# Lecture 10:

# Data Management
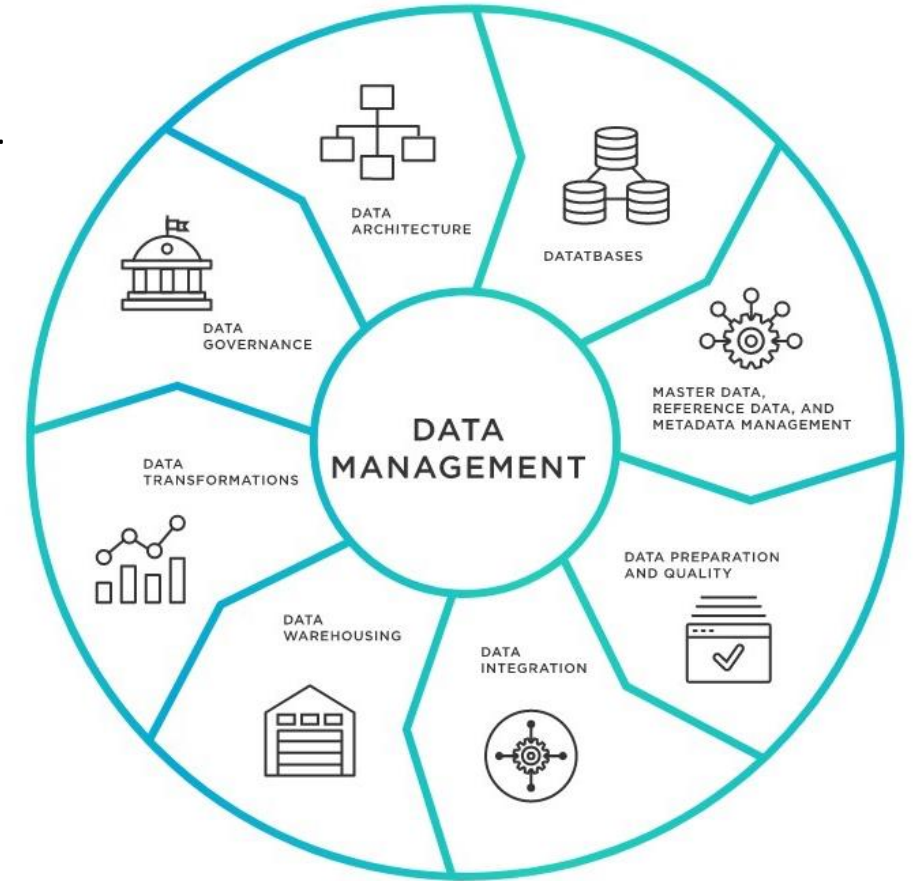
CS5481 Data Engineering

Instructor: Linqi Song

# Outline

1. What is Data Management

2. Data Quality

3. Data Security

4. Data Privacy

# What is data management?

- Data management is the practice of collecting, keeping, and using data securely, efficiently, and cost-effectively.

- As organizations create and consume data at unprecedented rates, data management solutions become essential for making sense of the vast quantities of data.

- Today's leading data management software ensures that reliable, up-to-date data is always used to drive decisions.

# Data management techniques (1)

**Types of data management:** Data management plays several roles in an organization's data environment, making essential functions easier and less time-intensive. These data management techniques include the following:

- **Data pipelines** enable the automated transfer of data from one system to another.
- **Data preparation** is used to clean and transform raw data into the right shape and format for analysis, including making corrections and combining data sets.
- **ETLs (Extract, Transform, Load)** are built to take the data from one system, transform it, and load it into the organization's data warehouse.
- **Data catalogs** help manage metadata to create a complete picture of the data, providing a summary of its changes, locations, and quality while also making the data easy to find.

# Data management techniques (2)

- **Data warehouses** are places to consolidate various data sources, contend with the many data types businesses store, and provide a clear route for data analysis.
- **Data governance** defines standards, processes, and policies to maintain data security and integrity.
- **Data architecture** provides a formal approach for creating and managing data flow.
- **Data security** protects data from unauthorized access and corruption.
- **Data modeling** documents the flow of data through an application or organization.

# Data quality – What is data quality?

- **Data quality** refers to the development and implementation of activities that apply <span style="color:red">quality management techniques to data</span> in order to ensure the data is fit to serve the specific needs of an organization in a particular context. Data that is deemed fit for its intended purpose is considered high quality data.

- Examples of data quality issues include duplicated data, incomplete data, inconsistent data, incorrect data, poorly defined data, poorly organized data, and poor data security.
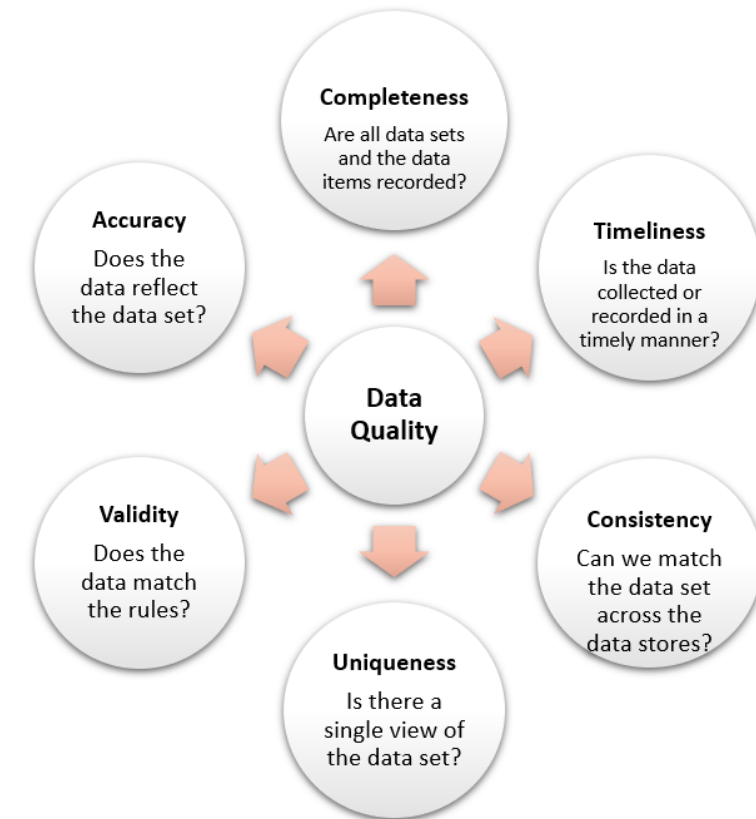
# Data quality dimensions (1)

There are six main dimensions of data quality: accuracy, completeness, consistency, validity, uniqueness, and timeliness.

**Accuracy:** The data should reflect actual, real-world scenarios; the measure of accuracy can be confirmed with a verifiable source.

**Completeness:** Completeness is a measure of the data's ability to effectively deliver all the required values that are available.

**Consistency:** Data consistency refers to the uniformity of data as it moves across networks and applications. The same data values stored in difference locations should not conflict with one another.



**Completeness**
Are all data sets and the data items recorded?

**Timeliness**
Is the data collected or recorded in a timely manner?

**Accuracy**
Does the data reflect the data set?

**Data Quality**

**Consistency**
Can we match the data set across the data stores?

**Validity**
Does the data match the rules?

**Uniqueness**
Is there a single view of the data set?

# Data quality dimensions (2)

**Validity:** Data should be collected according to defined business rules and parameters, and should conform to the right format and fall within the right range.

**Uniqueness:** Uniqueness ensures there are no duplications or overlapping of values across all data sets. Data cleansing and deduplication can help remedy a low uniqueness score.

**Timeliness:** Timely data is data that is available when it is required. Data may be updated in real time to ensure that it is readily available and accessible.

# How to improve data quality? (1)

Data quality measures can be accomplished with data quality tools, which typically provide data quality management capabilities such as:

- **Data preprocessing (Lecture 3)** – data cleaning, such as handling duplicates, outliers, and missing data, etc.

## DATA CLEANING CHECKLIST

**Up-to-date data**

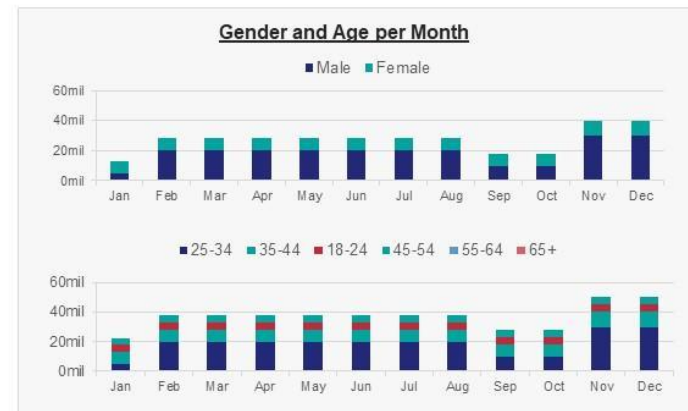Data should be up-to-date in order to obtain maximum value from the data analysis.

✔

**Missing values**

Count missing values and analyze where in the data they are missing. Missing values can disrupt some analyses and skew the results.

✔

**Duplicates**

Duplicate IDs indicate multiple records for one person, e.g. someone holds multiple functions at the same time.

✔

**Numerical outliers**

Numerical outliers are fairly easy to detect and remove. Define minimum and maximum to spot outliers easily.

✔

**Check IDs**

Check data labels of all the fields to see whether some categorical values are mislabeled.

✔

**Define valid output**

Define valid data labels for categorical data. Define data ranges for numerical variables. Non-matching data is presumably wrong.

✔

9

# How to improve data quality? (2)

- **Data profiling** - the process of examining, analyzing, and creating useful summaries of data.
  - The first step in the data quality improvement process is understanding your data.
  - Data profiling is the initial assessment of the current state of the data sets.
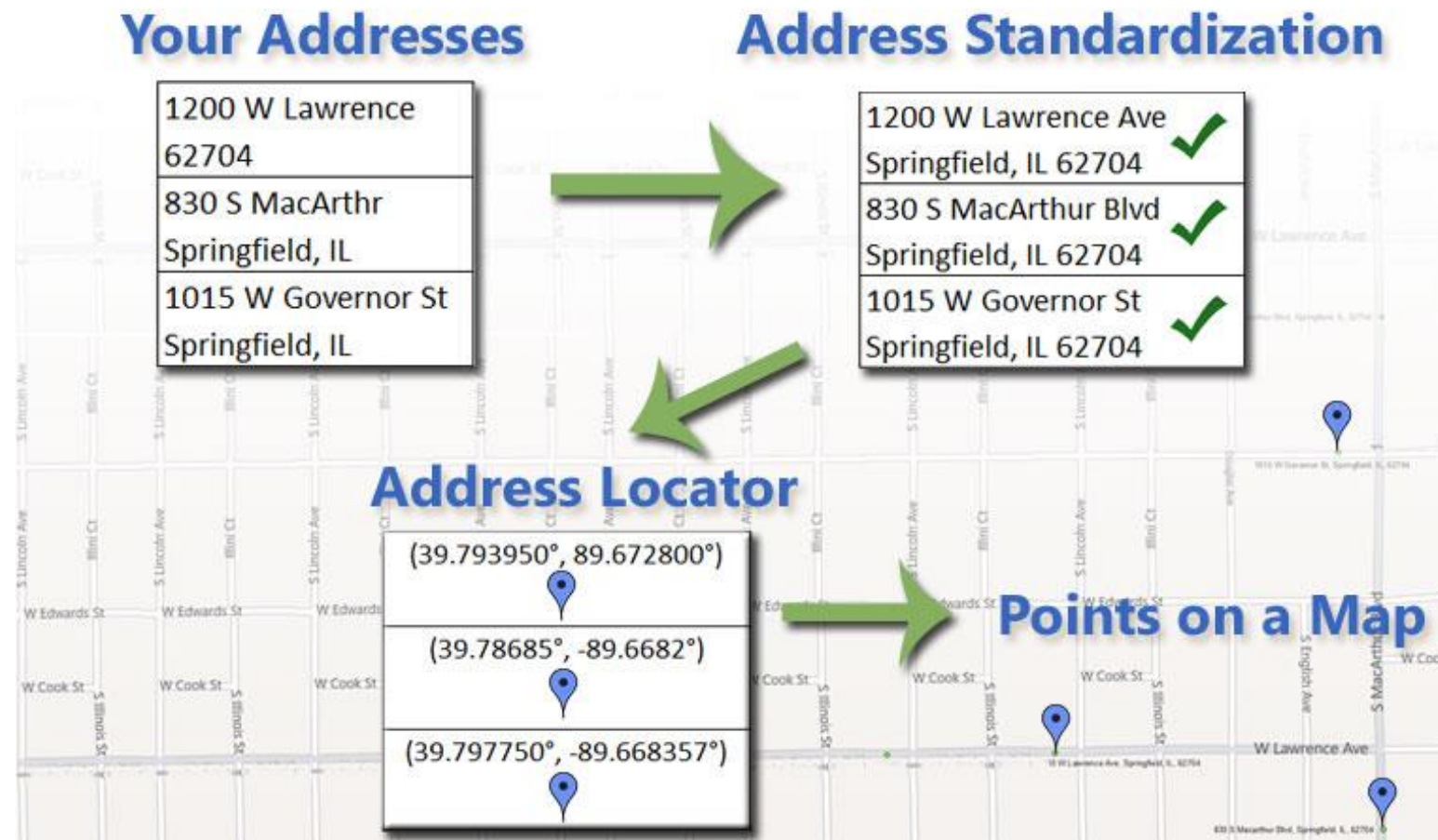


Original data

# How to improve data quality? (3)

- **Data Standardization** - disparate data sets are conformed to a common data format.

A set of rules and designs.



A healthcare domain data standardization example

# How to improve data quality? (4)

- **Geocoding** - The description of a location is transformed into coordinates that conform to U.S. and worldwide geographic standards.
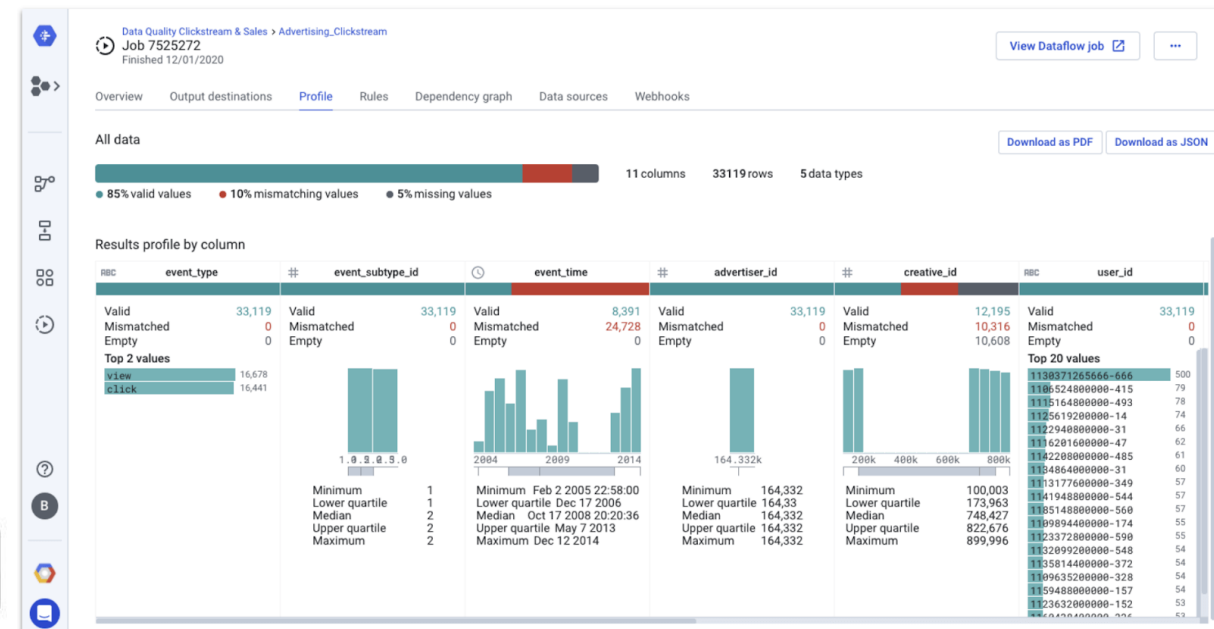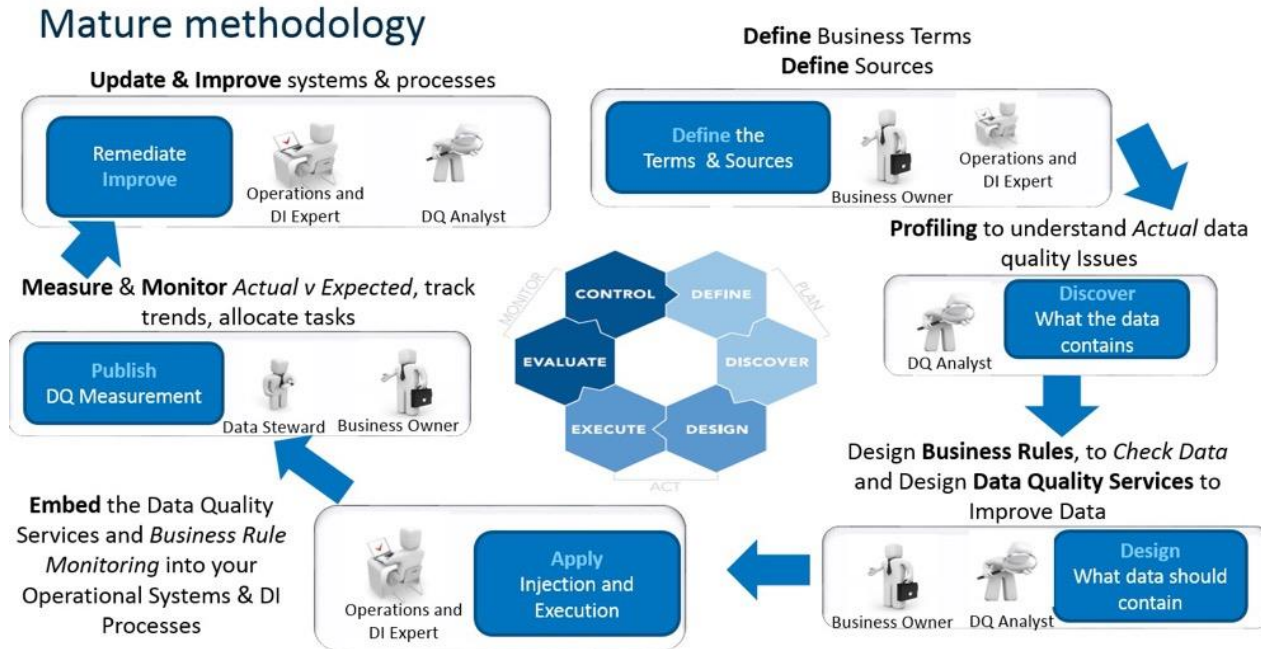


**Your Addresses**

| |
|---|
| 1200 W Lawrence 62704 |
| 830 S MacArthr Springfield, IL |
| 1015 W Governor St Springfield, IL |

**Address Standardization**

| | |
|---|---|
| 1200 W Lawrence Ave Springfield, IL 62704 | ✓ |
| 830 S MacArthur Blvd Springfield, IL 62704 | ✓ |
| 1015 W Governor St Springfield, IL 62704 | ✓ |

**Address Locator**

| |
|---|
| (39.793950°, 89.672800°) |
| (39.78685°, -89.6682°) |
| (39.797750°, -89.668357°) |

**Points on a Map**

# How to improve data quality? (5)

- **Matching or Linking** - identifies and merges matching pieces of information in big data sets. Finding records in a data set that refer to the same entity across different data sources (e.g., data files, books, websites, and databases)
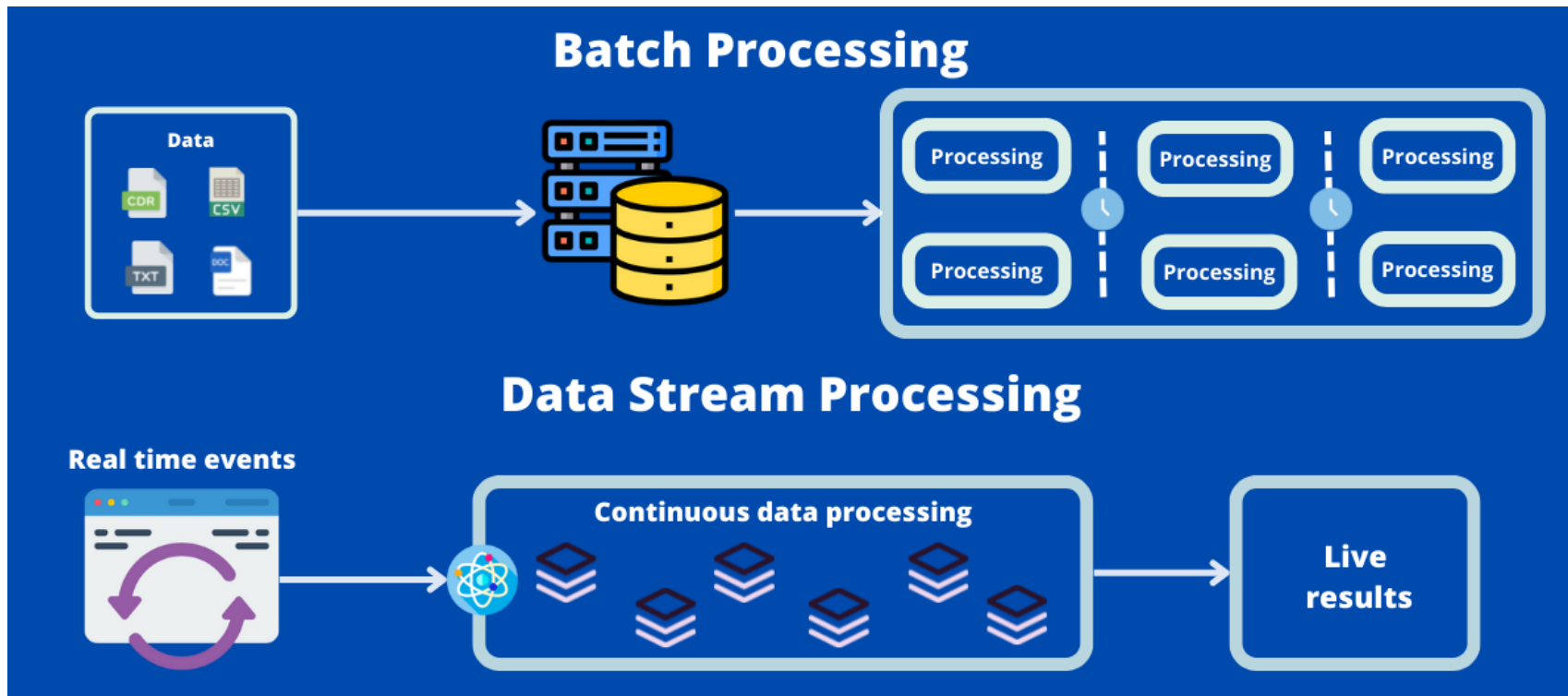


Contact Info
of Joe Smith

Demographic Info
of Joseph A. Smith

Unified Record: Joseph Smith

Purchase History
of J. smith

# How to improve data quality? (6)

- **Data quality monitoring** - frequent data quality checks are essential. Data quality software in combination with machine learning can automatically detect, report, and correct data variations based on predefined business rules and parameters.

# How to improve data quality? (7)

- **Batch and Real time** - Once the data is initially cleansed, an effective data quality framework should be able to deploy the same rules and processes across all applications and data types at scale over time.

# What is data security?

- **Data security is the process** of safeguarding digital information throughout its entire life cycle to protect it from corruption, theft, or unauthorized access.

    - It covers everything—hardware, software, storage devices, and user devices; access and administrative controls; and organizations' policies and procedures.

- **Data security uses tools and technologies** that enhance visibility of a company's data and how it is being used. These tools can protect data through processes like data masking, encryption, and redaction of sensitive information. The process also helps organizations streamline their auditing procedures and comply with increasingly stringent data protection regulations.

# Data security risks (1)

Organizations face an increasingly complex landscape of security threats with cyberattacks being launched by more sophisticated attackers. Some of the biggest risks to data security include:

- **Accidental Data Exposure**

Many data breaches are not a result of hacking but through employees accidentally or negligently exposing sensitive information. Employees can easily lose, share, or grant access to data with the wrong person, or mishandle or lose information because they are not aware of their company's security policies.

- **Insider Threats**

One of the biggest data security threats to any organization is its own employees. Insider threats are individuals who intentionally or inadvertently put their own organization's data at risk.

**Types of Insider Threats**

**Negligent**
Insiders who pose an unintentional threat due to human error and lack of security awareness

**Malicious**
Current or former employees who abused their access to steal intellectual property for personal gains

**Third - Party**
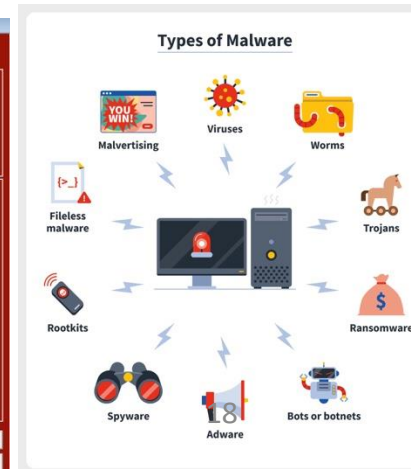Vendors who misuse their access and compromise the security of critical data
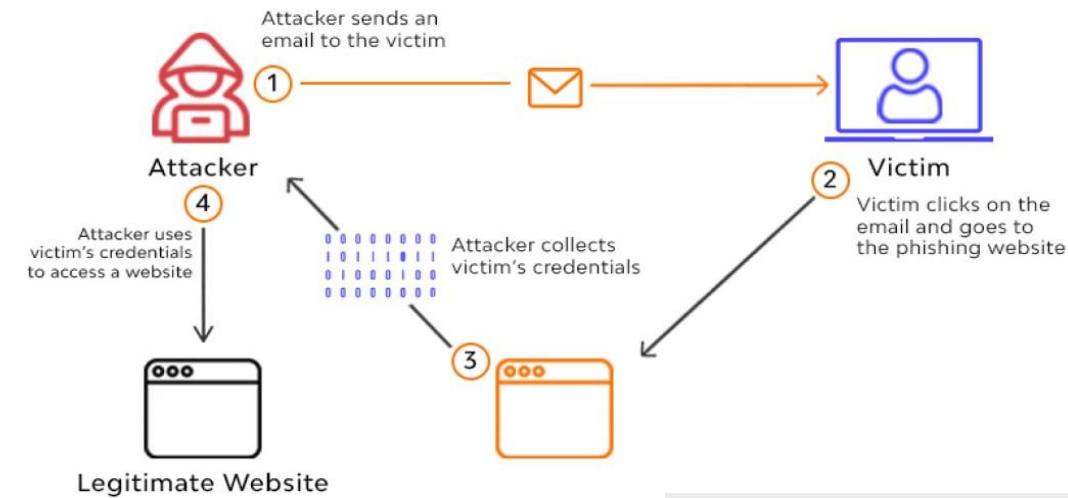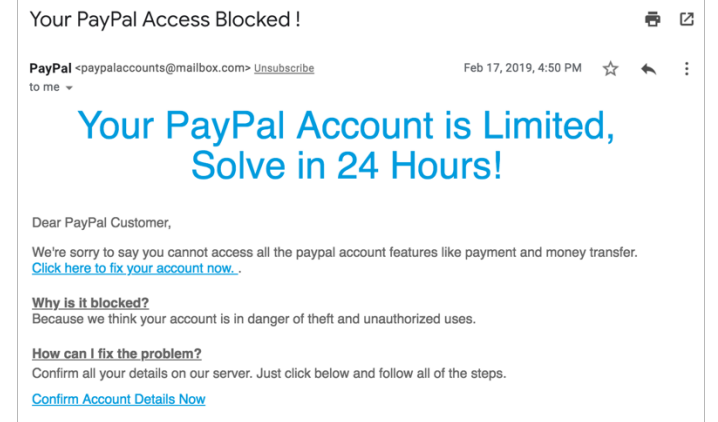
# Data security risks (2)

- **Phishing Attacks**

In a phishing attack, a cyber criminal sends messages, typically via email, short message service (SMS), or instant messaging services, that appear to be from a trusted sender. Messages include malicious links or attachments that lead recipients to either download malware or visit a spoofed website that enables the attacker to steal their login credentials or financial information.

- **Malware**

Malicious software is typically spread through email- and web-based attacks. Attackers use malware to infect computers and corporate networks by exploiting vulnerabilities in their software, such as web browsers or web applications. Malware can lead to serious data security events like data theft, extortion, and network damage.

# How to achieve data security? (1)

Organizations can use a wide range of data security methods to safeguard their data, devices, networks, systems, and users. Some of the most common types of data security, which organizations should look to combine to ensure they have the best possible strategy, include:
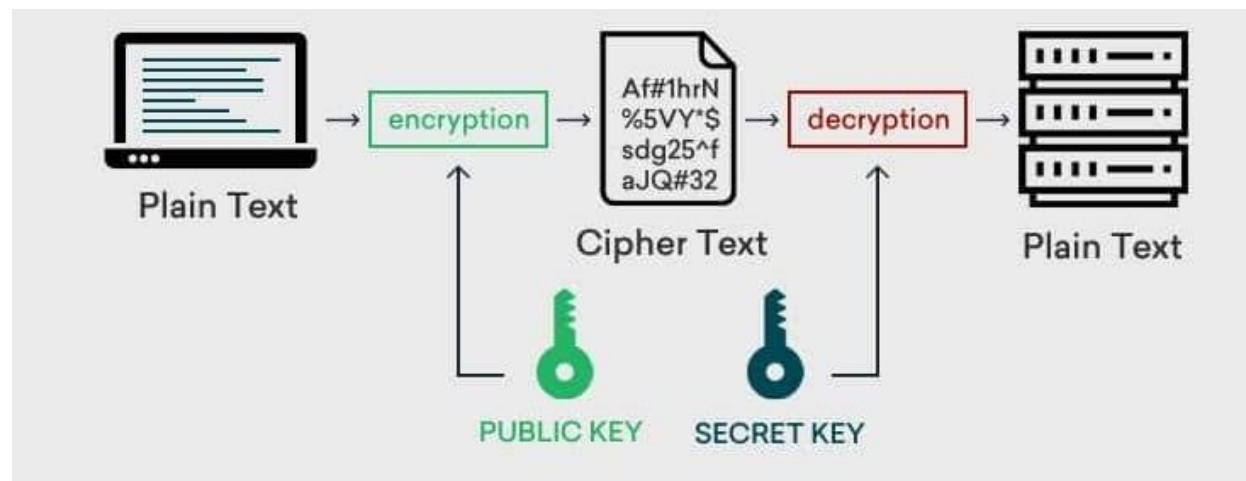
- **Regulation measures:**
  - Monitor user behavior and manage accounts
  - Enforce security policies
  - Provide security awareness training
  - Conduct proactive network monitoring
  - Consistently apply necessary software patches and updates on all systems
  - Implement malware protection to prevent attacks

# How to achieve data security? (2)

- **Encryption**

  Data encryption is the use of algorithms to scramble data and hide its true meaning. Encrypting data ensures messages can only be read by recipients with the appropriate decryption key. This is crucial, especially in the event of a data breach, because even if an attacker manages to gain access to the data, they will not be able to read it without the decryption key.

  Data encryption also involves the use of solutions like tokenization, which protects data as it moves through an organization's entire IT infrastructure.
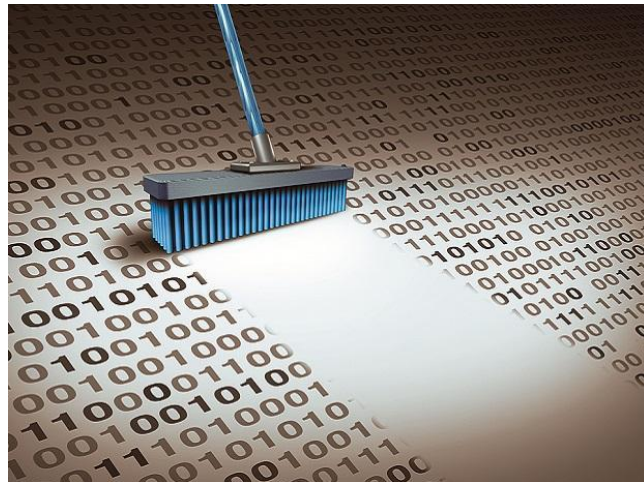
# How to achieve data security? (3)

- **Data Erasure**

    There will be occasions in which organizations no longer require data and need it permanently removed from their systems. Data erasure is an effective data security management technique that removes liability and the chance of a data breach occurring.

    This is a software-based technique for **effectively overwriting electronically stored data** with random binary information as per a given standard.
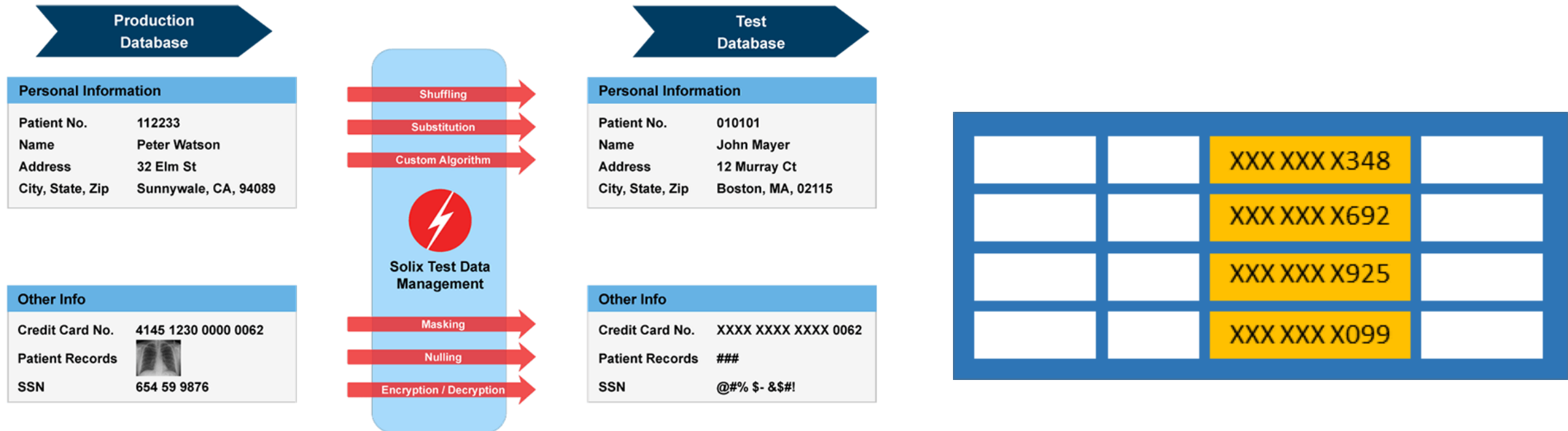
    This is a necessary part of any secure hardware disposal or decommissioning process.
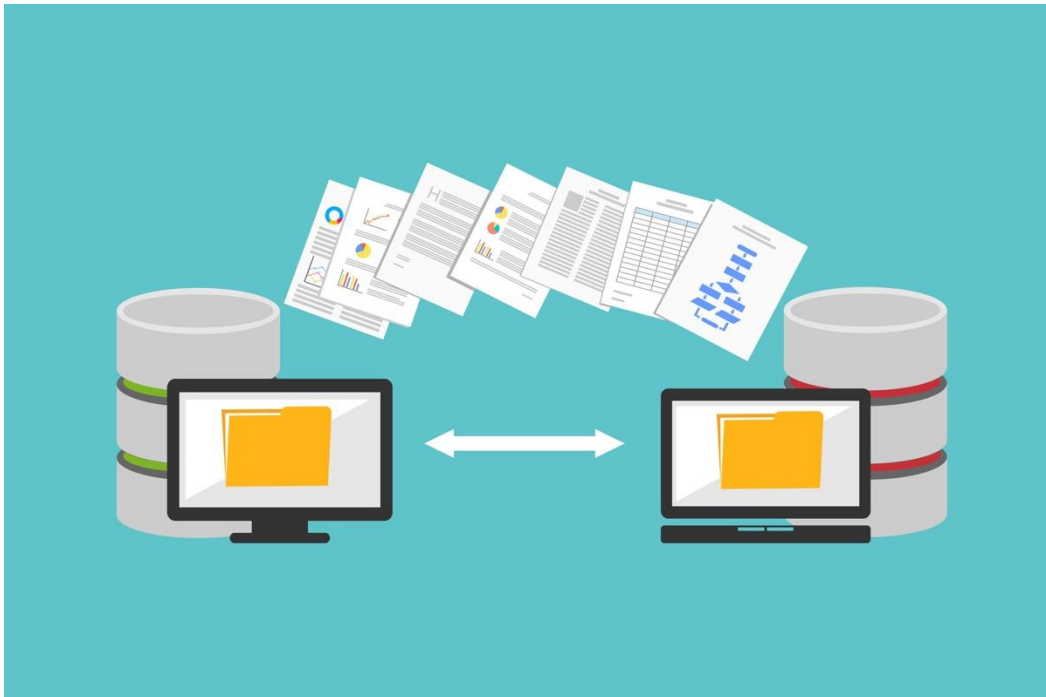
# How to achieve data security? (4)

- **Data Masking**

    Data masking enables an organization to hide data by obscuring and replacing specific letters or numbers. This process is a form of encryption that renders the data useless should a hacker intercept it. The original message can only be uncovered by someone who has the code to decrypt or replace the masked characters.

# How to achieve data security? (5)

- **Data Resiliency**

    Organizations can mitigate the risk of accidental destruction or loss of data by creating backups or copies of their data. Data backups are vital to protecting information and ensuring it is always available. This is particularly important during a data breach or ransomware attack, ensuring the organization can restore a previous backup.

# Data privacy

Data privacy generally means the ability of a person to determine for themselves when, how, and to what extent personal information about them is shared with or communicated to others, or being used.

This personal information can be one's name, location, contact information, or online or real-world behavior.

# Data privacy's focus

- Data Privacy or Information privacy is a part of the data protection area that deals with the proper handling of data focusing on compliance with data protection regulations.

- Data Privacy is centered around how data should be collected, stored, managed, and shared with any third parties, as well as compliance with the applicable privacy laws (such as California Consumer Privacy Act- CCPA or General Data Protection Regulation GDPR).
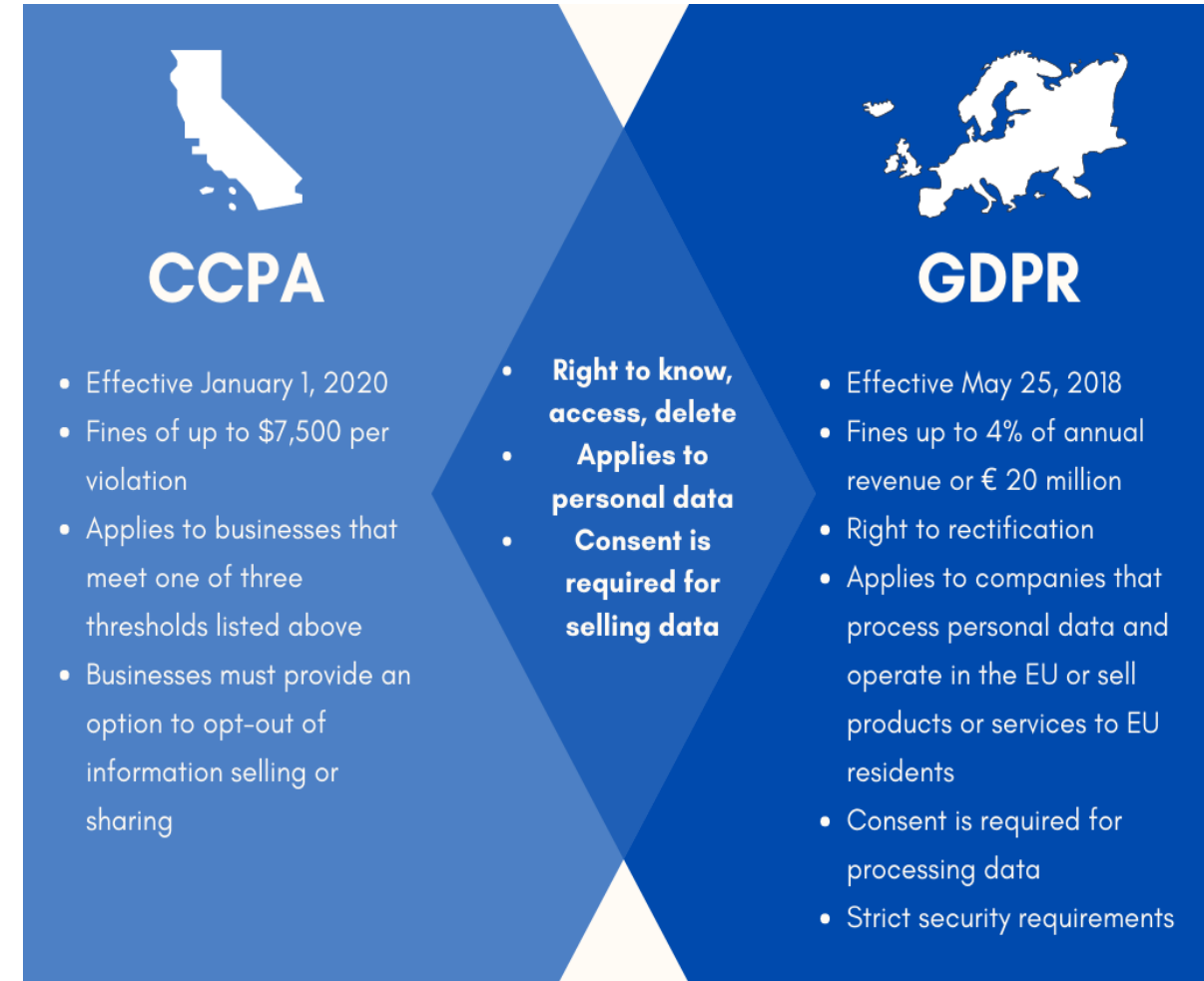
# Data privacy regulations

- **General Data Protection Regulations (GDPR) in Europe**

The 2018 GDPR legislation is a piece of law that protects the personal data of European citizens. It aims to increase people's control and privacy rights over their data and places strict controls on how organizations process that information. GDPR ensures that organizations process personal data securely and protect it from unauthorized processing, accidental loss, damage, and destruction.

- **California Consumer Privacy Act (CCPA) in California, US**

The CCPA aims to give California's consumers more control over how businesses collect their personal data. This includes the right to know what information a business has and how it is shared or used, the right to delete that information, the right to opt out of that data being sold to third parties, and the right to avoid discrimination for exercising these CCPA rights. Organizations must provide consumers with notice of their privacy practices.

## CCPA

- Effective January 1, 2020
- Fines of up to $7,500 per violation
- Applies to businesses that meet one of three thresholds listed above
- Businesses must provide an option to opt-out of information selling or sharing

- **Right to know, access, delete**
- **Applies to personal data**
- **Consent is required for selling data**

## GDPR

- Effective May 25, 2018
- Fines up to 4% of annual revenue or € 20 million
- Right to rectification
- Applies to companies that process personal data and operate in the EU or sell products or services to EU residents
- Consent is required for processing data
- Strict security requirements

# Data privacy regulations

- **Health Insurance Portability and Accountability Act (HIPAA) in US**

HIPAA of 1996 is a United States Act of Congress enacted federal law. It protects patients' health data from being exposed without their consent or knowledge. HIPAA contains a privacy rule, which addresses the disclosure and use of patient information and ensures that data is properly protected. It also has a security rule, which protects all individually identifiable health information that an organization creates, maintains, receives, or transmits electronically.

- **Personal Data (Privacy) Ordinance (PDPO) in HK**

The purpose of this 1995 ordinance is to protect the privacy rights of a person in regard to his personal data, i.e., the Data Subject: the information which relates to a living person and can be used to identify that person and it exists in a form in which access or processing is practicable. Examples of data subject protected by this ordinance include name, address, phone number, identity card number, photo, medical record and employment records. The data user, who collects, holds, or process this data is liable for any unlawful or wrongful use of this data.



TOP 10 CONSIDERATIONS FOR A
**HIPAA-COMPLIANT WEBSITE**

- HIPAA PRIVACY RULE
- HIPAA SECURITY RULE
- SSL ENCRYPTION
- HIPAA-COMPLIANT WEBSITE PLATFORM
- BUSINESS ASSOCIATE AGREEMENTS
- HEALTHCARE FOCUS OF INFRASTRUCTURE
- SECURITY OF DATA CENTER & AUDITING
- ONSITE AND OFFSITE BACKUPS
- MANAGED MULTI-FACTOR AUTHENTICATION
- MANAGED FIREWALL



- The six data protection principles form the base of the Ordinance.

- Data users must comply with the six data protection principles in the collection, holding, accuracy, retention period, security, privacy policy and access to and correction of personal data.

香港個人資料私隱專員公署
Office of the Privacy Commissioner for Personal Data, Hong Kong

27

# Why is data privacy regulation important?

- **For individuals**

Privacy laws around the world aim to give back individuals control over their data, empowering them to know how their data is being used, by whom and why, giving them control over how their personal data is being processed and used.

Organizations that collect personal data are obligated to respond to those questions and manage personal data in a compliant way. According to Gartner's predictions for the future of privacy, privacy is today what "organic" or "cruelty-free" was in the past decade.

- **From A Business Perspective**

Businesses can not operate without processing personal data in some way. However, in order to stay compliant, companies now have to manage personal data in a transparent and compliant way, be accountable for personal data they process, and adhere to privacy principles.
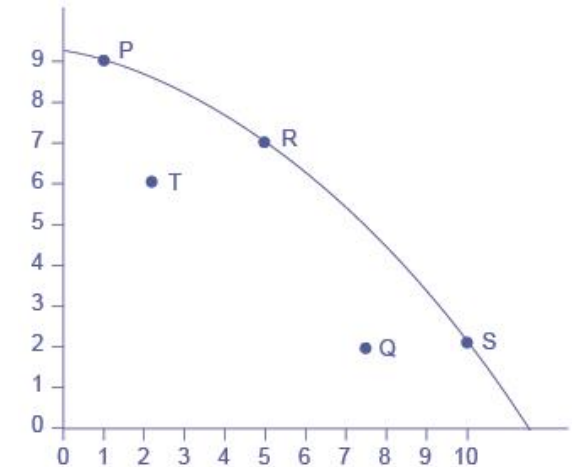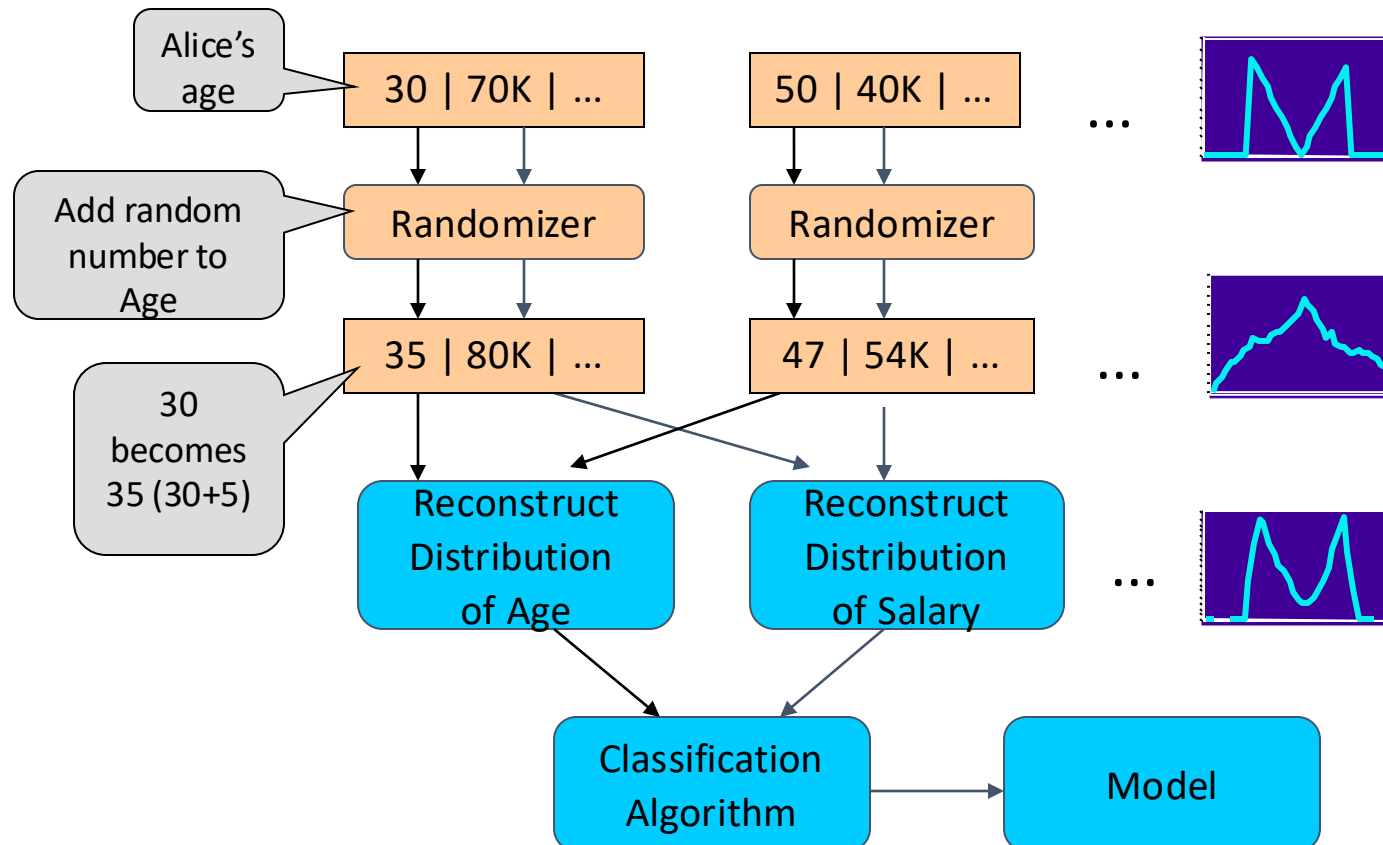
Otherwise, they risk huge regulatory fines, loss of customers' trust, investor appeal, and data breaches.

# Importance of transparency

- In this age of **data economy**, true company value lies in the collected personal data. This means data is an asset worthy of protecting and keeping.

- What companies keep forgetting is that the personal data of individuals processed by the companies are only borrowed.

- Privacy laws enable individuals to exercise their rights, such as the **right to be forgotten**, and in certain circumstances, **individuals can take back ownership of their data.**

- In order for companies to keep the data and keep the trust, they will have to demonstrate transparency by openly communicating how they process and manage personal data.

# Privacy preserving methods: data perturbation

- Perturb data with value distortion
  - User provides $x_i + r$ instead of $x_i$
  - $r$ is a random value
    - Uniform, uniform distribution between $[-\alpha, \alpha]$
    - Gaussian, normal distribution with $\mu = 0$, $\sigma$



Alice's age

30 | 70K | ...    50 | 40K | ...    ...

Add random number to Age

Randomizer    Randomizer

30 becomes 35 (30+5)

35 | 80K | ...    47 | 54K | ...    ...

Reconstruct Distribution of Age    Reconstruct Distribution of Salary    ...

Classification Algorithm → Model

Privacy-performance tradeoff
More noises
higher privacy, lower performance

# Privacy preserving methods: anonymity

- A release of data is said to have the k-anonymity property if the information for each person contained in the release cannot be distinguished from at least k−1 other individuals whose information also appear in the release.
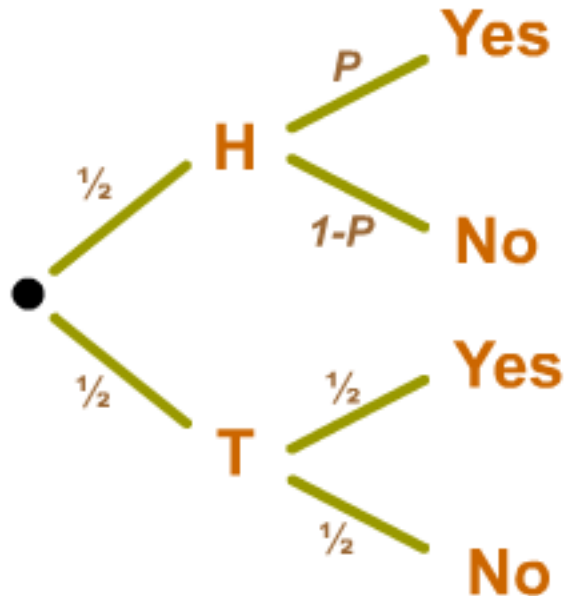
- Hiding a person among k of them.



| Bob | |
|-----|-----|
| *Zip* | *Age* |
| 47678 | 27 |

A 3-diverse patient table

| Zipcode | Age | Salary | Disease |
|---------|------|--------|---------------|
| 476** | 2* | 20K | Gastric Ulcer |
| 476** | 2* | 30K | Gastritis |
| 476** | 2* | 40K | Stomach Cancer |
| 4790* | ≥40 | 50K | Gastritis |
| 4790* | ≥40 | 100K | Flu |
| 4790* | ≥40 | 70K | Bronchitis |
| 476** | 3* | 60K | Bronchitis |
| 476** | 3* | 80K | Pneumonia |
| 476** | 3* | 90K | Stomach Cancer |

# Privacy preserving methods: randomized response

- Randomized response: allows respondents to respond to sensitive issues

  (such as criminal behavior) while maintaining confidentiality.

  - For example, in an interview, "did you go to the bar last weekend?".

  - Before they answer, they flip a coin. They are then instructed to answer "yes" or "no" randomly if the coin comes
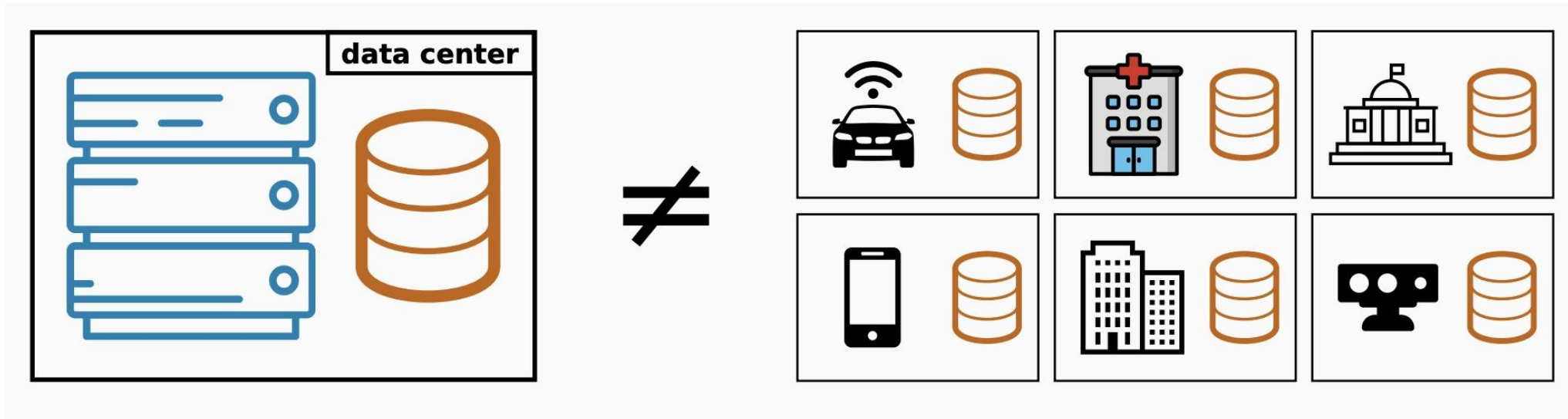
    up tails, and truthfully, if it comes up heads.

If for example, we get 60% of "yes" and 40% of "no", then we
kind of know the true "yes" rate.
Solving this equation: 0.25+0.5p = 0.6, 0.25+0.5(1-p) = 0.4
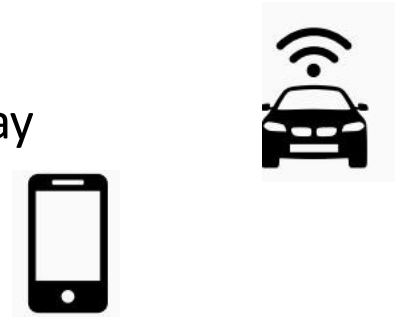
# Privacy preserving methods: federated learning (1)

- The standard setting in Machine Learning (ML) considers a centralized dataset processed in a tightly integrated system.

- But in the real world data is often decentralized across many parties.

# Privacy preserving methods: federated learning (2)

- Sending the data may be too costly

  - Self-driving cars are expected to generate several TBs of data a day

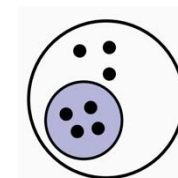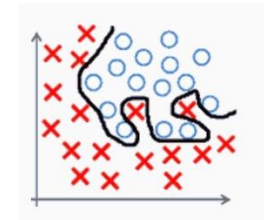  - Some wireless devices have limited bandwidth/power

- Data may be considered too sensitive

  - We see a growing public awareness and regulations on data privacy

  - Keeping control of data can give a competitive advantage in business and research
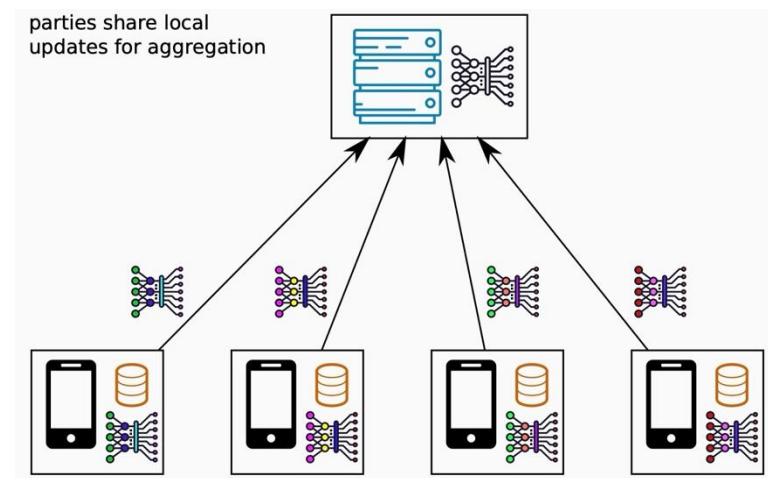
- The local dataset may be too small and biased

  - Bad predictive performance (e.g., due to overfitting)

  - Non-statistically significant results (e.g., medical studies)

  - Not representative of the target distribution

# Privacy preserving methods: federated learning (3)

- Federated learning is a machine learning technique that trains an algorithm across multiple decentralized edge devices or servers holding local data samples, without exchanging them.

- This approach stands in contrast to traditional centralized machine learning techniques where all the local datasets are uploaded to one server, as well as to more classical decentralized approaches which often assume that local data samples are identically distributed.

parties share local
updates for aggregation

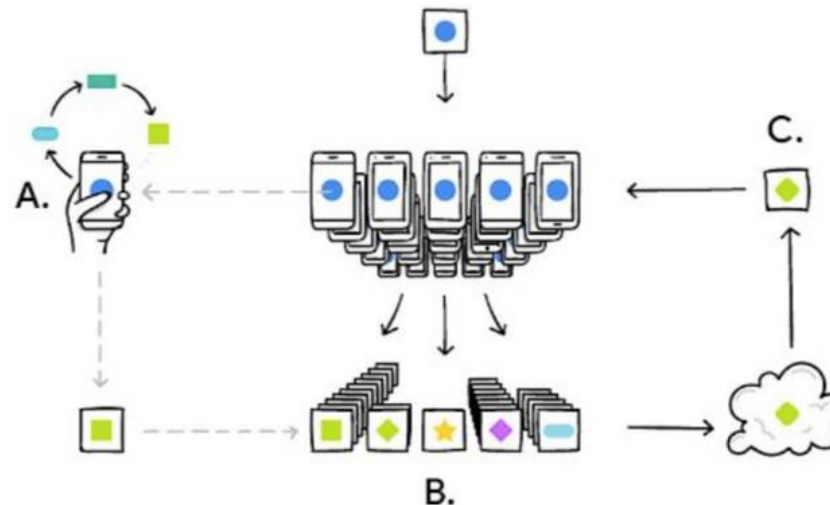In 2016, the term FL was first coined by Google researchers

# Privacy preserving methods: federated learning (4)

- Federated learning enables multiple agents to build a common, robust machine learning model without sharing data, thus allowing to address critical issues such as data privacy, data security, data access rights and access to heterogeneous data. Its applications are spread over a number of industries including defense, IoT, and pharmaceutics.

A: Your phone and your data are used to update a model (blue circle)
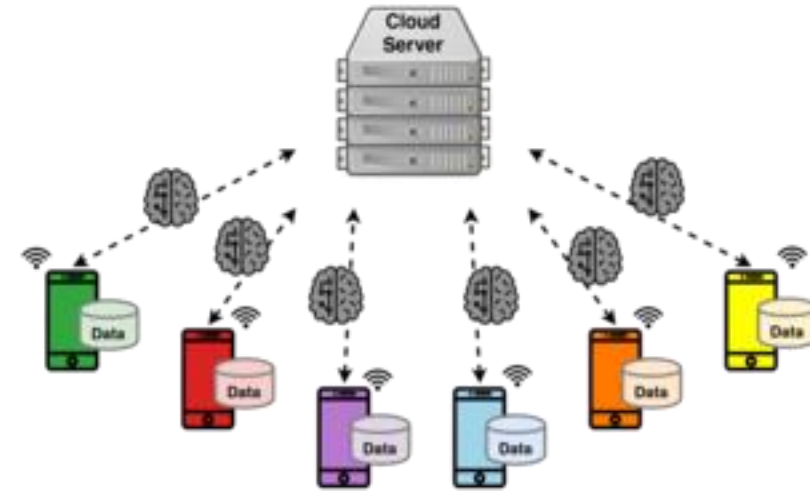
B: All updates from all participants are sorted and sent to the aggregator.

C: After aggregation, the global model updates are shared back to the participants when a new round can begin.
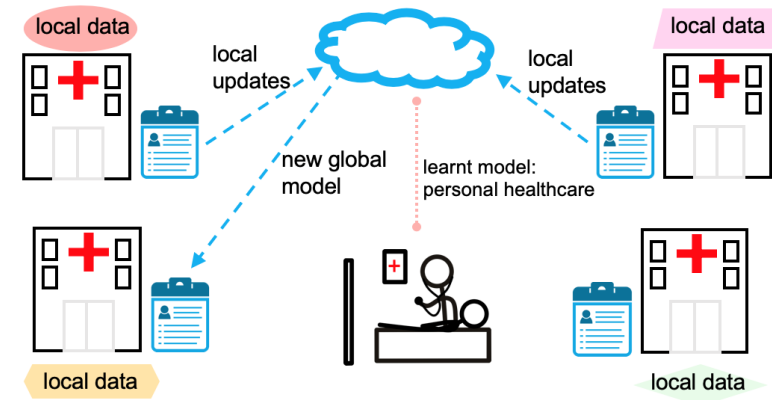
# Privacy preserving methods: federated learning (5)

- Cross-device federated learning
  - Smart phones, apps, IoT, edge devices
  - Massive number of parties (up to $10^{10}$)
  - Small dataset per party (could be size 1)
  - Limited availability and reliability
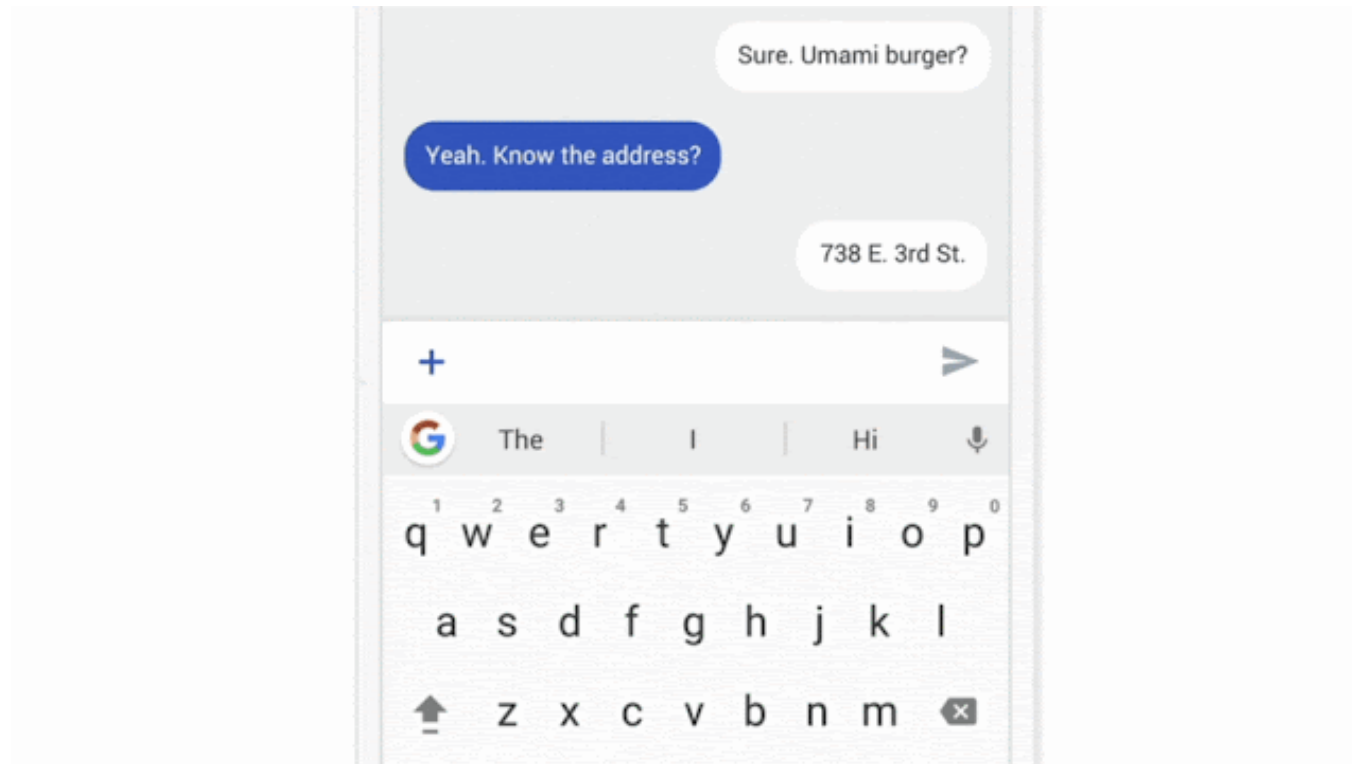  - Some parties may be malicious



- Cross-silo federated learning
  - Institutions, organizations, hospitals, etc.
  - Small number of parties (2-100)
  - Medium to large dataset per party
  - Reliable parties, almost always available
  - Parties are typically honest

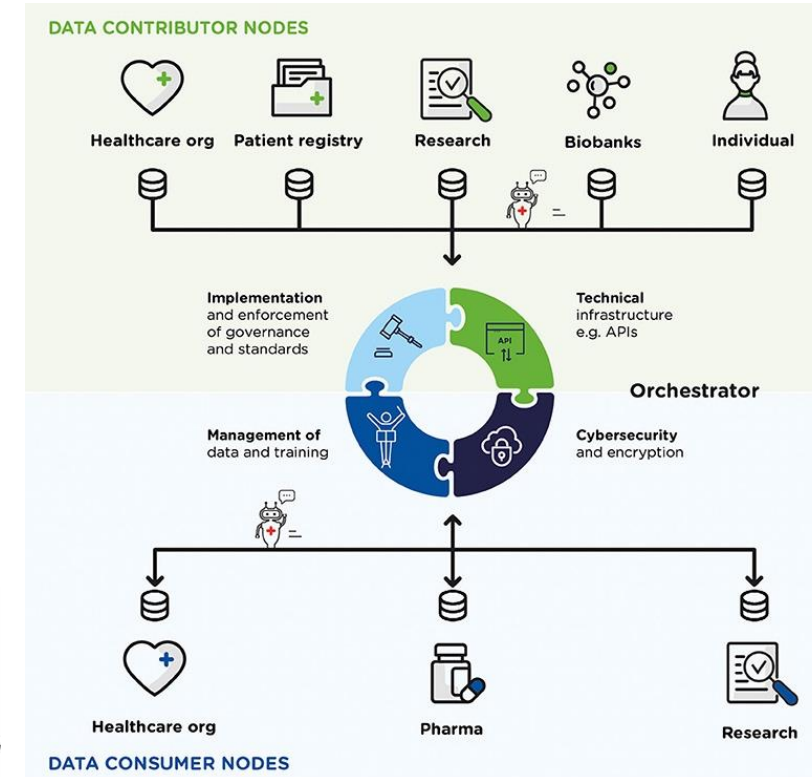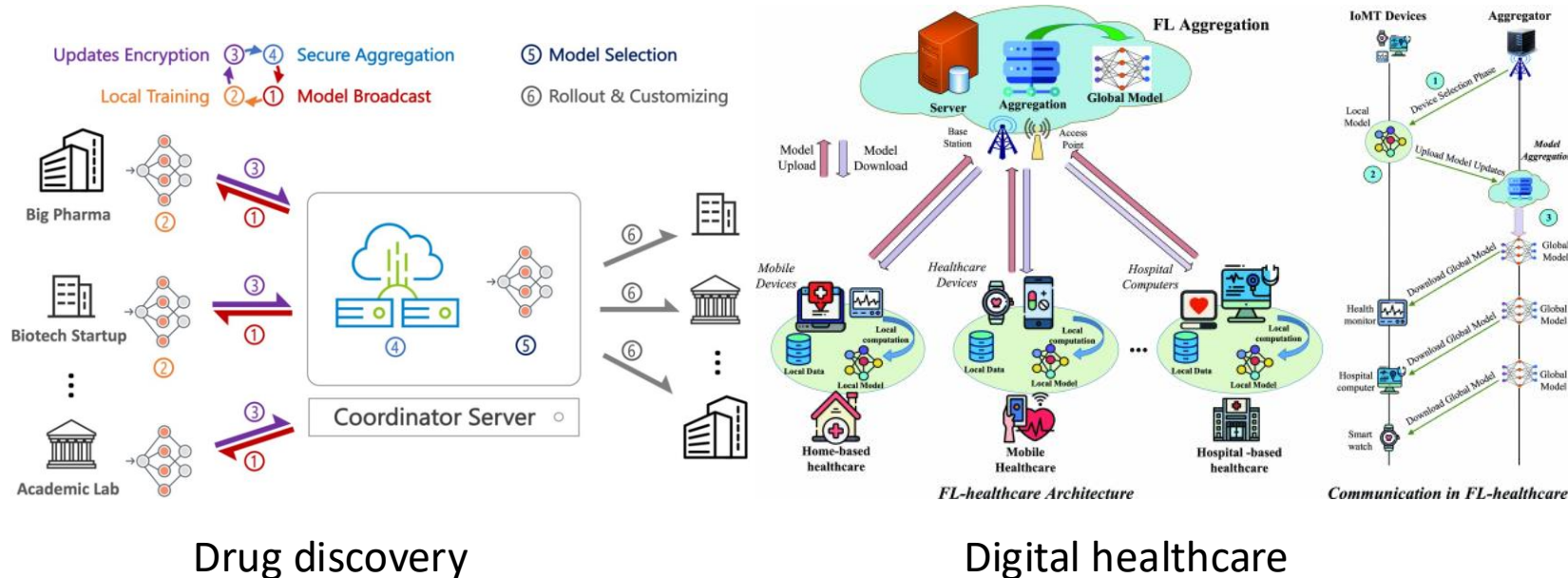# Privacy preserving methods: federated learning (6)

- Examples
  - Google's Gboard on Android, next typing word prediction federated learning model

# Privacy preserving methods: federated learning (7)

- Examples in healthcare
  - Healthcare data are scarce, sensitive, and distributed
  - Data-driven medicine (e.g., drug discovery: the process by which new candidate medications are discovered).
  - Digital healthcare (digital care programs, technologies with health, healthcare, living, and society to enhance the efficiency of healthcare delivery and to make medicine more personalized and precise).



Healthcare data



Drug discovery



Digital healthcare

# Privacy preserving methods: federated learning (8)

- Challenges

    - Statistical and resource heterogeneity

    - Personalization

    - Communication efficiency

    - Privacy preserving

# Thanks for your attention!

# Appendix

1. https://dataprivacymanager.net/5-things-you-need-to-know-about-data-privacy/

2. https://www.fortinet.com/resources/cyberglossary/data-security#:~:text=Data%20security%20is%20the%20process,and%20organizations'%20policies%20and%20procedures.

3. https://www.heavy.ai/technical-glossary/data-quality