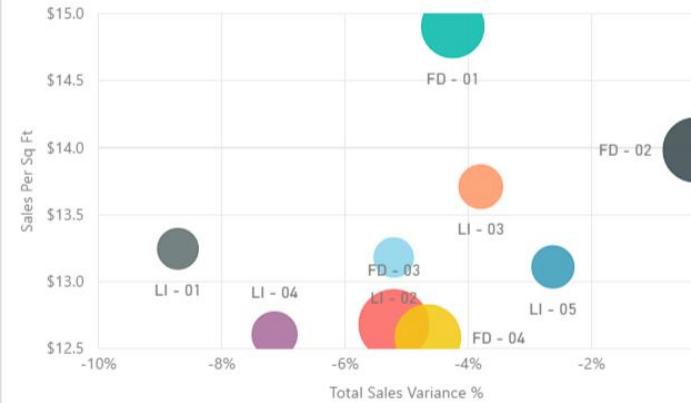
New Stores, New Stores Target YEAR TO DATE
● New Stores ● New Stores Target

This Year's Sales NEW STORES ONLY
\$2M

This Year's Sales, Last Year's Sales BY FISCAL MONTH
● This Year Sales ● Last Year Sales



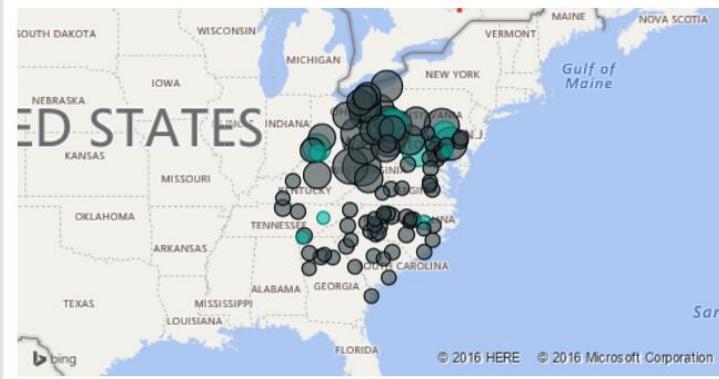
Total Sales Variance %, Sales Per Sq Ft, This Year's Sales BY DISTRICT



New Stores NEW STORES ONLY



This Year's Sales BY POSTAL CODE, STORE TYPE



This Year's Sales BY CITY, CHAIN



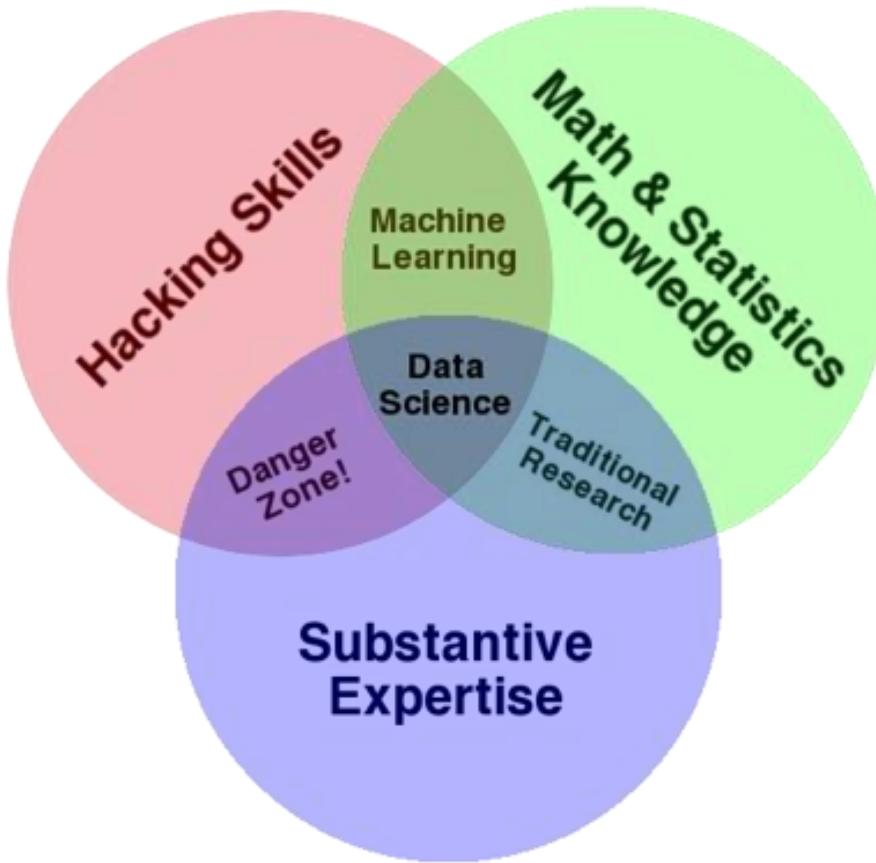
Sales Per Sq Ft BY NAME



02 - Data Science for Business

- Data Science aims to derive **knowledge** from **big data, efficiently and intelligently**.
- Data Science encompasses the **set of activity, tools, and methods** that enable **data-driven activities** in science, business, medicine, and government.
- Data science is a multidisciplinary blend of **data inference, algorithm development, and technology** in order to solve analytically complex problems.

What is Data Science?



Data Science Venn Diagram (Drew Conway)

Content Reference : <http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>

Element	Databases	Data Science
Data Value	Precious	Cheap
Data Volume	Modest	Massive
Example	Bank records, Census Information, Medical records,	Online clicks, GPS logs, Social Media (Facebook, Tweets), ...
Priorities	Consistency, Error Recovery, Auditability	Speed, Availability, Query richness
Structured	Strongly (Schema)	Weakly or None (Text)

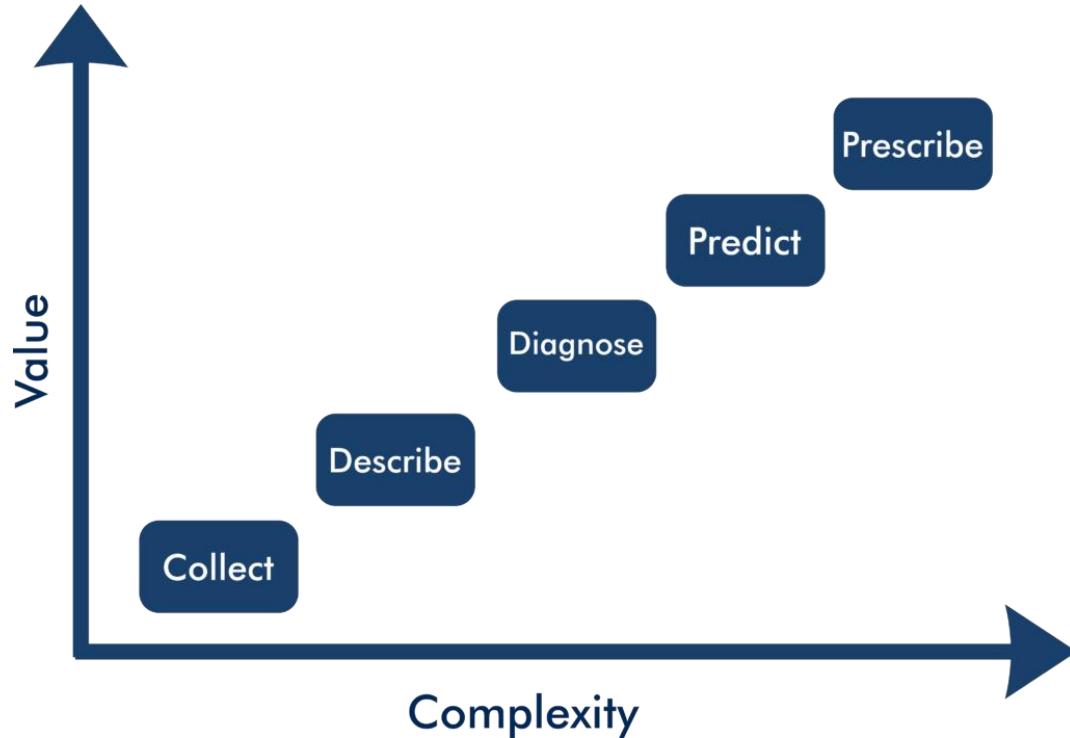
Data Science vs. Database

Element	Databases	Data Science
Properties	Transactions, Atomicity, Consistency, Isolation, Durability	Consistency, Availability, Partition Tolerance, theorem (2/3), eventual consistency
Realizations	SQL (Oracle, SQL Server, MariaDB, ...)	NoSQL (MongoDB, Apache Hbase, ...)
Query	Querying the past	Querying the future

Data Science vs. Database

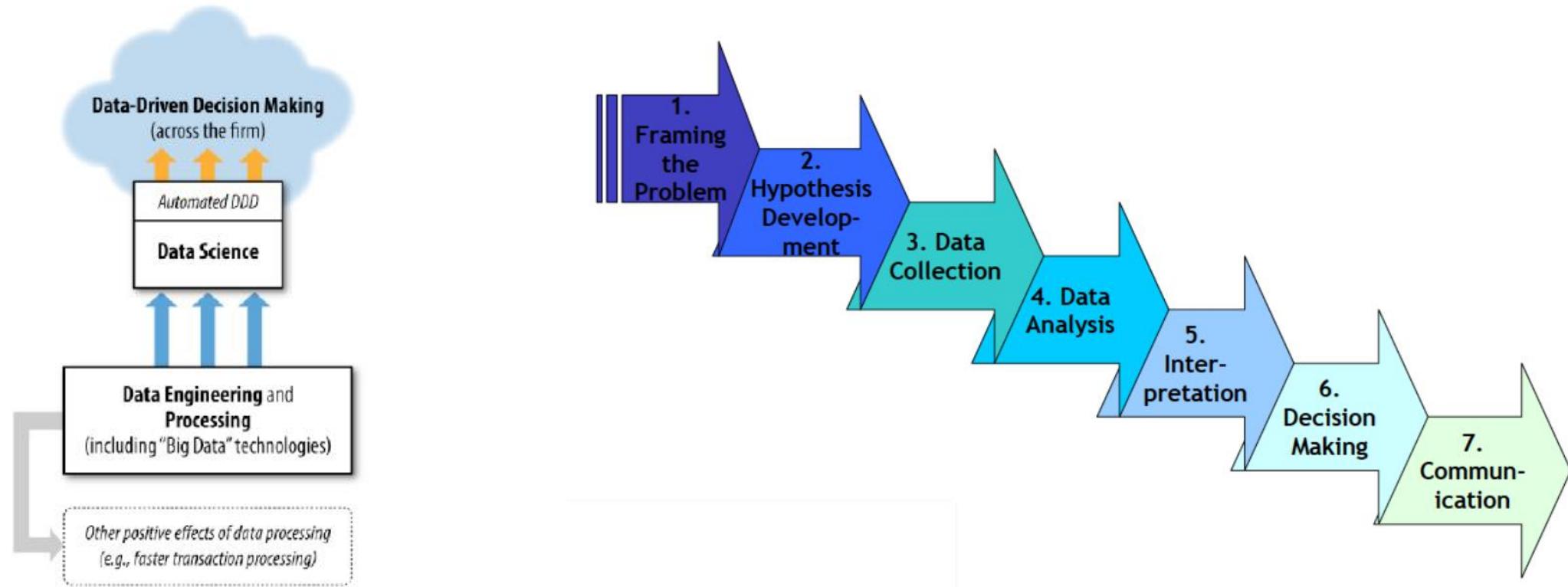
- Cloud Computing
- Big Data
- Machine Learning
- Statistics and Probability
- Programming Language (R, Python)

Basic Components of Data Science



- Describing the data is the first step in extracting value.
- Descriptive statistics are the core of most business reporting and are an essential first step in analysing the data.
- Diagnostics or analysis is the core activity of most professions.
- Predictive analysis seems to be the holy grail for many managers.
- Prescriptive analysis uses the knowledge created in the previous phases to automatically run a business process and even decide on future courses of action.

Strategic Data Science



Data-Driven Decision Making Process

- Guard against your biases
- Define objectives
- Gather data now
- Find the unresolved questions
- Find the data needed to solve these questions



Data-Driven Decision Making Guide

- Analyze and understand
- Don't be afraid to revisit and reevaluate
- Present the data in a meaningful way
- Set measurable goals for decision making
- Continue to evolve your data driven business decisions



Data-Driven Decision Making Guide

- Curiosity
- Creativity
- Focus
- Attention to Detail

Characteristics of Data Scientist

- Finding rich data sources.
- Working with large volumes of data despite hardware, software, and bandwidth constraints.
- Cleaning the data and making sure that data is consistent.
- Melding multiple datasets together.
- Visualizing that data.

Data Scientist

- Business User
- Project Sponsor
- Project Manager
- BI Analyst
- Database Administrator
- Data Engineer
- Data Scientist

Data Science Team

- Transforming
- Creating
- As a Service
- Crowdsourcing

Approaches to Developing Data Science Capabilities

-
- Problem Identification
 - Business Understanding
 - Data Acquisition - ETL
 - Data Understanding - EDA
-
- Feature Engineering
 - Model Training/Validation
 - Model Deployment
 - Model Monitoring

Data Science Life Cycles

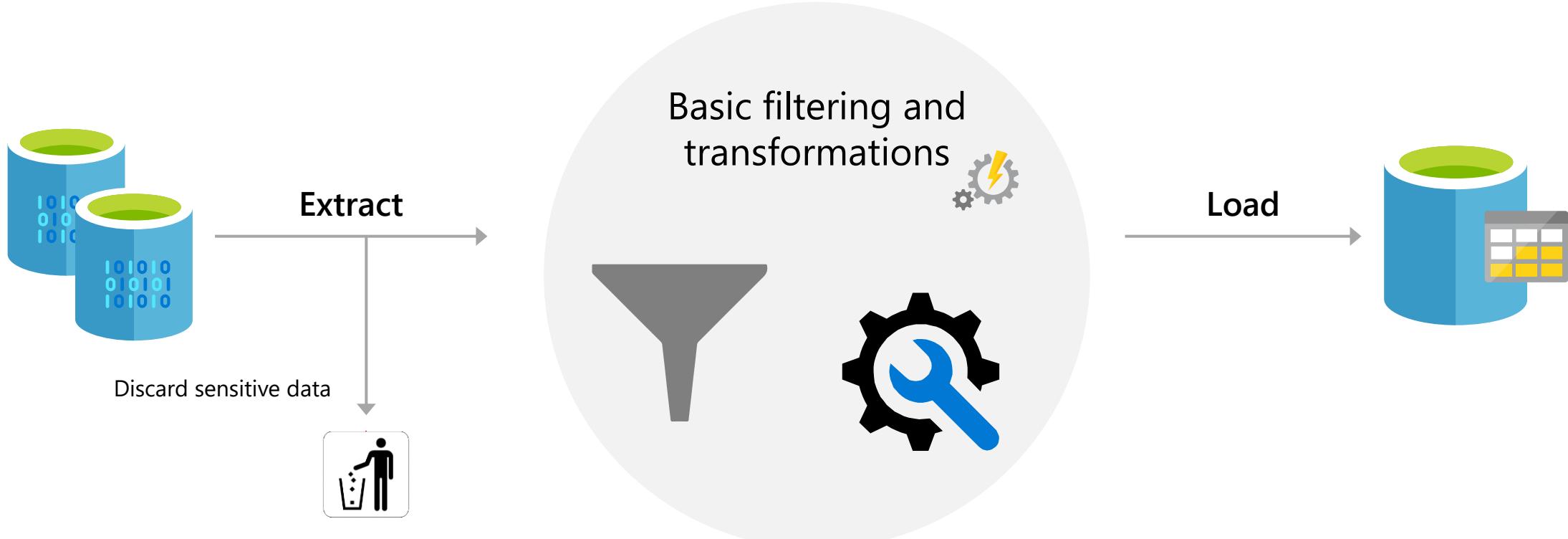
- Frame the Business Problem
- Previous know outcome from reliable and different source and industry experts
- Figure out whether cost of Failure does create reputation damage
- Facilitate Experimentation with Automation
- Decide where you need to scale

Problem Identification

- Understand the Business Objective
- Understand scalability need
- Defining your success Criteria and success measuring Metrics KPI & SLA
- Identify the key business variables that the model needs to predict
- Discuss Integration of model with business process

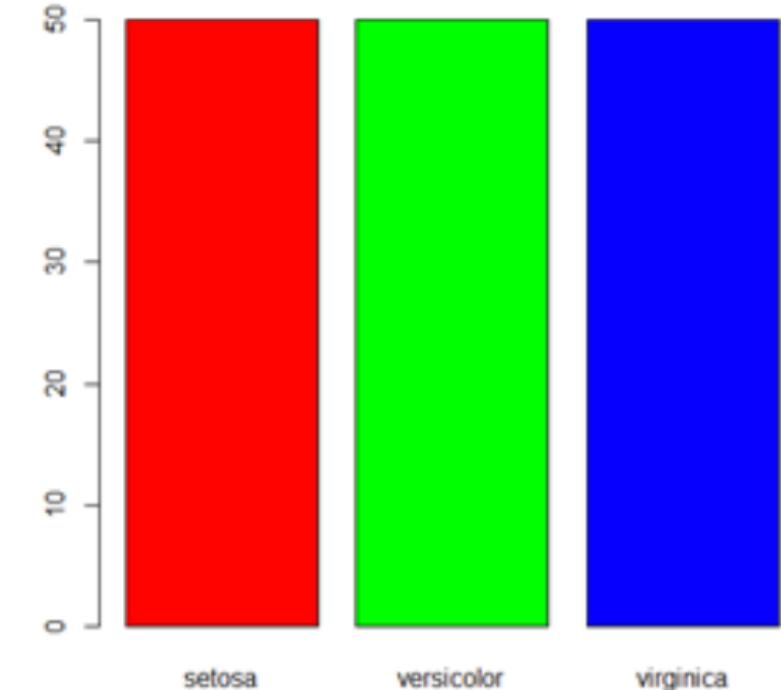
Business Understanding

Transform



Data Acquisition - ETL

- Descriptive Statistics
 - Measures of Frequency
 - Measures of Central Tendency
 - Measures of Dispersion or Variation
 - Measures of Position
- Inferential Statistics
 - The estimation of the parameter(s)
 - Testing of statistical hypotheses.

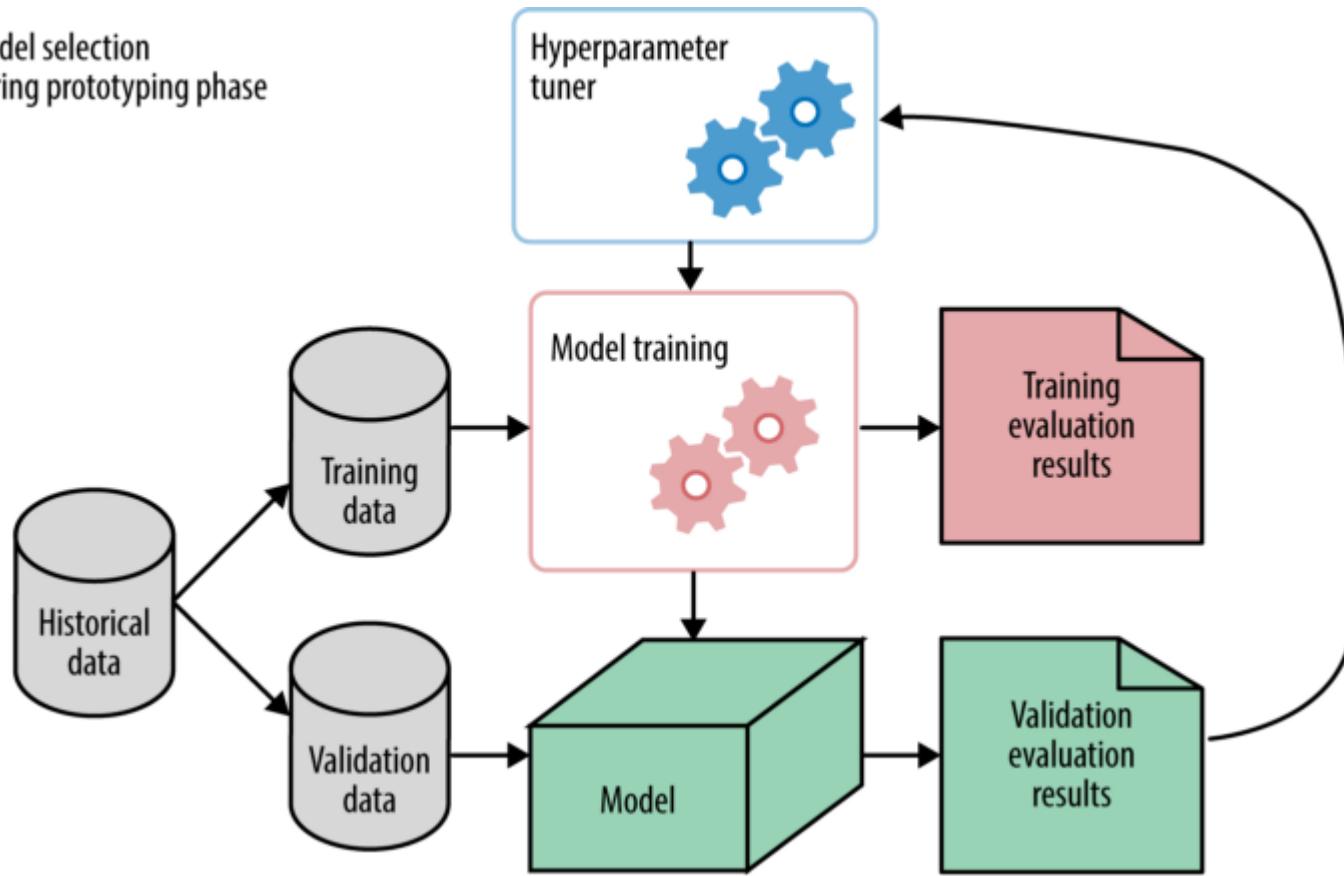


Data Understanding - EDA

- Feature engineering is the process of using domain knowledge to extract features from raw data via data mining techniques.
- These **features** can be used to strengthen the performance of machine learning models.
- Feature engineering can be considered as applied machine learning itself.

Feature Engineering

Model selection
during prototyping phase



Model Training/Validation

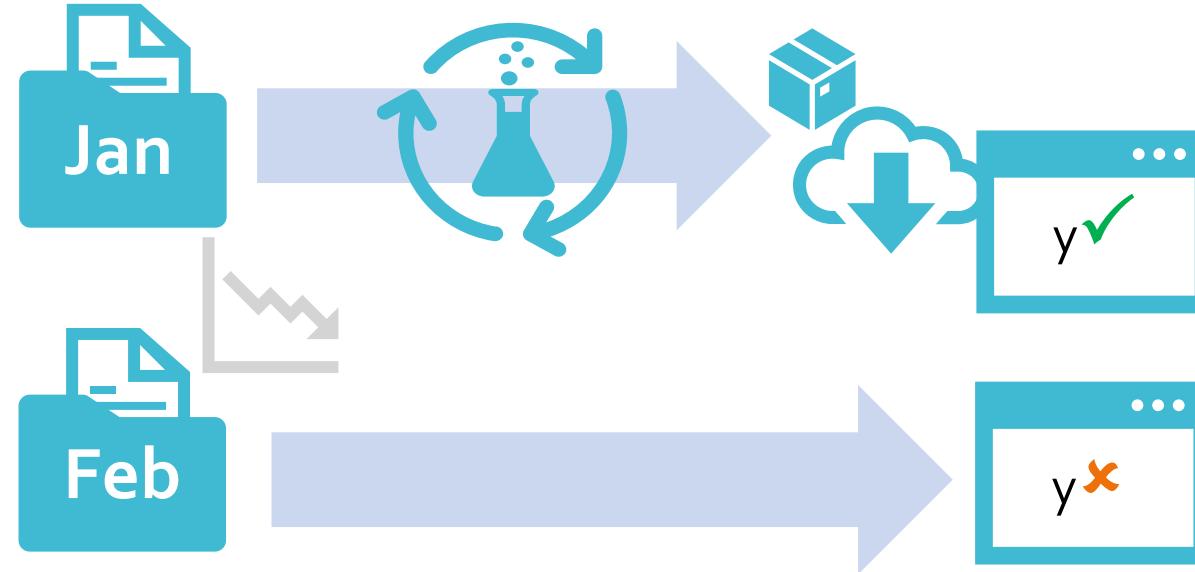


Deploy a Real-Time Pipeline



Publish a Batch Pipeline

Model Deployment



Model Monitoring



Thank you for stopping by.

Google Flu Trends and Google Dengue Trends are no longer publishing current estimates of Flu and Dengue fever based on search patterns. The historic estimates produced by Google Flu Trends and Google Dengue Trends are available below. It is still early days for nowcasting and similar tools for understanding the spread of diseases like flu and dengue – we're excited to see what comes next. Academic research groups interested in working with us should fill out this [form](#).

Sincerely,

The Google Flu and Dengue Trends Team.

Google Flu Trends Data:

You can also see this data in [Public Data Explorer](#)

- World
- Argentina
- Australia
- Austria
- Belgium
- Bolivia
- Brazil
- Bulgaria
- Canada
- Chile
- France
- Germany
- Hungary
- Japan

- **Google Flu Trends
(Nowcasting Example)**

- <https://www.google.org/flutrends/about/>

Nowcasting vs Forecasting

US World Environment Soccer **US politics** Business Tech Science Homelessness

Nate Silver

Numbers nerd Nate Silver's forecasts prove all right on election night

FiveThirtyEight blogger predicted the outcome in all 50 states assuming Barack Obama's Florida victory is confirmed



Luke Harding

Wed 7 Nov 2012 15.45 GMT



This article is over 5 years old

7

Advertisement

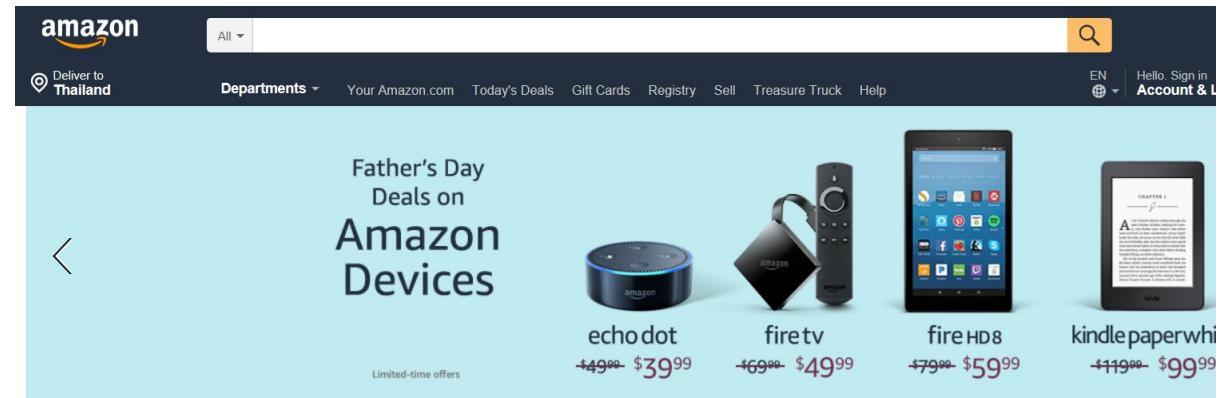
A week in the life of the world



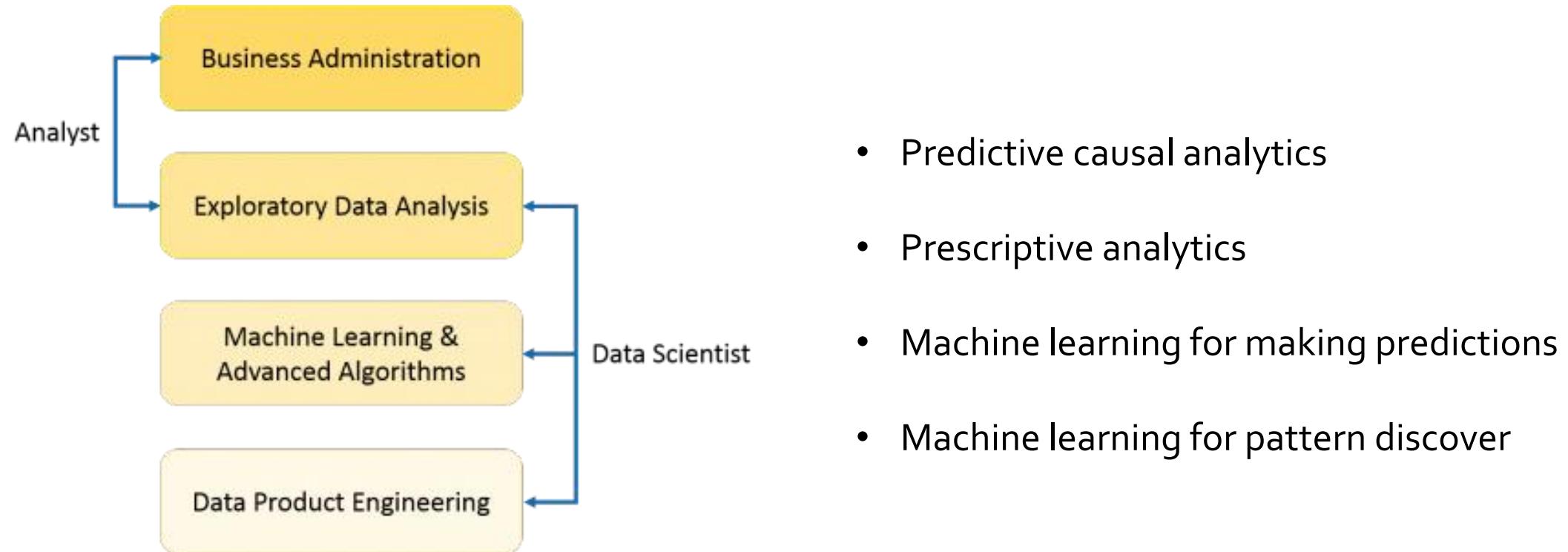
elections2012 (Forecasting Example)

Nowcasting vs Forecasting

- Transform Data Into ...
 - Valuable Insights
 - Data Productions
 - Interesting Stories



Summary of Data Science



Summary of Data Science

03 - Big Data Fundamentals

- Big data is a term that describes the large volume of data – both structured and unstructured – that inundates a business on a day-to-day basis. But it's not the amount of data that's important. It's what organizations do with the data that matters.
- Big data can be analyzed for insights that lead to better decisions and strategic business moves.

Overview - Big Data

- It's all happening online

- Click
- Ad impression
- Billing event
- Fast Forward, pause,...
- Server request
- Transaction
- Network message
- Fault

- User Generate Content on Web & Mobile



Where Does Big Data Come From?

- Health and Scientific Computing

Biology 2.0

A decade after the human-genome project, writes Geoffrey Carr (interviewed here), biological science is poised on the edge of something wonderful

Jun 17th 2010

Timekeeper

Like 391

Tweet

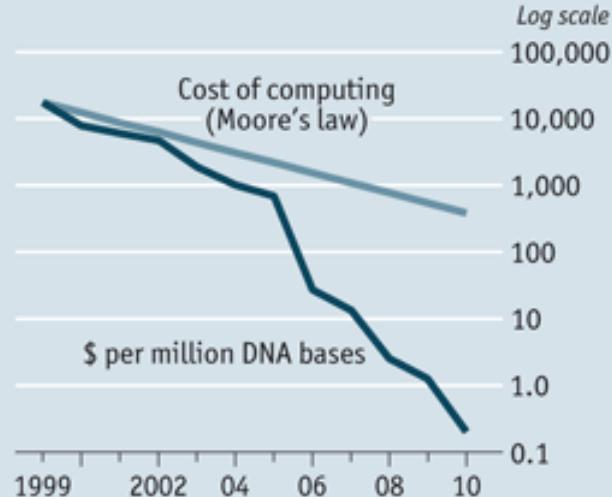


TEN years ago, on June 26th 2000, a race ended. The result was declared a dead heat and both runners won the prize of shaking the hand of America's then president, Bill

Baseline information

1

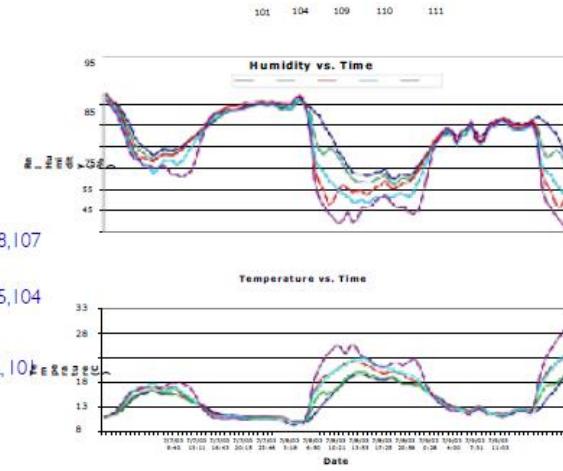
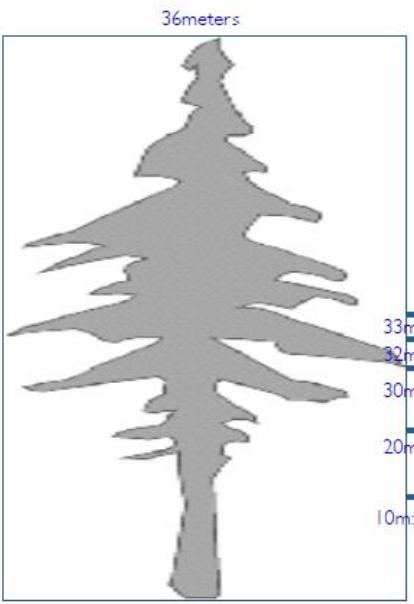
Cost of genome sequencing compared with
Moore's law for computers



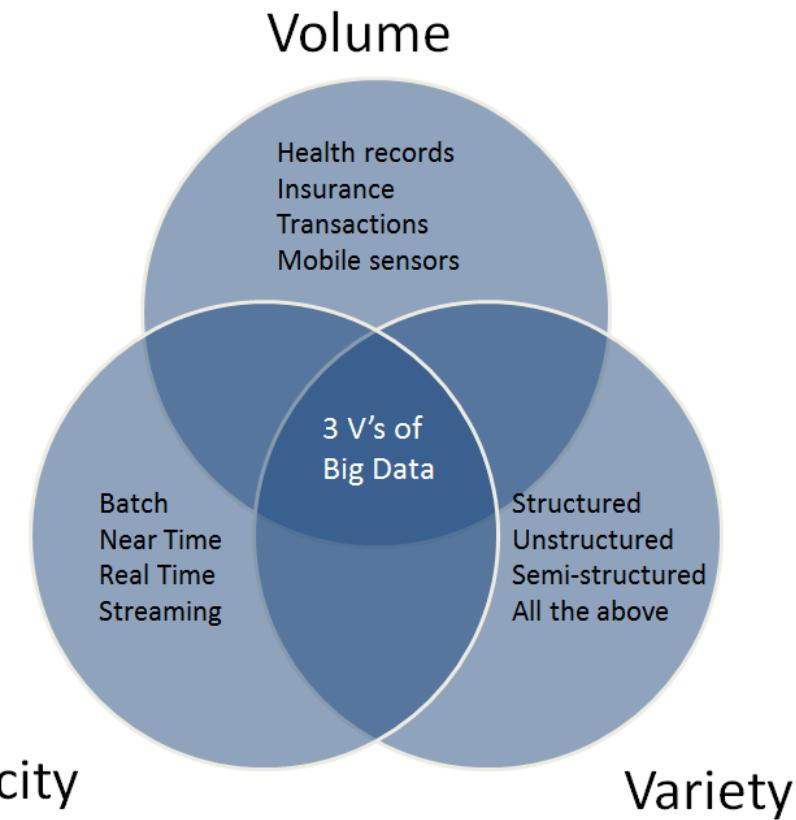
Source: Broad Institute

Where Does Big Data Come From?

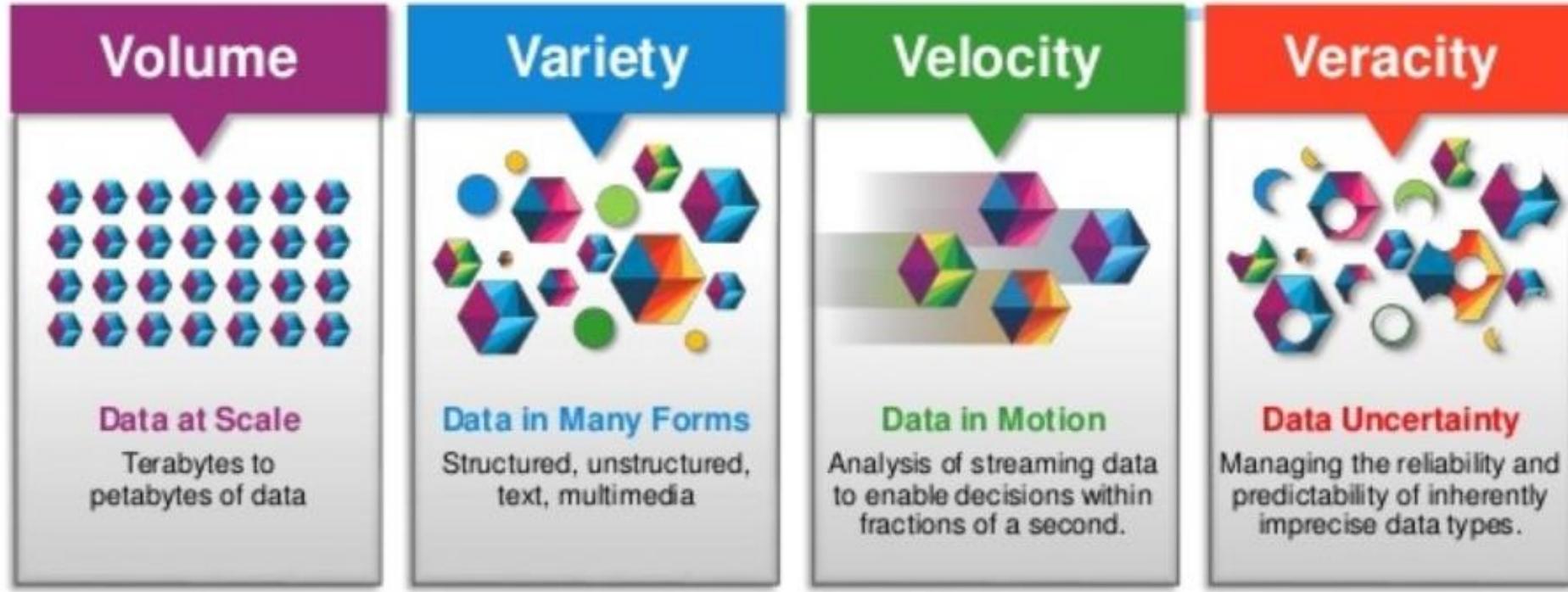
- Internet of Thing (IoT)



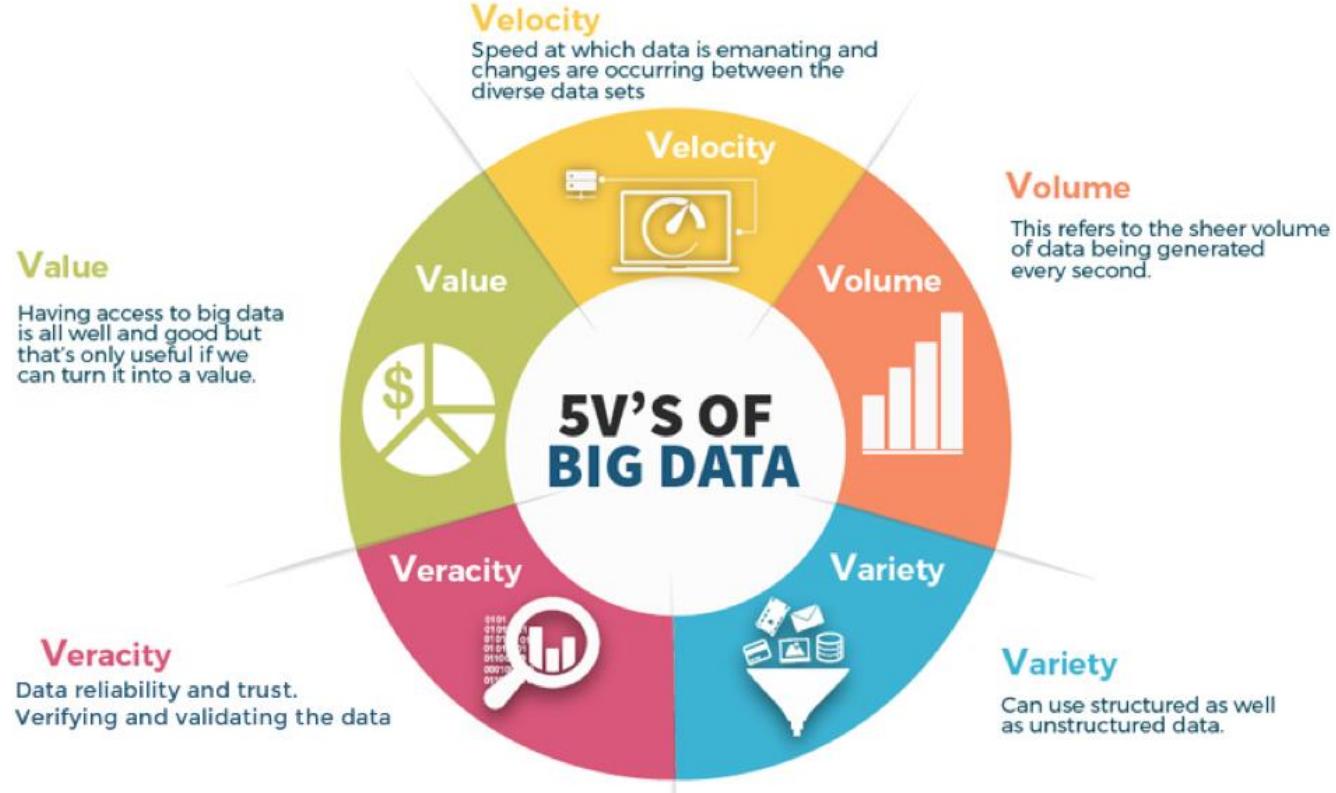
Where Does Big Data Come From?



The 3 V's of Big Data Characteristics

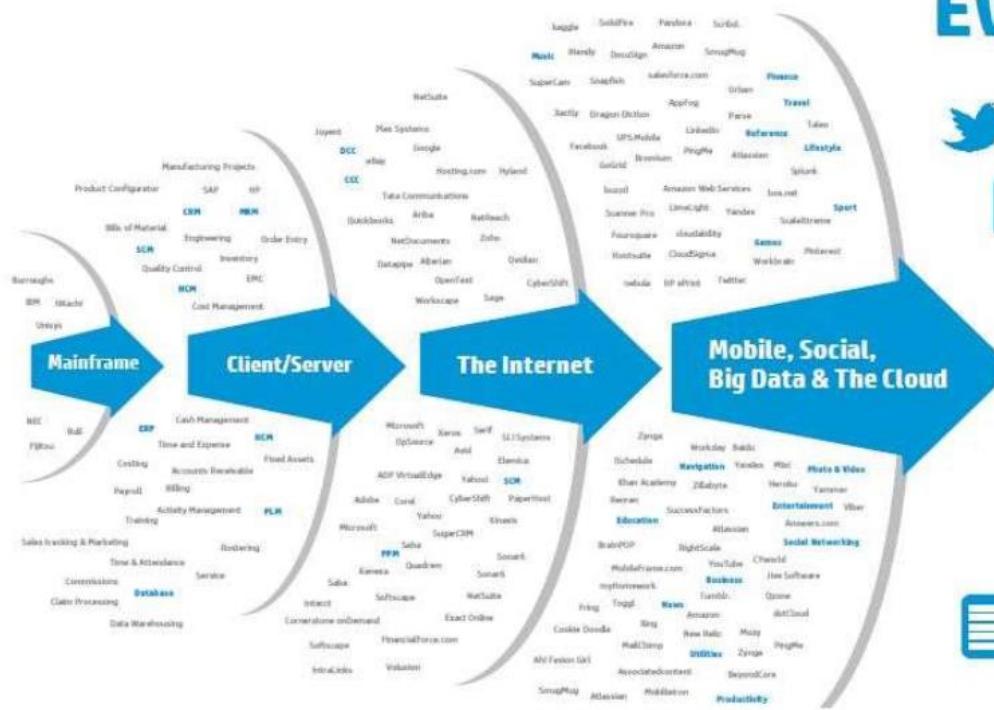


The 4V's of Big Data Characteristics (IBM)



The 5V's of Big Data

Content Reference : <https://www.techentice.com/the-data-veracity-big-data/>



Every 60 seconds

98,000+ tweets

695,000 status updates

11million instant messages

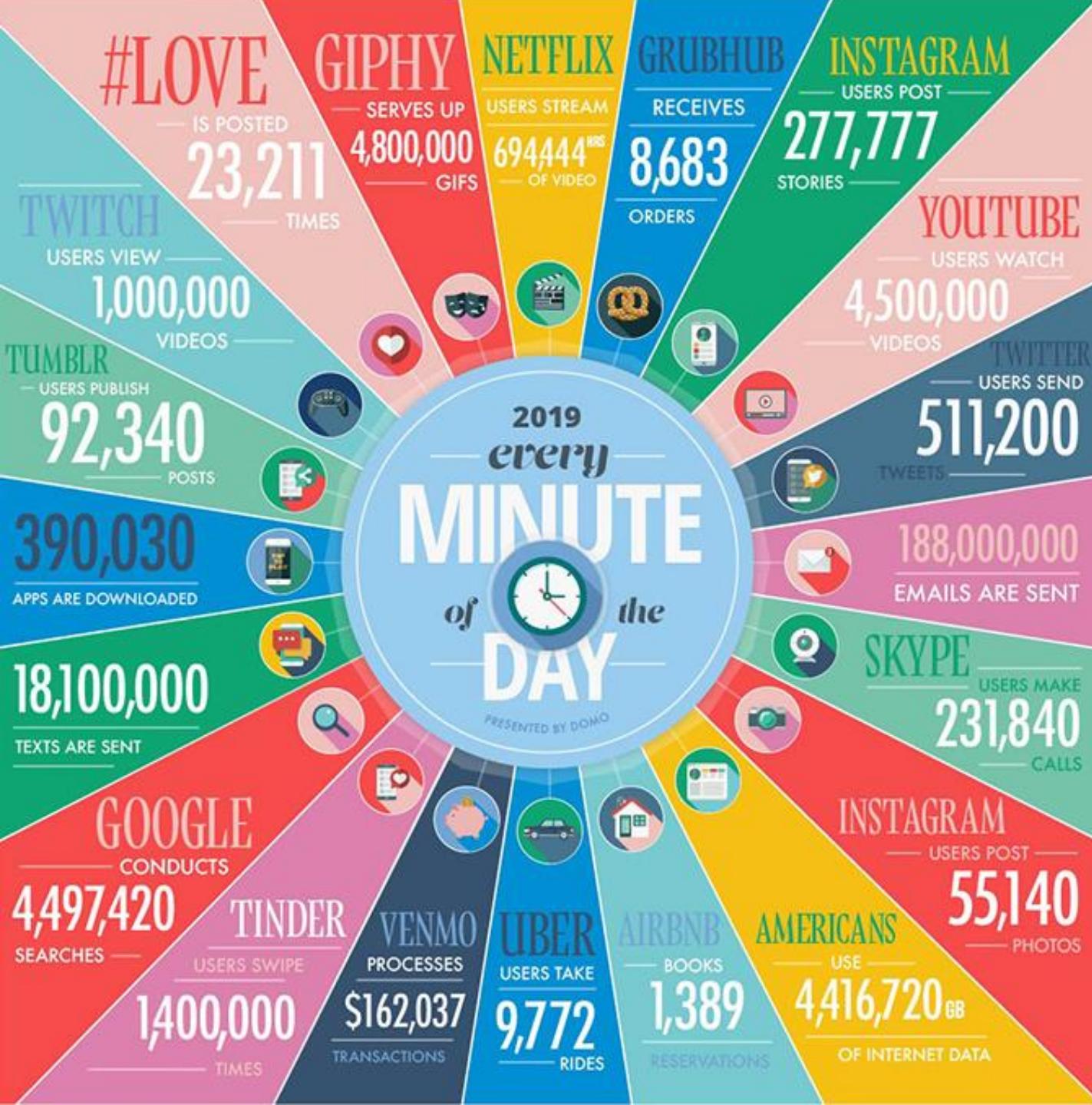
698,445 Google searches

168 million+ emails sent

1,820TB of data created

217 new mobile web users

Big Data Characteristics



Data Never Sleeps

Content Reference :
<https://web-assets.domo.com/blog/wp-content/uploads/2019/07/data-never-sleeps-7-896kb.jpg>



华为分别在德国和中国完成了车辆编队外场演示。
Huawei has also conducted platooning trial in Germany and in China.

- Hadoop – Big Data Platform
- MongoDB – Document Database [No-SQL]
- Cloud Computing



Big Data Technologies Overview

- Hadoop is an Apache open source framework written in java that allows **distributed processing of large datasets across clusters of computers** using simple programming models.
- A Hadoop framework application works in an environment that provides distributed storage and computation across clusters of computers.
- Hadoop is designed to scale up from single server to thousands of machines, each offering local computation and storage.



Introduction to Hadoop

- When to use Hadoop
 - For processing really big data
 - For storing a diverse set of data
 - For parallel data processing
- When NOT to use Hadoop
 - For real-time data analysis
 - For a relational database system
 - For a general network file system
 - For non-parallel data processing

When (not) to use Hadoop

HADOOP 1.0

MapReduce
(cluster resource management
& data processing)

HDFS
(redundant, reliable storage)

HADOOP 2.0

MapReduce
(data processing)

Others
(data processing)

YARN

(cluster resource management)

HDFS
(redundant, reliable storage)

Hadoop History



MapReduce
(Distributed Computation)

HDFS
(Distributed Storage)

YARN Framework

Common Utilities

- **Hadoop Common**

These are Java libraries and utilities required by other Hadoop modules. These libraries provides filesystem and OS level abstractions and contains the necessary Java files and scripts required to start Hadoop.

- **Hadoop YARN**

This is a framework for job scheduling and cluster resource management.

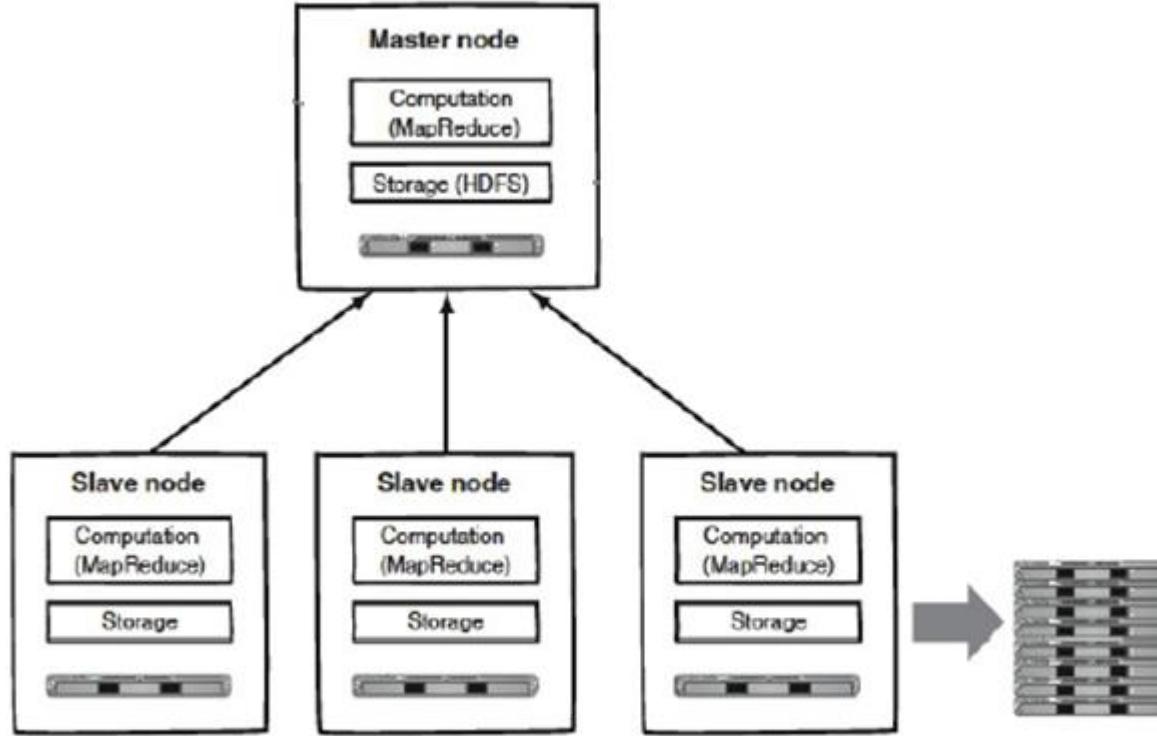
- **Hadoop Distributed File System (HDFS)**

A distributed file system that provides high-throughput access to application data.

- **Hadoop MapReduce :**

This is YARN-based system for parallel processing of large data sets.

Hadoop Architecture



Hadoop Architecture



Hadoop User Experience (HUE)



Data Exchange



Sqoop

Flume



Log Control



ZooKeeper

Coordination

Pig Scripting



Hive SQL



Mahout ML



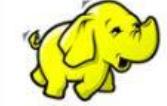
Oozie Workflow



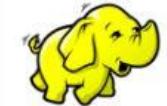
APACHE HBASE

Hbase

Columnar data store



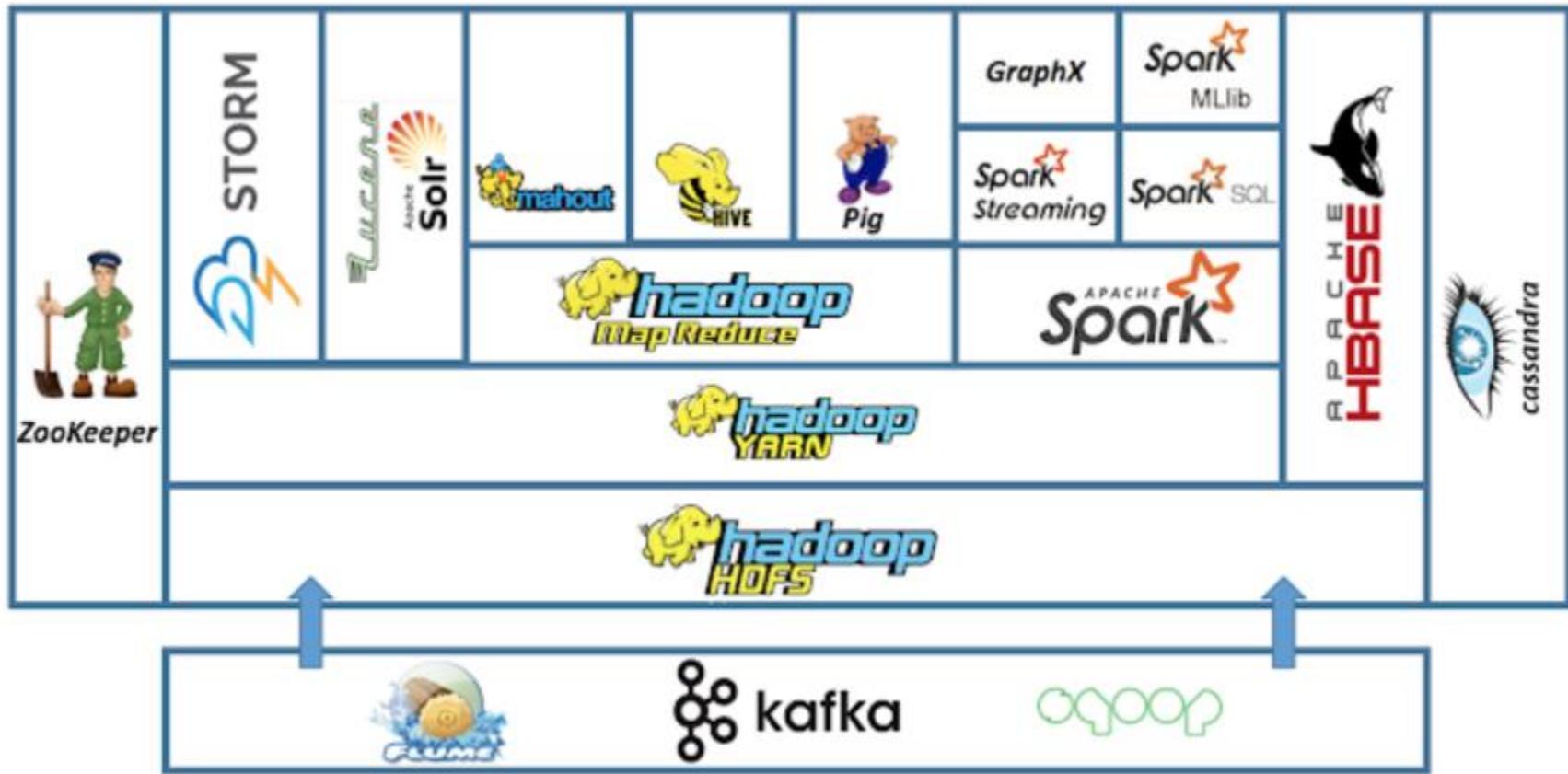
YARN/Map Reduce V2



Hadoop Distributed File System



Hadoop Ecosystem (Version 2)



Hadoop Ecosystem - Update

- Apache Hive is a data warehouse software project built on top of Apache Hadoop for providing data summarization, query and analysis.
- Hive gives an SQL-like interface to query data stored in various databases and file systems that integrate with Hadoop.
- SQL-like queries (HiveQL)



Hadoop Technology - Apache Hive

- HiveQL Sample

```
1 DROP TABLE IF EXISTS docs;
2 CREATE TABLE docs (line STRING);
3 LOAD DATA INPATH 'input_file' OVERWRITE INTO TABLE docs;
4 CREATE TABLE word_counts AS
5 SELECT word, count(1) AS count FROM
6 (SELECT explode(split(line, '\s')) AS word FROM docs) temp
7 GROUP BY word
8 ORDER BY word;
```

Hadoop Technology - Apache Hive

- Apache Pig is a high-level platform for creating programs that run on Apache Hadoop.
- The language for this platform is called Pig Latin

```
1 input_lines = LOAD '/tmp/word.txt' AS (line:chararray);
2 words = FOREACH input_lines GENERATE FLATTEN(TOKENIZE(line)) AS word;
3 filtered_words = FILTER words BY word MATCHES '\\w+';
4 word_groups = GROUP filtered_words BY word;
5 word_count = FOREACH word_groups GENERATE COUNT(filtered_words) AS count, group AS word;
6 ordered_word_count = ORDER word_count BY count DESC;
7 STORE ordered_word_count INTO '/tmp/results.txt';
```



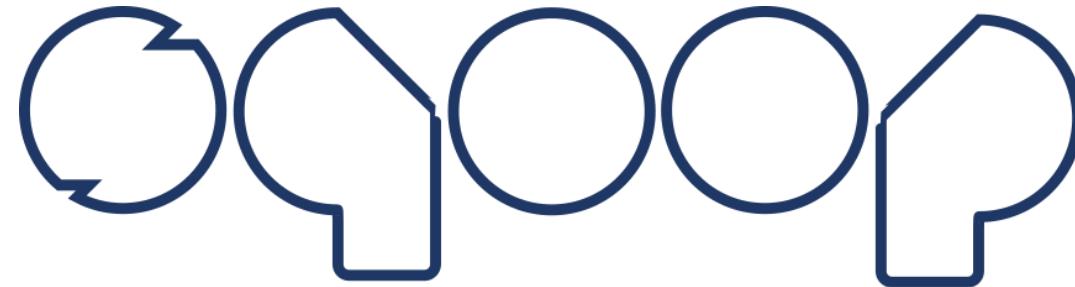
Hadoop Technology - Apache Pig

- HBase is an open-source, non-relational, distributed database modeled after Google's big table and is written in Java.
- HBase is not a direct replacement for a classic SQL database



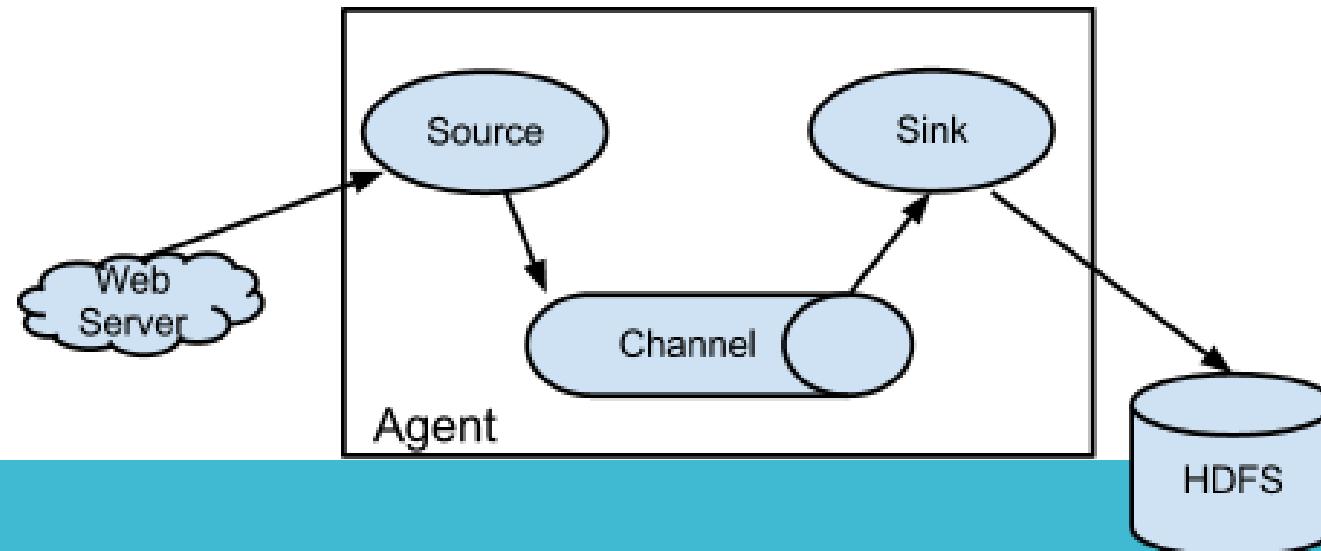
Hadoop Technology - Apache HBase

- Sqoop is a command-line interface application for transferring data between relational databases and Hadoop



Hadoop Technology - Apache Sqoop

- Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of log data.
- It has a simple and flexible architecture based on streaming data flows.

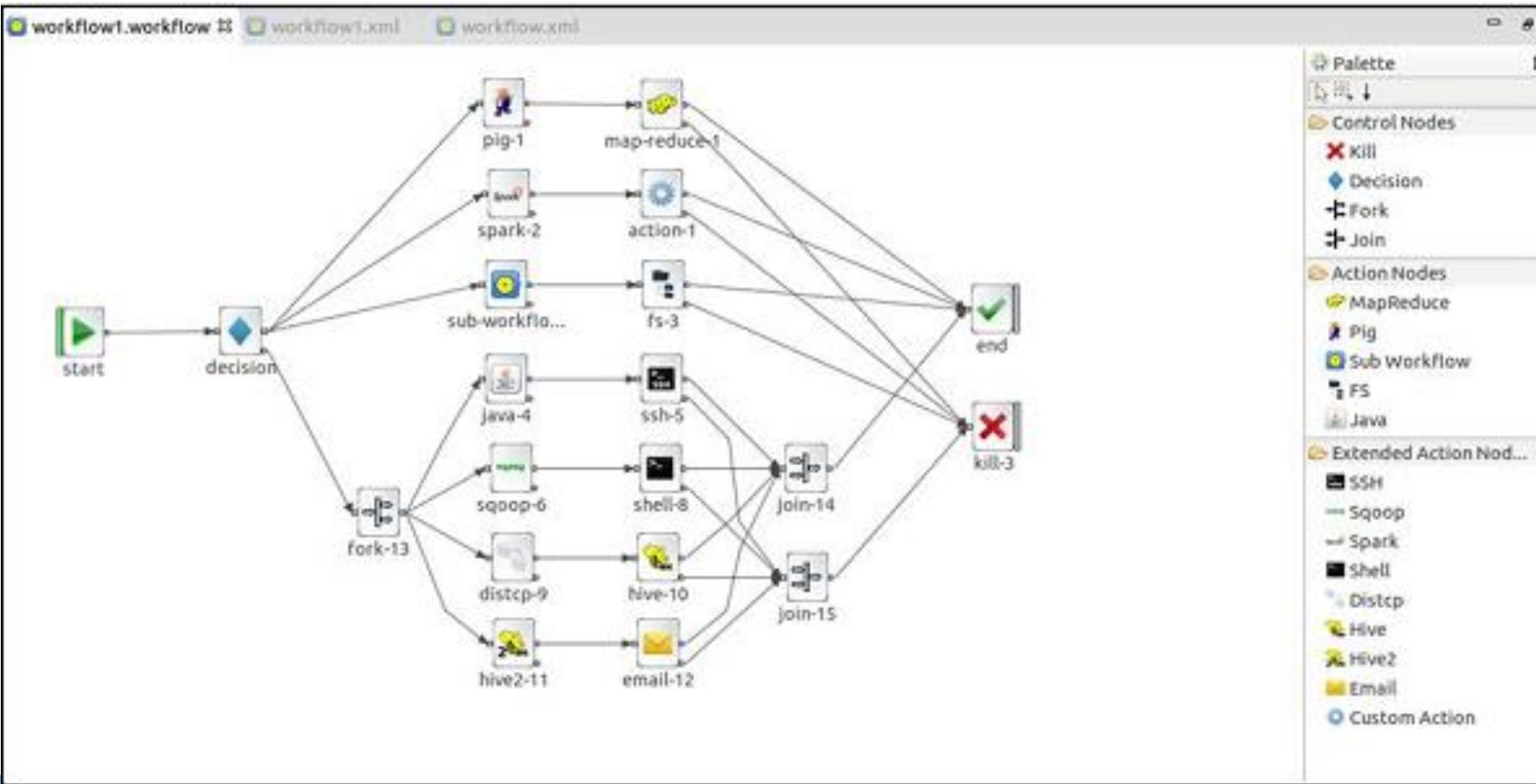


Hadoop Technology - Apache Flume

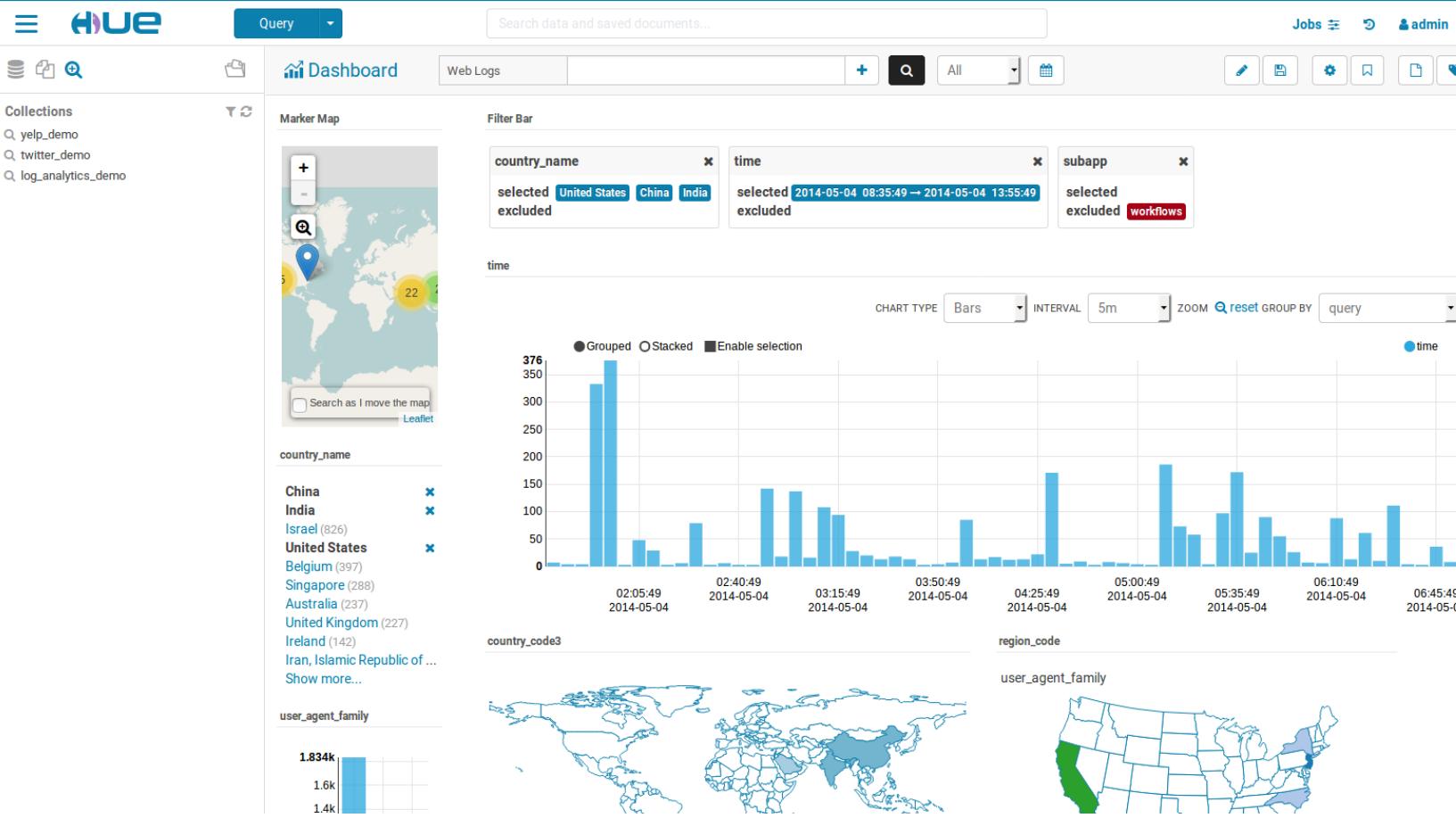
- Apache Oozie is a server-based workflow scheduling system to manage Hadoop jobs.
- Workflows in Oozie are defined as a collection of control flow and action nodes in a directed acyclic graph.
- Control flow nodes define the beginning and the end of a workflow (start, end, and failure nodes) as well as a mechanism to control the workflow execution path (decision, fork, and join nodes).



Hadoop Technology - Apache Oozie



Hadoop Technology - Apache Oozie



- Hue is an open source Analytics Workbench for browsing, querying and visualizing data.

HUE

Hadoop Technology - Apache Hue

- Apache Mahout is a project of the Apache Software Foundation to produce free implementations of distributed or otherwise scalable machine learning algorithms focused primarily in the areas of collaborative filtering, clustering and classification.
- Mahout also provides Java libraries for common maths operations (focused on linear algebra and statistics) and primitive Java collections.



Hadoop Technology - Apache Mahout

MongoDB is a database using document-oriented storage. Under this model, data is stored in documents and documents are combined into collections.

- Database
- Collection
- Document



Big Data Technology - MongoDB

RDBMS	MongoDB
Database	Database
Table	Collection
Tuple/Row	Document
column	Field
Table Join	Embedded Documents
Primary Key	Primary Key (Default key <code>_id</code> provided by mongodb itself)

Big Data Technology - MongoDB

- Schema less
- Structure of a single object is clear
- No complex joins
- Deep query-ability (document-based query language)
- Tuning
- Ease of scale-out

Advantages of MongoDB over RDBM

- Big Data
- Content Management and Delivery
- Mobile and Social Infrastructure
- User Data Management
- Data Hub

Where to Use MongoDB?

```
var MongoClient = require('mongodb').MongoClient;

//Create a database named "myDatabase":
var url = "mongodb://localhost:27017/myDatabase";

MongoClient.connect(url, function(err, db) {
  if (err) throw err;
  console.log("Database created!");
  db.close();
});
```

Create Database using NodeJS

```
var MongoClient = require('mongodb').MongoClient;
var url = "mongodb://localhost:27017/";

MongoClient.connect(url, function(err, db) {
  if (err) throw err;
  var dbo = db.db("myDatabase");

  //Create a collection name "customers":
  dbo.createCollection("customers",
    function(err, res) {
      if (err) throw err;
      console.log("Collection created!");
      db.close();
    });
});
```

Create Collection using NodeJS

```
var MongoClient = require('mongodb').MongoClient;
var url = "mongodb://localhost:27017/";

MongoClient.connect(url, function(err, db) {
  if (err) throw err;
  var dbo = db.db("myDatabase");

  var myobj = {
    name: "Bank of Thailand", address: "Tewet" };
  dbo.collection("customers").insertOne(myobj, function(err, res) {
    if (err) throw err;
    console.log("1 document inserted");
    db.close();
  });
});
```

Create Document using NodeJS

```
var MongoClient = require('mongodb').MongoClient;
var url = "mongodb://localhost:27017/";

MongoClient.connect(url, function(err, db) {
  if (err) throw err;
  var dbo = db.db("myDatabase");

  var myquery = { address: 'Bangkok' };
  dbo.collection("customers").deleteOne(myquery, function(err, obj)
  {
    if (err) throw err;
    console.log("1 document deleted");
    db.close();
  });
});
```

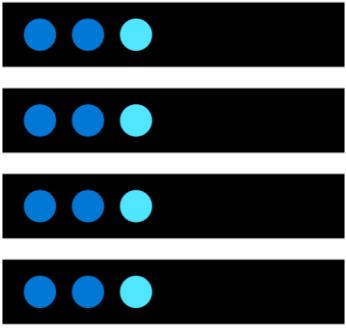
Delete Document using NodeJS

```
var MongoClient = require('mongodb').MongoClient;
var url = "mongodb://localhost:27017/";

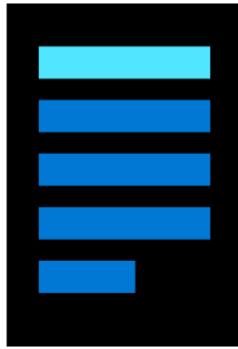
MongoClient.connect(url, function(err, db) {
  if (err) throw err;
  var dbo = db.db("myDatabase");
  dbo.collection("customers")
    .drop(function(err, delOK) {
      if (err) throw err;
      if (delOK) console.log("Collection deleted");
      db.close();
    });
});
```

Drop Collection using NodeJS

Cloud Computing



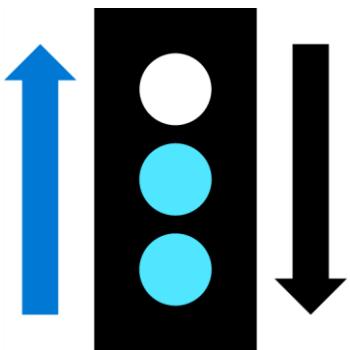
Compute



Storage



Cloud providers include
Microsoft, Amazon, and Google



Networking



Analytics

Explore key cloud concepts

High availability

Fault tolerance

Scalability

Elasticity

Global reach

Customer latency capabilities

Agility

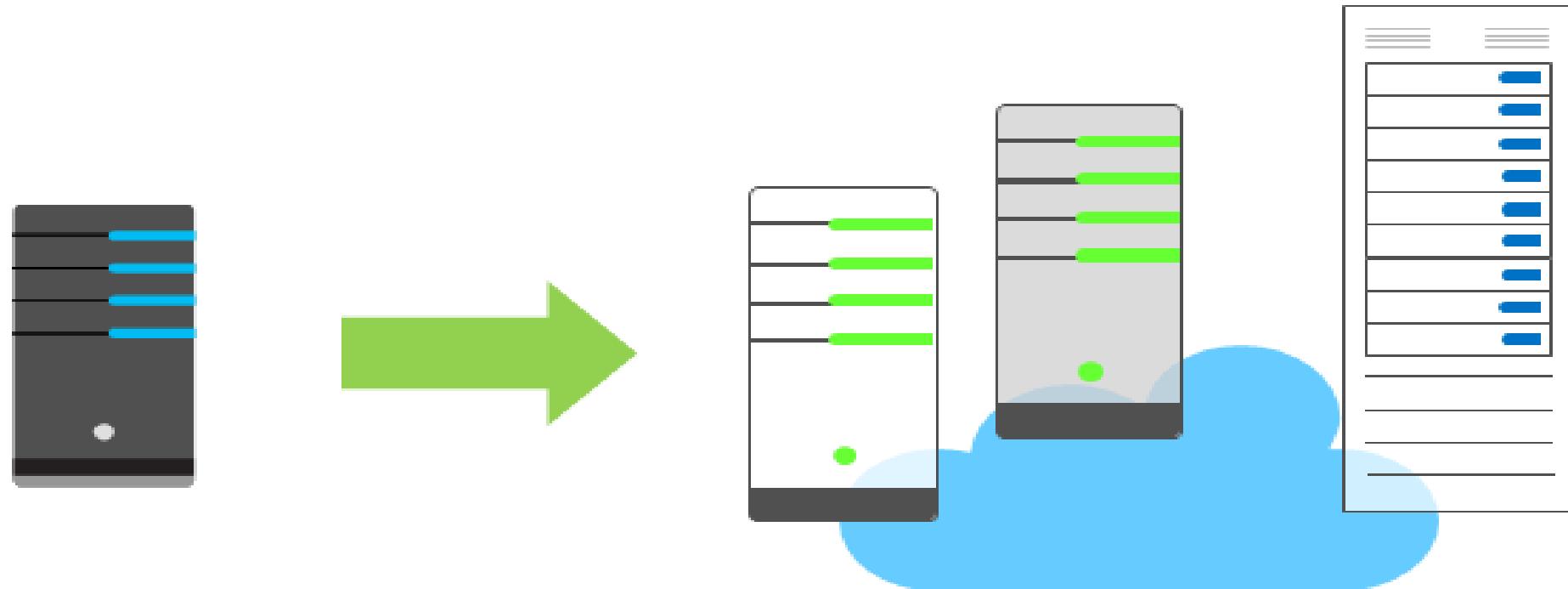
Predictive cost considerations

Disaster recovery

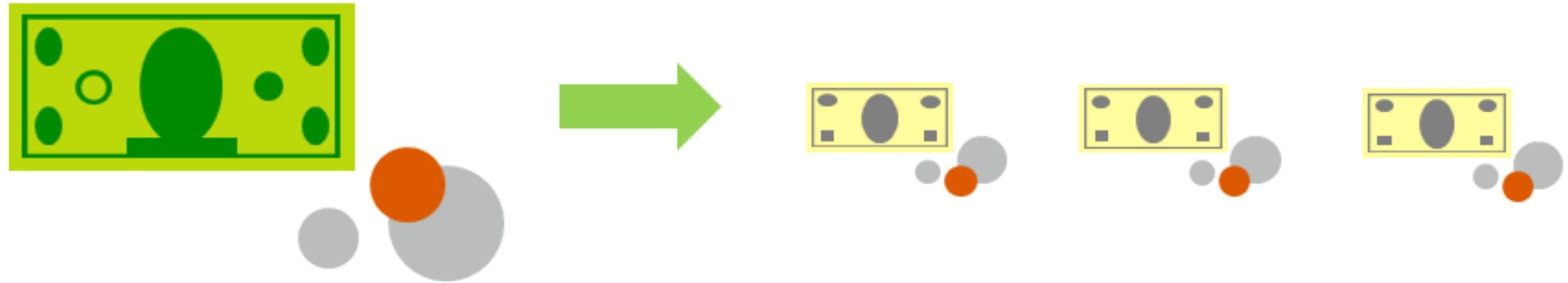
Security

Discuss economies of scale

Economies of scale – Cloud providers can reduce costs and gain efficiency when operating at a large scale.



Compare CapEx vs. OpEx



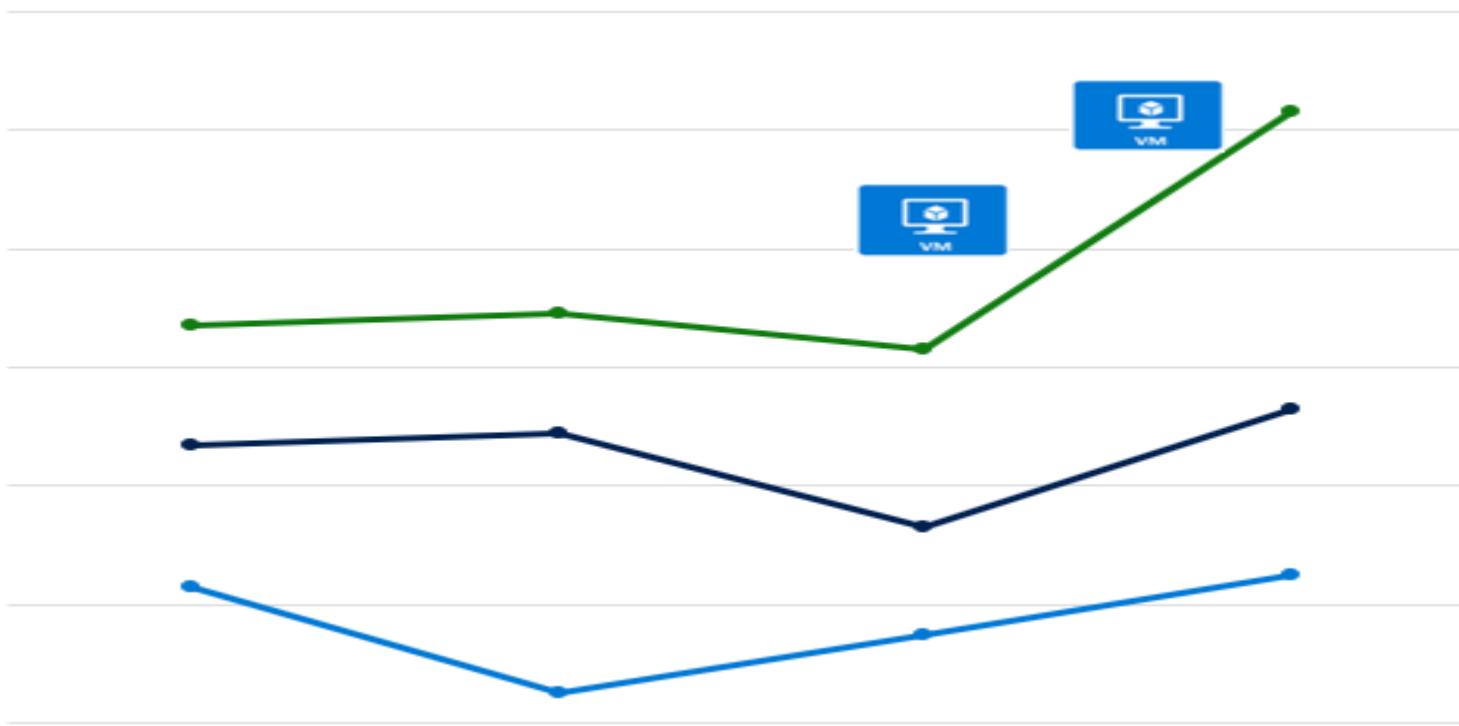
Capital Expenditure (CapEx)

- High upfront cost, value of investment reduces over time.

Operational Expenditure (OpEx)

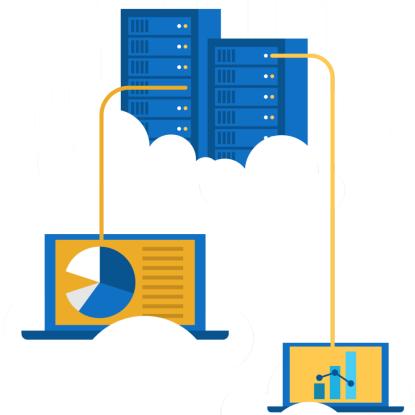
- Spend on services or products as needed.
- No upfront cost, pay-as-you use.

Define consumption-based model



Consumption-based model = Pay only for the resources you use

Distinguish types of cloud models



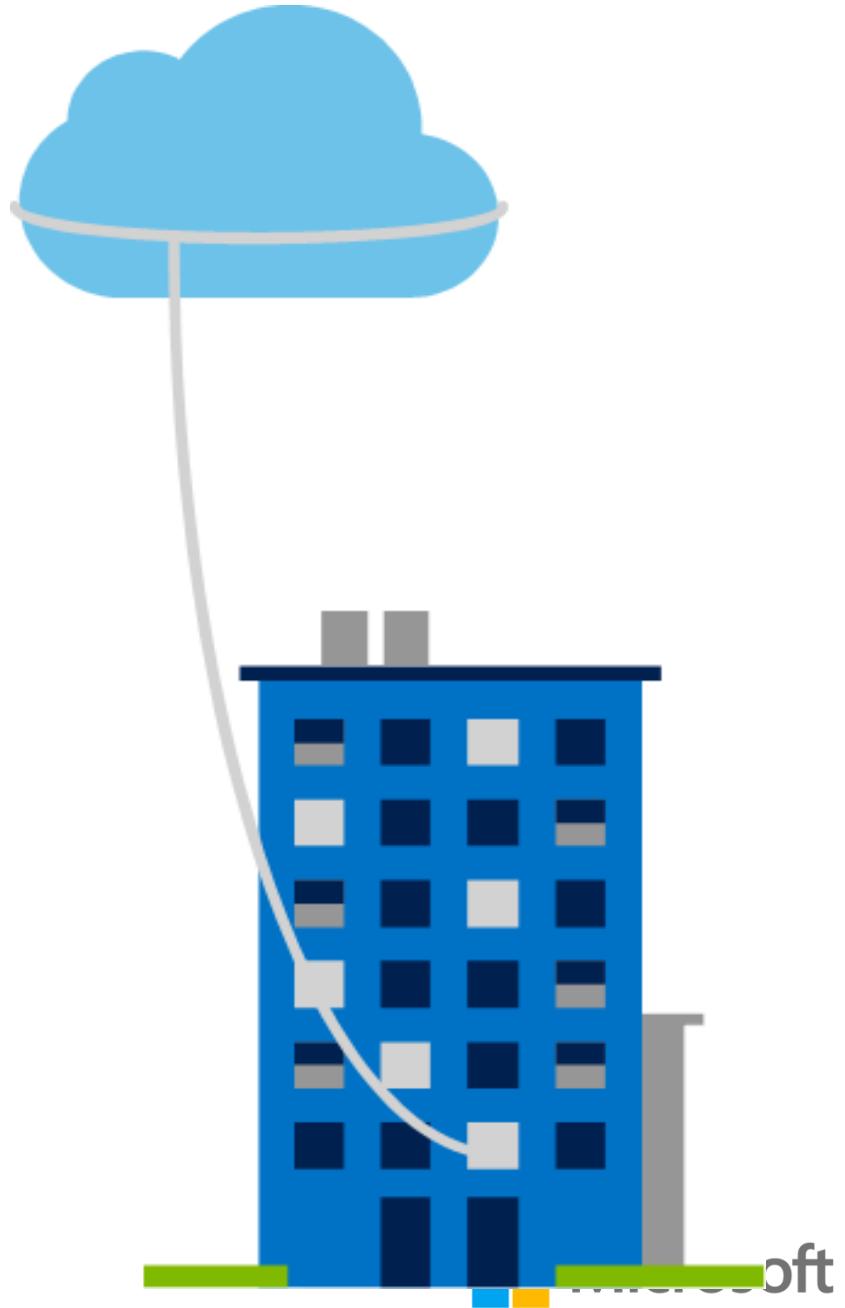
public cloud



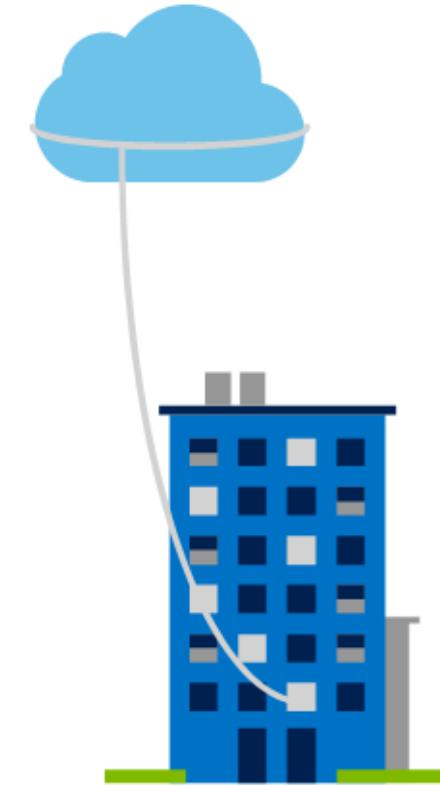
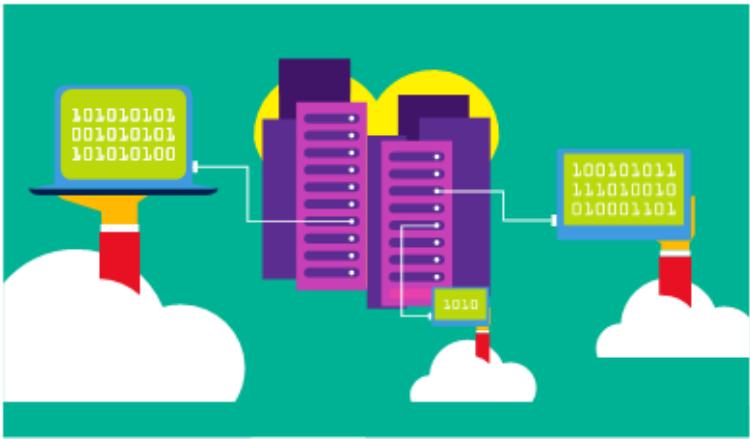
- Owned by cloud services or *hosting* provider.
- Provides resources and services to multiple organizations and users.
- Accessed via secure network connection (typically over the internet).

private cloud

- Organizations create a cloud environment in their datacenter.
- Organizations responsible for operating the services they provide.



hybrid cloud



Combines *Public* and *Private* clouds to allow applications to run in the most appropriate location.

Compare cloud models

Public cloud:

- No capital expenditures to scale up.
- Applications can be quickly provisioned and deprovisioned.
- Organizations pay only for what they use.

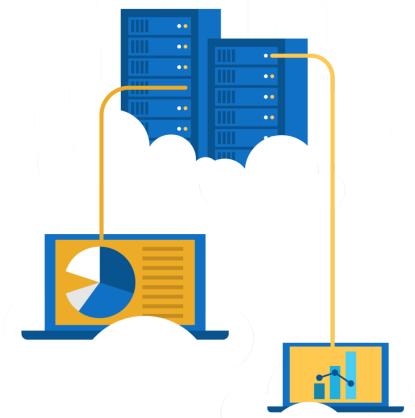
Private cloud:

- Organizations have complete control over resources.
- Organizations have complete control over security.

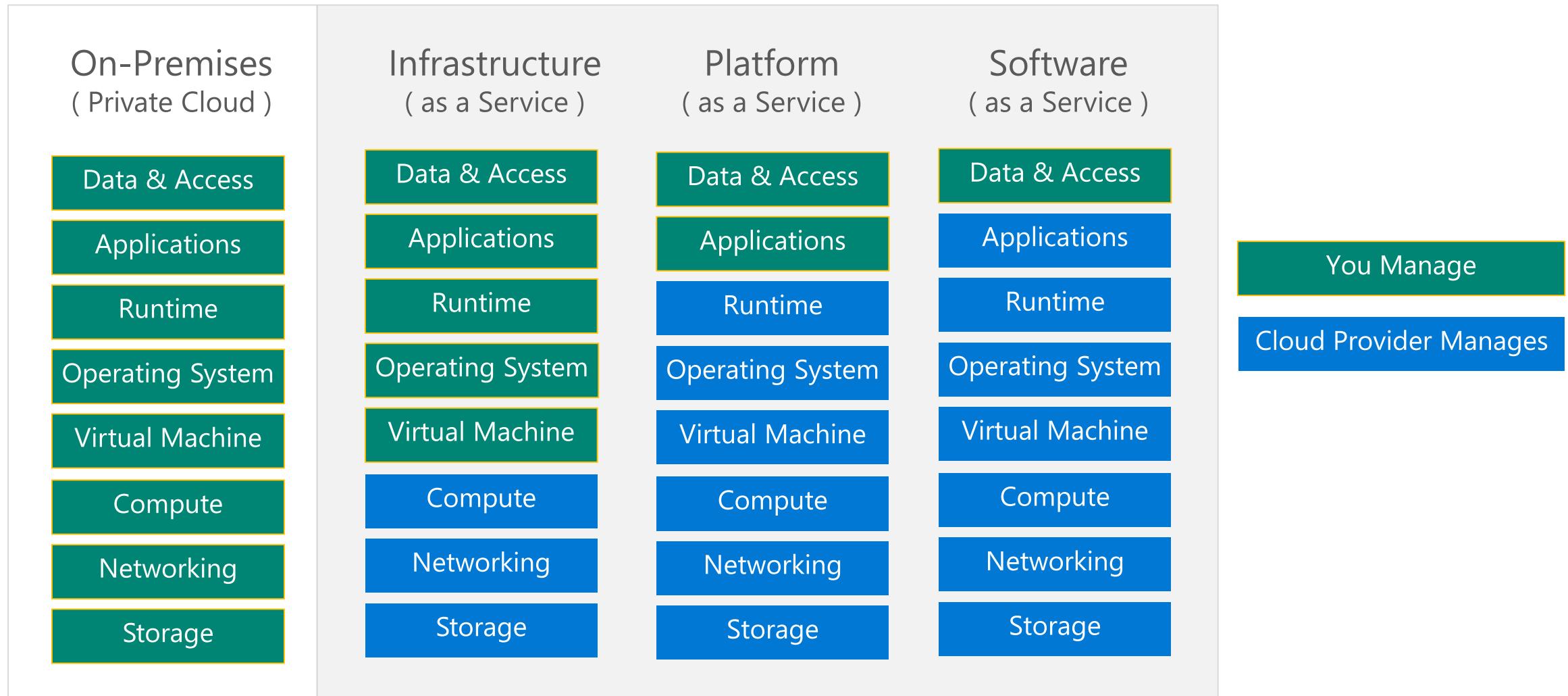
Hybrid cloud:

- Most flexibility.
- Organizations determine where to run their applications.
- Organizations control security, compliance, or legal requirements.

Explore types of cloud services

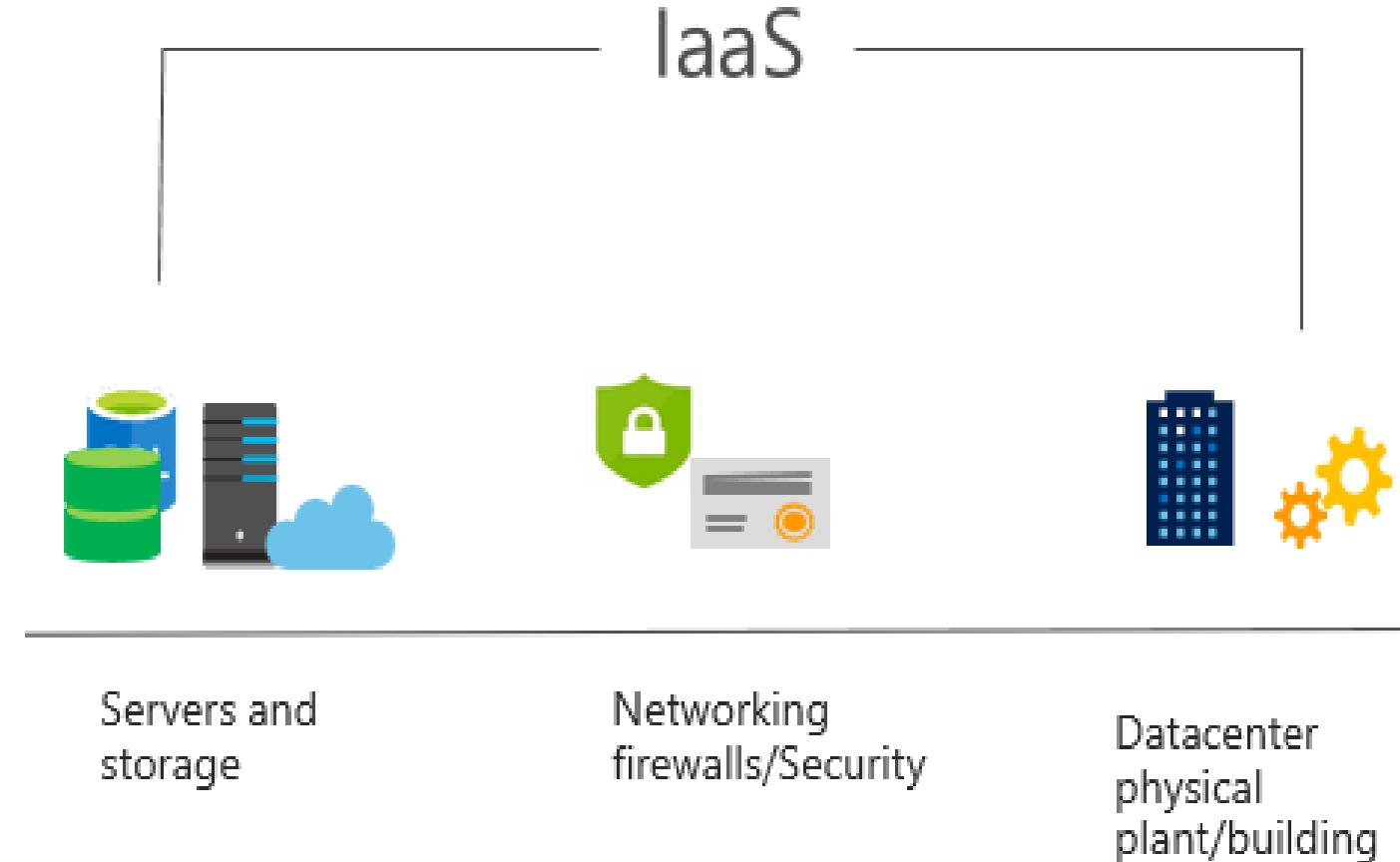


Discuss shared responsibility model

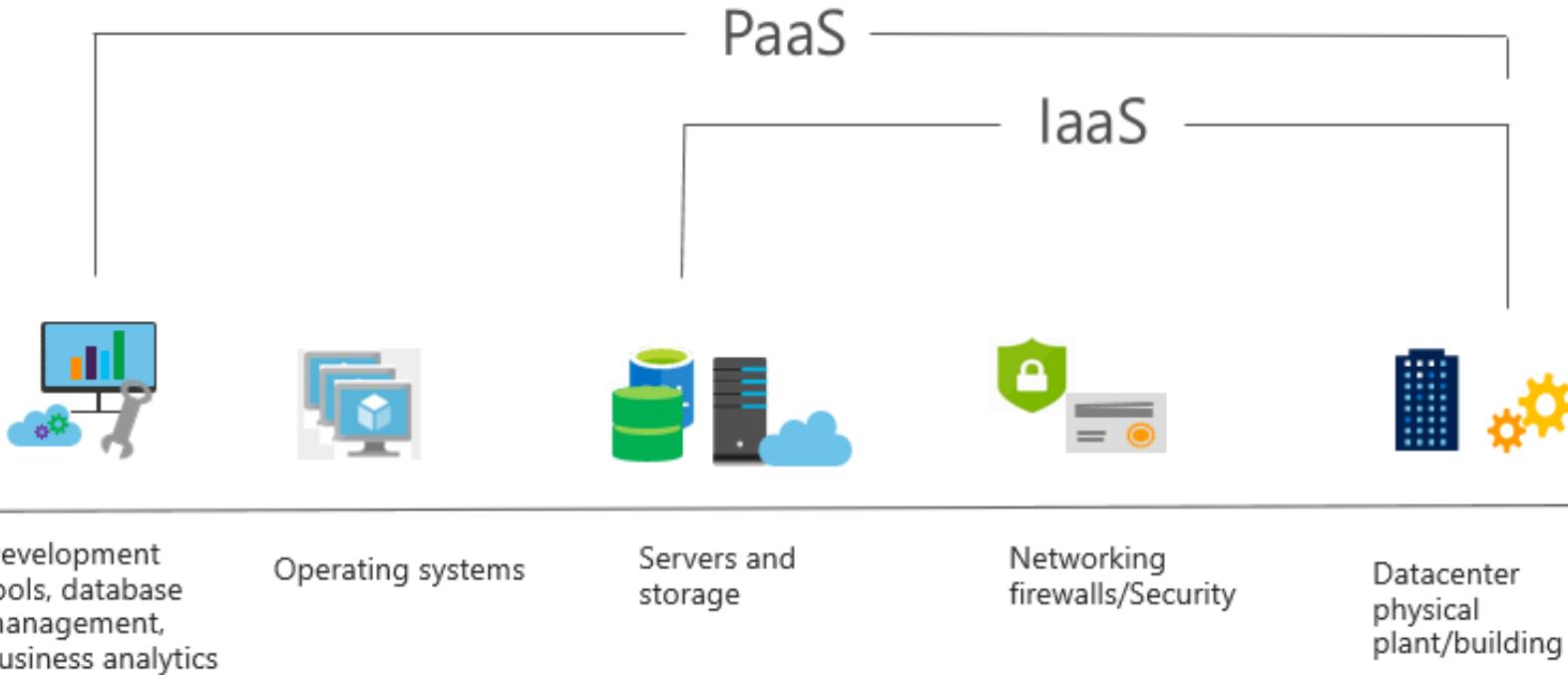


Define Infrastructure as a Service (IaaS)

Build pay-as-you-go IT infrastructure by renting servers, virtual machines, storage, networks, and operating systems from a cloud provider.

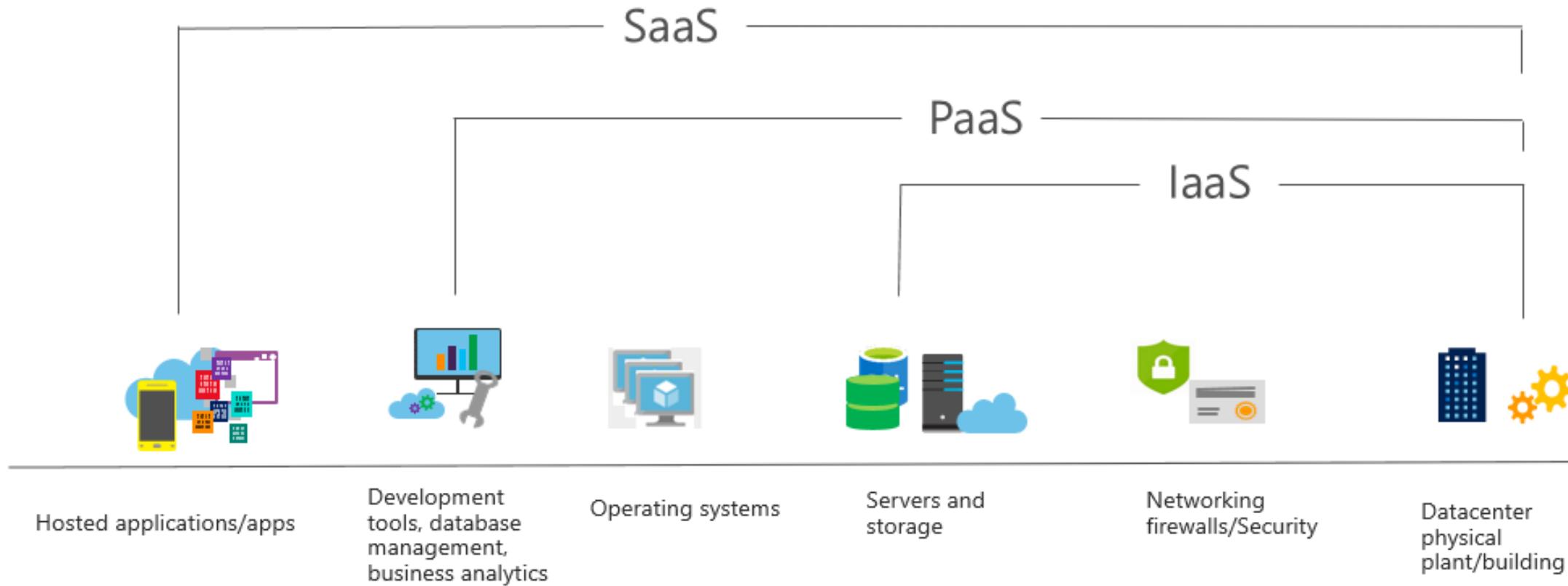


Define Platform as a Service (PaaS)



Provides environment for building, testing, and deploying software applications; without focusing on managing underlying infrastructure.

Define Software as a Service (SaaS)



Users connect to and use cloud-based apps over the internet: for example, Microsoft Office 365, email, and calendars.

Compare cloud services

IaaS

- The most flexible cloud service.
- You configure and manage the hardware for your application.

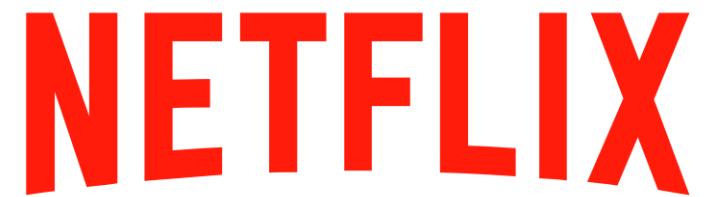
PaaS

- Focus on application development.
- Platform management is handled by the cloud provider.

SaaS

- Pay-as-you-go pricing model.
- Users pay for the software they use on a subscription model.

- Finance and Banking
- Telecommunication
- Retail
- Healthcare
- Agriculture
- Entertainment
- Real Estate

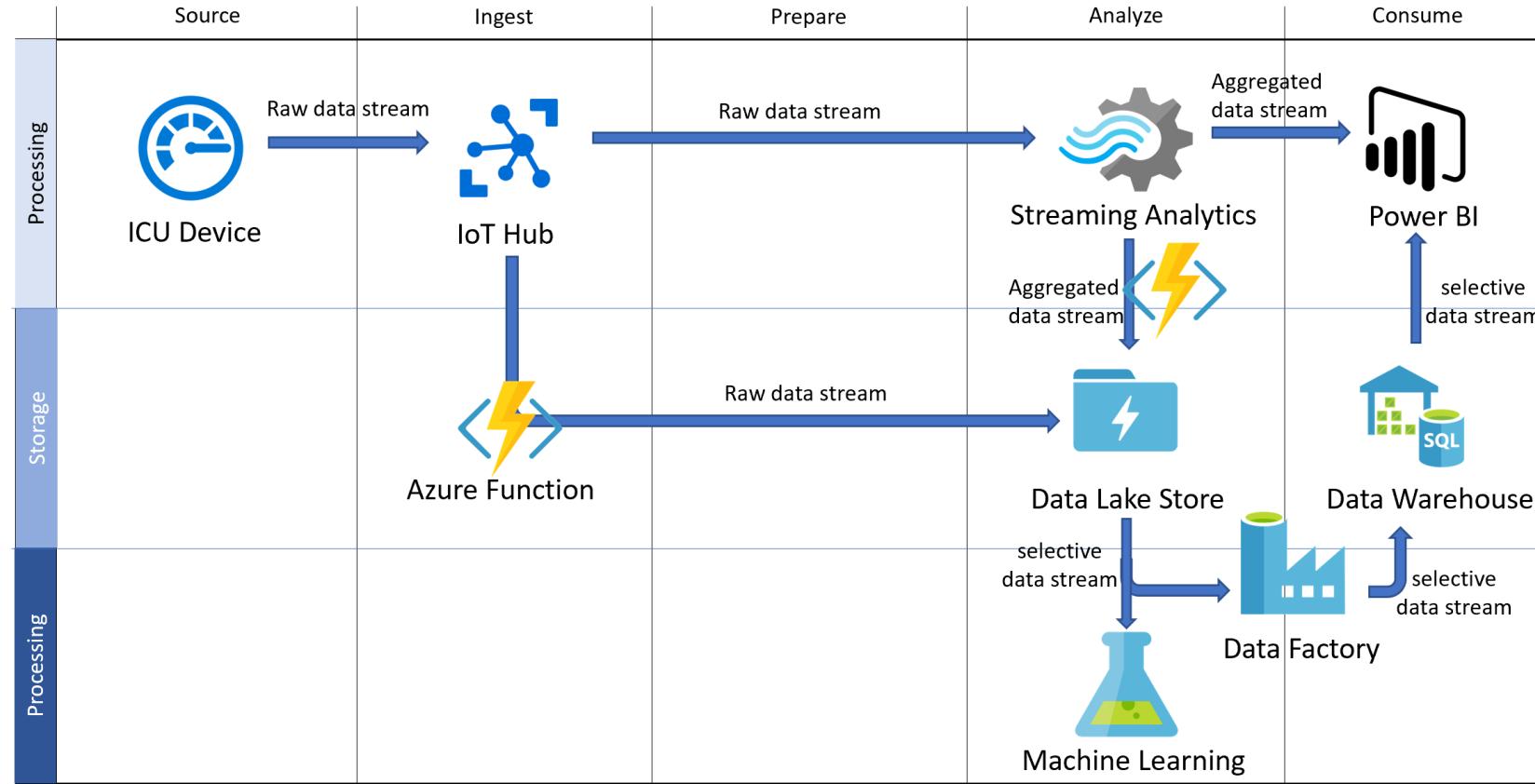


Big Data Use Case & Success Stories

- Data privacy is responsibly collecting, using and storing data about people, in line with the expectations of those people, your customers, regulations and laws.
 - Data ethics is doing the right thing with data, considering the human impact from all sides, and making decisions based on your brand values.

Privacy and Ethics of Big Data

Content Reference : <https://looker.com/blog/big-data-ethics-privacy>



Big Data Project – Case Example (Azure)

What to use for Data



Storage Account



- When you need a **low cost, high throughput** data store.
- When you need to store **No-SQL** data.
- When you **do not need to query** the data directly. **No ad hoc query** support.
- Suits the storage of archive or **relatively static data**.
- Suits acting as a **HDInsight Hadoop** data store.



Data Lake Store



- When you need a **low cost, high throughput** data store.
- **Unlimited storage** for **No-SQL** data
- When you **do not need to query** the data directly. **No ad hoc query** support.
- Suits the storage of archive or **relatively static data**.
- Suits acting as a **Databricks**, **HDInsight** and **IoT** data store.



Azure Databricks



- **Eases the deployment** of a Spark based cluster.
- Enables the **fastest processing** of Machine Learning solutions.
- **Enables collaboration** between data engineers and data scientists.
- Provides **tight enterprise security integration** with Azure Active Directory
- **Integration with other Azure Services** and **Power BI**.



Azure CosmosDB



- Provides **global distribution** for both structured and unstructured data stores.
- **Millisecond query response time**.
- **99.999% availability** of data.
- **Worldwide elastic scale** of both the storage and throughput
- **Multiple consistency levels** to control data integrity with concurrency



Azure SQL Database



- When you require a **relational** data store.
- When you need to manage **transactional workloads**
- When you need to manage a **high volume on inserts and reads**
- When you need a service that **requires high concurrency**
- When you require a solution that can scale **elastically**

What to use for Data



Azure Synapse Analytics



- When you require an integrated **relational** and **big data store**.
- When you need to manage **data warehouse** and **analytical workloads**.
- When you need **low cost storage**.
- When you require the ability to **pause and restart the compute**.
- When you require a solution that can scale **elastically**



Azure Stream Analytics



- When you require a **fully managed event processing** engine.
- When you require **temporal analysis of streaming** data.
- Support for analyzing **IoT streaming** data.
- Support for analyzing application data through **Event Hubs**.
- Ease of use with a **Stream Analytics Query Language**.



Azure Data Factory



- When you want to **orchestrate the batch movement** of data.
- When you want to connect to a **wide range of data platforms**.
- When you want to **transform or enrich** the data in movement.
- When you want to **integrate with SSIS packages**.
- Enables **verbose logging** of data processing activities.



Azure HDInsight



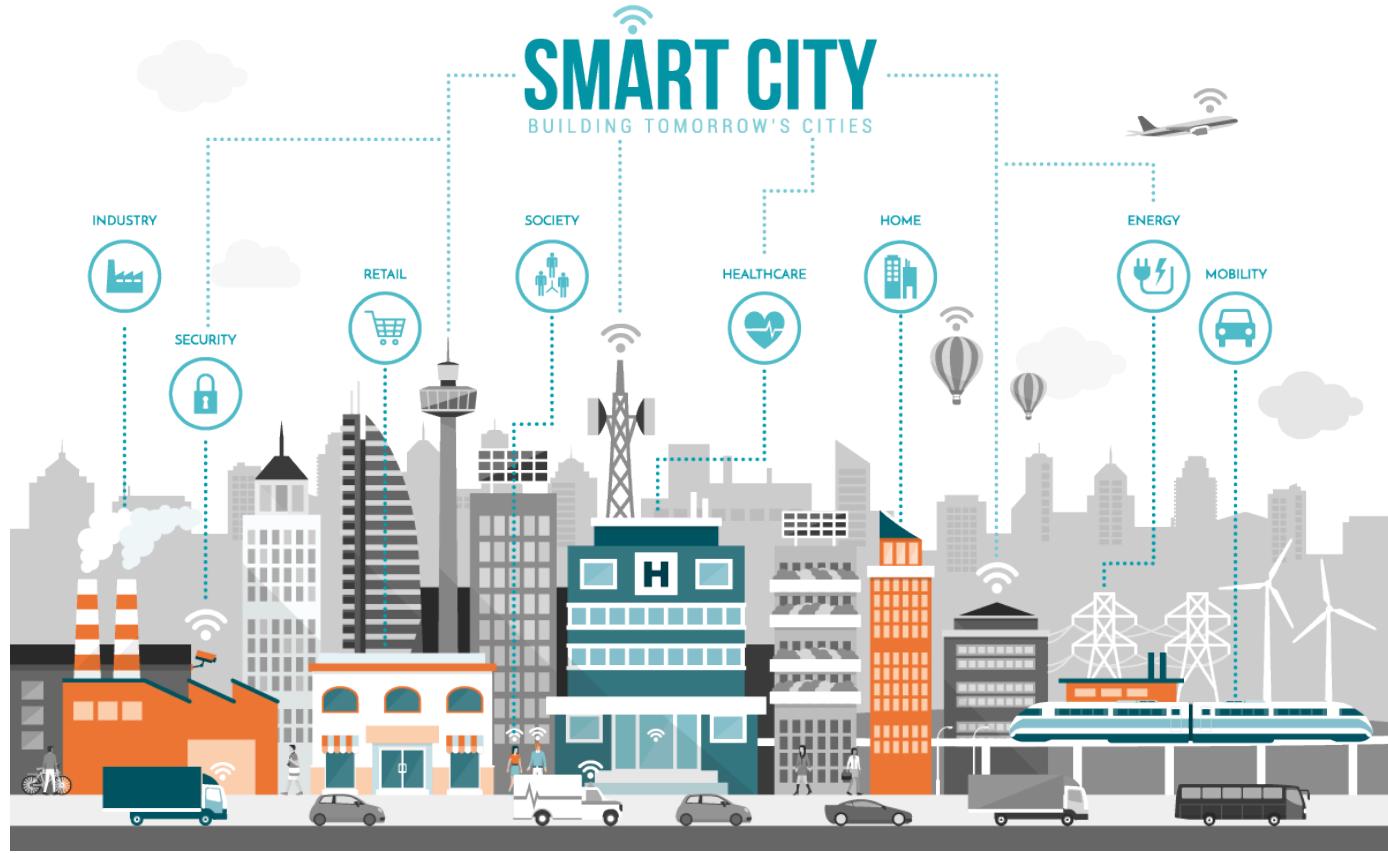
- When you need a **low cost, high throughput** data store.
- When you need to store **No-SQL** data.
- Provides a Hadoop **Platform as a Service** approach
- Suits acting as a **Hadoop, Hbase, Storm or Kafka** data store.
- **Eases the deployment and management** of clusters.



Azure Data Catalog



- When you require **documentation** of your data stores.
- When you require a **multi user** approach to documentation.
- When you need to **annotate data sources** with descriptive metadata.
- A **fully managed cloud service** whose users can discover the data sources.
- When you require a **solution that can help business users** understand their data.



Smart City

Content Reference : <https://www.arcweb.com/industries/smart-cities>



Smart Farm

Content Reference : <https://www.luda.farm/>

04 - Big Data Sources

- Enterprise Data are functionally different application supporting the business operations in various departments.
- Enterprise Data Sources within an organization integrate and retrieve data for both internal applications and external communication.
- Enterprise Data Sources ensure trust and confidence in data assets.

Introduction to Enterprise Data Source

- The vast majority of contemporary Enterprise Systems are based on relational database technology and, consequently, provide well-structured datasets that can be easily accessed.
- Enterprise Systems are Oracle, SQL Server, DB2, SAP,
- Use ODBC, JDBC, OLEDE, ... For access to dataset.



Enterprise System

ORACLE® Sales Results versus Forecast

The screenshot displays the Oracle Sales Management interface. At the top, there's a navigation bar with 'Sales Management > Expense Mana' and a date range 'Q2 FY05 Day -6 24-Jun-20'. Below this are sections for 'Report Views' (Jacques' EMEA View, Jacques view 1, Key Acct View, My Personal View), 'Period' (Quarter), 'Compare' (Current vs. Last Year), and 'Actions' (Printable Page, Send an Email, Export to a File, Delegate, Save View, Edit View, Delete View, Reset Parameter Default Values, Personalize Links). The main area shows three charts: 'Forecast, Won' (bar chart for Industry Accounts, Key Accounts, Mid Market Accounts, Partner Accounts), 'Forecast, Won Change' (bar chart for Industry Accounts, Key Accounts, Mid Market Accounts, Partner Accounts), and a table titled 'Sales Group Forecast by Product Category' with columns like Won % of Forecast, Net Booked (K), Revenue (K), and Revenue % of Forecast. A sidebar on the left lists various sales-related links such as Sales Group Forecast by Product Category, Leads, Opportunities and Backlog, Forecast Overview, Opportunity Win/Loss, Opportunity Win/Loss (with Counts), Opportunity Activity, Weighted Pipeline, Forecast versus Won Trend, Forecast, Pipeline, Won Trend, Extended Forecast versus Won Trend, Extended Forecast versus Pipeline Trend, Pipeline Trend, and Win/Loss Trend. A bottom link 'Sales Forecast Management' is also present.

Flexible Time Periods & Comparisons

Report Views
Jacques' EMEA View
Jacques view 1
Key Acct View
My Personal View

Period Quarter

Compare Current vs. Last Year

Actions

- Printable Page
- Send an Email
- Export to a File
- Delegate
- Save View
- Edit View
- Delete View
- Reset Parameter Default Values
- Personalize Links

Forecast, Won

Forecast, Won Change

Personalize and Save Favorite Queries

Graphs and Trends

Flexible, Pre-Built Analytics

Sales Group Forecast by Product Category

View By	Won % of Forecast	Net Booked (K)	Revenue (K)	Revenue % of Forecast
Sales Group	72.4%	134.7%	4,233	-22.2%
Product Category	1,077	5,118	529.7%	
%	75	91	217.1%	
	384	222		
	1,536	204	567.1%	

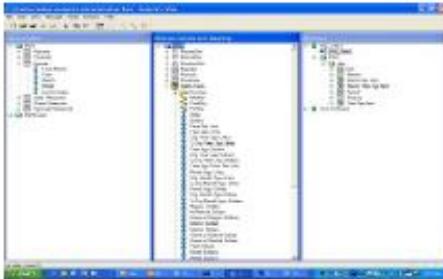
Drill & Pivot to explore data

Links for Guided Analysis

- Oracle provides a range of “Enterprise Application” including HR/Payroll, Finance, Customer Relationship Management (CRM).
- Oracle suite of application is Oracle E-Business Suite, PeopleSoft and JD Edwards.

Enterprise System - Oracle

Define



Analyze



Report



Consistent
Multi-Source
Intelligence

BI EE Analytic Server

- Enterprise Semantic Model
- Multi-Source Data Warehouse Model
- Pre-Packaged Dashboards
- Pre-Packaged ETL Maps

ORACLE®

SIEBEL



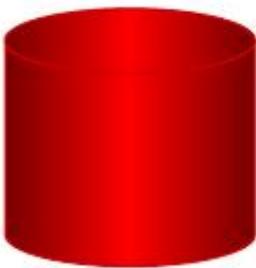
ORACLE®

PEOPLESOFT ENTERPRISE



ORACLE®

E-BUSINESS SUITE



ORACLE®



- SAP providers enterprise application that include modules such as HR/Payroll, Finance, and SAP Suite of offering also includes a number of vertical application.
- Extracting data from SAP is more difficult than other enterprise system since SAP users a proprietary language, ABAP, for data manipulation. Additionally, SAP database use tables which are different from standard database.
- In order to extract data from SAP we need to use additional tools.

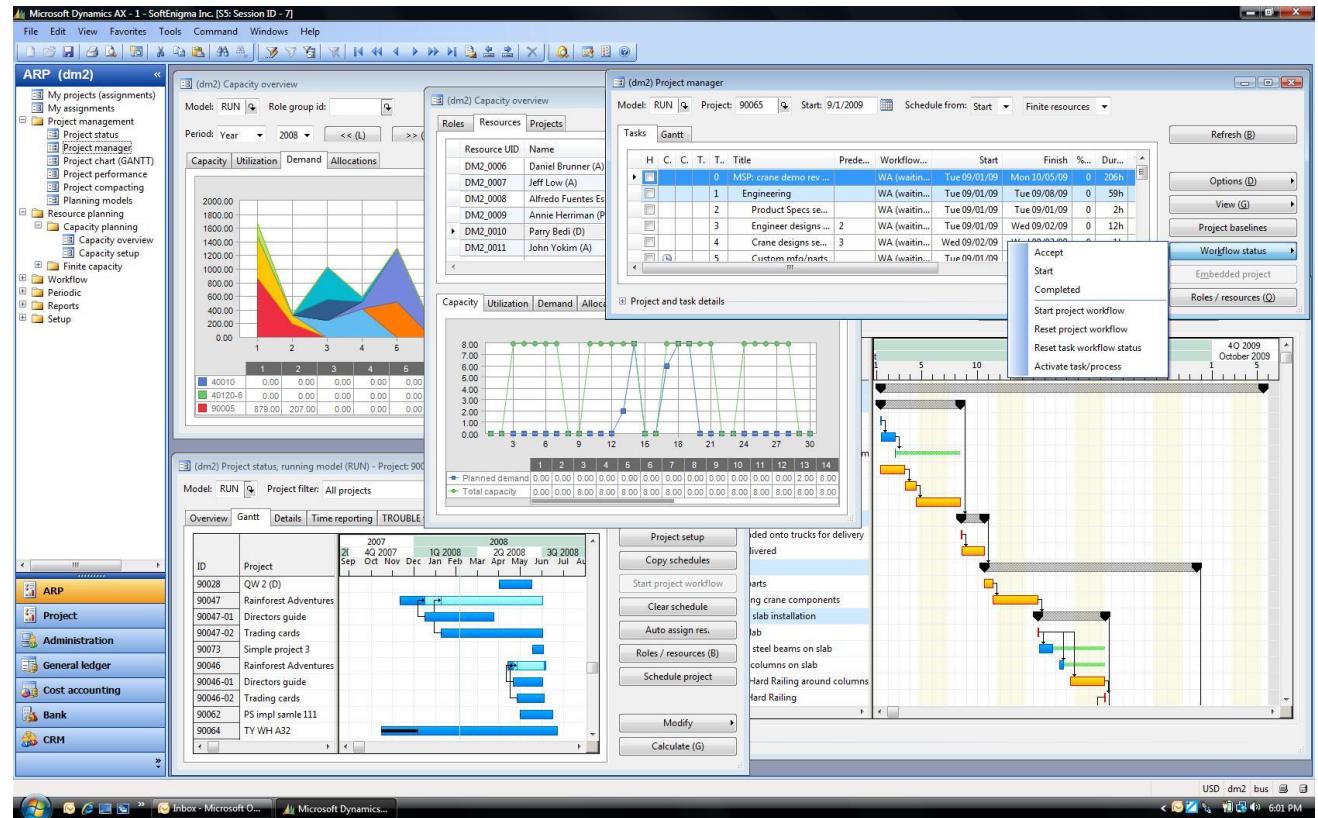
- SAP Connector from Oracle
- SAP Business Connector
- SAP Java Connector
- SAP .NET Connector from SAP



Enterprise System - SAP

- Microsoft provides a set of enterprise application through its Dynamic suite. The suite includes CRM and ERP modules and is targeted at small and medium enterprises.

- The Microsoft Dynamics Suite is based on relational database technology (SQL Server)



Enterprise System - Microsoft

- Defined in many different ways, but not rigorously.
 - A decision support database that is maintained separately from the organization's operational database
 - Support information processing by providing a solid platform of consolidated, historical data for analysis.
- "A data warehouse is a subject-oriented, integrated, time-variant, and nonvolatile collection of data in support of management's decision-making process."—W. H. Inmon
- Data warehousing: The process of constructing and using data warehouses

Data Warehouse

- Organized around major subjects, such as customer, product, sales
- Focusing on the modeling and analysis of data for decision makers, not on daily operations or transaction processing
- Provide a simple and concise view around particular subject issues by excluding data that are not useful in the decision support process

Data Warehouse : Subject-Oriented

- Constructed by integrating multiple, heterogeneous data sources
 - relational databases, flat files, on-line transaction records
- Data cleaning and data integration techniques are applied.
 - Ensure consistency in naming conventions, encoding structures, attribute measures, etc. among different data sources E.g., Hotel price: currency, tax, breakfast covered, etc.
 - When data is moved to the warehouse, it is converted.

Data Warehouse : Integrated

- The time horizon for the data warehouse is significantly longer than that of operational systems
 - Operational database: current value data
 - Data warehouse data: provide information from a historical perspective (e.g., past 5-10 years)
- Every key structure in the data warehouse
 - Contains an element of time, explicitly or implicitly
 - But the key of operational data may or may not contain “time element”

Data Warehouse : Time Variant

- A physically separate store of data transformed from the operational environment
- Operational update of data does not occur in the data warehouse environment
 - Does not require transaction processing, recovery, and concurrency control mechanisms
 - Requires only two operations in data accessing: initial loading of data and access of data

Data Warehouse : Nonvolatile

	OLTP	OLAP
users	clerk, IT professional	knowledge worker
function	day to day operations	decision support
DB design	application-oriented	subject-oriented
data	current, up-to-date detailed, flat relational isolated	historical, summarized, multidimensional integrated, consolidated
usage	repetitive	ad-hoc
access	read/write index/hash on prim. key	lots of scans
unit of work	short, simple transaction	complex query
# records accessed	tens	millions
#users	thousands	hundreds
DB size	100MB-GB	100GB-TB
metric	transaction throughput	query throughput, response

Online Transactional Processing (OLTP)

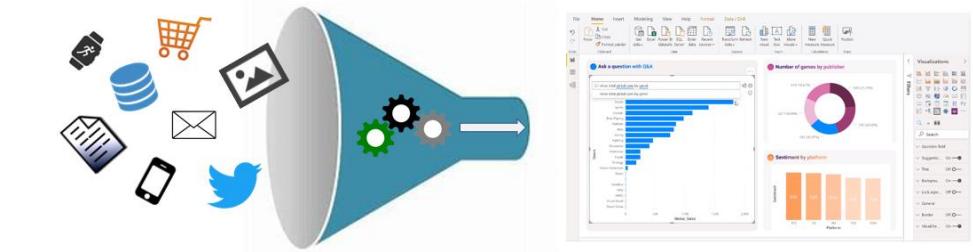
Customer

CustomerID	CustomerName	CustomerPhone

Orders

OrderID	CustomerID	OrderDate

Online Analytical Processing (OLAP)

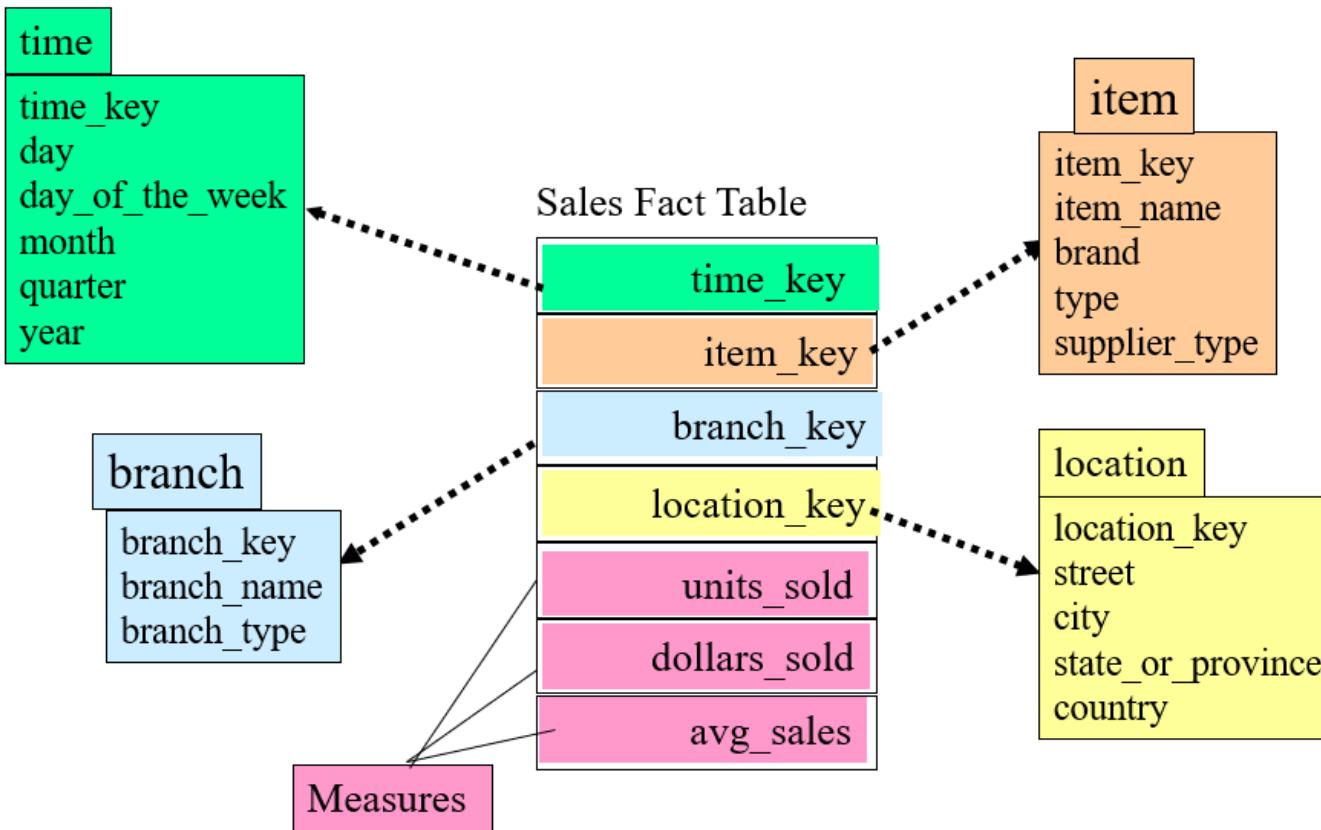


Transactional vs analytical data stores

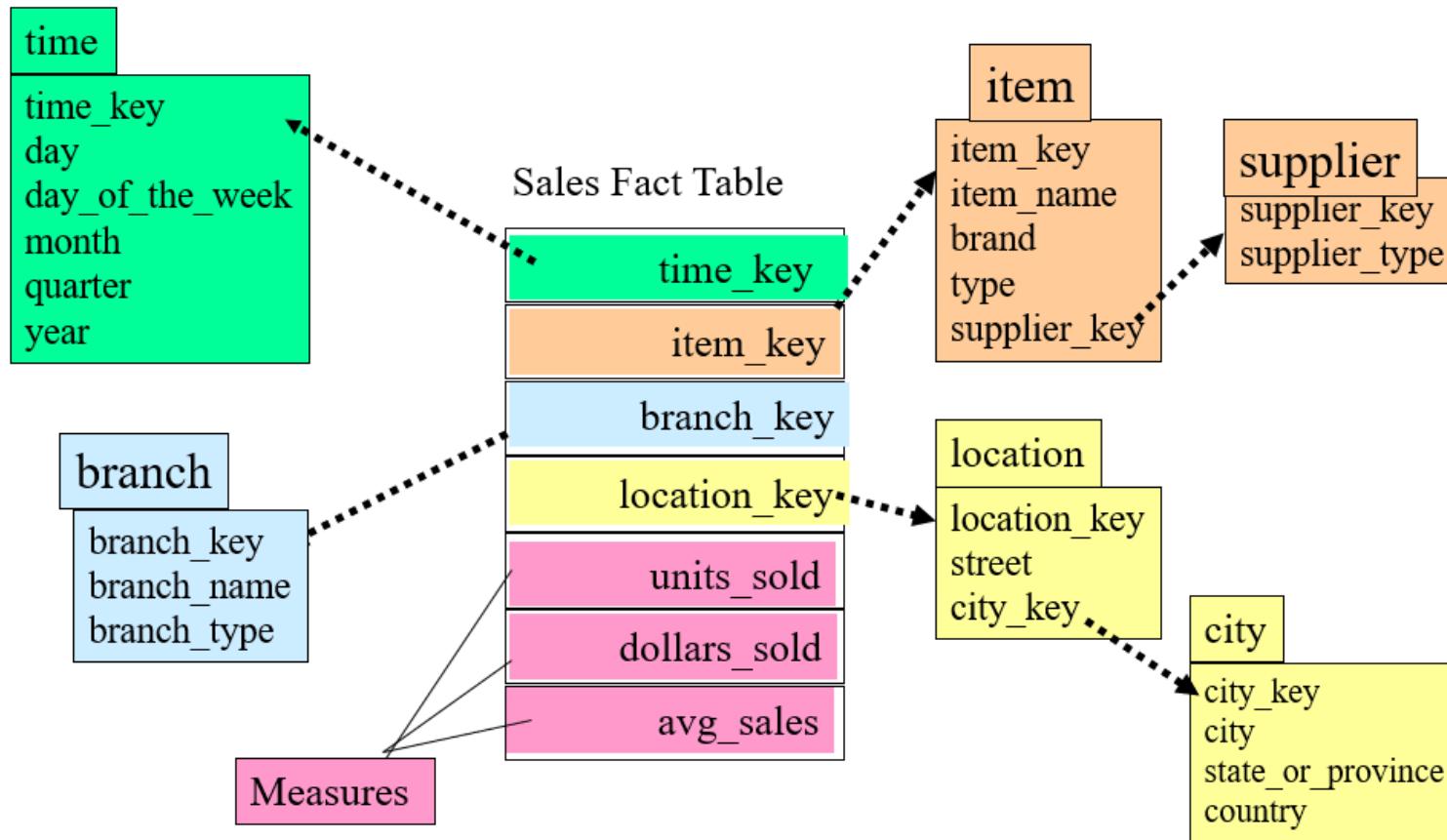
Modeling data warehouses: dimensions & measures

- **Star schema:** A fact table in the middle connected to a set of dimension tables
- **Snowflake schema:** A refinement of star schema where some dimensional hierarchy is normalized into a set of smaller dimension tables, forming a shape similar to snowflake
- **Fact constellations:** Multiple fact tables share dimension tables, viewed as a collection of stars, therefore called galaxy schema or fact constellation

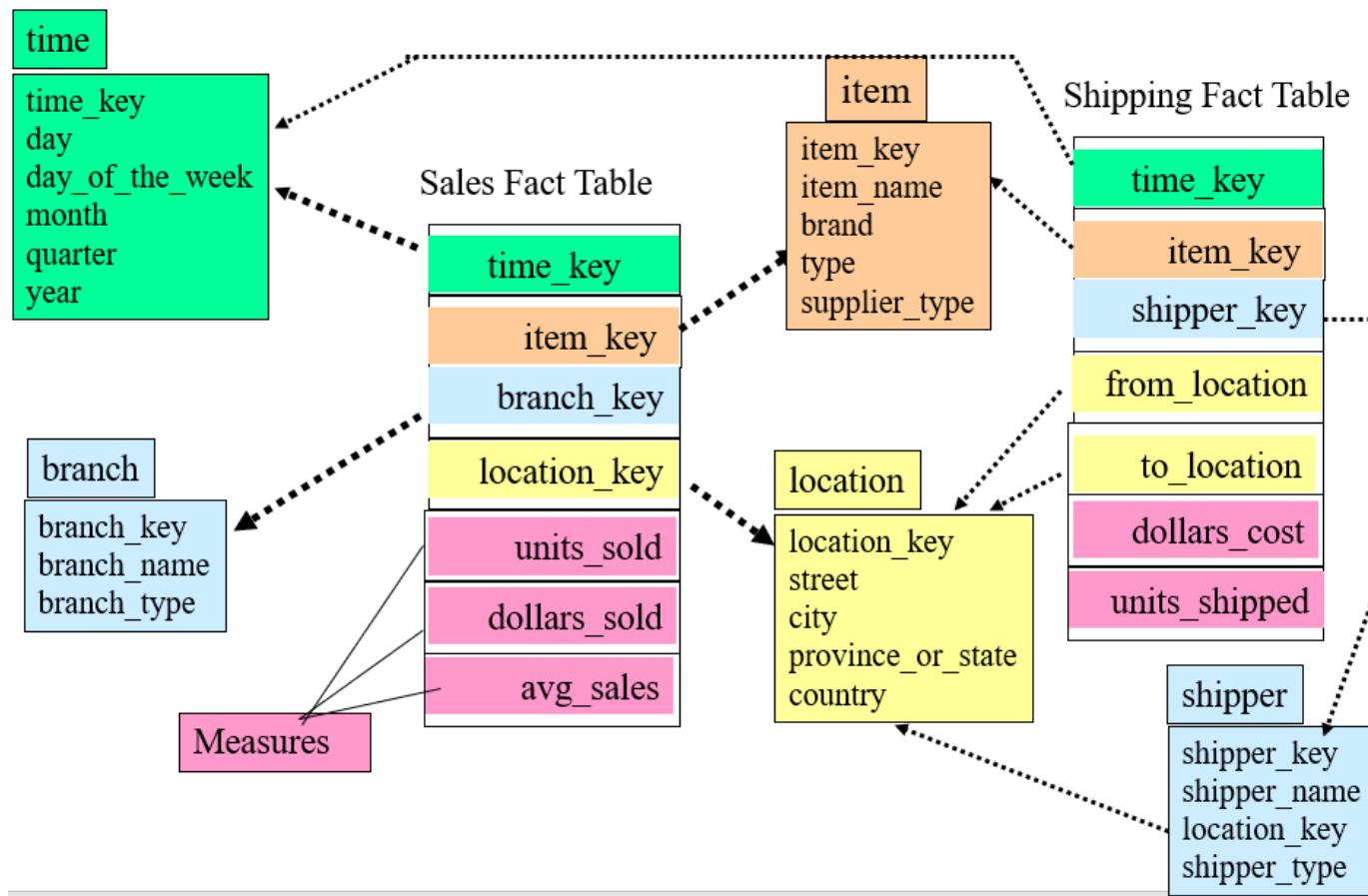
Conceptual Modeling of Data Warehouses



Example of Star Schema

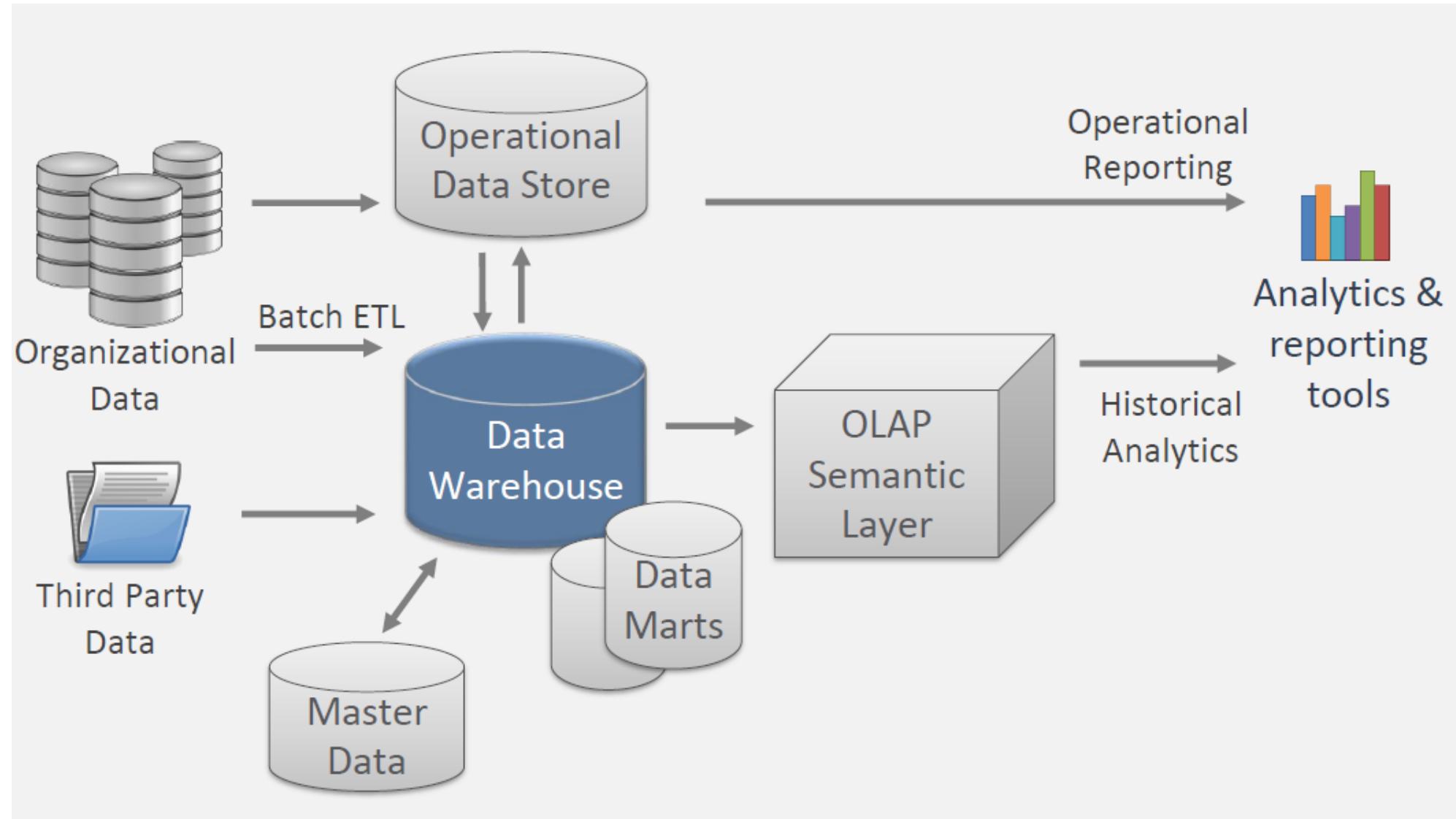


Example of Snowflake Schema

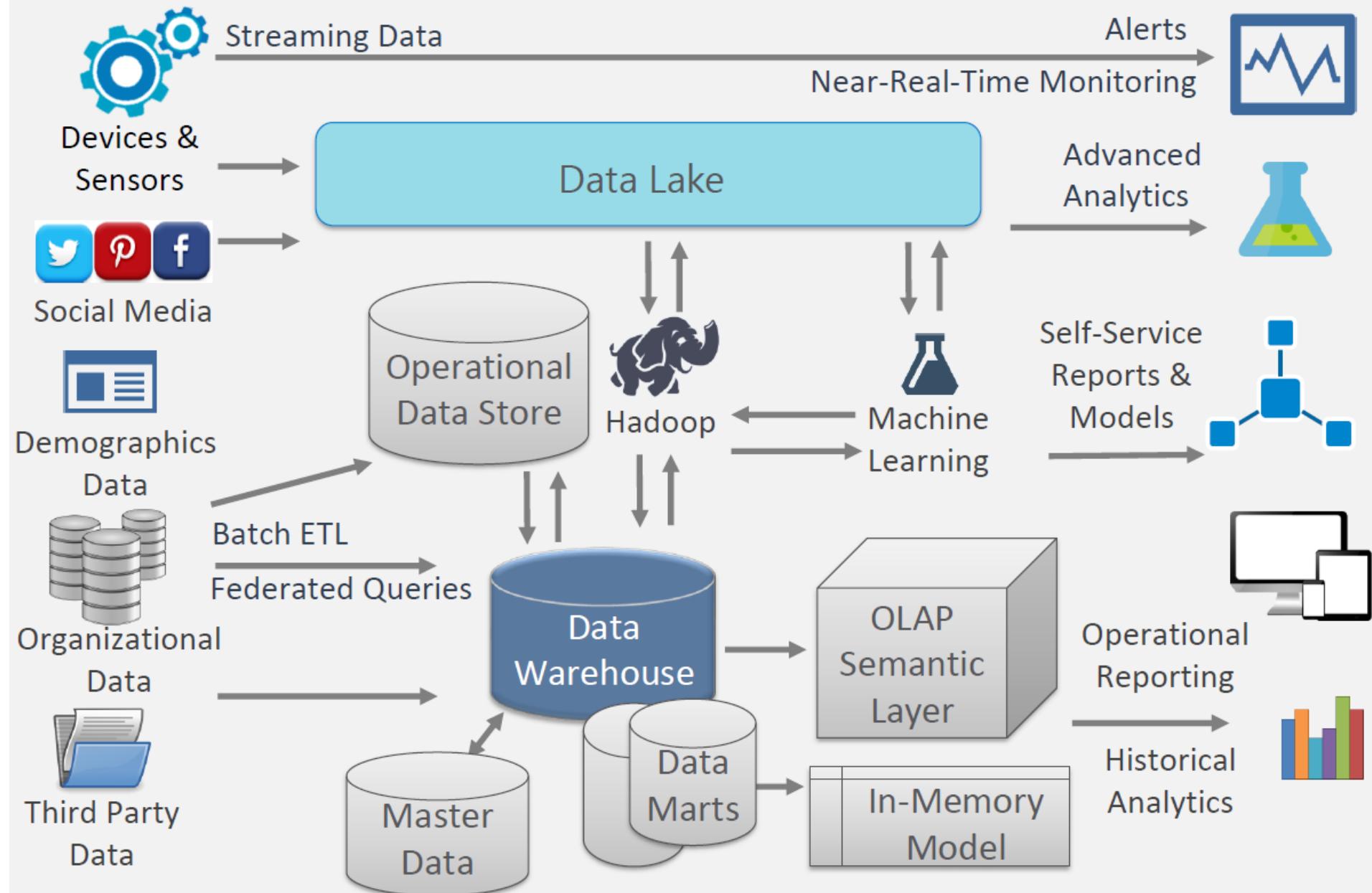


Example of Fact Constellation

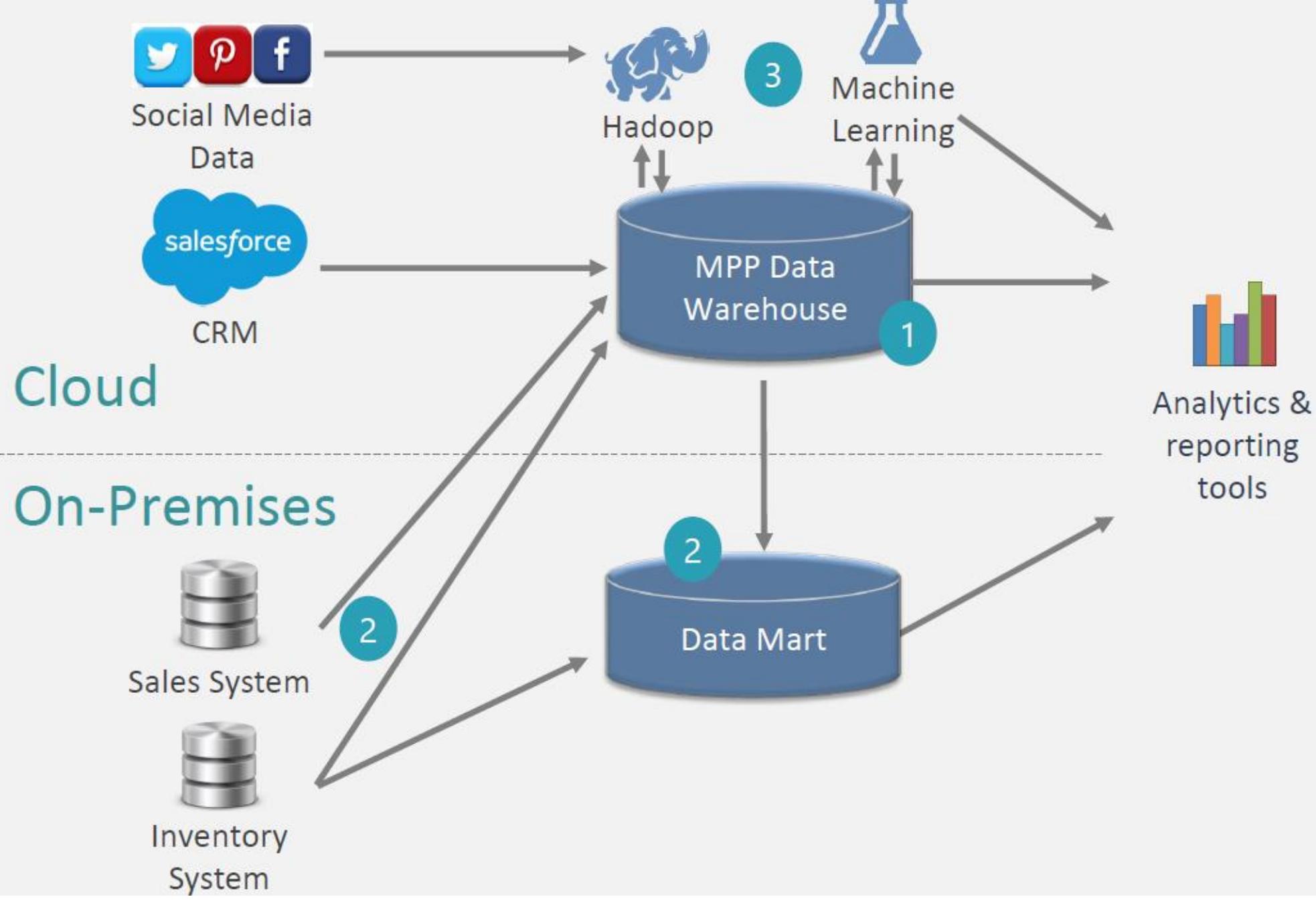
Traditional Data Warehousing



Modern Data Warehouses

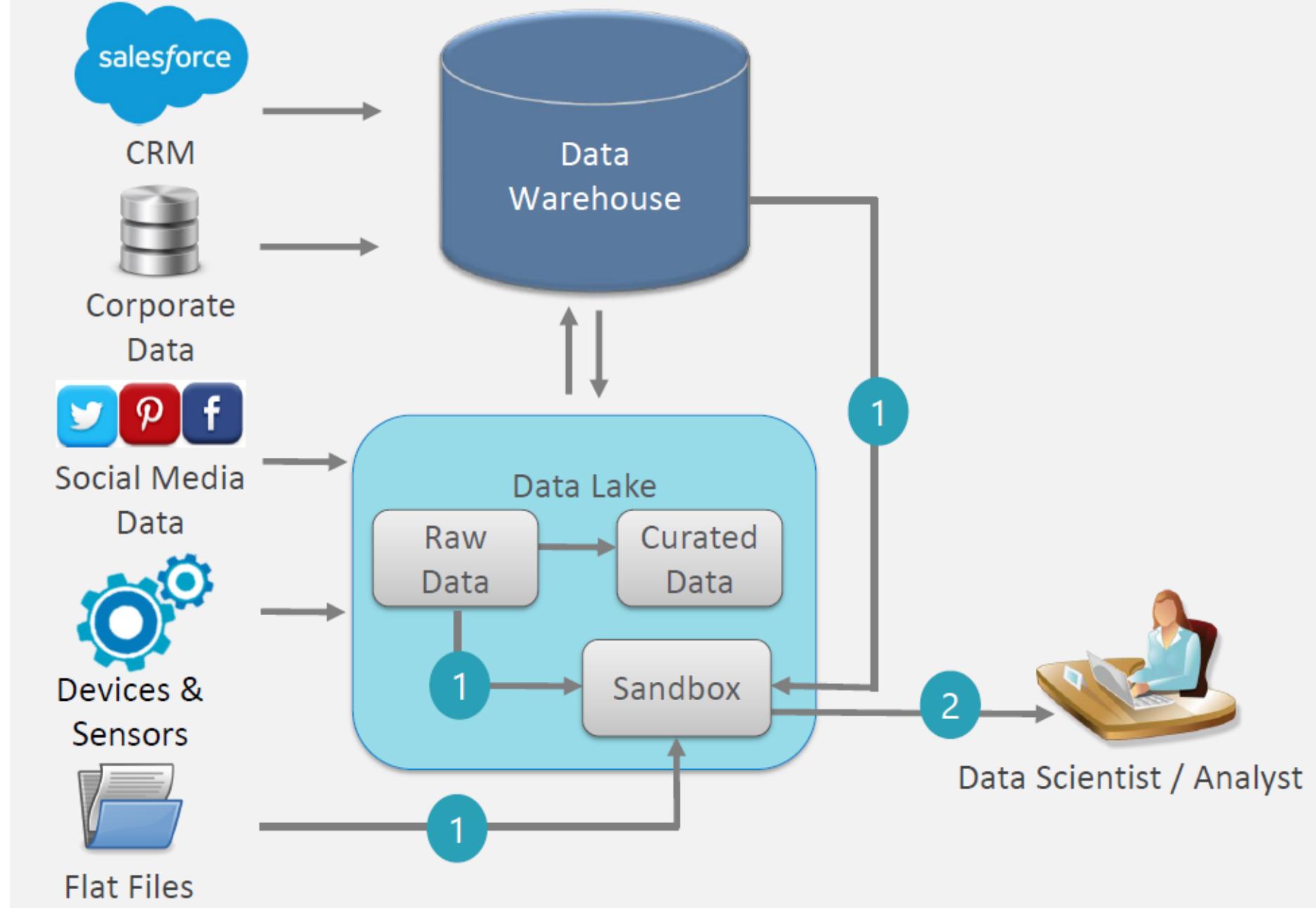


Hybrid Architecture



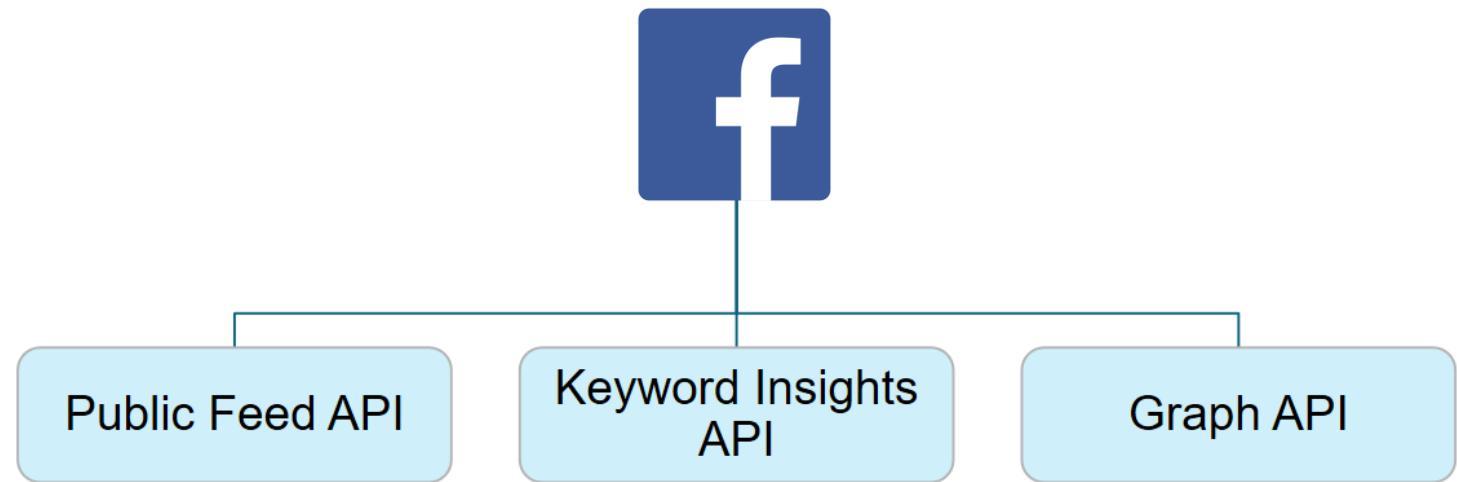
Content Reference : Melissa Coates [Blog: sqlchick.com/Twitter: @sqlchick]

Sandbox Solutions



Content Reference : Melissa Coates [Blog: sqlchick.com/Twitter: @sqlchick]

- Facebook provides several ways to access its data streams. The access is done through separate Application Programming Interface (APIs)
- Some APIs are available only to certain organizations and currently cannot be accessed by the general public.



Social Media Data Sources - Facebook

API ກາຣຕລາດ[ກາຣໃຊ້ API](#)[ຫັ້ນໝາໂພ້ນາ](#)[ກາຮປະມູລແລກປາກປົກປິກ](#)[ກາຮການແດກລຸ່ມເປົ້າທາຍາ](#)[ຫຼັມມູລເຊີງລຶກຂອງໂພ້ນາ](#)**API ຫຼັມມູລເຊີງລຶກ**[ພາຣາມີເຕେର](#)[ຫຼັມມູລແກຍຍ່ອຍ](#)[Estimated & In-Development](#)[ຫຼັດຈຳກັດແລກປົງບັດທີ່ສຸດ](#)[ກາຣຕິດຕາມກາຣຄລິກແລກສົກຕິ](#)[ຄື່ງໄວ້ຮົດ](#)[ຄອນເວຼອຮັ້ນແບບອຳໂລ່ນ](#)

ຫຼັມມູລເຊີງລຶກ

ໃຫ້ບັນດາອິນເທົ່ານີ້ທີ່ສອດຄລັງກັນສໍາຫັບເຮັດວຽກດູສົກຕິຂອງໂພ້ນາ

- ພາຣາມີເຕେର - ພາຣາມີເຕେରທີ່ມີໃຫ້ຈຳນັກໃນປະລາຍທາງນີ້
- ພິລົດ - ດ້ວຍເລືອກໃນ [fields](#)
- ຫຼັມມູລແກຍຍ່ອຍ - ພລັກພົກຂອງກຸ່ມ
- ຫຼັມມູລແກຍຍ່ອຍການດໍາເນີນການ - ທ່ານຈະເຫັນໃຈການຕອບສອນຈາກຫຼັມມູລແກຍຍ່ອຍການດໍາເນີນການ
- ຂາຍທີ່ມີເຊີງລຶກ - ສໍາຫັບຄໍາຂອງທີ່ມີຜົລັກພົກຈ້າງນາກ ໃຫ້ໃຫ້ຈຳນັກທີ່ມີເຊີງລຶກໃນໆ
- ຂີ່ຈຳຈັດແລກປົງບັດທີ່ສຸດ - ຂີ່ຈຳຈັດຂອງການເຮັກ ກາຣກຮອງ ແລກປົງບັດທີ່ສຸດ
- ຂົ່ວເກັນທີ່ຫຼັດແລກປົງບັດ - ຂົ່ວເກັນທີ່ຫຼັດໃນເຄື່ອງມືໂພ້ນານັ້ນ Facebook ພຮັດກັນຂອງ API ທີ່ສອດຄລັງກັນ

ກາຮເຮີມຕັ້ນໃຈໝາຍວ່າງຈ່າຍ

ກາຮເຂົ້າຄົ່ງການຮ່າງຈາກແລກປົງບັດທີ່ສຸດ ແລກປົງບັດທີ່ສຸດທີ່ມີສິນໃຈການຮ່າງຈາກແລກປົງບັດທີ່ສຸດ

ບານເພຈນີ້[ກາຮເຮີມຕັ້ນໃຈໝາຍວ່າງຈ່າຍ](#)[1. ສ່ວນແພ](#)[2. ສົກຕິແຄມເປັນ](#)[ທ່ານເຮັກ](#)[ຈະດັບ](#)[ຈຳການຮັບທີ່ມາ](#)[ກາຮບ່າຍຢ່ອງ](#)[ກາຮເຮີມຕັ້ນ](#)[ປ້າຍໂພ້ນາ](#)[ຄໍານີ້ຍຳມາເກີ່ວກັບຈຳນານຄລິກ](#)[ອືບເຈິດທີ່ຄົນແລ້ວແລກທີ່ເກີນກາວ](#)[ກາຮແກ້ໄຂປົງຫາ](#)[ແຫລ່ງຫຼັມມູລທີ່ເກີ່ວຂຶ້ນ](#)

Social Media Data Sources - Facebook



Formerly the National Climatic Data Center (NCDC)... [more about NCEI](#) »

[Home](#) [Climate Information](#) **Data Access** [Customer Support](#) [Contact](#) [About](#)

Search



Home > Data Access

Quick Links

- [Land-Based Station](#) ▾
- [Satellite](#) ▾
- [Radar](#) ▾
- [Model](#) ▾
- [Weather Balloon](#) ▾
- [Marine / Ocean](#) ▾
- [Paleoclimatology](#) ▾
- [Severe Weather](#) ▾
- [Blended & Global](#)

Data Access

For API data access use the NCEI suite of API services:

- Access Data Service API
- Access Search Service API
- Access Order Service API
- Access Support Service API

NCEI is the world's largest provider of weather and climate data. Land-based, marine, model, radar, weather balloon, satellite, and paleoclimatic are just a few of the types of datasets available. Detailed descriptions of the available products and platforms are below.

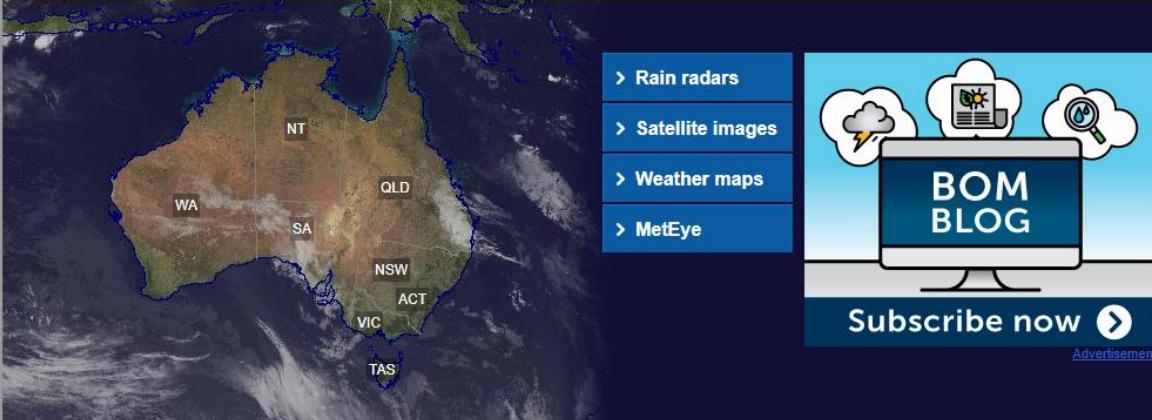
Public Data Sources - Weather



Warnings current

NSW VIC QLD WA SA TAS ACT NT

Warning services ▾



Forecast for Monday 12 October

City observations ▾

Sydney Now	Melbourne Now	Brisbane Now	Perth Now	Adelaide Now	Hobart Now	Canberra Now	Darwin Now
17.1° SSW 7km/h	6.5° NE 7km/h	17.6° W 2km/h	12.0° CALM 0km/h	18.4° NE 13km/h	6.3° NW 19km/h	11.6° SSE 6km/h	25.8° NNE 7km/h
16° 24° Partly cloudy.	7° 25° Mostly sunny.	17° 28° Partly cloudy.	11° 28° Mostly sunny.	15° 28° Partly cloudy.	6° 20° Possible shower.	10° 22° Partly cloudy.	25° 35° Sunny.

Public Data Sources - Weather

Content Reference : <http://www.bom.gov.au/>

New to this site? [Start Here](#)

[DataBank](#) [Microdata](#) [Data Catalog](#) 

World Bank Open Data

Free and open access to global development data

Search data e.g. GDP, population, Indonesia

Browse by [Country](#) or [Indicator](#)

MOST RECENT

If development data is so important, why is it chronically underfinanced? 

Michael M. Lokshin, Jun 11, 2018

Beyond Proof of Concept: do we have the right structure to take disruptive technologies to production? 

Michael M. Lokshin, May 30, 2018

The 2018 Atlas of Sustainable Development Goals: an all-new visual guide to data and development 

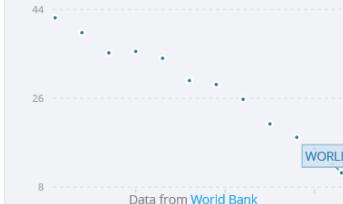
World Bank Data Team, May 24, 2018

[View all news](#) 

[View all blogs](#) 

WHAT YOU CAN LEARN WITH OPEN DATA

Poverty headcount ratio at \$1.90 a day (2011 PPP) (% of population)



Data from [World Bank](#)

Extreme Poverty

The proportion of the world's population living in extreme poverty has dropped significantly

Atlas of Sustainable Development Goals 2018 From World Development Indicators



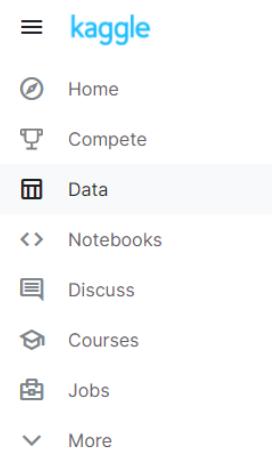
SDG Atlas 2018

May 25, 2018

Public Data Sources – Economics & Finance

Content Reference : <https://data.worldbank.org/>

kaggle



Search

Datasets

Find and use datasets or complete tasks. [Learn more.](#)

Engage With Dataset Tasks

You can now actively engage with datasets with thousands of tasks! Help the community by creating and solving Tasks on datasets!

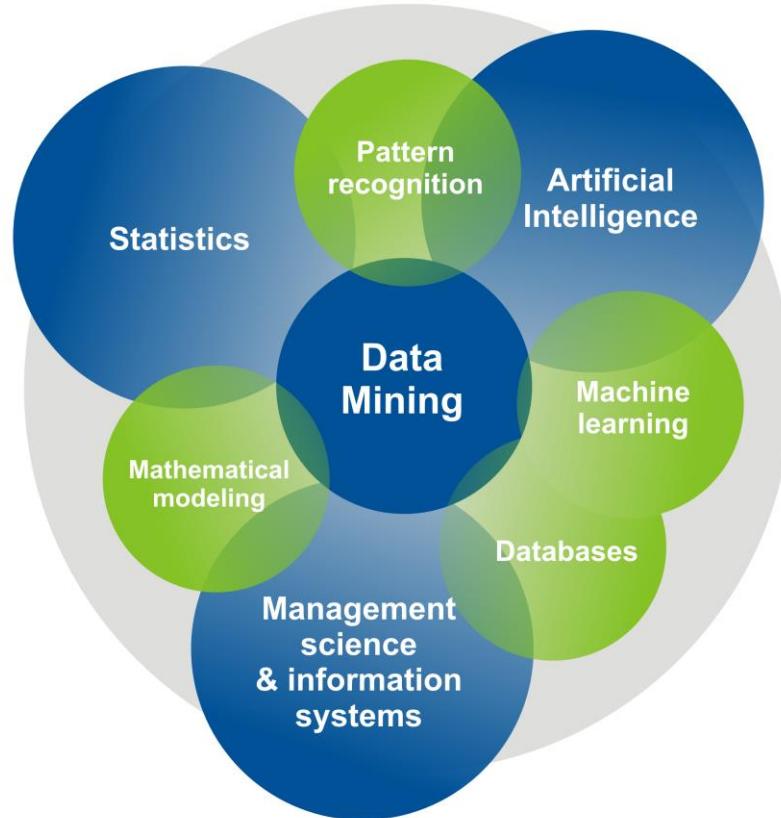
[Tackle a new task](#) [See Details](#)

Search 56,293 datasets [Feedback](#) [Filter](#)

Public Sort by: Hottest

Dataset	Uploader	Created	Size	Rating	Tasks			
COVID-19 Open Research Dataset Challenge (CORD-19)	Allen Institute For AI	Link	2 days	5 GB	8.8	200514 Files (JSON, CSV, other)	17 Tasks	8410
Novel Corona Virus 2019 Dataset	SRK	17 days	3 MB	9.7	8 Files (CSV)	6 Tasks	4849	

Content Reference : <https://www.kaggle.com/datasets>



- Data mining is the process of discovering patterns in large data sets involving methods at the intersection of machine learning, statistics, and database systems.

Introduction to Data Mining

- The explosive growth in data collection
 - The storing of data in data warehouses.
 - The availability of increased access to data from Web navigation and intranet.

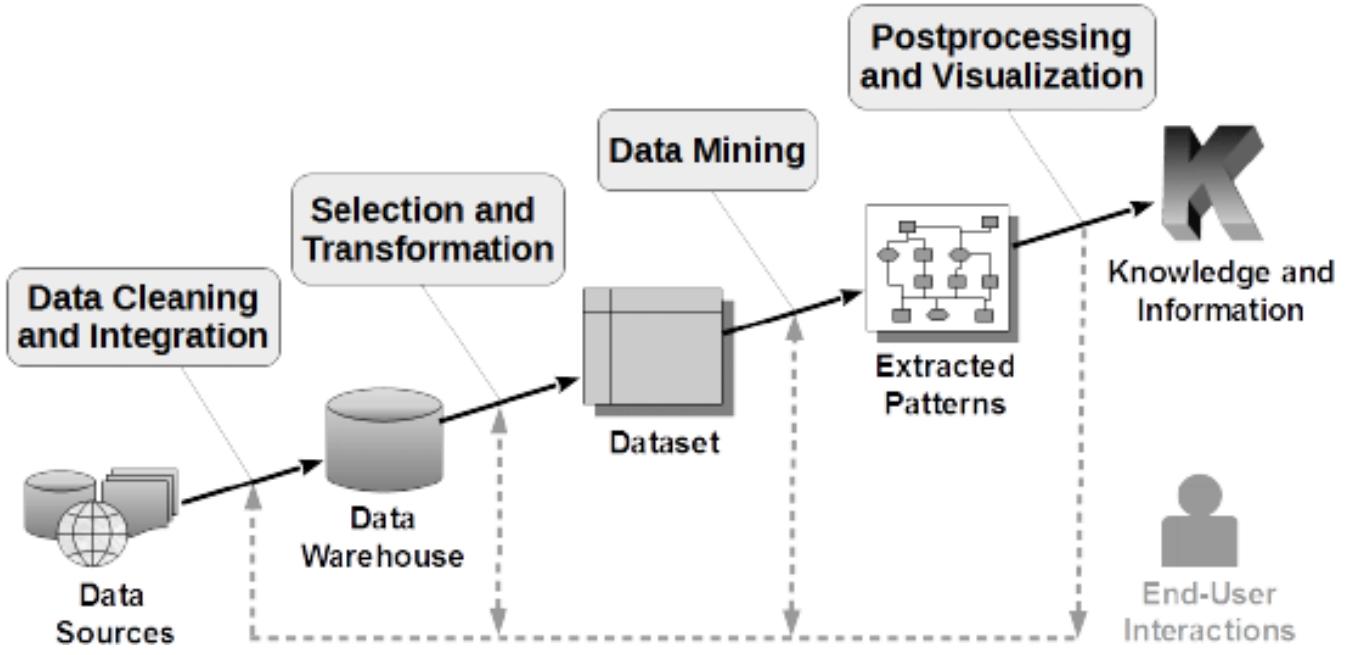
***** We have to find a more effective way to use these data in decision support process than just using traditional query languages *****

Why Data Mining?

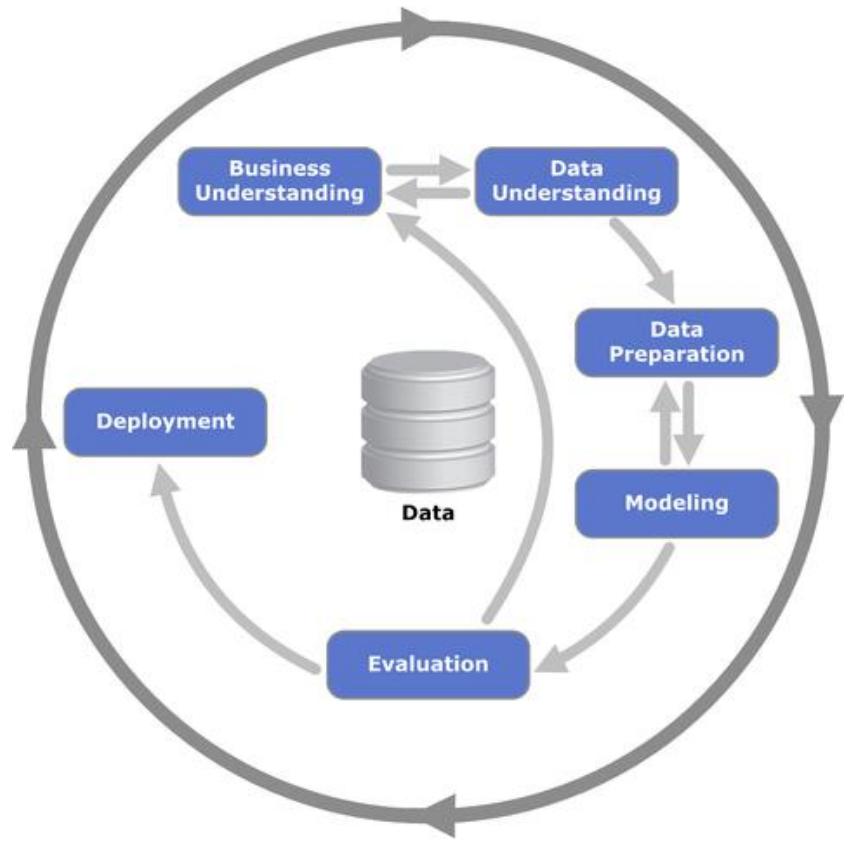
- Knowledge discovery in databases
- Knowledge extraction
- Data archeology
- Data exploration
- Data pattern processing
- Data dredging

Data Mining - Objectives

- Data Cleaning
- Data Integration
- Data Selection
- Data Transformation
- Data Mining
- Pattern Evaluation
- Knowledge Representation



Data Mining Process



Cross-Industry process for data mining

Data Mining Process - CRISP-DM

- Prediction
 - how certain attributes within the data will behave in the future.
- Identification
 - identify the existence of an item, an event, an activity.
- Classification
 - partition the data into categories.
- Optimization
 - optimize the use of limited resources.

Data Mining - Goals

- **Association**
 - Providing the rules correlate the presence of a set of items with another set of item.
- **Classification**
 - Classification is the process of learning a model that describes different classes of data, the classes are predetermined.
 - The model that is produced is usually in the form of a decision tree or a set of rules.
- **Clustering**
 - The previous data mining task of classification deals with partitioning data based on a pre-classified training sample.

Basic Types of Data Mining

Market-basket model.

- Look for combinations of products.
- Put the SHOES near the SOCKS so that if a customer buys one they will buy the other.

Transactions: is the fact the person buys some items in the itemset at supermarket.



Association

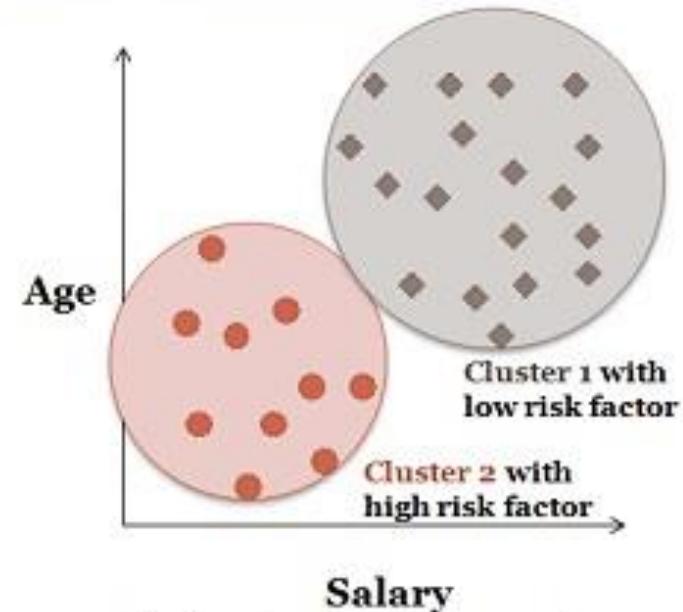
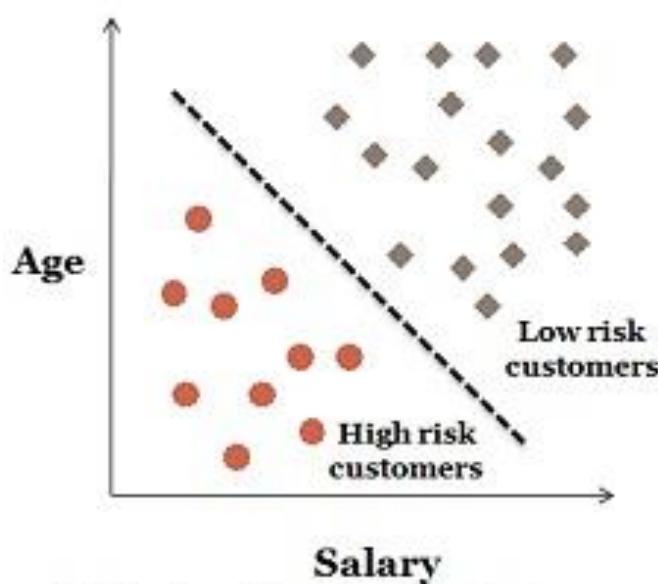
	predictors					
case ID	CUST_ID	CUST_GENDER	EDUCATION	OCCUPATION	AGE	target
	101501	F	Masters	Prof.	41	0
	101502	M	Bach.	Sales	27	0
	101503	F	HS-grad	Cleric.	20	0
	101504	M	Bach.	Exec.	45	1
	101505	M	Masters	Sales	34	1
	101506	M	HS-grad	Other	38	0
	101507	M	< Bach.	Sales	28	0
	101508	M	HS-grad	Sales	19	0
	101509	M	Bach.	Other	52	0
	101510	M	Bach.	Sales	27	1

Classification

Classification

VS

Clustering



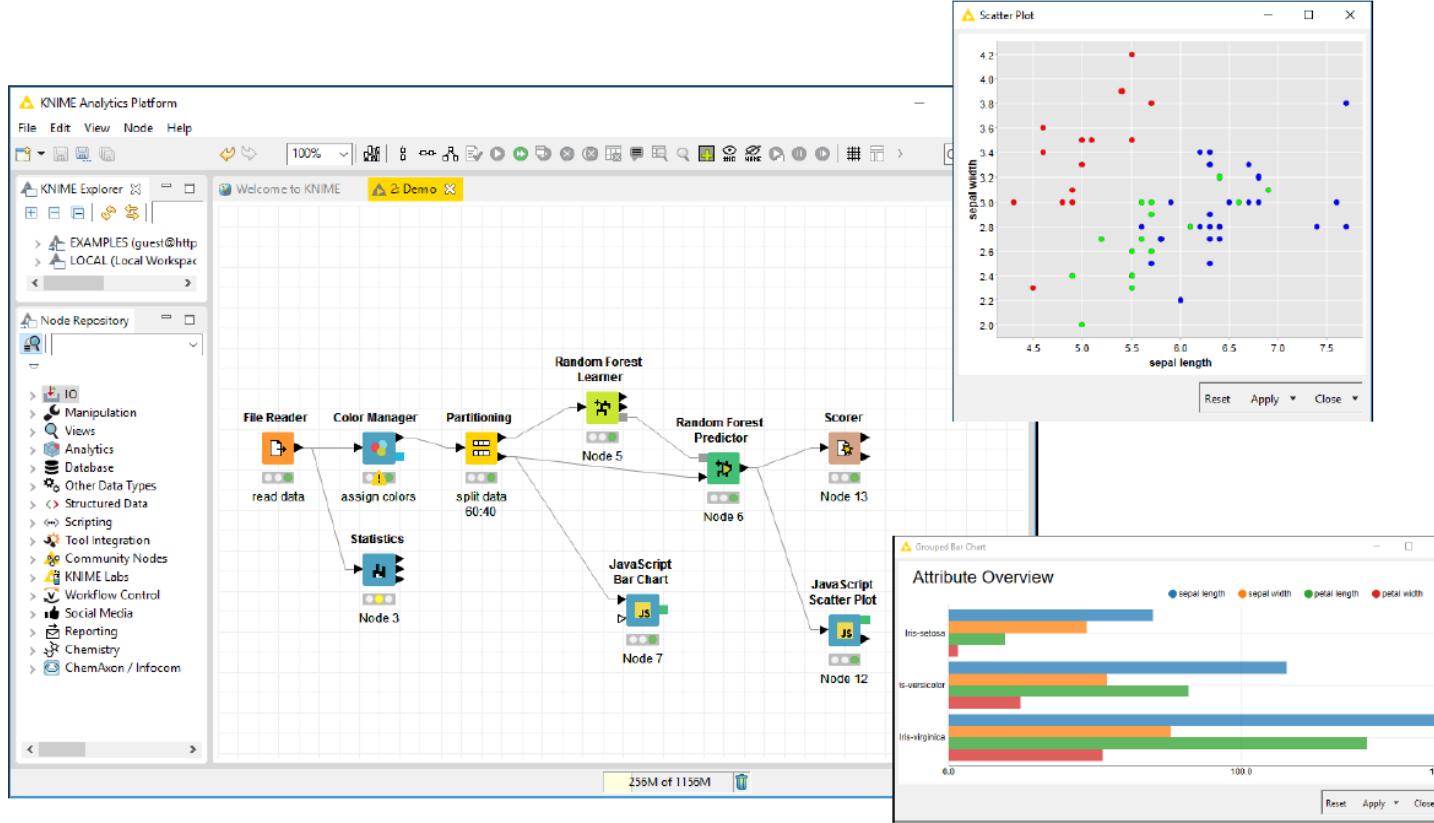
Risk classification for the loan payees on the basis of customer salary

Classification vs Clustering

- Weka is data mining software that uses a collection of machine learning algorithms.
- These algorithms can be applied directly to the data or called from the Java code.
- Weka is a collection of tools for:
 - Regression
 - Clustering
 - Association
 - Data pre-processing
 - Classification
 - Visualization



Data Mining Tools - Weka



- KNIME is a free and open-source data analytics, reporting and integration platform.
- KNIME integrates various components for machine learning and data mining through its modular data pipelining concept.

Data Mining Tools - KNIME

Content Reference : <https://www.knime.com>

Business Applications

- Predict behavior
- Deliver personalized services
- Measure profitability

Challenges

- Noisy data
- Scalability
- Incomplete data

Data Mining Summary

Best Practices

- Preserve data
- Have a good idea of the insights you seek
- Strive for data quality
- Recognize outliers

Data Mining Summary

05 - Introduction to Data Analytics

Data Science - Analytic Type

- Analysis
- Business Analytics

Data Science - Build Model

- Build
- AI & Machine Learning

Two Type of Data Science

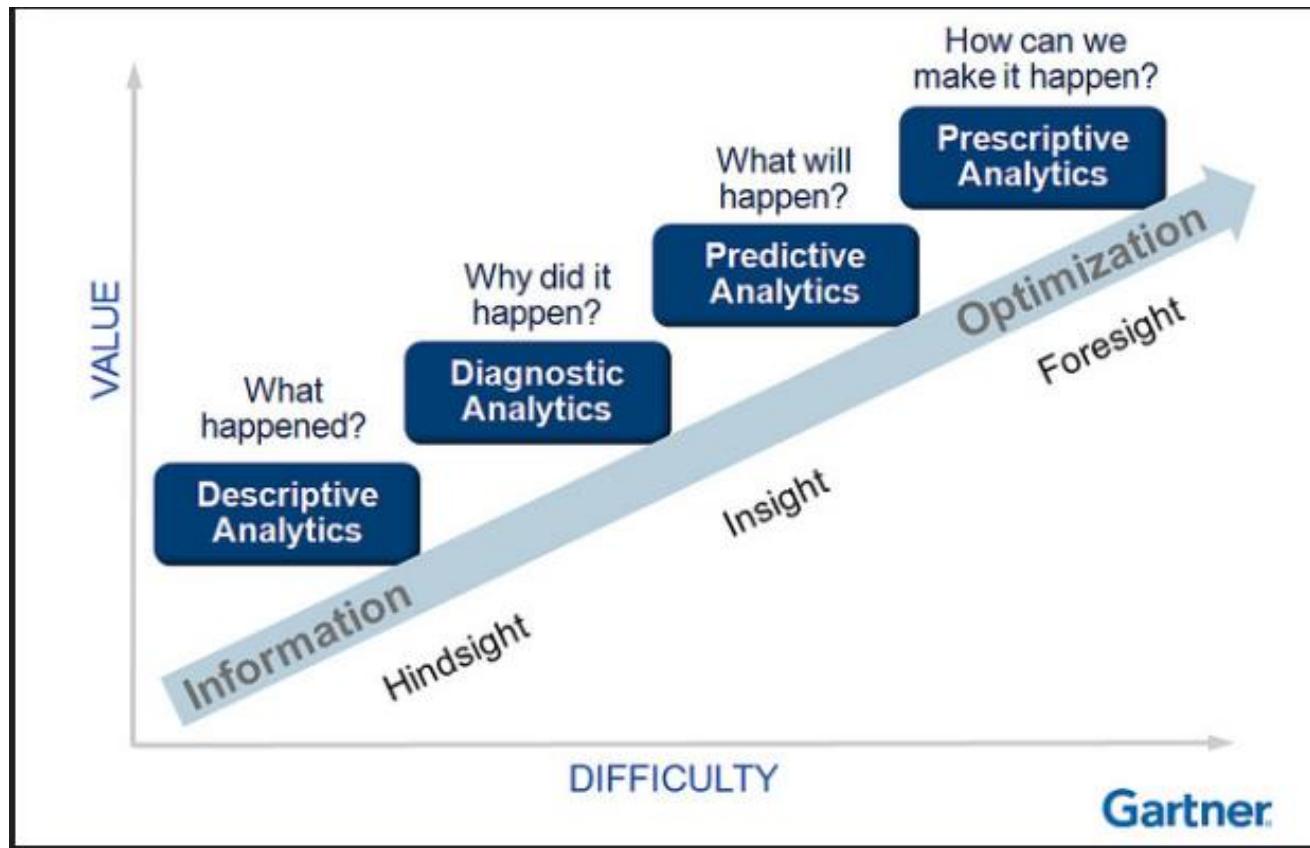
Overview of Data Analysis

Data Analysis is telling a story with data.

Five categories of analytics:

- Descriptive
- Diagnostic
- Predictive
- Prescriptive
- Cognitive



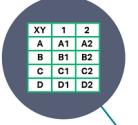


Data Analytics Maturity

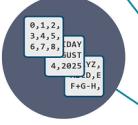
Content Reference : Gartner

Structured Data vs Unstructured Data

Can be displayed in rows, columns and relational databases



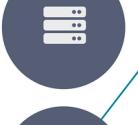
Numbers, dates and strings



Estimated 20% of enterprise data (Gartner)



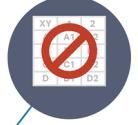
Requires less storage



Easier to manage and protect with legacy solutions



Cannot be displayed in rows, columns and relational databases



Images, audio, video, word processing files, e-mails, spreadsheets



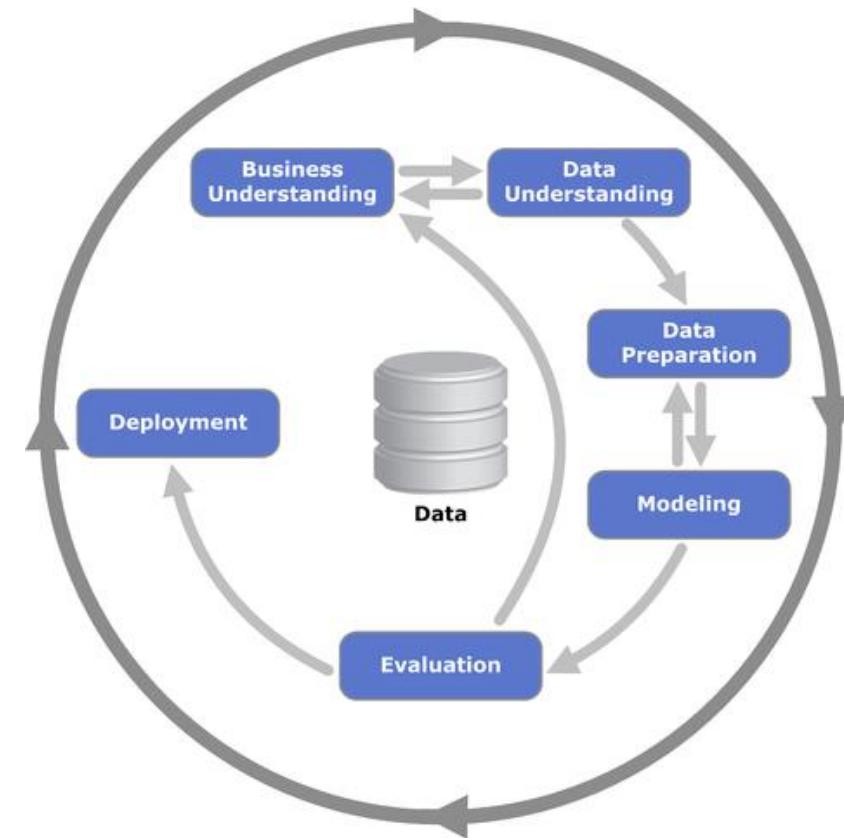
Estimated 80% of enterprise data (Gartner)



Requires more storage



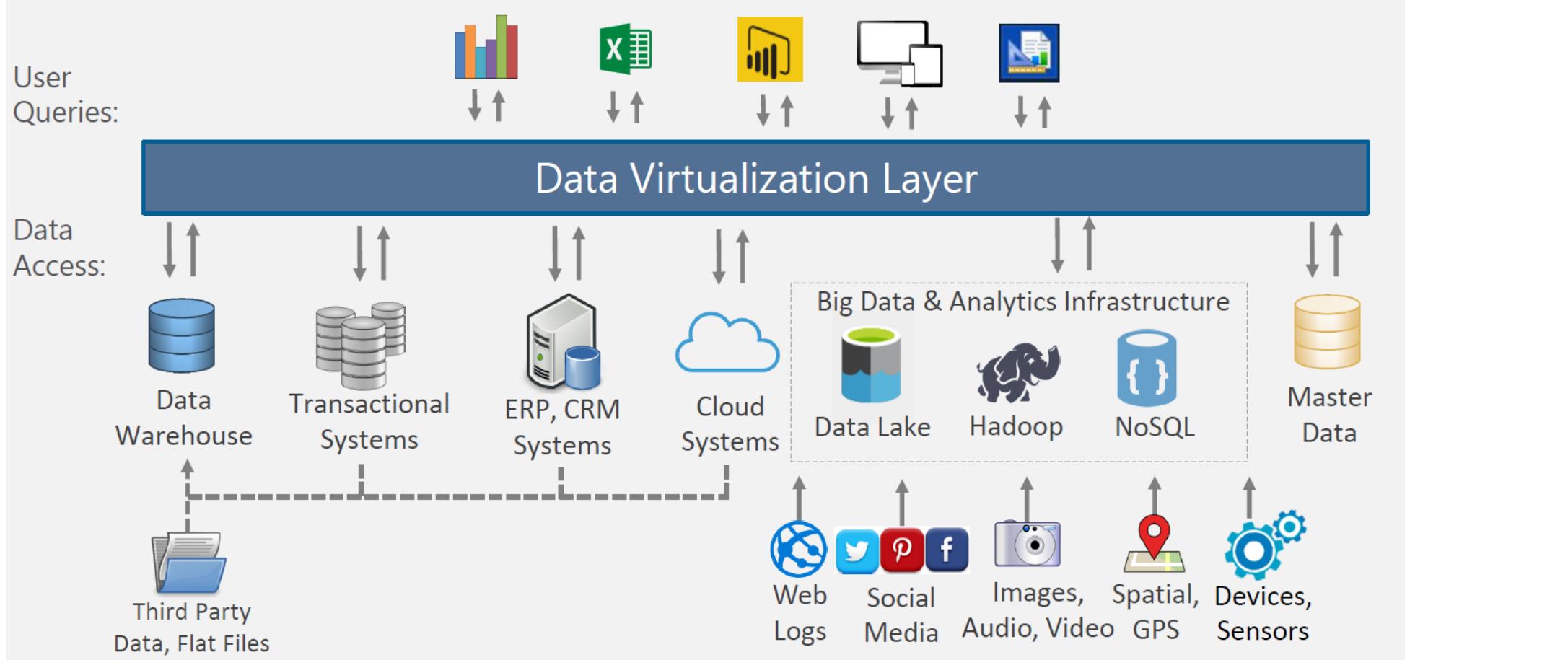
More difficult to manage and protect with legacy solutions



Understanding the Nature of the Data and The Data Analysis Process

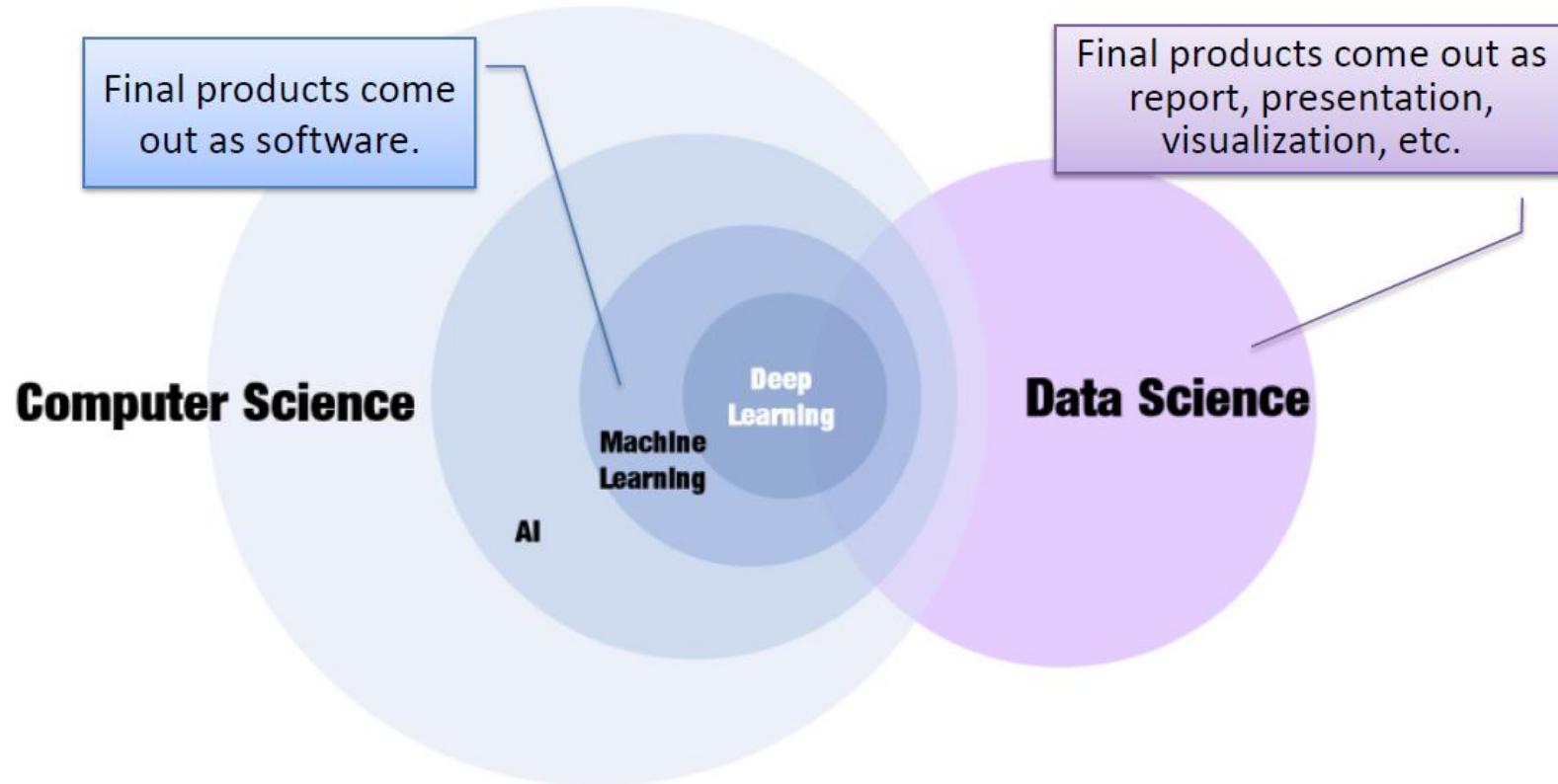
- Customer Lifetime Value
- Customer Segmentation
- Up and Cross Selling
- Next Best Action
- Propensity to buy
- Chun Prevention
- Fraud Detection
- Risk Management
- Demand Forecast
- Price Optimization
- Quality Assurance
- Predictive Maintenance

Data Analytics in Business



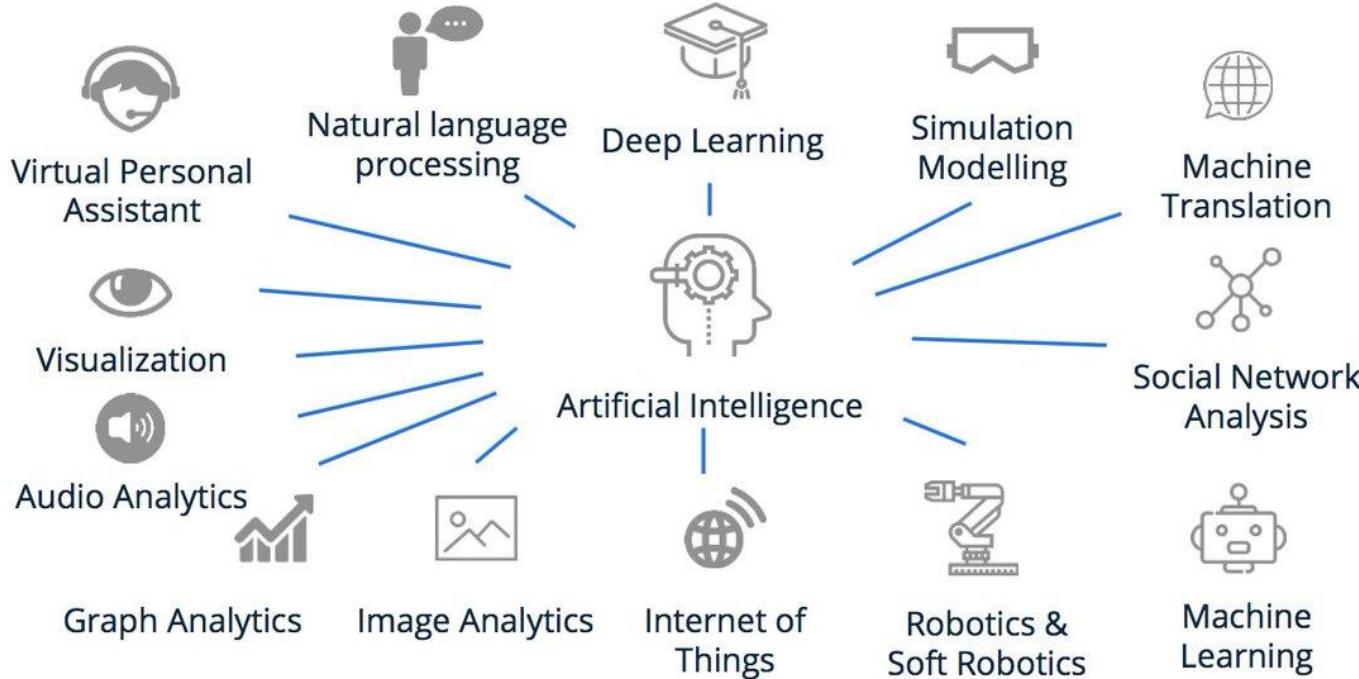
Data Virtualization

06 - Understanding Machine Learning



Terminology

Possible applications for Artificial Intelligence



source statista via @mikequindazzi

Artificial Intelligence (AI)

Content Reference : statista via @mikequindazzi

Software that imitates human capabilities

- Making decisions based on data and past experience
- Recognizing abnormal events
- Interpreting visual input
- Understanding written and spoken language
- Engaging in dialogs and conversations

What is Artificial Intelligence?

	Machine Learning	Predictive models based on data and statistics – the foundation for AI
	Anomaly Detection	Systems that detect unusual patterns or events, enabling pre-emptive action
	Computer Vision	Applications that interpret visual input from cameras, images, or videos
	Natural Language Processing	Applications that can interpret written or spoken language
	Conversational AI	AI agents, (or <i>bots</i>), that can engage in dialogs with human users

Common Artificial Intelligence Workloads

Challenge or Risk	Example
Bias can affect results	A loan-approval model discriminates by gender due to bias in the data with which it was trained
Errors may cause harm	An autonomous vehicle experiences a system failure and causes a collision
Data could be exposed	A medical diagnostic bot is trained using sensitive patient data, which is stored insecurely
Solutions may not work for everyone	A predictive app provides no audio output for visually impaired users
Users must trust a complex system	An AI-based financial tool makes investment recommendations - what are they based on?
Who's liable for AI-driven decisions?	An innocent person is convicted of a crime based on evidence from facial recognition – who's responsible?

Challenges and Risks with AI



Fairness



Reliability & Safety



Privacy & Security



Inclusiveness



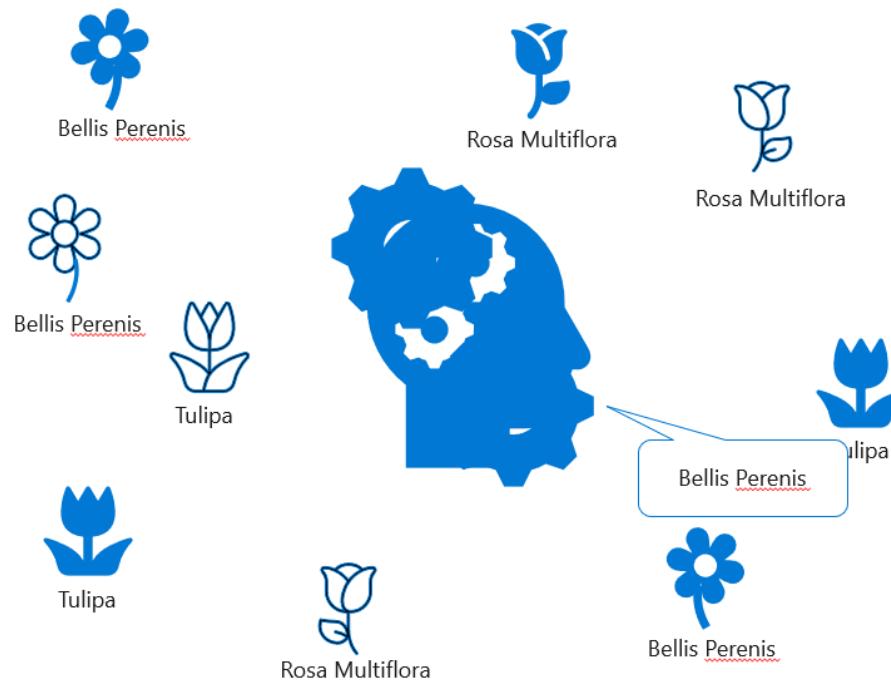
Transparency



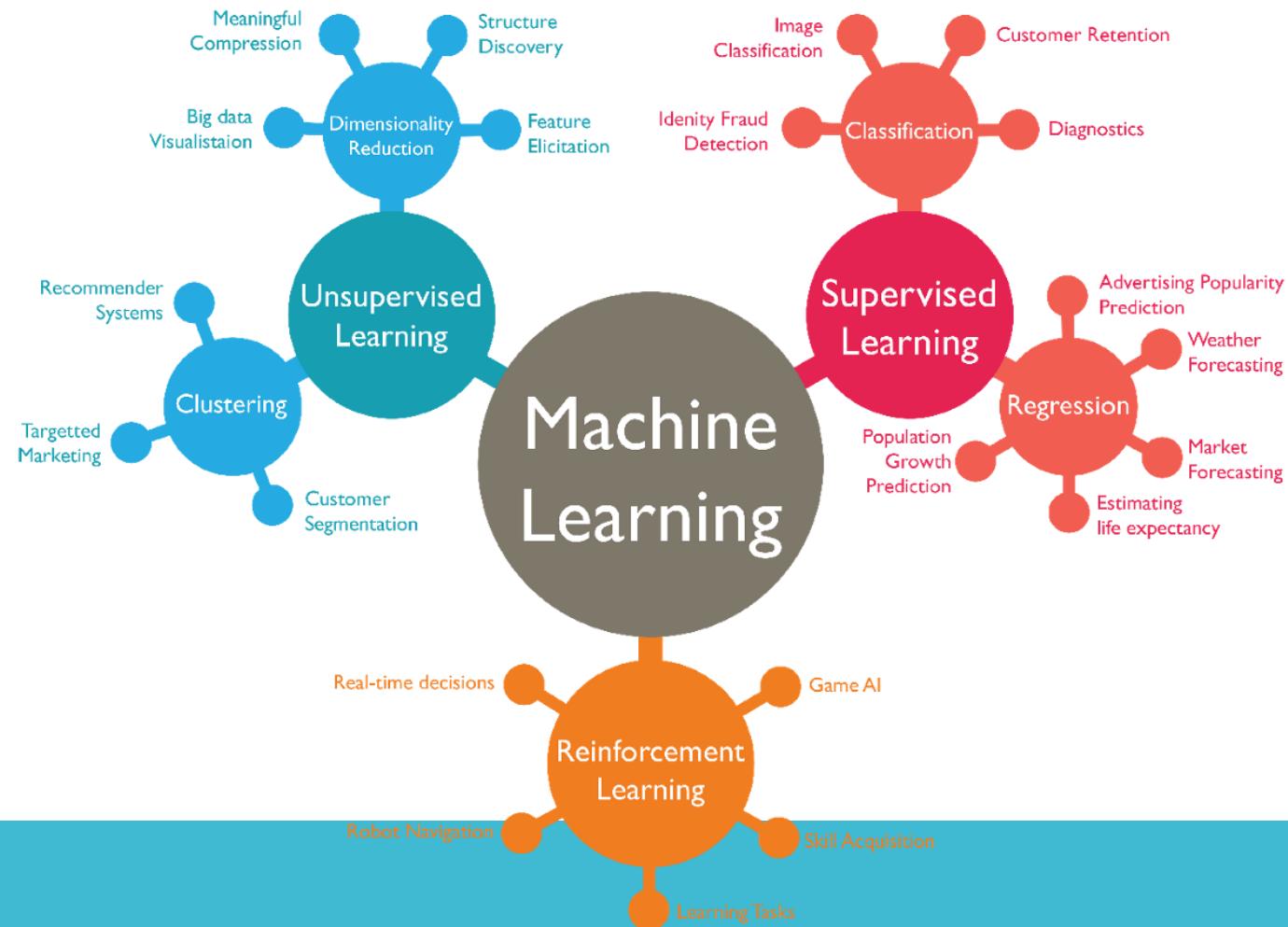
Accountability

Principles of Responsible AI

Creating predictive models by finding relationships in data



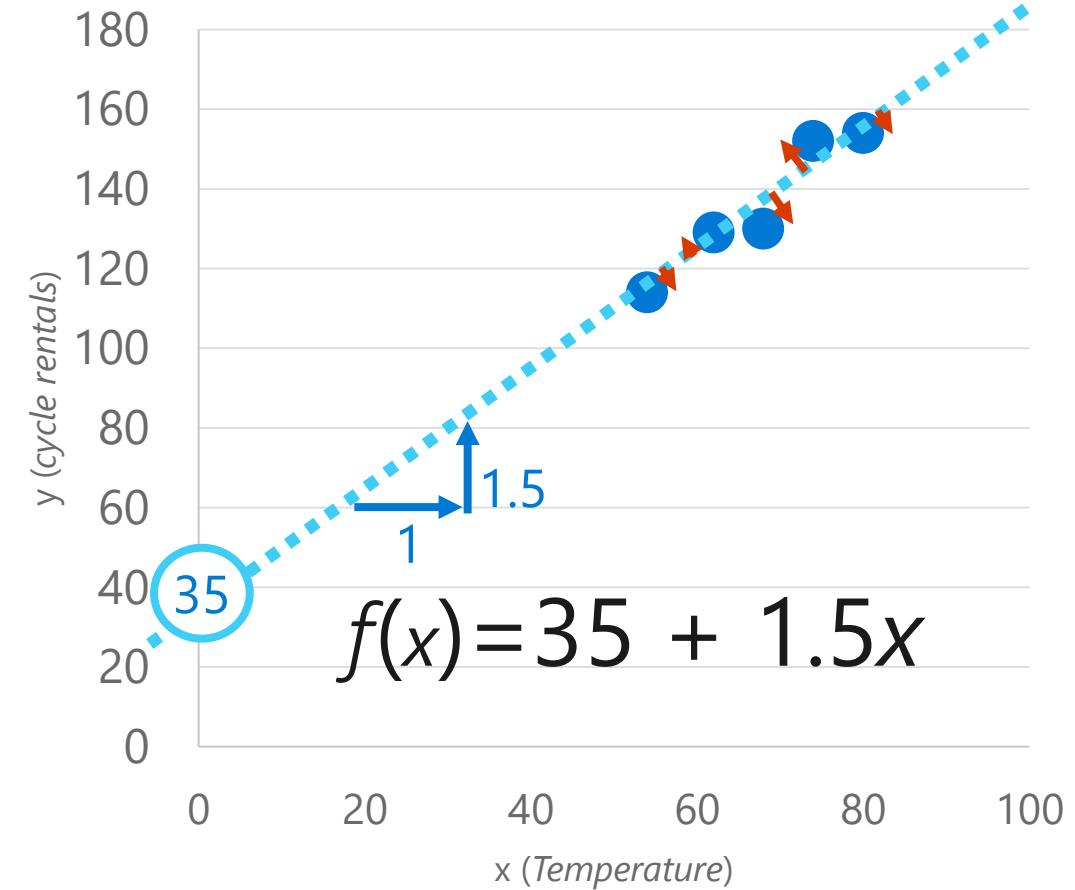
What is Machine Learning



Type of Machine Learning

Regression

		$f(x)$	\hat{y}
		x	y
Training	56		115
Validation	61		126
Training	67		137
Validation	72		140
Training	76		152
Validation	82		156
Training	54		114
Validation	62		129
Training	68		130
Validation	74		152
Validation	80		154



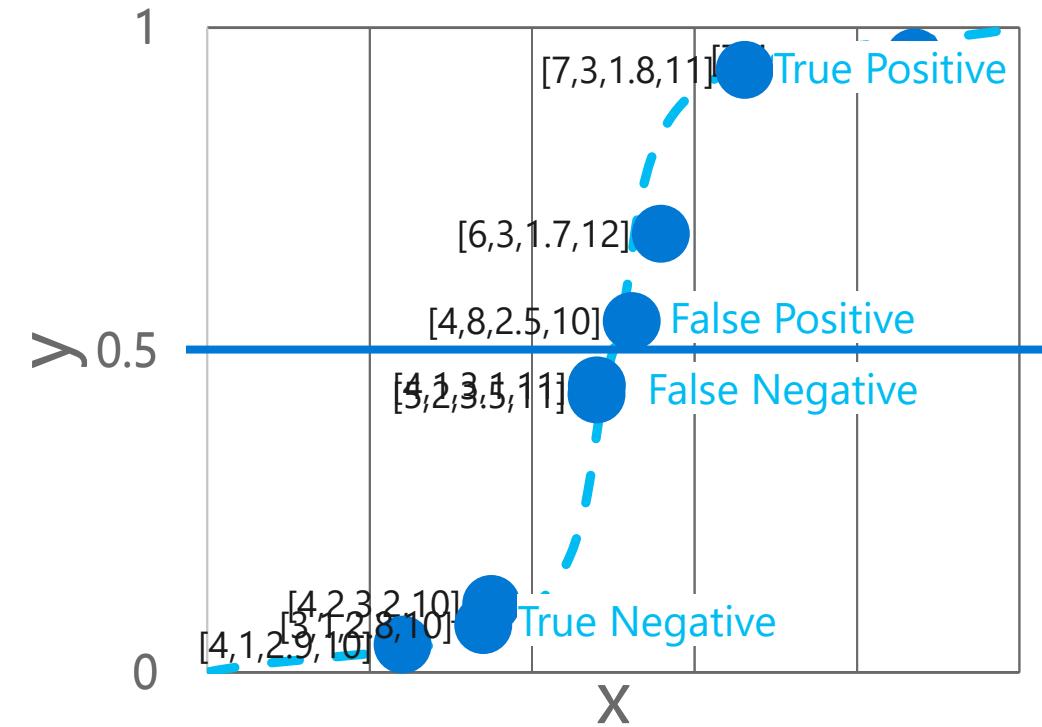
Content Reference : Microsoft Corporation

Classification

		X	y
Training	[4,2,3,2,10]	0	0
	[6,3,1,7,12]	1	1
	[5,2,3,5,11]	0	0
	[4,1,2,9,10]	0	0
	[7,4,2,1,11]	1	1
	[3,1,2,8,10]	0	0
	[7,3,1,8,11]	1	1
	[4,8,2,5,10]	0	0
	[4,1,3,1,11]	1	1
Validation			

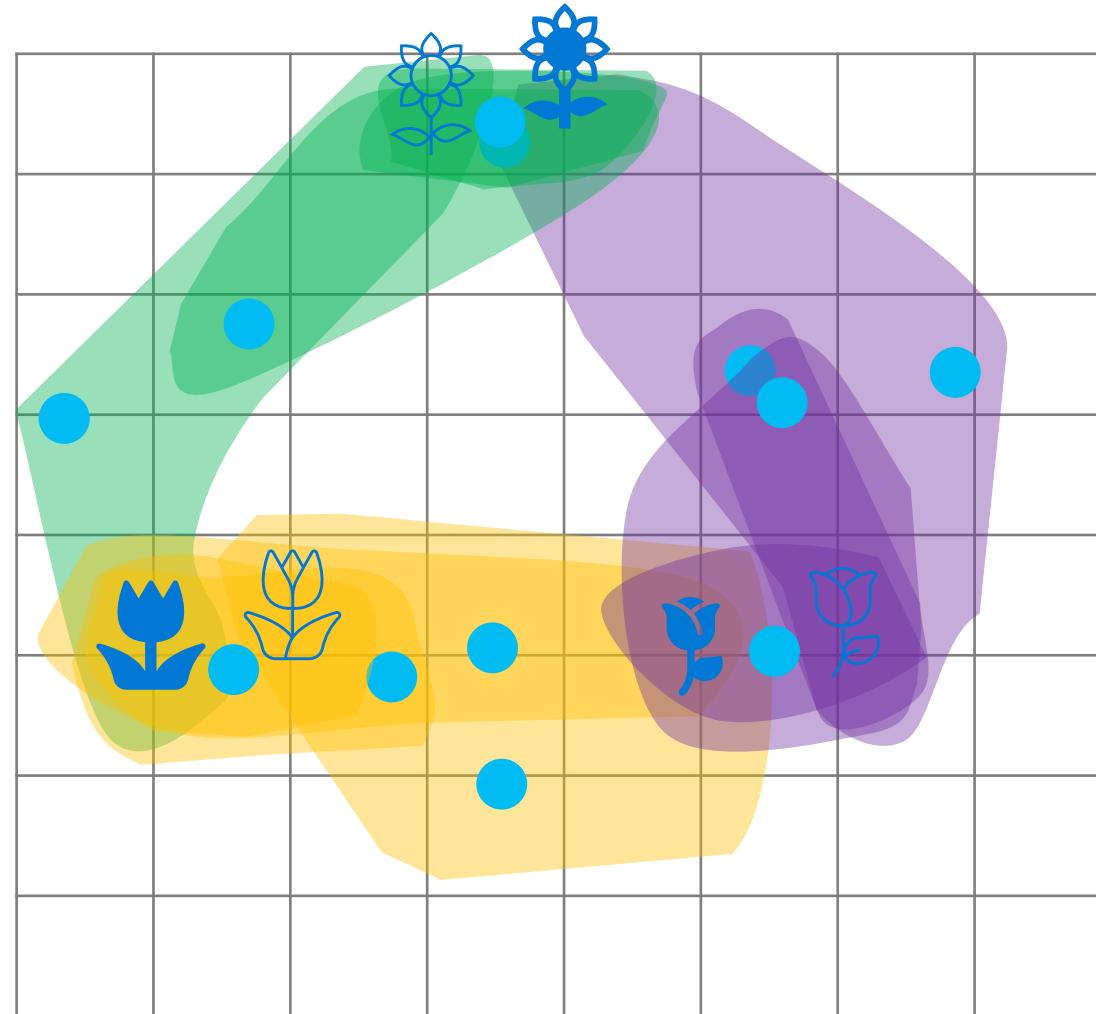
		Predicted
Actual	1	0
1	126	21
0	7	119

P(1)	P(0)	\hat{y}	
0.2	0.8	0	✓
0.9	0.1	1	✓
0.6	0.4	1	✗
0.3	0.7	0	✗



Clustering

	↔	
	6	3
	5	3
	2	3
	1	3
	3	8
	4	8

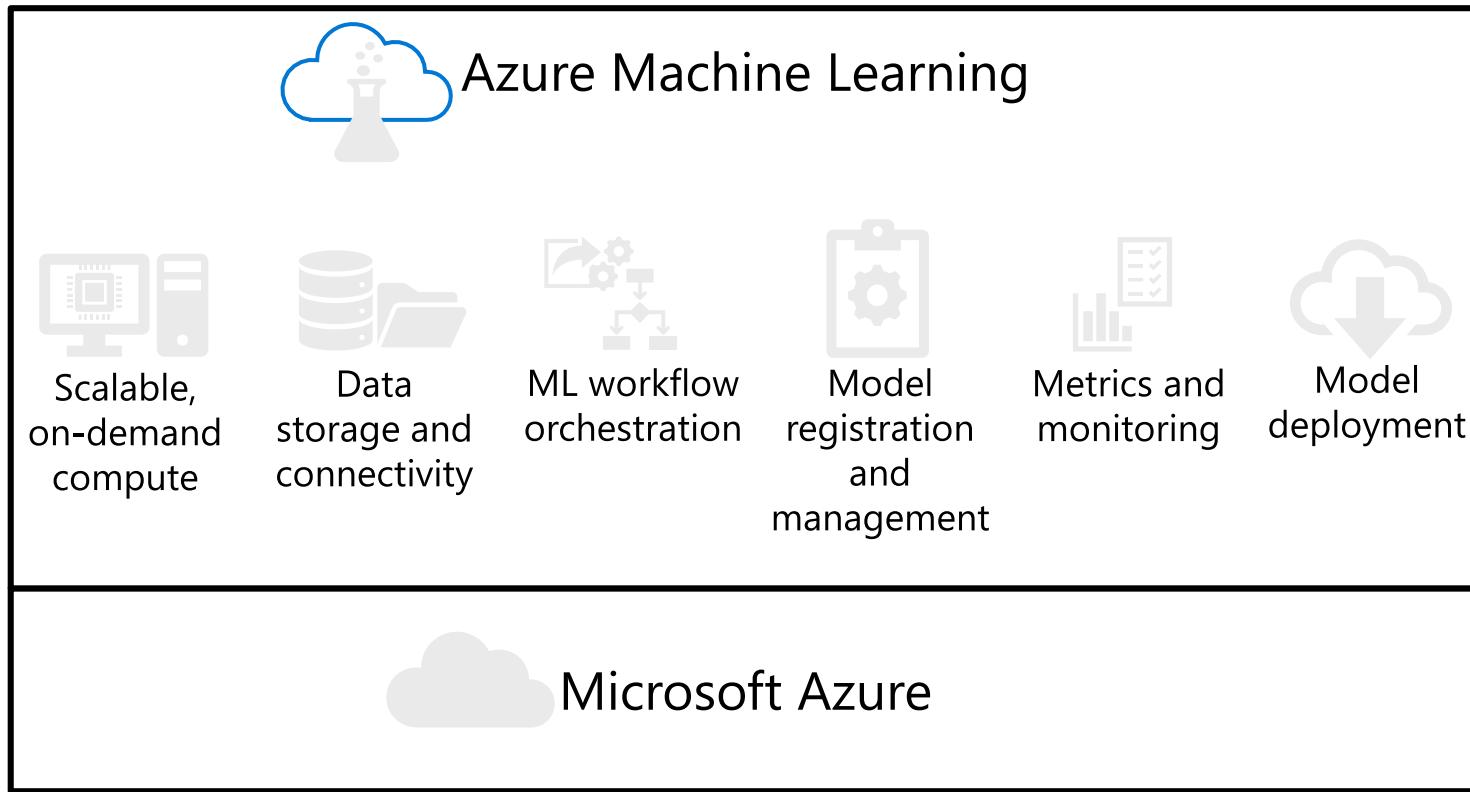


Azure Machine Learning No-Code with Designer

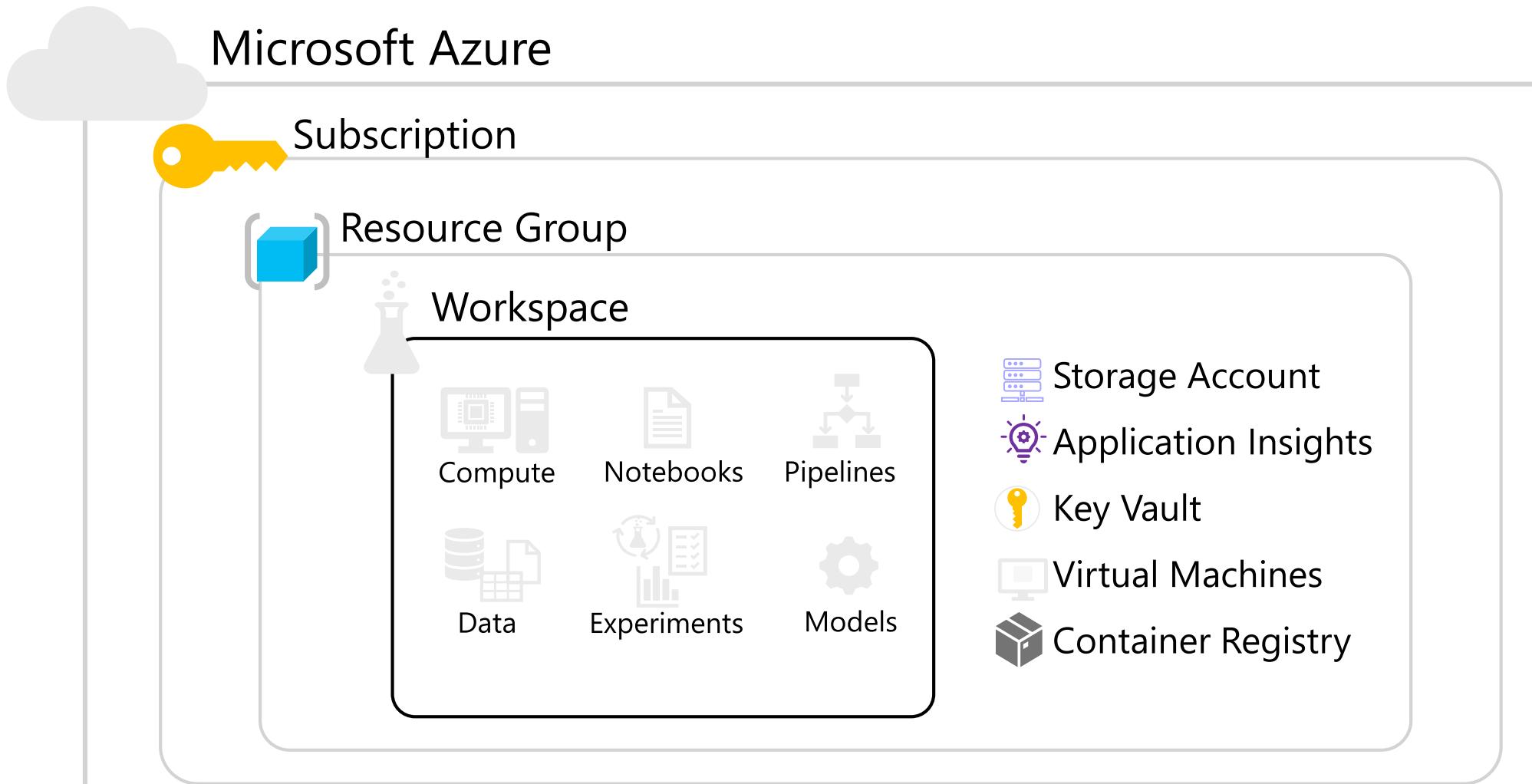


What is Azure Machine Learning?

A platform for operating machine learning workloads in the cloud



Azure Machine Learning Workspaces



Considerations for Creating a Workspace



Region

Check Azure Resource availability

For example, NC-Series Virtual Machines for GPU processing



Edition

Enterprise

- All features

Basic

- No Visual Designer
- No Automated ML user interface
- No Data Drift user interface

Azure Machine Learning studio

Manage compute and data

Run experiments

View metrics

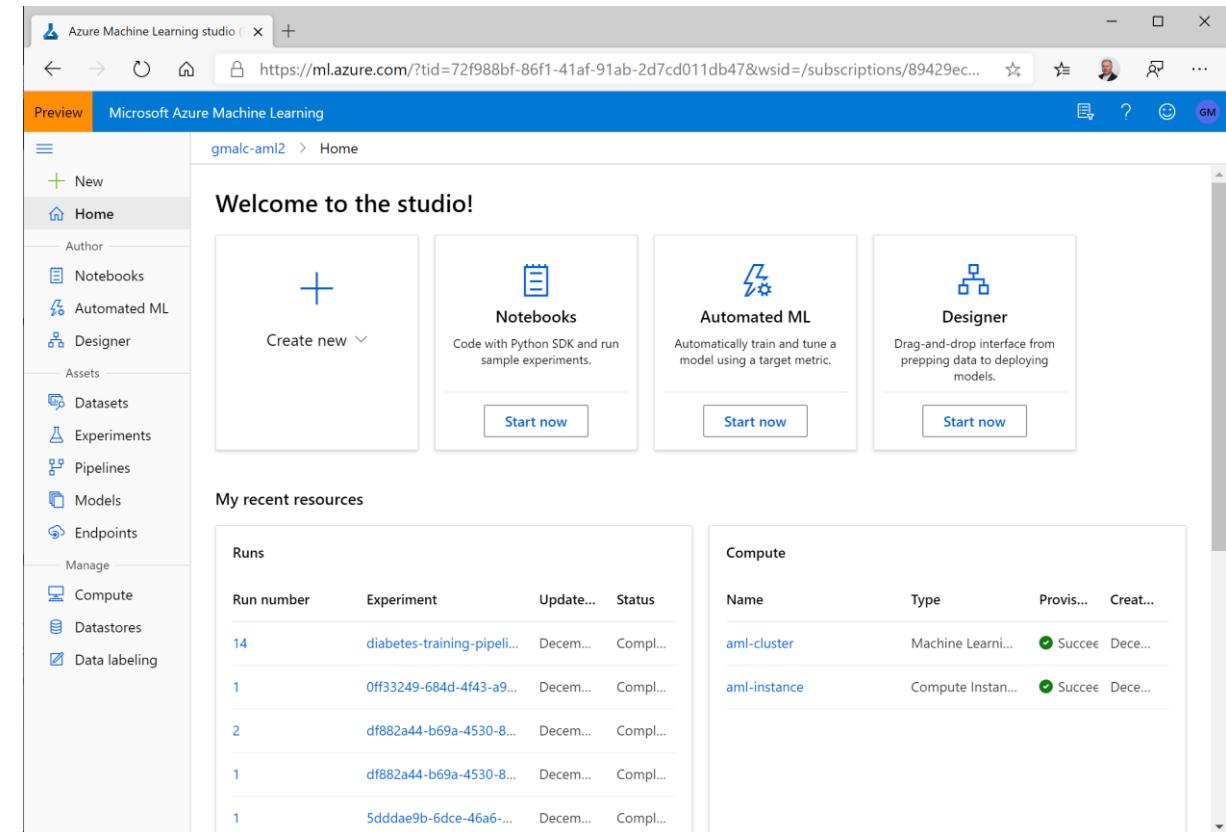
Manage and deploy models

Manage endpoints

Use graphical modeling tools:

Designer - "no-code" model development

Automated Machine Learning - find the best model for your data



The Azure Machine Learning SDK for Python

Code-based configuration for machine learning assets:

Automate repeatable asset creation

Ensure consistency across development, test, and production environments

Incorporate machine learning asset configuration into DevOps

```
pip install azureml-sdk
```

```
from azureml.core import Workspace

ws = Workspace.from_config()
for compute_name in ws.compute_targets:
    compute = ws.compute_targets[compute_name]
    print(compute.name, ":", compute.type)
```

Compute Instances

Jupyter Notebook and JupyterLab servers in your workspace
Choose the compute specifications you need

The screenshot shows the Azure Machine Learning studio interface. On the left, the navigation pane is open with the 'Compute' section selected. In the center, the 'Compute Instances' tab is active, displaying a table with one row: 'aml-instance' (Status: Running, Application URI: JupyterLab, Virtual Machine: STANDARD_DS). A green arrow points from the 'Compute Instances' table towards the Jupyter notebook window on the right. The Jupyter notebook window displays Python code for interacting with the workspace:

```
from azureml.core import Workspace
ws = Workspace.from_config()
print(ws.name, "loaded")
```

Below the code, a section titled 'View Azure ML Resources' shows more code for listing compute targets, datastores, and datasets:

```
from azureml.core import ComputeTarget, Datastore, Dataset
print("Compute Targets:")
for compute_name in ws.compute_targets:
    compute = ws.compute_targets[compute_name]
    print("\t", compute.name, ":", compute.type)

print("Datastores:")
for datastore_name in ws.datastores:
    datastore = Datastore.get(ws, datastore_name)
    print("\t", datastore.name, ":", datastore.datastore_type)

print("Datasets:")
for dataset_name in list(ws.datasets.keys()):
    dataset = Dataset.get_by_name(ws, dataset_name)
    print("\t", dataset.name)
```

At the bottom of the notebook, there is a note: 'Now you've seen how to use the Azure ML SDK to view the resources in your workspace. The SDK provides a great way to script the creation and configuration of the resources you need to operate machine learning workloads using Azure ML. For more details, see the [Azure ML SDK documentation](#). On the File menu, click Close and Halt to close this notebook. Then return to the lab instructions.'

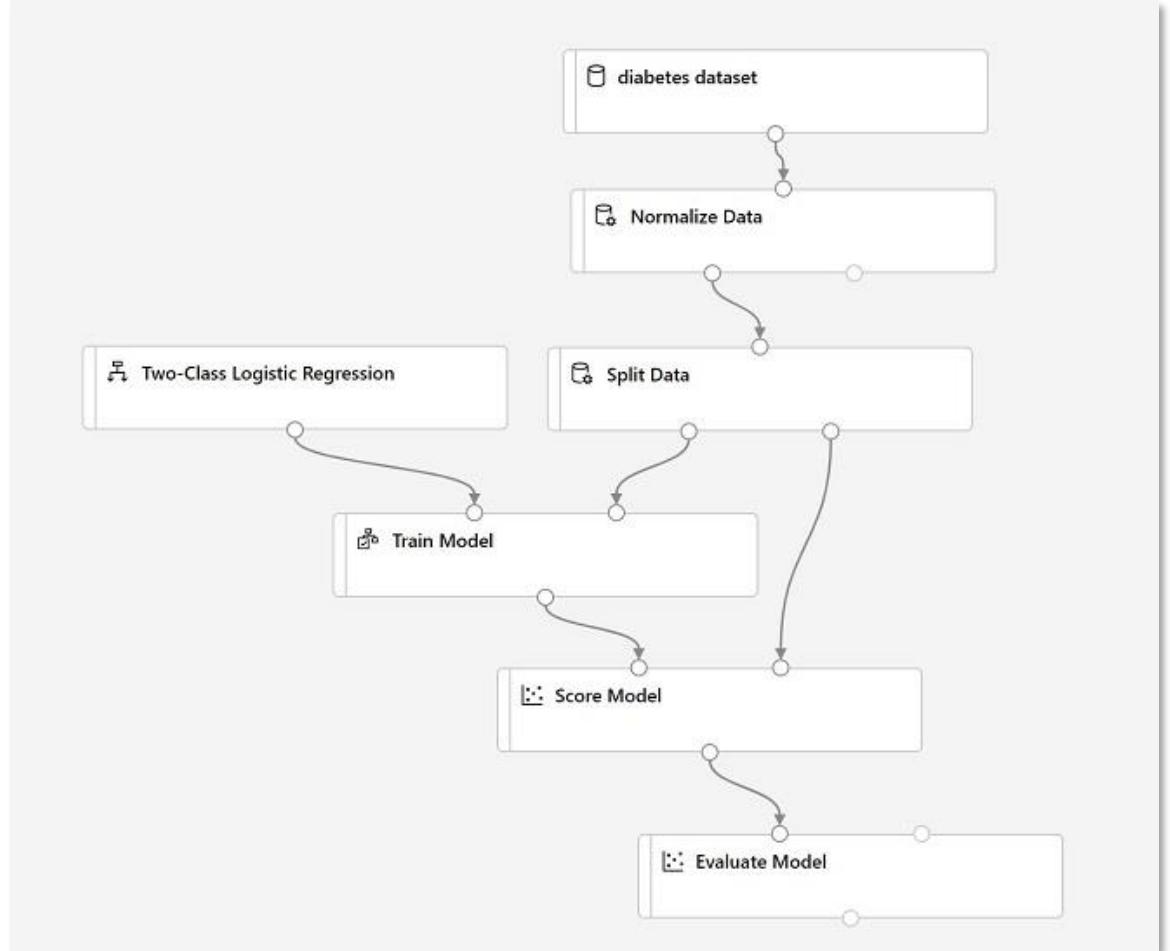
Microsoft

Walkthrough :

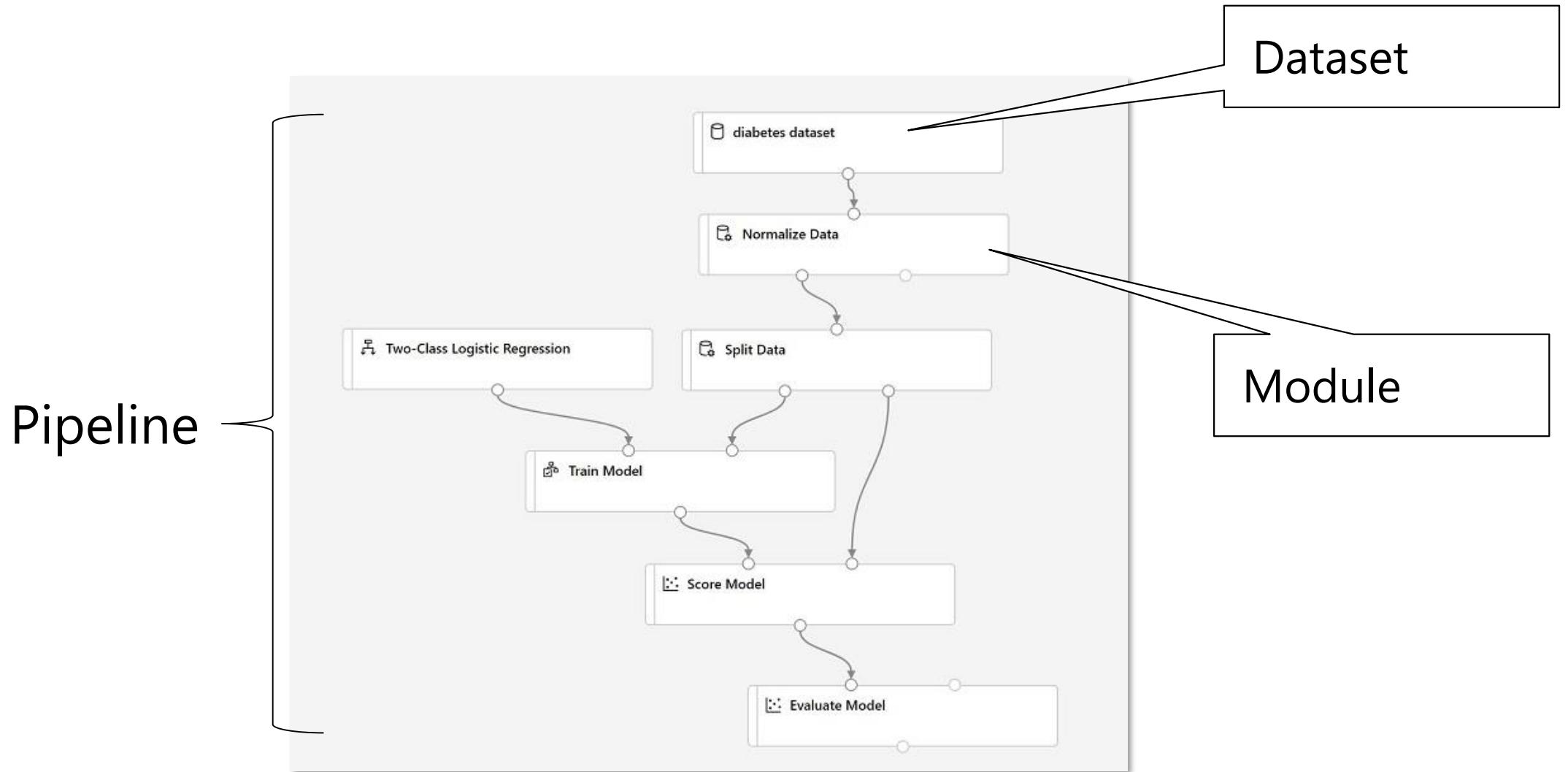
Creating an Azure Machine Learning Workspace

What is Azure Machine Learning Designer?

Drag-and-Drop Interface for:
Preparing data and training models
Publishing models as services



Designer Pipelines and Modules



Training, Scoring, and Evaluating Models

1. Algorithm

Classification, Regression, Clustering

Two-Class Logistic Regression

diabetes dataset

Normalize Data

Split Data

Train Model

Score Model

Evaluate Model

3. Score

Predict from test data

2. Train

Fit model for specified label

4. Evaluate

model-specific metrics

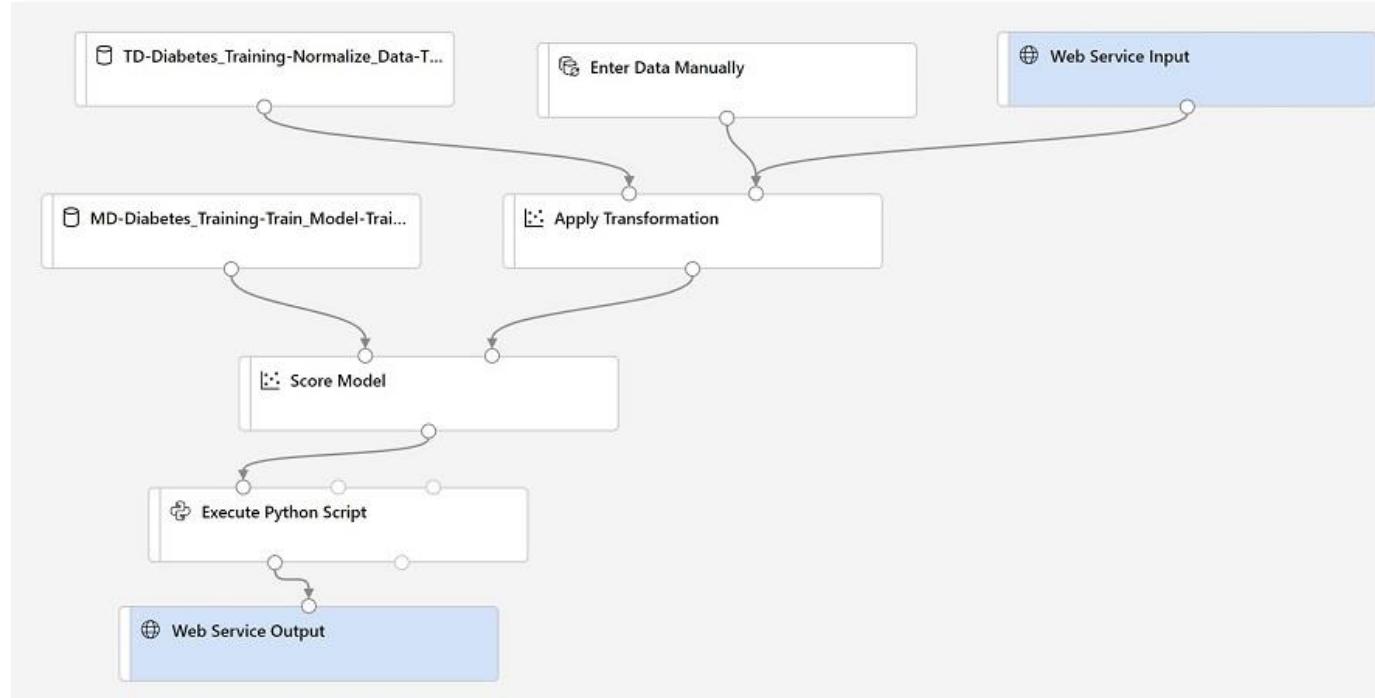
Custom Code Modules

Apply SQL Transformation	Use a SQL statement to process up to three input tables
Execute Python Script	Implement a custom Python function to process up to two dataframes
Create Python Model	Implement a custom Python model in place of a built-in algorithm
Execute R Script	Implement a custom R function to process up to two dataframes

Walkthrough :

Creating a Training Pipeline with the Azure ML Designer

What is an Inference Pipeline?



A data flow defining a web service for using the trained model

A **Web Service Input** defines the input data schema

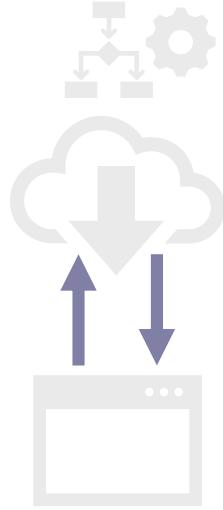
Transformations based on training data are encapsulated in datasets

The trained model is encapsulated in a dataset

A **Web Service Output** defines the output data schema

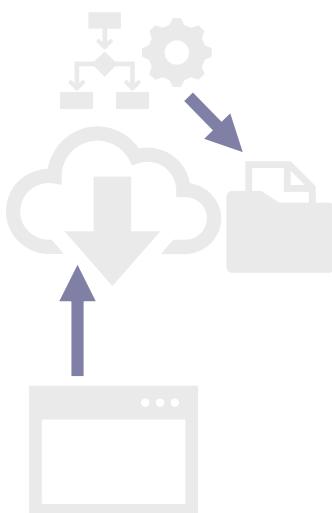
You may want to modify the pipeline before deploying its as a web service

Publishing a Service Endpoint



Deploy a Real-Time Pipeline:

Requires Azure Kubernetes Services Inference Compute
Submit new data to HTTP endpoint for immediate results



Publish a Batch Pipeline

Requires Azure Machine Learning Training Compute
Initiate pipeline experiment run through HTTP endpoint
Results saved in run output

Consuming a Service Endpoint

**View endpoints in Azure Machine Learning studio
Use starter code to build client applications**

```
data = {"Inputs": {"input0": [{"feature1": "123", "feature2": "99"}]},  
        "GlobalParameters": {}}  
body = str.encode(json.dumps(data))  
  
url = 'http://10.0.0.1:80/api/v1/service/diabetes_predictor/score'  
api_key = 'a1234567890x'  
headers = {'Content-Type': 'application/json',  
           'Authorization': ('Bearer ' + api_key)}  
  
req = urllib.request.Request(url, body, headers)  
response = urllib.request.urlopen(req)  
result = response.read()
```

Walkthrough :

Deploying a Service with the Azure ML Designer

Complete the Course