



**UNIMORE**  
UNIVERSITÀ DEGLI STUDI DI  
MODENA E REGGIO EMILIA

# Big Data Integration for E-Health

**Data4SmartHealth 2021 October 27 2021**  
**NOI Techpark Bozen-Bolzano**

**Prof. Sonia Bergamaschi**

Department of Engineering "Enzo Ferrari"

[sonia.bergamaschi@unimore.it](mailto:sonia.bergamaschi@unimore.it)

[www.dbgroup.unimore.it](http://www.dbgroup.unimore.it)



- Internet of Medical Things (IoMT)
- Big Data
- Big Data Technologies: Relational or NoSQL systems?
- Big Data Integration
- Data Science & Data Driven AI for Big Data
- Big Data Integration & Data Science in E-Health
- GDPR and Ealth Data

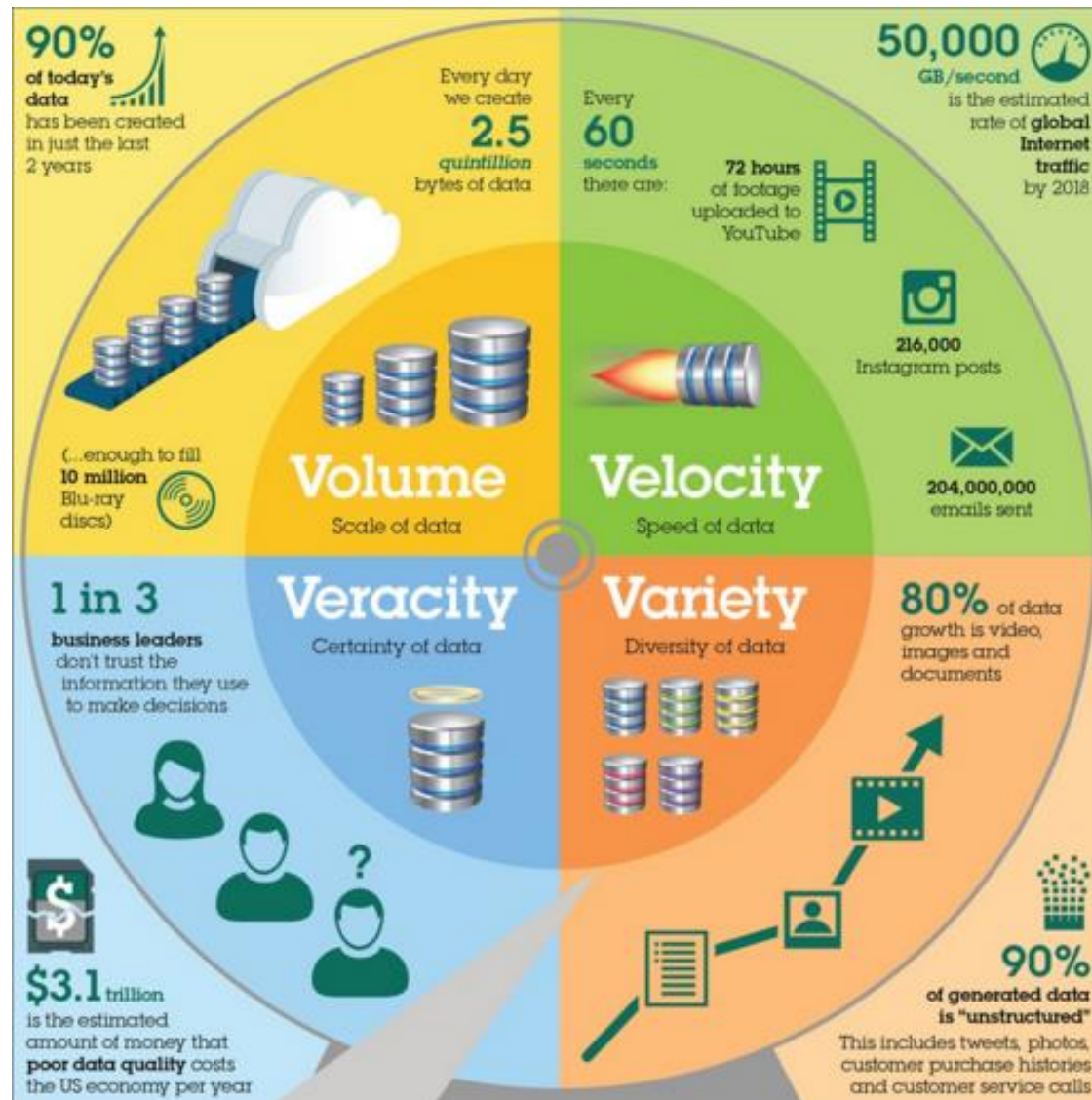


# Big Data - Full faith in the power of data

The quest for knowledge used to begin with grand theories. Now it begins with massive amounts of data. **Welcome to the Petabyte Age!**



# The FOUR V's of Big Data





The production of data is expanding at an astonishing pace. Experts now point to a 4300% increase in annual data generation by 2020. Drivers include the switch from analog to digital technologies and the rapid increase in data generation by individuals and corporations alike.

**2020: MORE THAN 1/3 OF THE DATA PRODUCED WILL LIVE IN OR PASS THROUGH THE CLOUD.**

Size of Total Data  
Enterprise Created Data  
Enterprise Managed Data

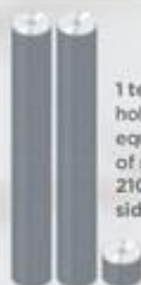
**Only 0.5% to 1% of the data is used for analysis.**

**2012: CUSTOMERS WILL START STORING 1 EB OF INFORMATION.**



### WHAT IS A ZETTABYTE?

1,000,000,000,000	gigabytes
1,000,000,000,000	terabytes
1,000,000,000,000	petabytes
1,000,000,000,000	exabytes
1,000,000,000,000	zettabyte



1 terabyte holds the equivalent of roughly 210 single-sided DVDs.

It took roughly 1 petabyte of local storage to render the 3D CGI effects in Avatar.



In 2007, the estimated information content of all human knowledge was 295 exabytes.

### DATA PRODUCTION WILL BE 44 TIMES GREATER IN 2020 THAN IT WAS IN 2009

More than 70% of the digital universe is generated by individuals. But enterprises have responsibility for the storage, protection and management of 80% of it.\*

# Value – The most important V of all!



- Then there is another V to take into account when looking at *Big Data: Value!*
- Having access to big data is no good unless we can turn it into value.
- What technologies?

## Technologies for Big Data

- Big Data Management
- Big Data Integration
- Big Data Science



God made integers,  
all else is the work of man.

*(Leopold Kronecker, 19<sup>th</sup> Century Mathematician)*

**Codd made relations,  
all else is the work of man.**

*(Raghu Ramakrishnan, DB text book author)*

THE POWER OF INFINITE POSSIBILITIES

Stonebraker Says  
Turing award 2014

**One Size Fits None**  
**“The elephants are toast”**

## At This Point, RDBMS is “long in the tooth”

There are at least 6 (non trivial) markets where a row store can be clobbered by a specialized architecture !

- Warehouse (Vertica, Red Shift, Sybase IQ, DW Appliances)

- OLTP (VoltDB, HANA, Hekaton)

- RDF (Vertica, et. al.)

- Text (Google, Yahoo, ...)

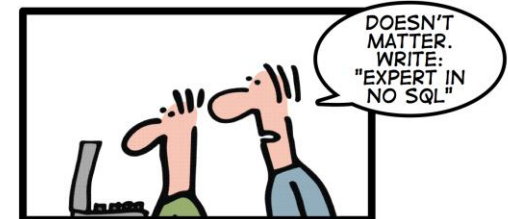
- Scientific data (R, MatLab, SciDB)

- Data Streaming (Storm, Spark Streaming, InfoSphere)

## An emerging “movement” around non-relational software for Big Data

- NOSQL stands for “Not Only SQL” (but is not entirely agreed upon), where SQL doesn’t really mean the query language, but instead it denotes the traditional relational DBMS.
- Google **Bigtable** & **Mapreduce**, **Memcached**, and Amazon’s **Dynamo** are the “proof of concept” that inspired many of the NOSQL systems:
  - Memcached demonstrated that in-memory indexes can be highly scalable, distributing and replicating objects over multiple nodes
  - Dynamo pioneered the idea of *eventual consistency* as a way to achieve higher availability and scalability
  - BigTable demonstrated that persistent record storage could be scaled to thousands of nodes & Mapreduce introduces parallel computation for distributed data platforms.

### HOW TO WRITE A CV



Leverage the NoSQL boom

# Challenges (1) – Selection of the Big Data Technology

- **Volume, Velocity**

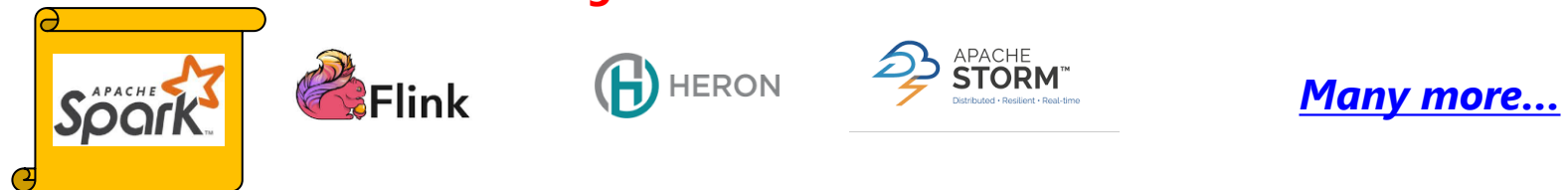
Calling for new **Big Data** systems:

- **Big Data Management Systems: NOSQL & more**



- **Big Data Analysis Systems:**

- **Batch + Streaming**

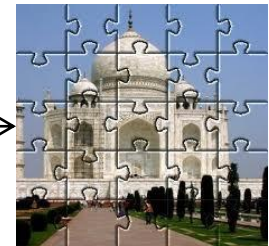
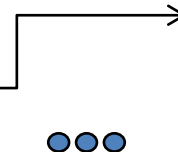


*Not only Relational Database Management Systems  
and Business Intelligence*



# Data integration as a new technological solution

- The discipline of data integration comprises the practices, architectural techniques and tools that ingest, transform, combine and provision data across the spectrum of information types in the enterprise and beyond in order to meet the data consumption requirements of all applications and business processes.
- Applications of Data Integration
  - Business, Science, Government, The Web, Health, Pretty much everywhere
- Small Data integration = solving lots of puzzles
  - Each puzzle (e.g., Taj Mahal) is an **integrated entity**
  - Each piece of a puzzle comes from some **source**
  - Small data integration → solving small puzzles



# Data integration as a new commercial software

## ✓ According to Gartner:

Gartner estimates that the data integration tool market generated more than \$2.7 billion in software revenue (in constant currency) at the end of 2016.

- ✓ A projected five-year compound annual growth rate of 6.32% will bring the total market revenue to around \$4 billion in 2021 (see "Forecast: Enterprise Software Markets, Worldwide, 2014-2021, 2Q17 Update" )
- ✓ ***\$3.3 billion software revenue in 2020.***

## Market Overview:

- ✓ The biggest change in the market from 2016 is the pervasive yet elusive demand for metadata-driven solutions.
- ✓ Consumers are asking for hybrid deployment not just in the cloud and on-premises but also across multiple data tiers throughout broad deployment models, plus the ability to blend data integration with application integration platforms (which is metadata driven in combination with workflow management and process orchestration) and a supplier focus on product and delivery initiatives to support these demands.

# Data integration in the research community

- The research community has been investigating data integration for more than 30 years: different research communities (database, artificial intelligence, semantic web) have been developing and addressing issues related to data integration:
  - Definitions, architectures, classification of the problems to be addressed
  - Different approaches have been proposed and benchmarks developed
- **Open issues**
  - Uncertainty, **Provenance**, and Cleaning
  - Lightweight Integration
  - Visualizing Integrated Data
  - Integrating Social Media
  - **Big Data Integration**



- Data integration = solving lots of puzzles
  - Big data integration → **big messy** puzzles
  - E.g., missing, duplicate, damaged pieces



**Big Data Science**  
**Data Analysis (Business Intelligence, Statistics, Data Mining, Math)**  
**+**  
**DataDriven Artificial Intelligence)**



- From the Big Data era people **do not focus on improvement the quality of data**, but just add more data to overcome errors from noisy and poor-quality information;
- From a recent talk Andrew Ng ([https://it.wikipedia.org/wiki/Andrew\\_Ng](https://it.wikipedia.org/wiki/Andrew_Ng)) states that 99% of the papers were model-centric;
- As results many models do not work well on real data;
- A recent paper from Google researchers analyzes the work of 53 AI practitioners reports that *“data cascades—compounding events causing negative, downstream effects from data issues—triggered by conventional AI/ML practices that undervalue data quality... are pervasive (92% prevalence), invisible, delayed, but often avoidable.”*

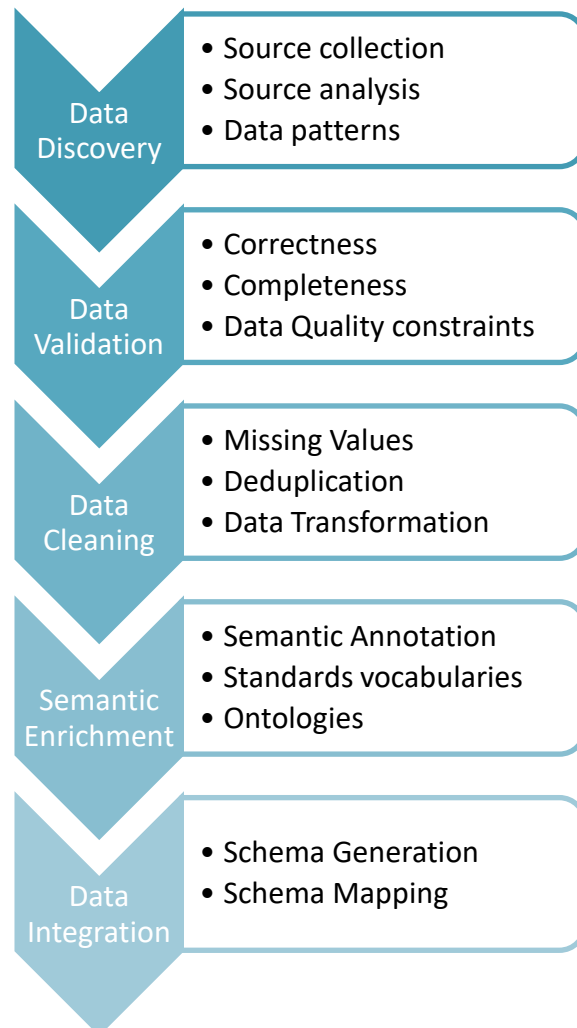
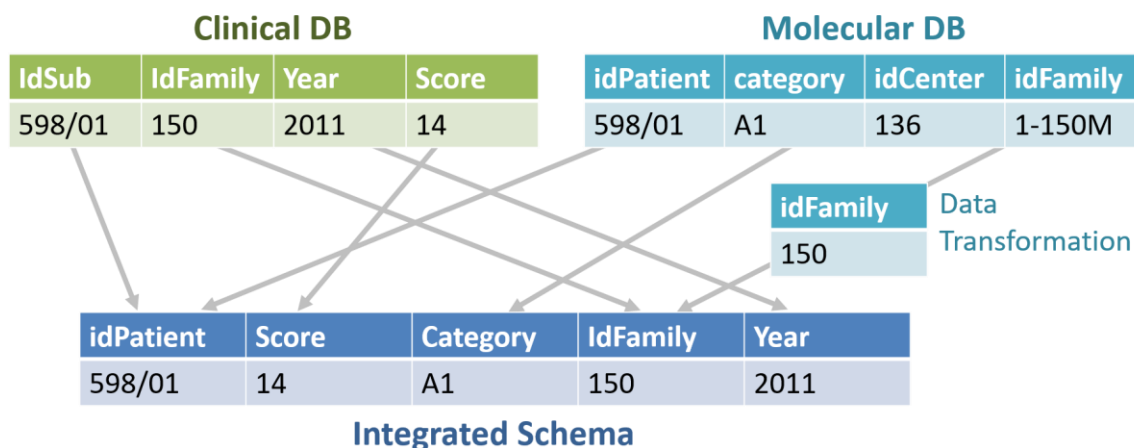
Model-Centric	Data-Centric
<ul style="list-style-type: none"><li>- Collect as much data as possible</li><li>- Iteratively <b>improve the model</b> to deal with the noise in the data</li></ul>	<ul style="list-style-type: none"><li>- Hold the model fixed</li><li>- Iteratively <b>improve the quality of the data</b> to obtain good results</li></ul>

- **Data education lack** of adequate training on AI data quality, collection, and ethics. AI courses focuses on toy datasets with clean values, but AI in practice required the creation of data pipelines, often from scratch, going from groundtruth to model maintenance.
- We have to define a **systematic pipeline** to improve the quality of data;
- Systematic **improvement of data quality** on a basic model is better than using the state-of-the-art models with low-quality data;
- In his recent talk Andrew Ng states that **good data** for ML/AI is:
  - Defined consistently (the label definition is unambiguous)
  - Cover important cases (good coverage of inputs)
  - Has a feedback from the production data
  - Sized appropriately

# Big Data Integration in Healthcare



**MOMIS** is a data integration system able to aggregate data from heterogeneous (structured and semi-structured) and distributed sources (EHR, PHR, eCRF, ePRO, medical devices) in a semi-automatic way, exploiting the **semantic relationships** existing in the data sources (available as open source by Datariver [www.datariver.it](http://www.datariver.it)).



- 2009 – Founded as Spin-Off of the University of Modena and Reggio Emilia
- 2011 – Accreditation as **Industrial Research Lab of the High Technology** Network of the Emilia-Romagna Region
- 2017 Self-certification as **CRO** (Contract Research Organization) by Italian Medicines Agency (AIFA) for **Data Management** and **Statistical Analysis**
- 2017 **Digital Innovation Award in Health** - Politecnico of Milano Observatory
- 2019 **EU INNOLABS Acceleration Programme Award**
- 2021 **EU INNO4COV-19 Funding**
- 2021 Certification as **EHDEN Certified-SME**





## Mission

- Provide innovative solutions for Clinical Trials, Patient Support Programs (PSP), pathology and rare disease registries to Pharma and Biotech companies, Clinical research institutes and Hospitals
- Specialized in designing and developing Web and Mobile software solutions for **Clinical Data Management, Big Data Integration & Analytics, Internet of Medical Things (IoMT).**





**MyHealth** is a Web and Mobile IoTM platform for monitoring and **improving quality of life of individuals**

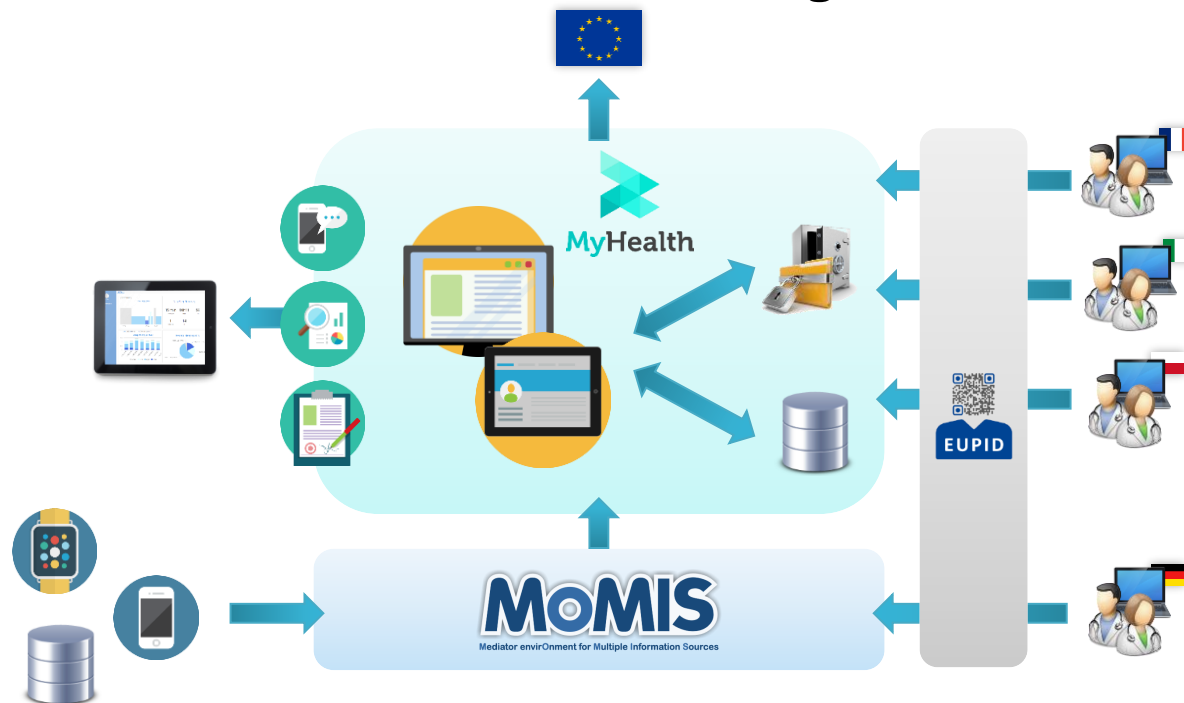
- MyHealth active projects:
  - Patients with **specific pathologies**  
12 Clinical trials, 4 Patient Support Programs (PSP)
  - Monitoring **elderly people**  
16 residential care & nursing homes
  - **Promoting multi-aging physical activity**  
1 research project on open air recreational activities
  - MyHealth-COV **EU INNO4COV-19 Funding**



# PARTNER – Paediatric Rare Tumors Network

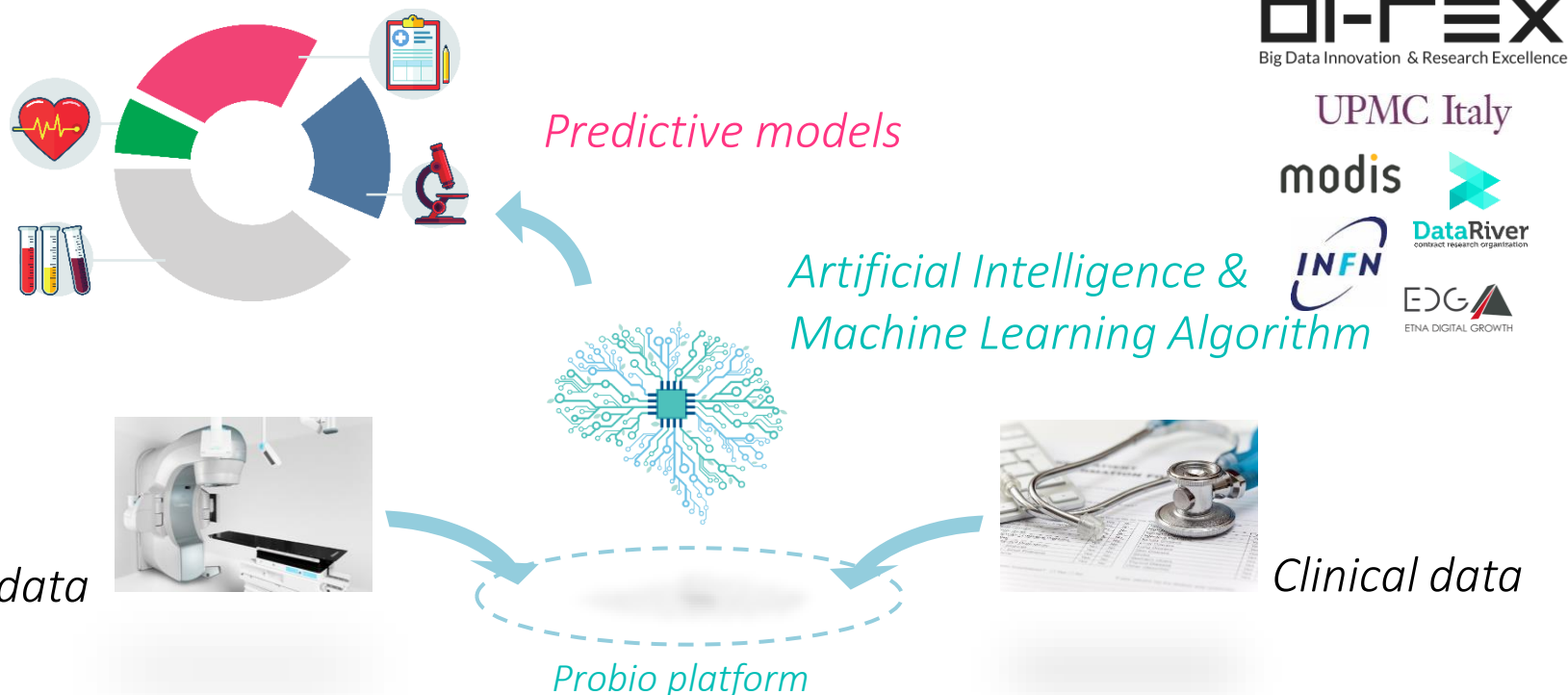
**MOMIS system for virtual data integration** of heterogeneous, fragmented and not standardized data from **Biobank databases, existing registries and other distributed data sources**.

- ✓ secure access to the information **without migration of data**;
- ✓ Elaboration of **indexes** that process the patient outcomes;
- ✓ increase **information** and **knowledge** about diseases and patients.

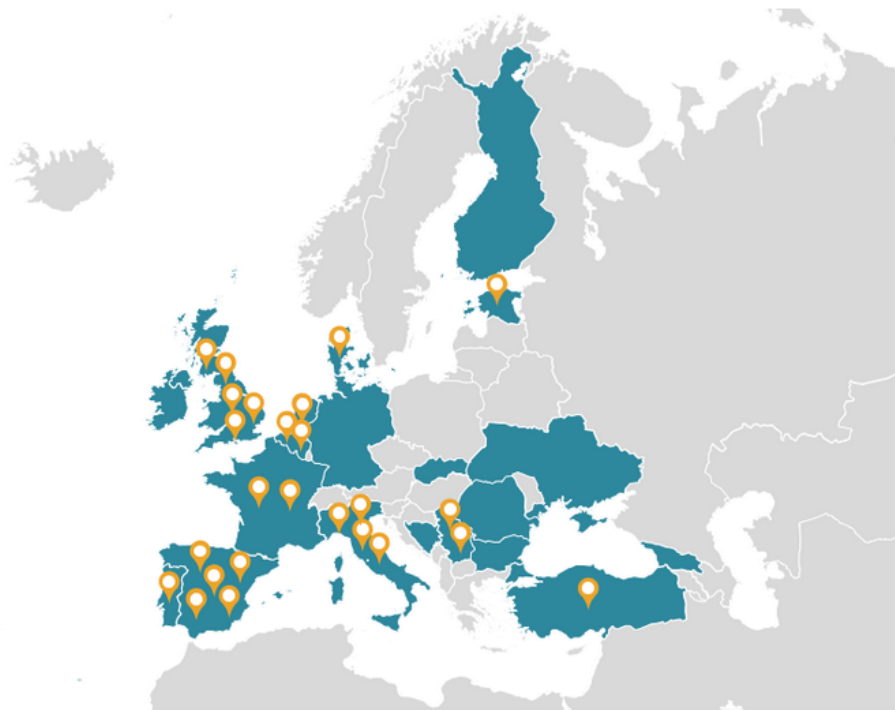


## Decision support tool

AI and ML algorithm applied to the integration of radiological images with clinical and laboratory data to **identify new information and relationships between data** that can subsequently be refined and extended to further sources to expand and increase the **prognostic value** of the platform.



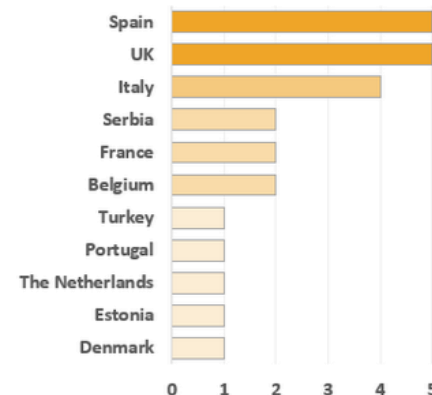
# EHDEN – COVID19 Rapid Collaboration Call



Awarded applicants



Applicant countries



- ✓ 25 Data Partners
- ✓ Over 1 million SARS-COV-2-tested patients
- ✓ 228,000 of whom tested positive

## Mapping

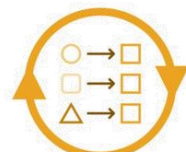


EHDEN COVID-19 Taskforce

Intense collaboration



Data Partner



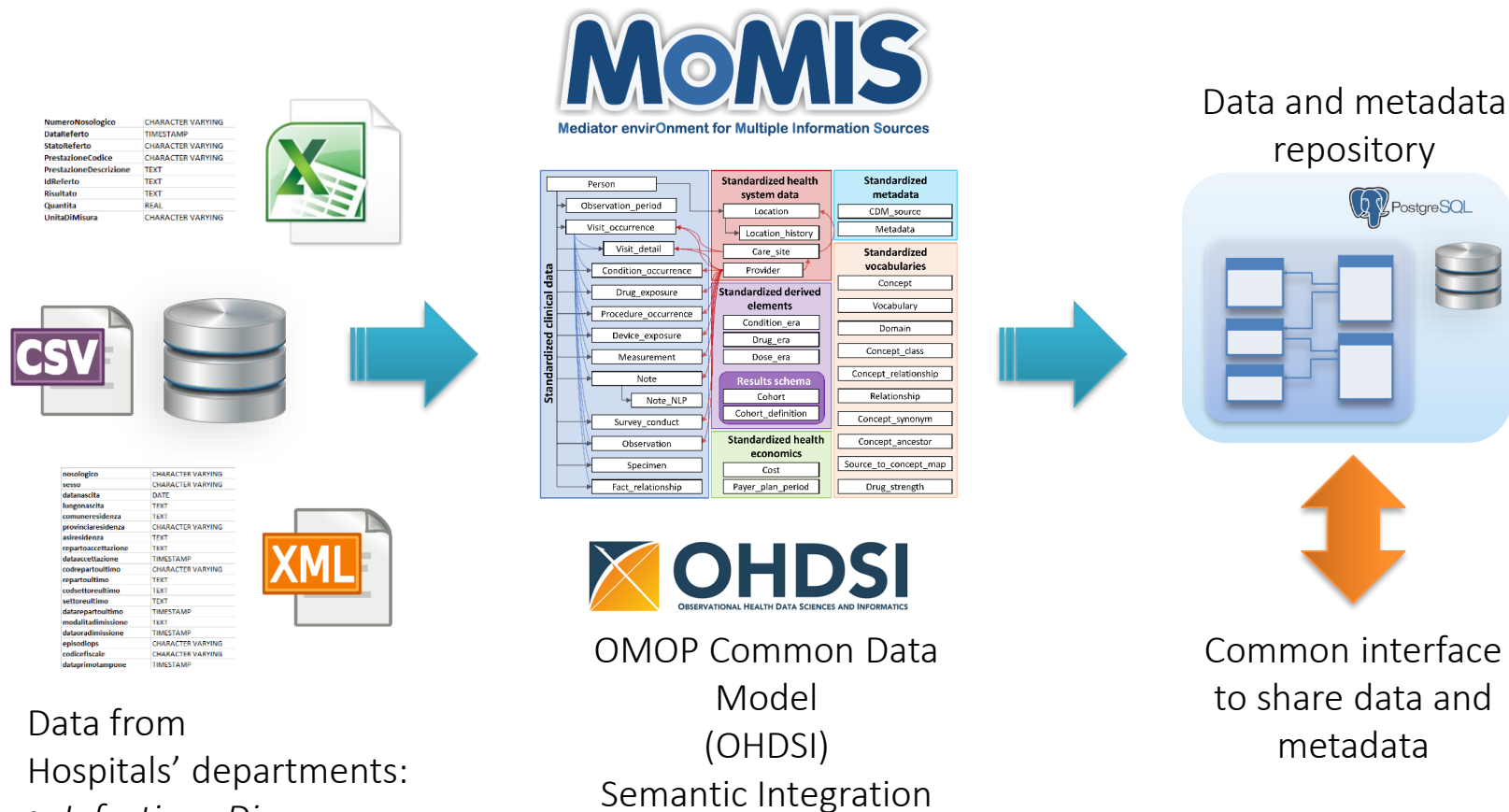
ETL development & iterative mapping review



Engagement in the EHDEN research network



# EHDEN – COVID19 Rapid Collaboration Call



Data from Hospitals' departments:

- *Infectious Diseases*
- *COVID*
- *Intensive and Sub-intensive Care Units*

Assist home care patients with COVID-19 and during the follow-up providing a remote and continuous monitoring of their parameters and tele support, improving the intensive care and fight the pandemic.

- **Telemonitoring system:** patient remote monitoring using wearables devices to collect patient physiological parameters and a mobile App for ePRO questionnaires;
- **AI data analysis** for multidimensional monitoring to engage and support patients at home at the right time;
- **Vocal assistant tool** for patient support and engagement;
- **Televisit system** for medical staff-patient video visits;
- **Web and mobile App** for doctor-patient communication.



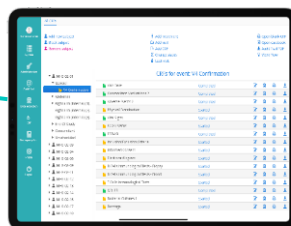
Videovisit and  
Tele support



Patients Mobile  
App  
Patient Reported  
Outcomes (ePRO)



Medical devices to  
collect physiological  
data



Electronic Health Records  
(EHR)



Doctor-patient  
communication Vocal  
Assistant & ChatBot

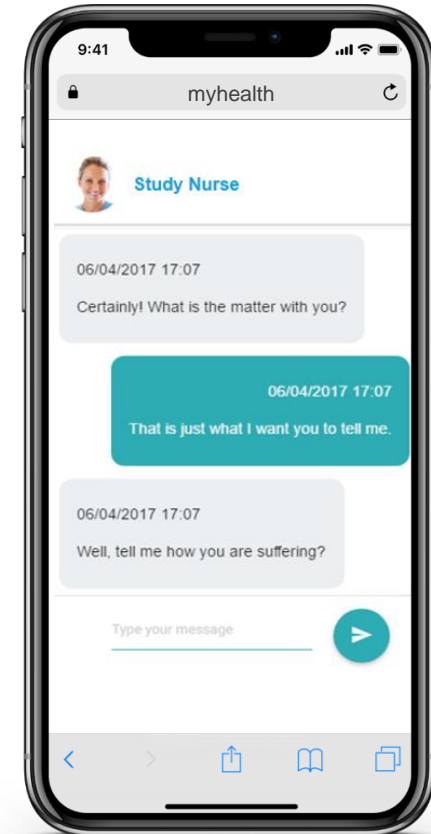
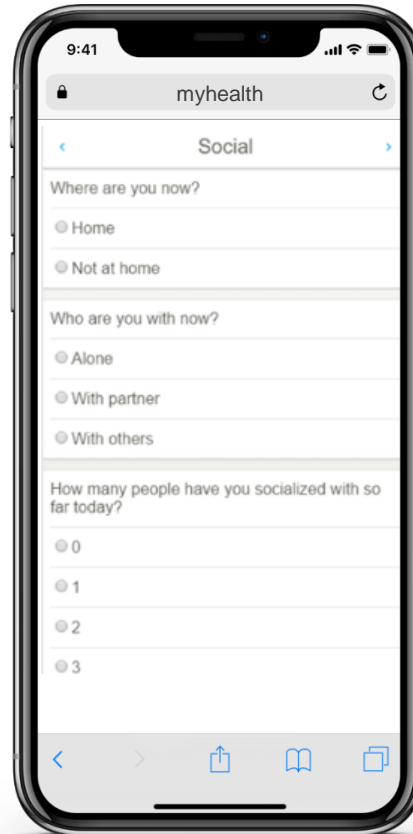


## Devices data

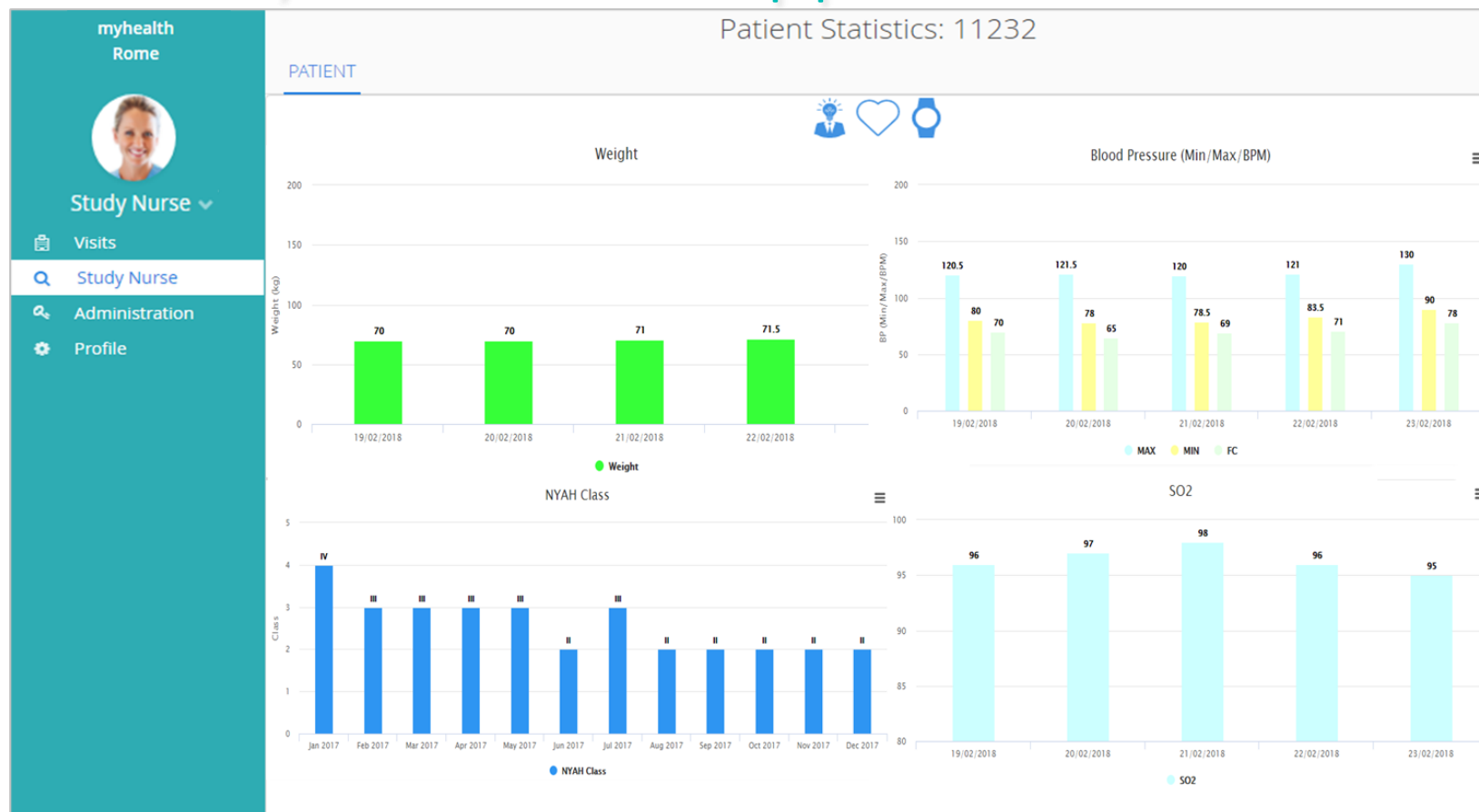
- Oxygen Saturation (SpO2), Breathing
- Diastolic and systolic blood pressure (mmHg)
- ECG and Detection of atrial fibrillation (AFib)
- Heart rate (BPM avg, max, min)
- Duration and quality of sleep, sleep apnea (n. of awakenings, deep and light sleep)
- Glycemia
- Physical Activity (walk, run, sleep, cycling), number of steps, Calories burned



# User Mobile App

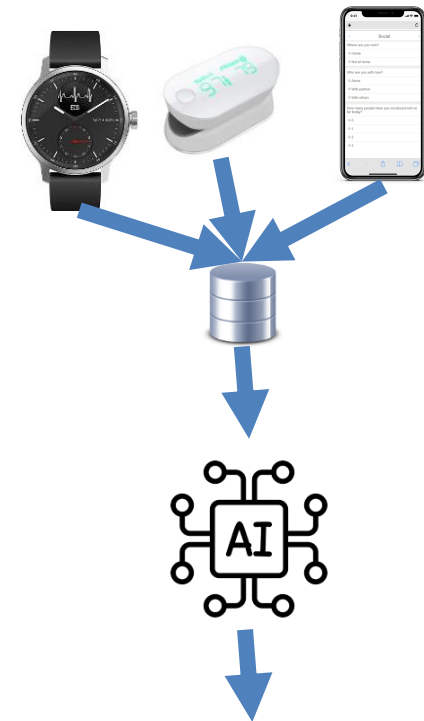


# Clinician/Nurse Web Application



## a) <sup>AI</sup> Personalized support for each patient

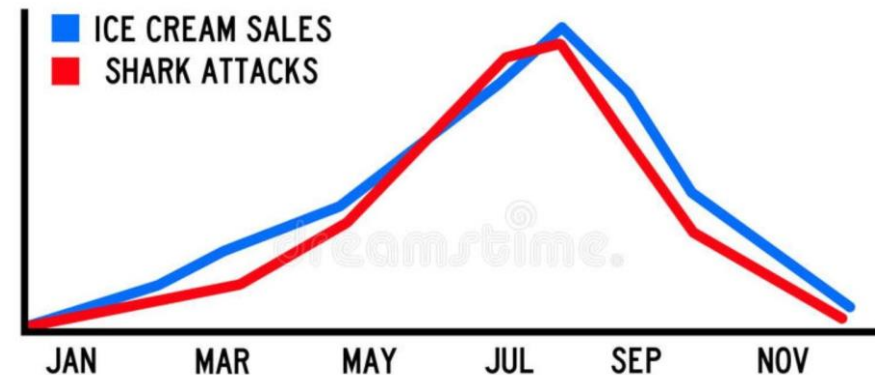
- **Input**: integrated data collected by wearable devices and mobile app's forms
- **Output**: personalized feedback for each user:
  - i. Personalized contents:
    - » Level of details based on the patient's profile
    - » Action Recommendation based on her status
  - ii. Personalized notifications:
    - » customized time, frequency





# AI is not the silver bullet

- Some tasks require human-in-the-loop to read the data correctly



## *b) Prioritization of the patients for clinical users*

- Input: integrated data about all the monitored patients
- Output: priority of patients and best channel of interaction
  - i. Priority for each clinical user: what patient's data to review so to take action

## GDPR & Ealth Data

- Personal data **concerning health** should include all data regarding to the health status of a subject such as:
  - **Information about the person** collected for provisioning health case services;
  - **Clinical** and **biometric data** derived from tests or examination of the patient, including biological samples;
  - Any information on e.g., disease, disability, medical history, clinical treatment, disease risks, etc.
- These kind of data are considered **sensitive data** by the General Data Protection Regulation (GDPR art. 51-52).



# GDPR and Secondary Use of Health Data

- Access to health data is crucial for medical research
  - Large volume of data is more likely to **provide compelling and robust evidence**;
  - The **integration** of more data sources helps the **improvement** of **disease diagnostic** techniques and treatments;
  - Exchange of data in the context of multi-national clinical trials ensures that **conclusions are valid** for different groups of people, avoid misleading conclusions;
- GDPR impose **limitations** on sharing and reuse of health data for research purpose;
- **Processing personal health data is lawful** only when:
  - The data subject given the **explicit consent** to the processing of personal data for one or more specified purposes;
  - **vital interests** of the data subject are protected;
  - Reasons of **public interest in the area of public health exist** (for instance ensuring high standards of quality and safety of health care and of medicinal products or medical devices).

- **Pseudonymization** means that personal data is processed in a manner that it can no longer be attributed to a specific data subject without the use of **additional information**;
- Personal data which have undergone pseudonymization is still **considered personal data**, meaning that are still protected by GDPR;
- **Pseudonymization ≠ Anonymization**
  - An information is anonymized when does not relate to an identified or identifiable person or regards to personal data rendered anonymous in such a manner that the **data subject is not or no longer identifiable**.
  - **GDPR does not concern anonymized information** which can be used for statistical research purposes.



- **Data masking:** hiding data with altered values. E.g. by using shuffling, encryption, character substitution. Data masking makes reverse engineering impossible.
- **Generalization:** removes some of the data to make it not identifiable. E.g. removes the house number in an address. The purpose is to eliminate identifiers while retaining data accuracy.
- **Shuffling and permutation:** rearrange the attribute values so they do not correspond with the original record.
- **Data perturbation:** applying techniques that round numbers and add a random noise. A small noise may lead to weak anonymization while a large one can reduce the utility of the dataset.
- **Synthetic data:** algorithmically manufactured information that has no connection to real events. The process involves creating statistical models based on patterns found in the original dataset. E.g. by using standard deviations, medians, linear regression or other statistical techniques.

# GDPR – Data Retention and Transparency

- Personal data should be kept (in a form which permits identification) for **no longer than is necessary for the purposes of the processing**
  - A time limit should be established for **erasure/ anonymization** of the data
- GDPR states that personal data should be processed **lawfully, fairly and in a transparent way**. Transparency means that any information regarding the processing of personal data must be:
  - Easily **accessible**;
  - Easy to **understand**;
  - **Clear** and written in **plain language**.





- Right of **access** to personal data which have been collected
  - For example, medical records such as examinations results, diagnoses, treatments, etc.
- Right to **rectification**
  - Wrong or missing data must be fixed without delays;
- Right to be **forgotten**
  - Personal data must be deleted, for example if data subject withdrawn her consent to data processing.



# How to protect privacy while doing Data Integration?

## • European Patient Identity Management



- prevent duplicate registration of patients
- avoid creating a transparent universal patient ID but
- provide distinct pseudonyms for patients in different contexts
- preserve the possibility for re-identification by a trusted third party
- keep a protected link between the different pseudonyms in the background
- which supports creating merged, datasets for secondary use

