GDSD 2025

# DATA-DRIVEN MEDICINE

## PROF. ANNE SCHWERK

- INTERNATIONAL UNIVERSITY OF APPLIED SCIENCES
- CHARITÉ / BERLIN INSTITUTE OF HEALTH

# AI-DRIVEN MEDICINE

iu | INTERNATIONALE HOCHSCHULE

# Overview: AI and Healthcare Data

# INCREASED RELEVANCE OF MEDICAL AI
## PUBLICATIONS

153.647

Total: 2,252,722

2428

1969

**Search terms**: AI OR artificial intelligence AND Medicine OR healthcare

2025

# ADOPTION OF AI PER INDUSTRY



Levels of AI maturity by industry, 2021 and 2024*

● 2021   ● 2024

| Industry | |
|---|---|
| Tech | |
| Automotive | |
| Life Science | |
| Retail | |
| Energy | |
| Communications & Media | |
| Insurance | |
| Travel | |
| Consumer Goods & Services | |
| Healthcare | |

30    40    50    60

Notes: * 2024 = estimated scores. Industries' AI maturity scores represent the arithmetic average of their respective Foundational and Differentiation index.
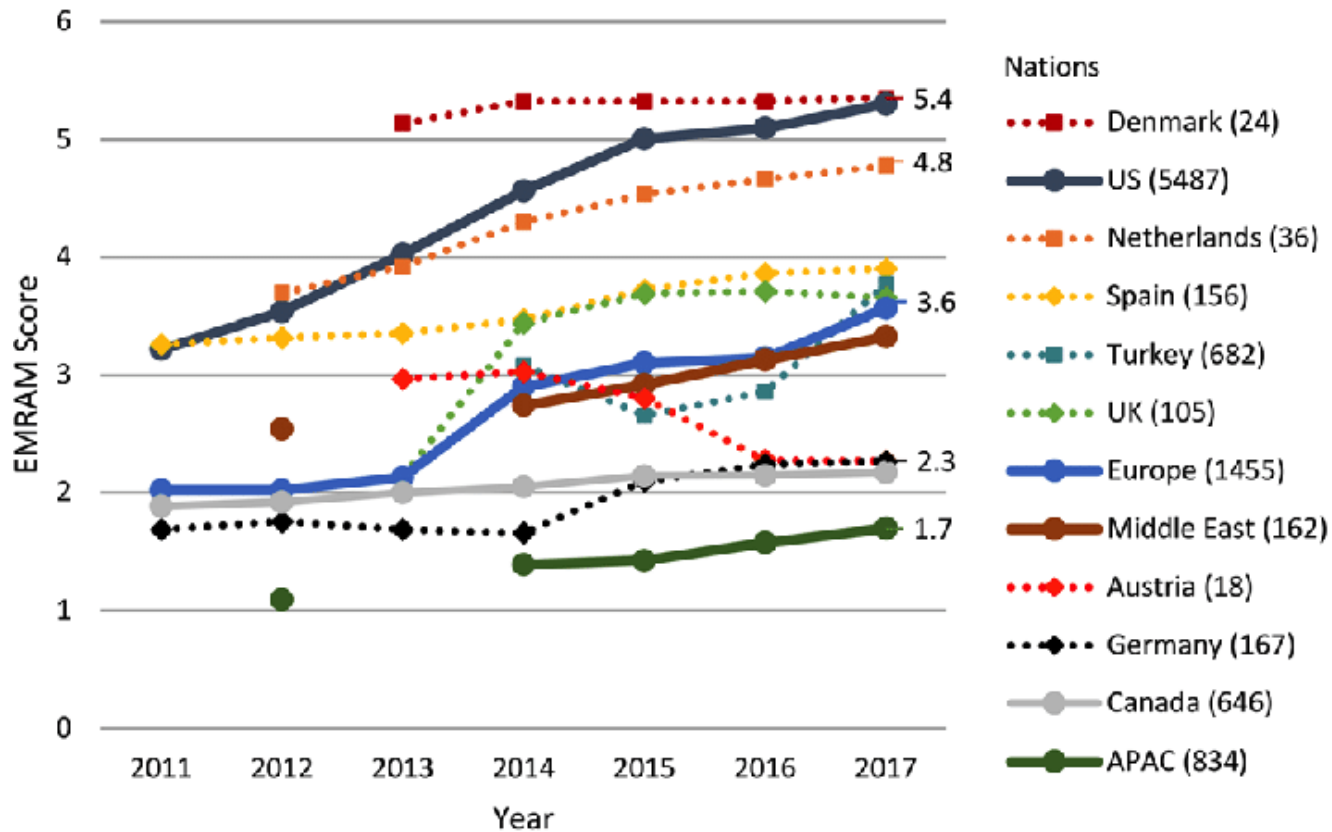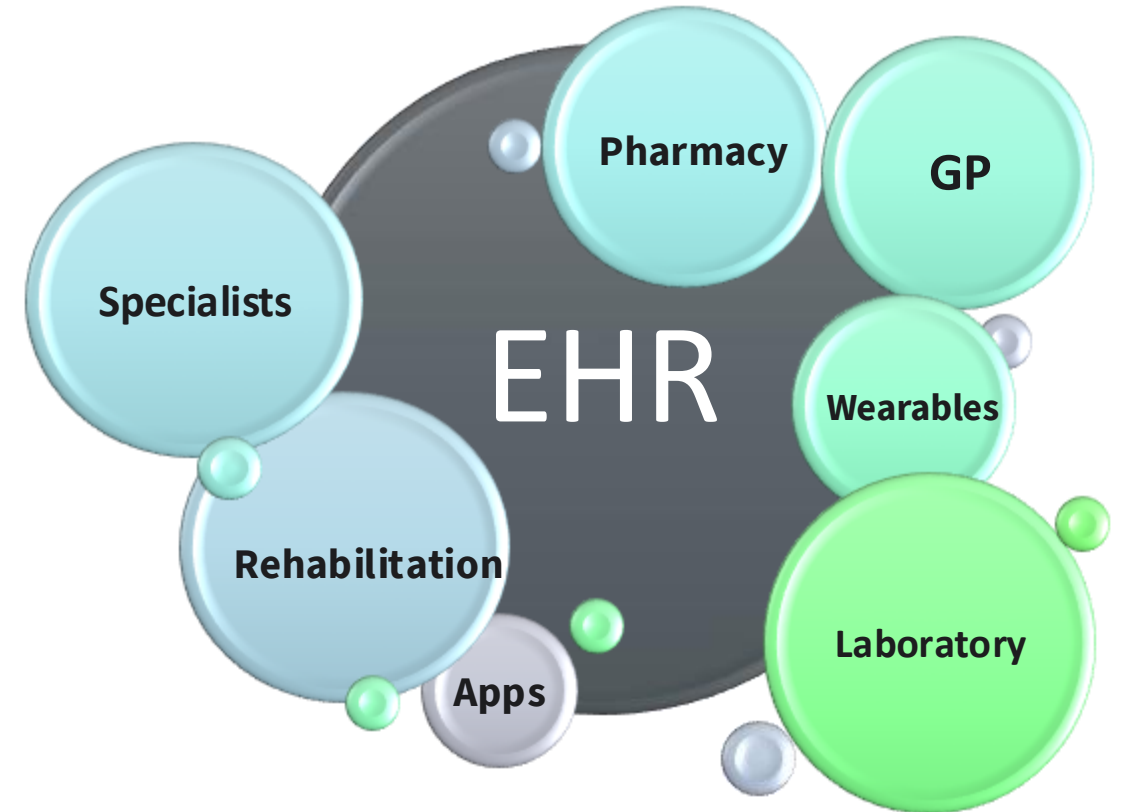
Source: Accenture Research

statista 𝄜

# LACK OF DIGITALISATION



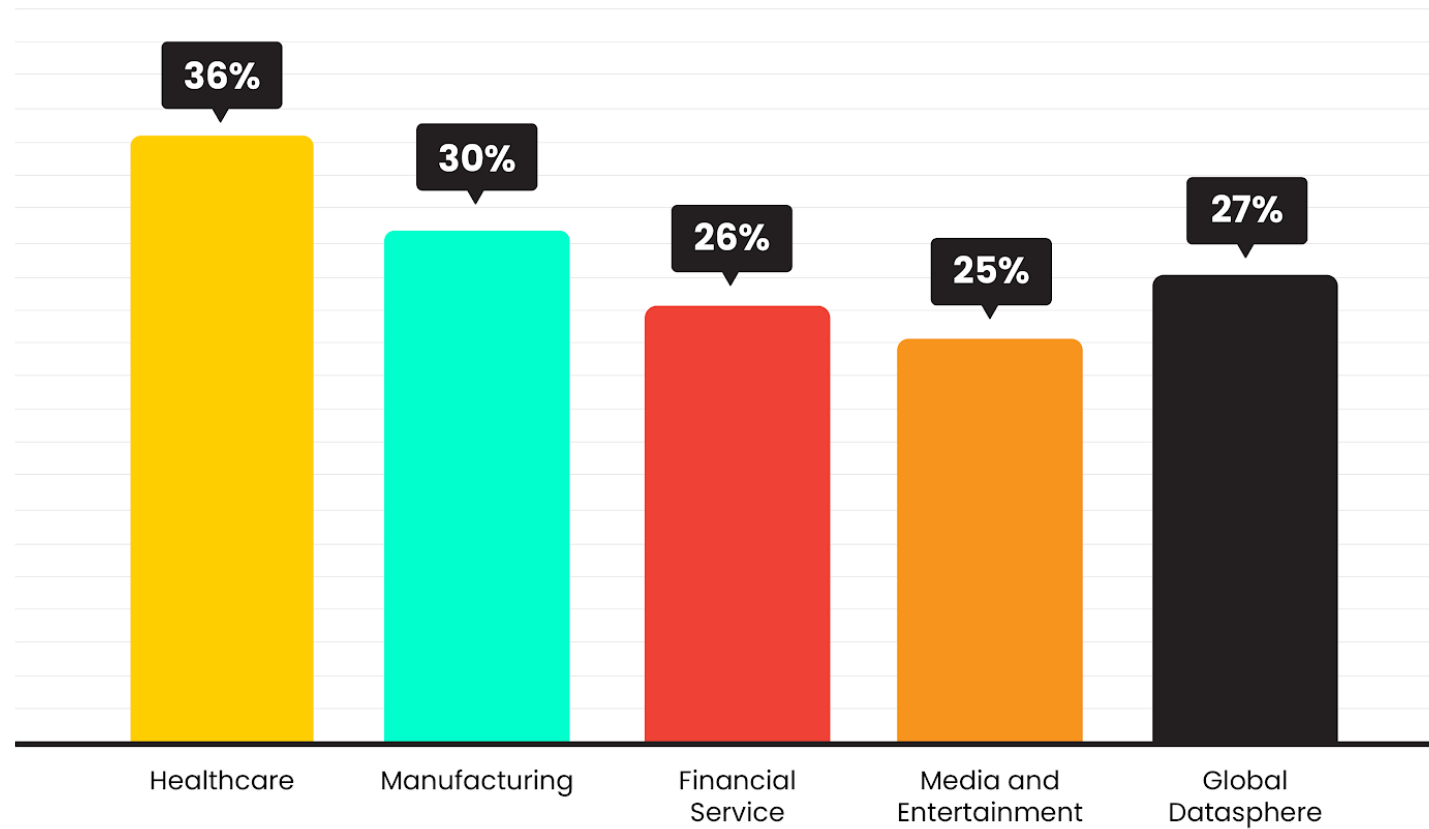| STAGE | HiMSS Analytics EMRAM EMR Adoption Model Cumulative Capabilities |
|---|---|
| 7 | Complete EMR; External HIE; Data Analytics, Governance, Disaster Recovery, Privacy and Security |
| 6 | Technology Enabled Medication, Blood Products, and Human Milk Administration; Risk Reporting; Full CDS |
| 5 | Physician documentation using structured templates; Intrusion/Device Protection |
| 4 | CPOE with CDS; Nursing and Allied Health Documentation; Basic Business Continuity |
| 3 | Nursing and Allied Health Documentation; eMAR; Role-Based Security |
| 2 | CDR; Internal Interoperability; Basic Security |
| 1 | Ancillaries - Laboratory, Pharmacy, and Radiology/Cardiology information systems; PACS; Digital non-DICOM image management |
| 0 | All three ancillaries not installed |

Bukowski  et al. (2020); https://www.nexus-marabu.de/nachricht/neue-himss-emram-kriterien

— Lack of data & inhibited sharing:

- Labeled data
- Benchmark data sets & open data
- Data silos
- Lack of cross-validation options
- Lack of standardization/ interoperability

— Data quality:

- Imbalanced data
- Missing data
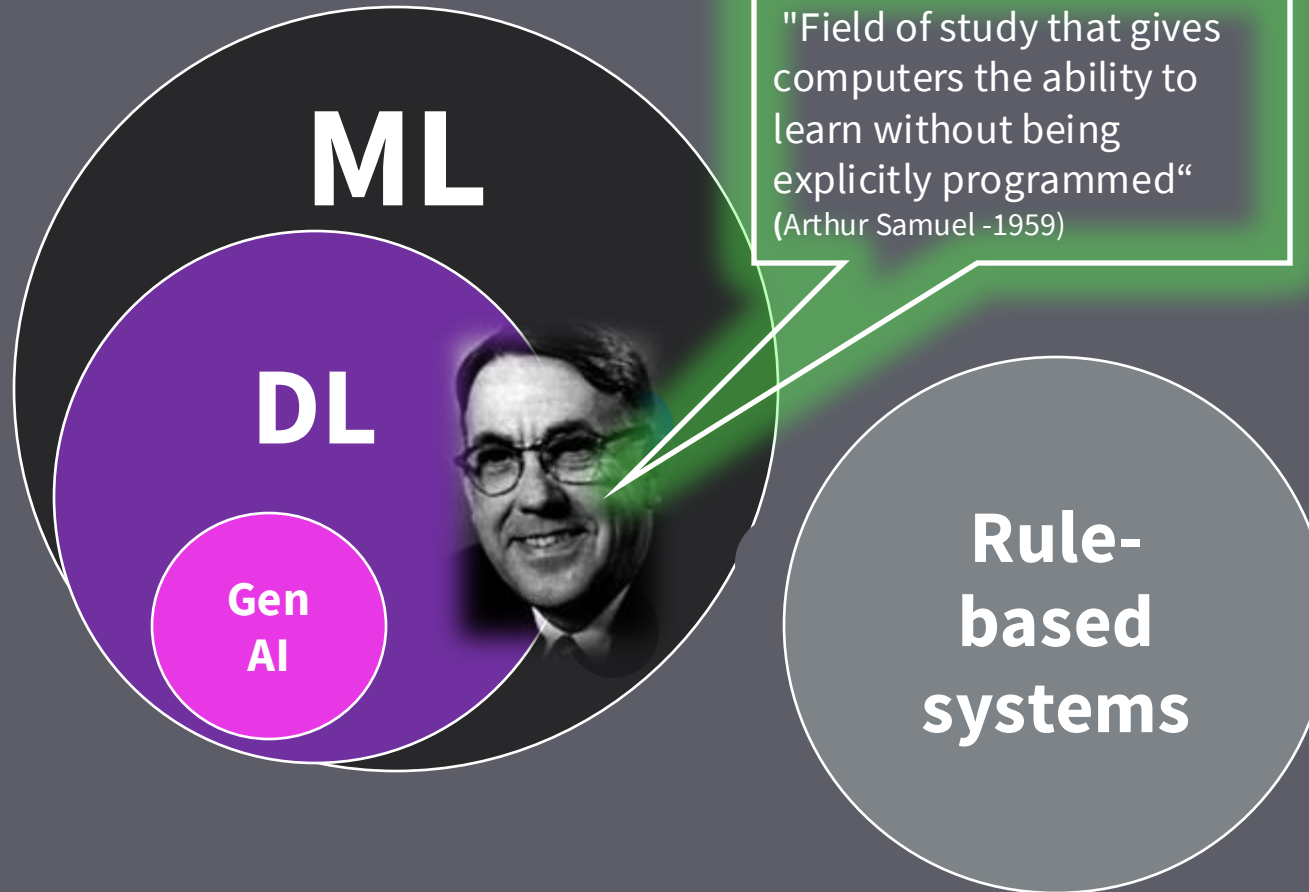- Incomplete data
- Standardization



Source of the image: Schwerk (2023)

# HEALTHCARE DATA INCREASE



## 2018-2025 Data
## Compound Annual Growth Rate (CAGR)

| | 36% | 30% | 26% | 25% | 27% |
|---|---|---|---|---|---|
| | Healthcare | Manufacturing | Financial Service | Media and Entertainment | Global Datasphere |

# THE END OF THEORY - DEDUCTION VERSUS INDUCTION



**Deduction**

1. Hypothesis

2. Experiment

3. Data collection

4. Data analysis

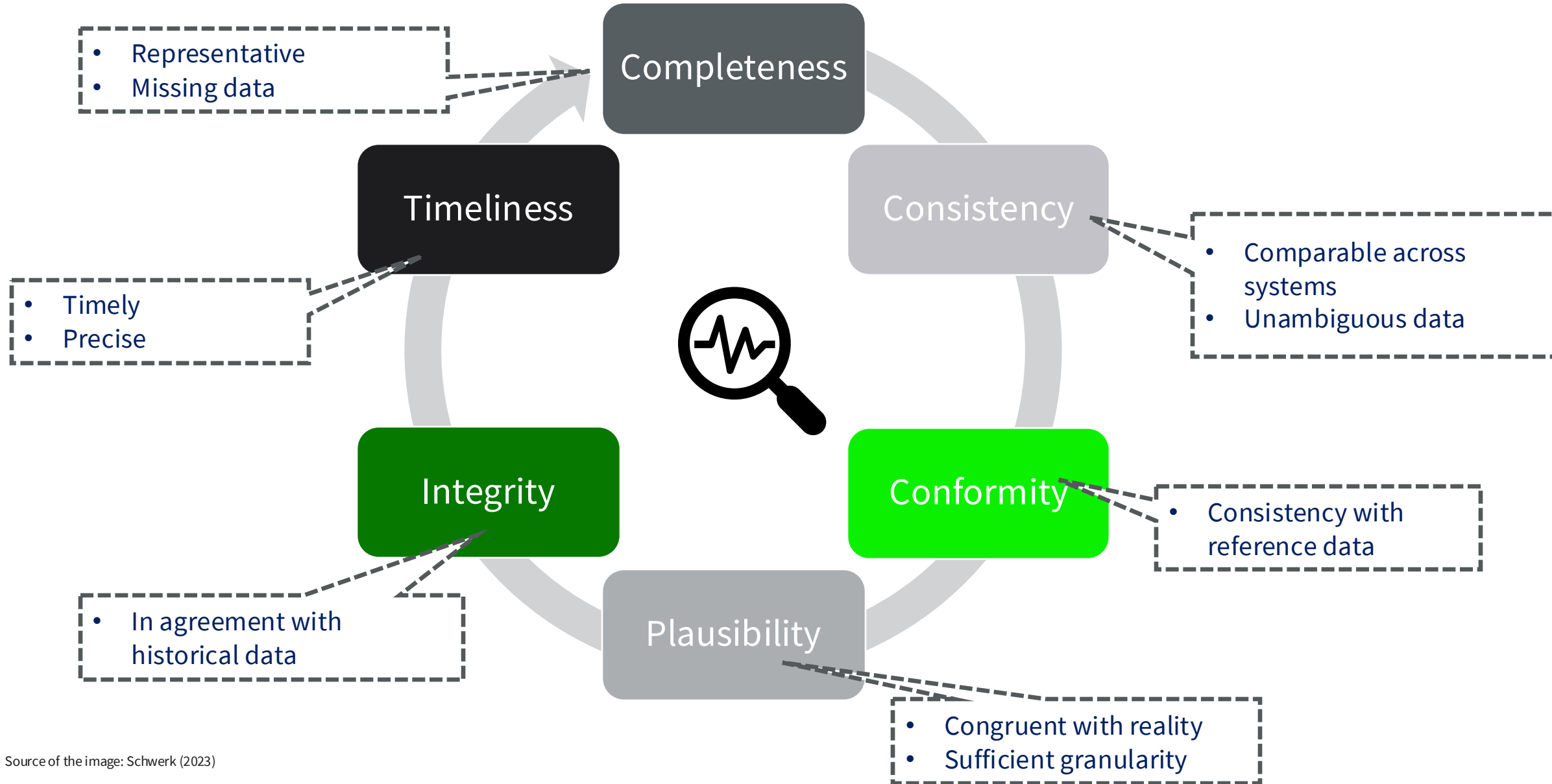5. Validation

**Induction**

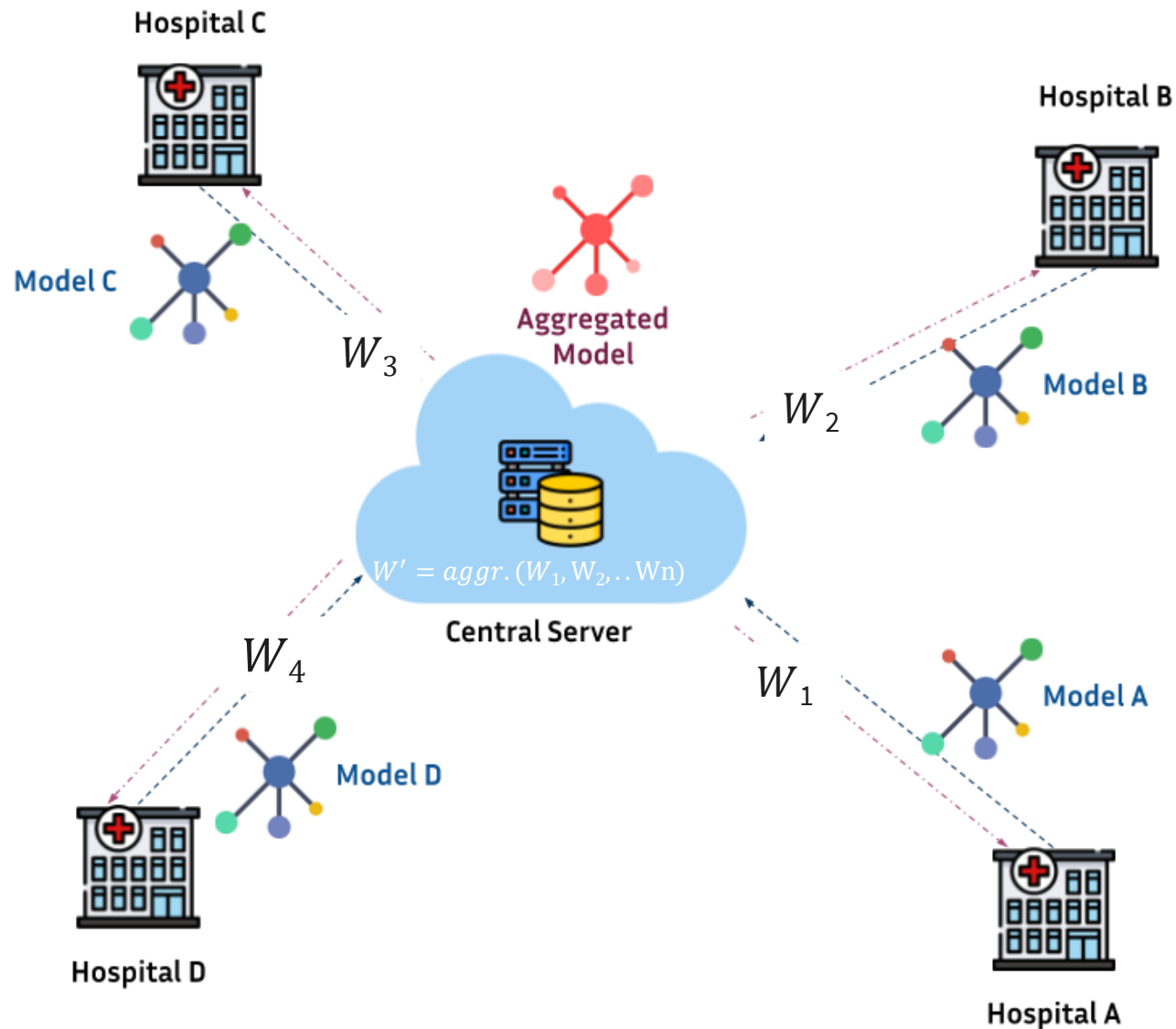1. Big Data Integration

2. Data mining

3. Pattern recognition

4. Hypothesis generation

5. Validation

# IMPORTANCE OF DATA: DATA QUALITY



- Representative
- Missing data

Completeness

Consistency

- Comparable across systems
- Unambiguous data

Timeliness

- Timely
- Precise

Integrity

Conformity

- Consistency with reference data

- In agreement with historical data

Plausibility

- Congruent with reality
- Sufficient granularity

Hospital C

Model C

$W_3$

Aggregated Model

Hospital B

Model B

$W_2$

$W' = aggr.(W_1, W_2, ..Wn)$

Central Server

$W_4$

Model D

Model A
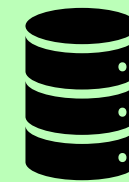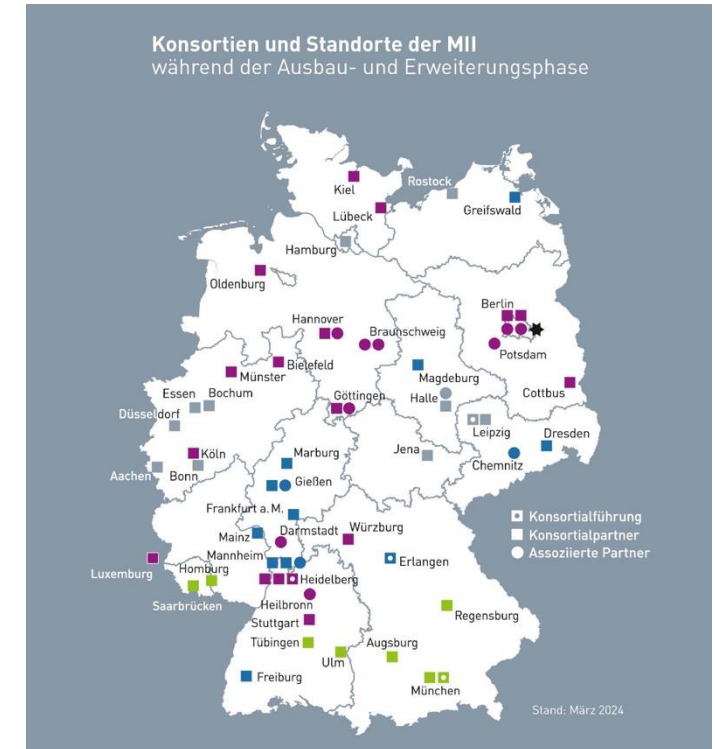
$W_1$

Hospital D

Hospital A

Source: https://fedbiomed.org/

- **Data protection:** Data remains at the original location - only model parameters are shared

- **Decentralization:** Hospitals train local models and send encrypted parameters to central coordinators

- **Iterative process:** The coordinator aggregates local models into a global model and shares it. The process is repeated until the model converges
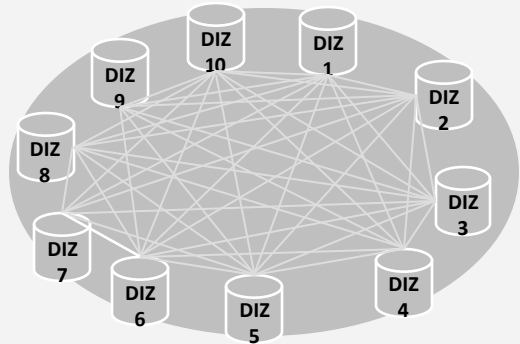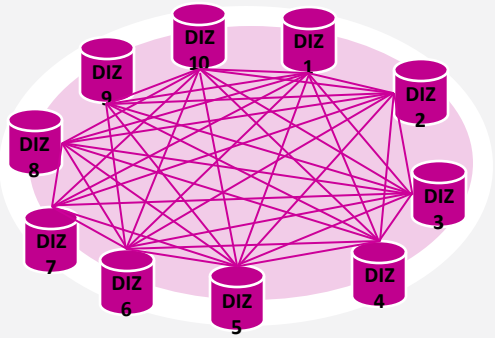
- Initiated by the Medical Informatics Initiative (MII) launched in 2016 by the BMBF
- A federated network for centralized and decentralized data access
- Core data set from primary IT systems of the universities →local data integration centers
- Standardization to MII format
- Depending on the basis of use:
  - Distributed evaluations
  - Central evaluations (with broad consent)



Konsortien und Standorte der MII
während der Ausbau- und Erweiterungsphase

Stand: März 2024

- 15 M EHRs
- 160 M Diagnoses
- 1.5 Billion Lab values
- → **2-3 M uncoded rare diseases**

# Examples

# RARE DISEASES: DIAGNOSTIC ODYSSEY



**Rare Disease Dilemma**

Most rare diseases are genetic — 72%

Most rare diseases affect children — 75%

Many are misdiagnosed — 70%

On average it takes 8 years to diagnose — 8 years

# DATA AVAILABILITY

- FDPG Data
- Orpha / ICD Codes

**L1:**

FDPG Data

Diagnosis codes

- Addition through clinics
- **Goal**: Europe-wide analysis

**L2:**

1) Diagnosis codes
2) Minimal core data

**L3:**
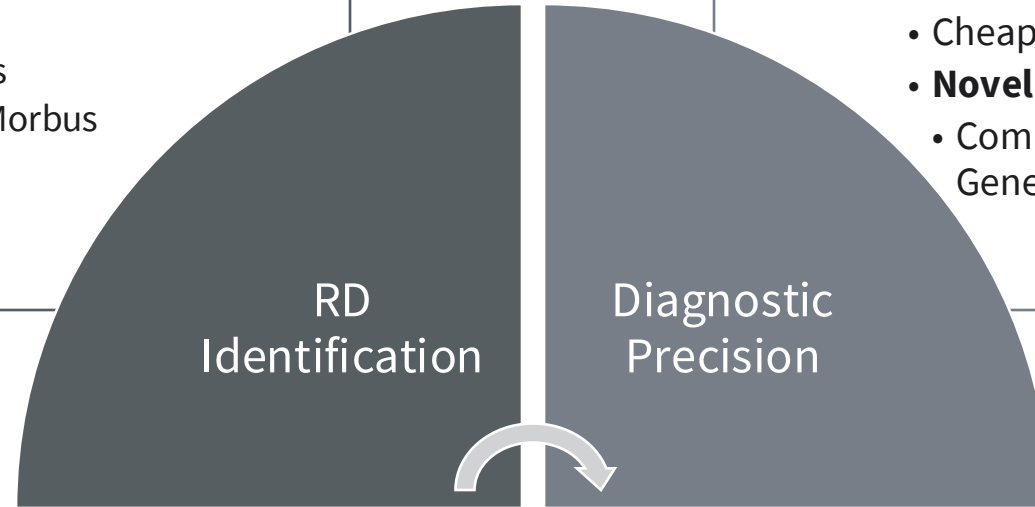
1) Diagnosis codes
2) Minimal core data
3) Stand. details

- Special data:
  - Progress data
  - Intervention data
  - PROMs
  - Omics
- European Registers

# ML FOR RARE DISEASE (RD) DIAGNOSTICS

iu INTERNATIONALE HOCHSCHULE

Morbus Osler

- Undected RDs
- EHR-based Phenotypes
- **Example**: Morbus Osler

- Combined markers
- Cheaper Diagn.
- **Novel Diagnoses**:
- Combined: Genetic (PRS) Epigenetic

Hypophosphatasia

**RD Identification**

**Diagnostic Precision**

**Subgroup Identification & therapy**

**Novel Biomarker**

Parkinson's disease

- Precision medicine
- Treatment

Monitoring
- Interventions
- Progression
- Diagnostics

Gaucher disease

- Hereditary disease: ALPL gene mutation

- 400 + disease-causing ALPL variants

- Symptoms:

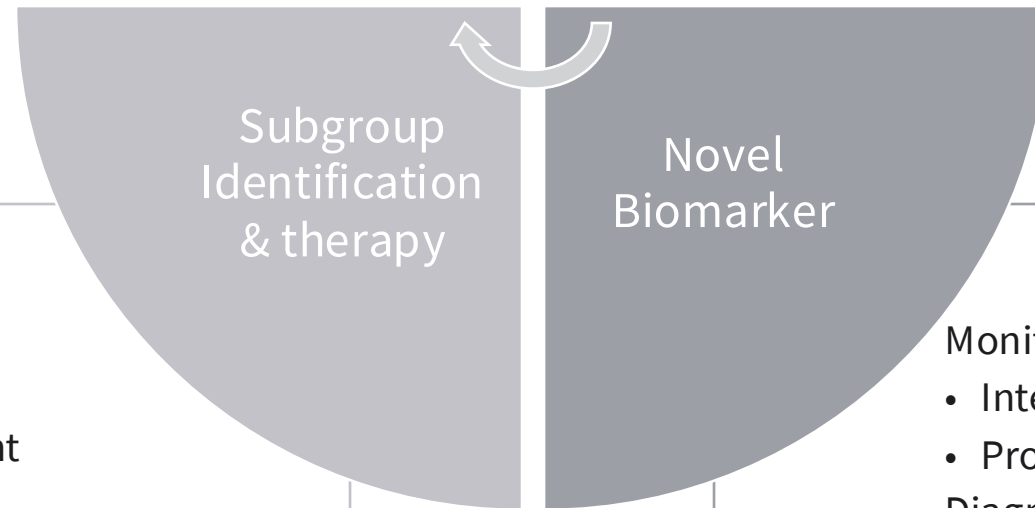  - **Severe**: Bone demineralization, respiratory failure, seizures

  - **Mild**: Tooth loss, periodontal disease

- Diagnosis: Ø 5.7 years delayed

  - Frequent misdiagnoses

  - ALP value + symptoms + genetics

  - Specific orpha code

- Incorrect treatment: bisphosphonates → Bone damage

→ **Early diagnosis is crucial**

**L1 Analysis**: Use ALP (+PLP) biomarkers + phenotypes to determine more specific biomarker thresholds
→ identify new patients and improve diagnostic precision



— Khan et al. 2024

INTERNATIONALE
HOCHSCHULE

- Hereditary disease: : ENG, ACVRL1, or SMAD4 gene mutation

- 600 + disease-causing variants

- Symptoms:

  - **Severe:** Arteriovenous malformations → Bleeding, strokes, cardiac stress

  - **Mild:** Nosebleeds (epistaxis), telangiectasias on skin and mucous membranes

- Diagnosis: Ø 26 years delayed

  - Frequent misdiagnoses

  - Curacao criteria + genetics

- Incorrect treatment: anticoagulants → Bleeding

→ Early diagnosis is crucial



HEALTHY

INTERNAL ORGO...
BLEEDING

INTEBEL
BLEEDING

**L1 Analysis**: Using the ICD-10 code for identification & extended phenotyping
→ RBS on Curacao Criteria → New identifications

— Pierucci et al. 2012

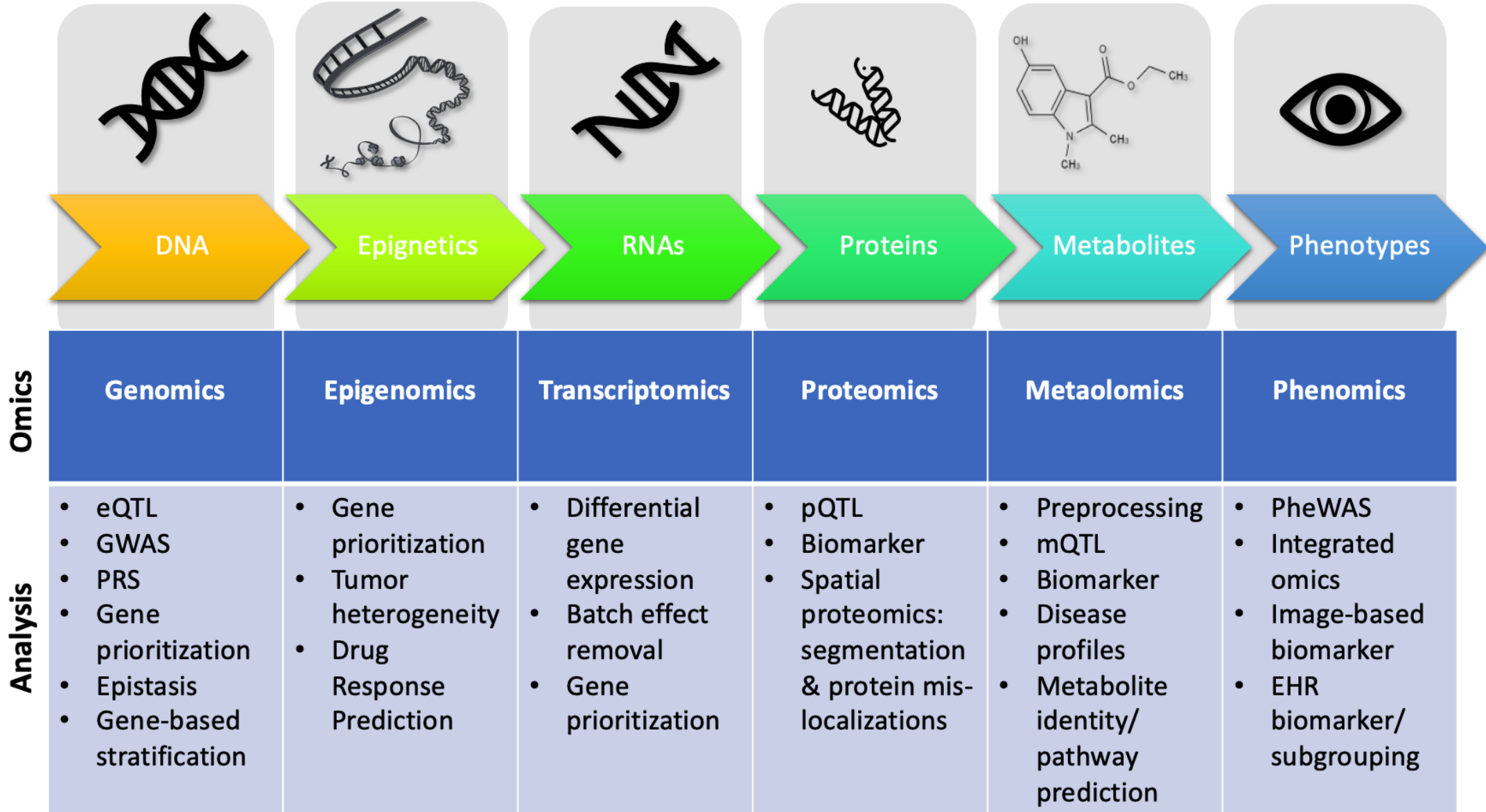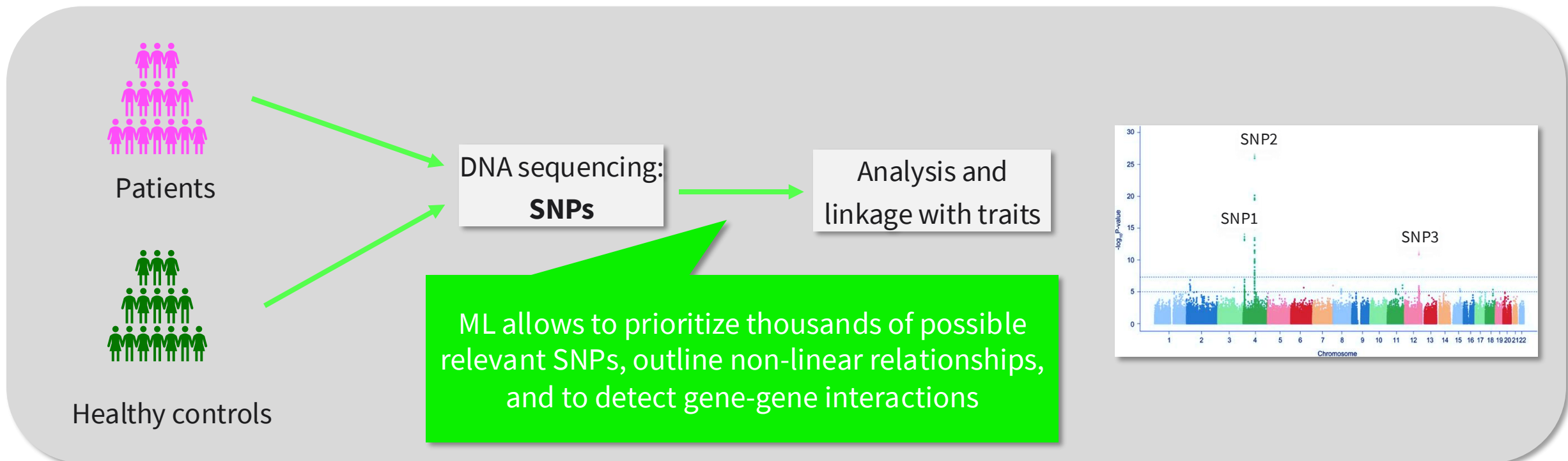The complete set of genetic data of a cell, including variations in DNA sequence or DNA structure.

The complete set of RNA molecules transcribed from the DNA of an organism or a cell population, including messenger RNA (mRNA), and non-coding RNA

The complete set of physical and biochemical data sources, including images, blood tests, and data from EHRs

The complete set of epigenetic changes, incl. changes in chromatin structure, DNA methylation, histone protein modifications, non-coding RNA

The complete set of the small-molecule metabolites present in a cell, tissue, or organism

The complete set of proteins expressed by a cell, tissue, or organism, e.g. amount, identity, function, and interaction.

Gen-omics

Transcriptomics

Epigen-omics

Prote-omics

Metabolomics

Phen-omics

Omics Data

Source of the image: Schwerk (2023)

# OMICS DATA AND ANALYSES

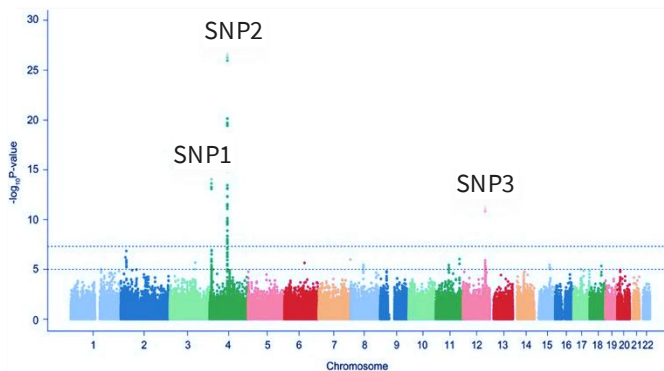| | DNA | Epignetics | RNAs | Proteins | Metabolites | Phenotypes |
|---|---|---|---|---|---|---|
| **Omics** | **Genomics** | **Epigenomics** | **Transcriptomics** | **Proteomics** | **Metaolomics** | **Phenomics** |
| **Analysis** | • eQTL<br>• GWAS<br>• PRS<br>• Gene prioritization<br>• Epistasis<br>• Gene-based stratification | • Gene prioritization<br>• Tumor heterogeneity<br>• Drug Response Prediction | • Differential gene expression<br>• Batch effect removal<br>• Gene prioritization | • pQTL<br>• Biomarker<br>• Spatial proteomics: segmentation & protein mis-localizations | • Preprocessing<br>• mQTL<br>• Biomarker<br>• Disease profiles<br>• Metabolite identity/ pathway prediction | • PheWAS<br>• Integrated omics<br>• Image-based biomarker<br>• EHR biomarker/ subgrouping |

Source of the image: Schwerk (2023)

**GWAS:**

- Analyzes the genomes of large groups of people with and without a disease / trait

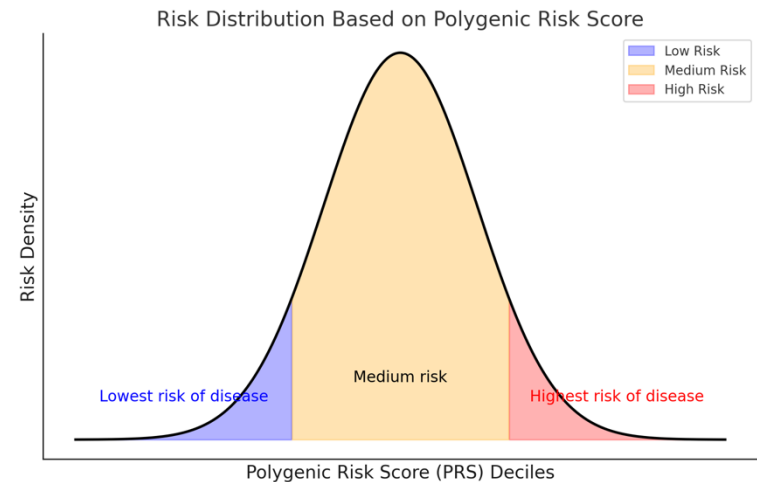- Identifyies genomic variants (SNPs) - by assessing linkage disequilibrium



Patients

DNA sequencing:
**SNPs**

Analysis and
linkage with traits

ML allows to prioritize thousands of possible relevant SNPs, outline non-linear relationships, and to detect gene-gene interactions

Healthy controls

— Source of the image : Schwerk (2023) and adapted from Matsuo et al. (2016;)

- PRS sum weighted effect sizes of risk variants from GWAS to estimate disease susceptibility
- Linear approach: does not account for complex genetic interactions (e.g., epistasis)
- **ML helps to:**
  - Account for non-linear interactions
  - Combine multiple existing PRS



GWAS

Top genetic variants
e.g. SNPs
(largest effect size)

Risk scoring

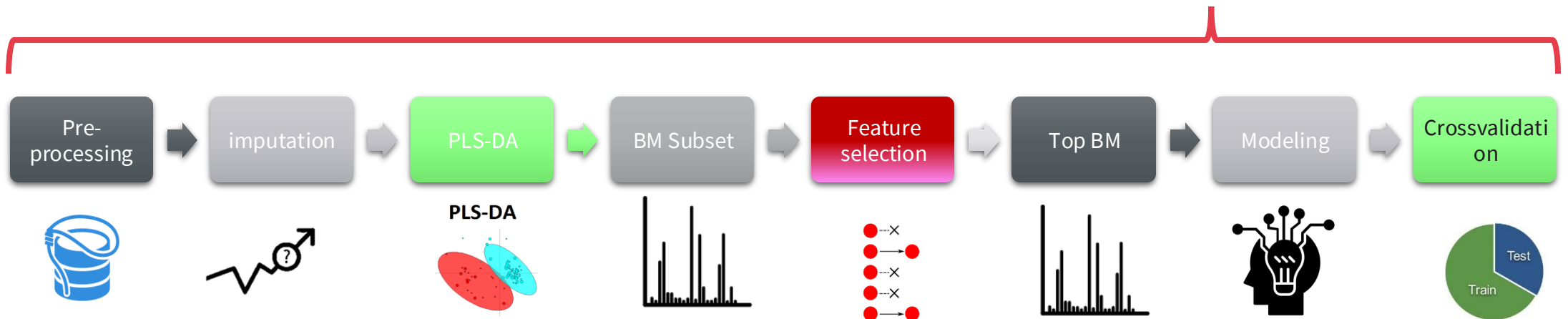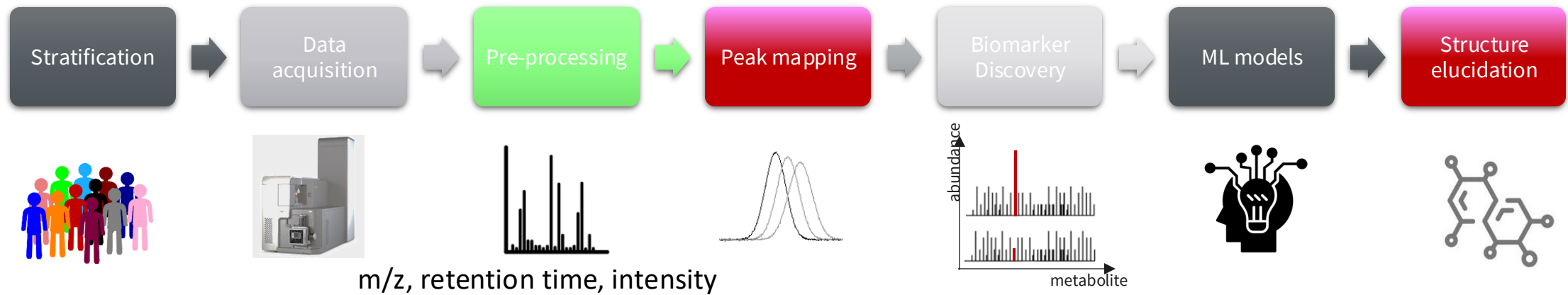# COST PER GENOME AND PRECISION MEDICINE PUBLICATIONS



Source of the image: NIH

# UNTARGETED METABOLOMICS FOR FINDING NOVEL BIOMARKERS



m/z, retention time, intensity

- **m/z (Mass-to-Charge Ratio):** the mass of a molecule divided by its charge
- **Retention Time:** the duration a metabolite takes to pass through the chromatography column before being detected
- **Intensity/ abundance:** concentration of metabolite in sample
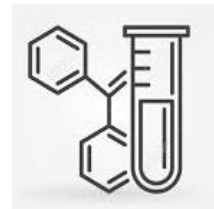
# EXPLAINING METABOLITES



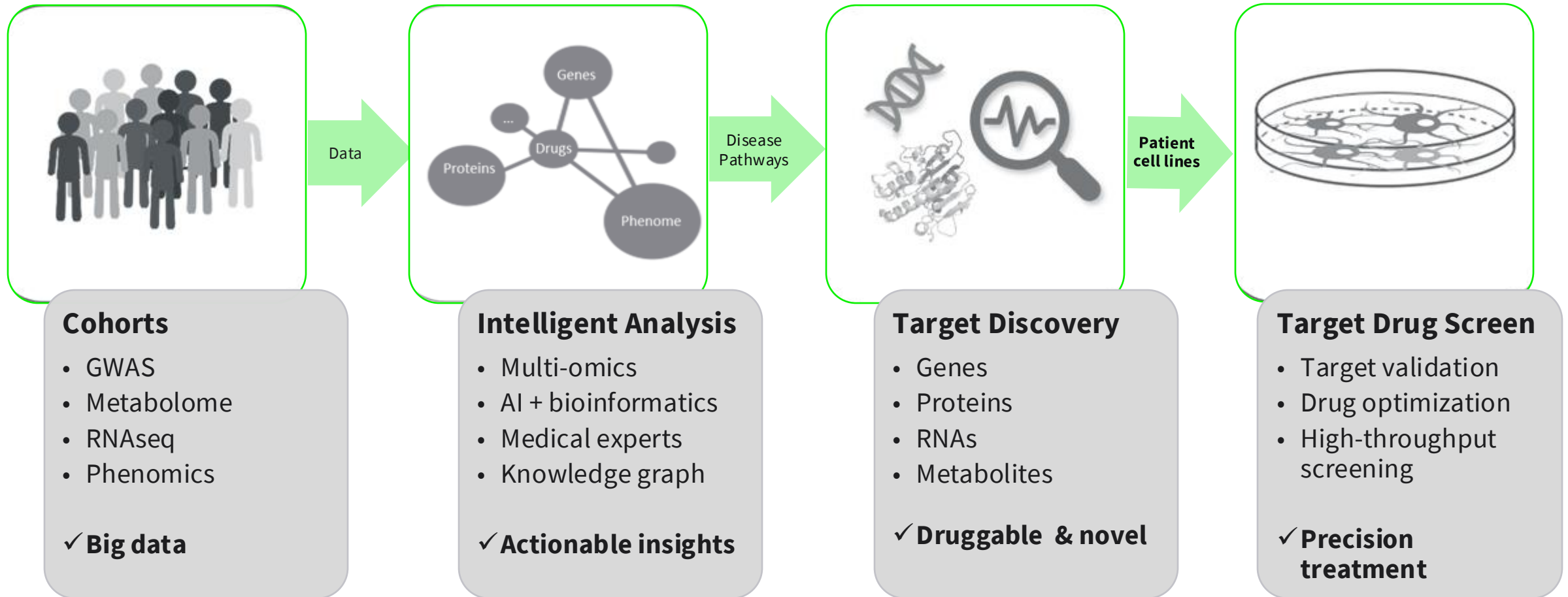Database Matching → Structure Elucidation → Chemical Synthesis → Biological Validation

MS FINDER
CFM-ID
CSI: FingerID
MetFrag

**Database search
In silico fragmentation**

# SCREEN4CARE

— **Genetic newborn screening:**
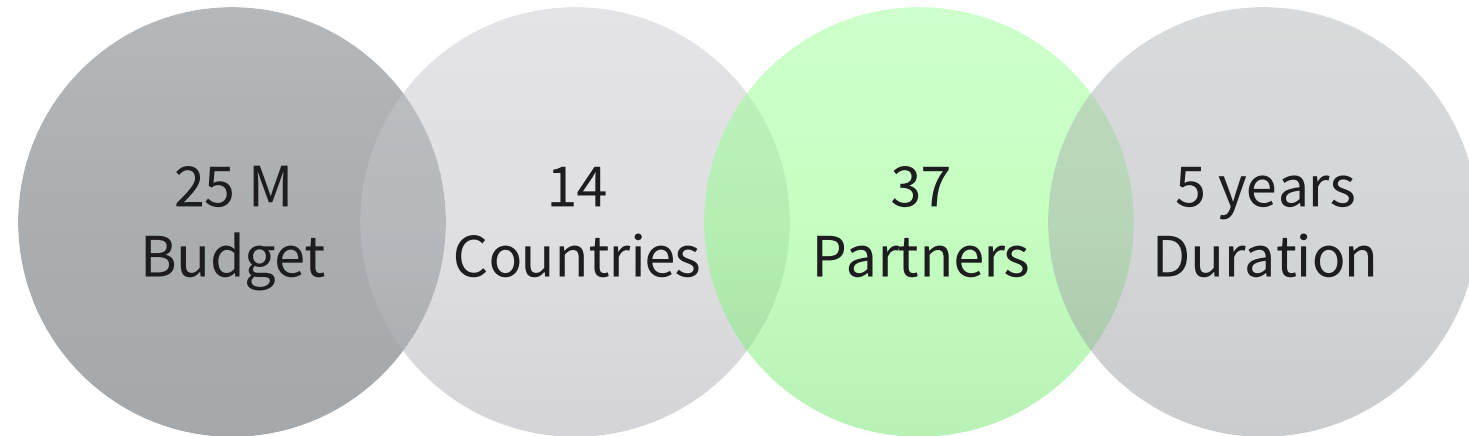
— Early diagnosis of genetic RDs

— **Digital tools:**

— Symptom checker

— EHR algorithms

— **Infrastructure**:

— Federated ML

SHORTENING THE PATH TO RARE DISEASE DIAGNOSIS BY USING NEWBORN GENETIC SCREENING AND DIGITAL TECHNOLOGIES

25 M Budget   14 Countries   37 Partners   5 years Duration

**Improved the accuracy and speed of diagnosis**

https://www.screen4care.eu/

# Take-Home-Message & Discussion

— **AI as a medical game changer:**

—Access to digital data

—High quality data

—**Expected AI revolution:**

—Fastest growing healthcare data source

—Text data and LLMs

# THANK YOU

Prof. Dr. Anne Schwerk

anne.schwerk@iu.org