

---

# SYNTHETIC DATA FOR SUPPORTING PRIVACY PRESERVING DATA SHARING

---



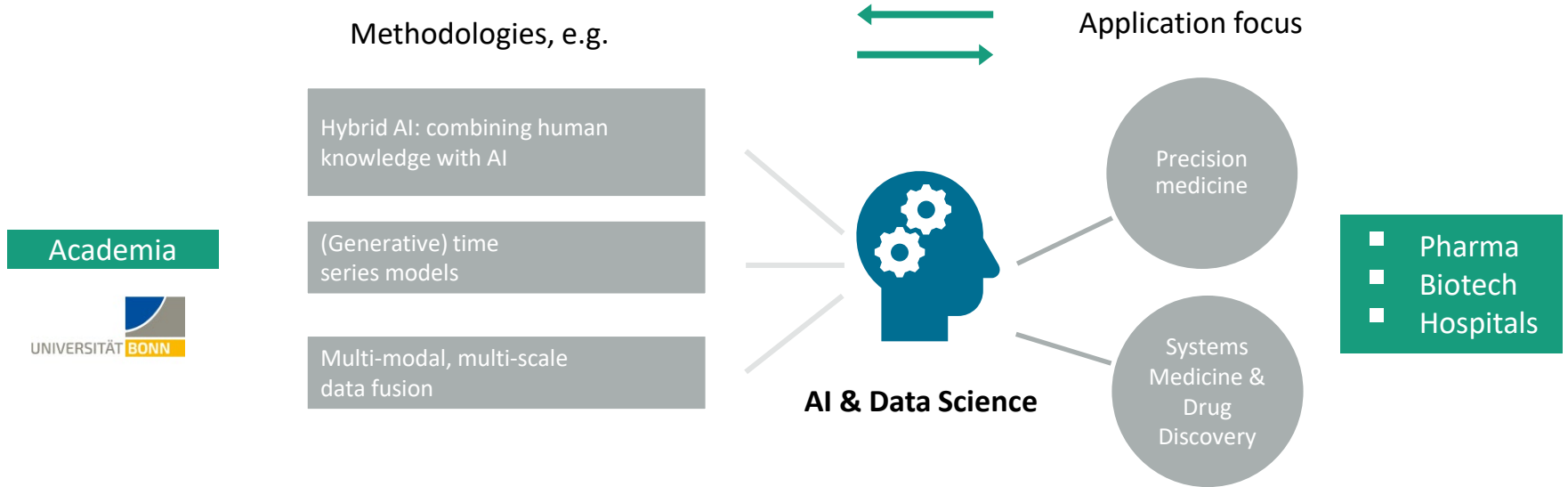
**Prof. Dr. Holger Fröhlich**

Head of AI & Data Science Group, Deputy Head of Department of Bioinformatics  
Fraunhofer Institute for Algorithms and Scientific Computing (SCAI)



# AI & Data Science Group @Fraunhofer SCAI

## Mission: Bringing Better Treatments to the Right Patients

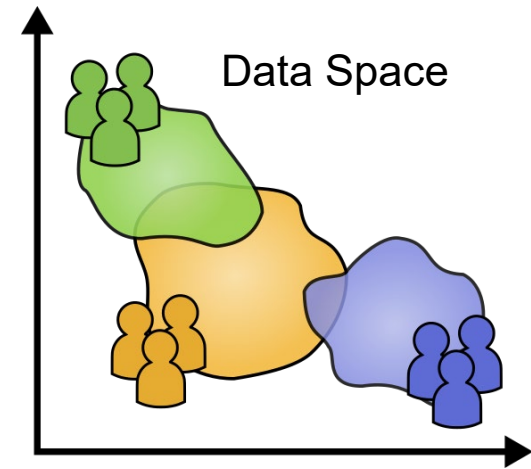


### Our value chain:



## Access to Data as a Key Challenge

- Access to patient-level data is highly restricted by EU and German law
- Data within a data holding organization is
  - Biased / not necessarily representative of the entire disease population
  - Often too small to develop generalizable ML models
- Is there a way to facilitate data sharing in a legally compliant manner?



# AI for Simulation of Synthetic Patient Data



Real World Clinical Study



AI-technology

VAMBN, MultiNODEs

A mathematical representation of the real world cohort based on an AI model



A synthetic, „virtual“ cohort

## Risk assessment for re-identification

Controlled risk algorithms such as differential privacy



## Quality assurance:

Good fit between observations in the real world and representation in the virtual cohort



## Potential use cases

### Simulation

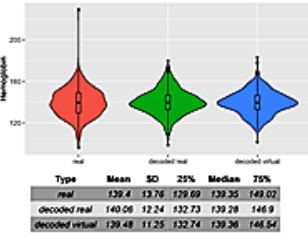
- synthetic data
- "what, if" scenarios
- „Patients-like me“

### Exploration

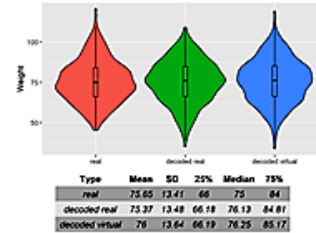
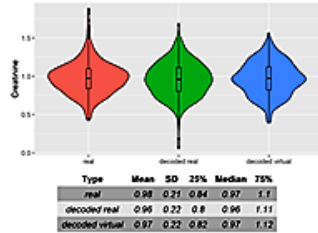
- inclusion / exclusion criteria
- extrapolation
- Statistical power

# Synthetic Data are Realistic While Respecting ( $\epsilon$ , $\delta$ )-Differential Privacy

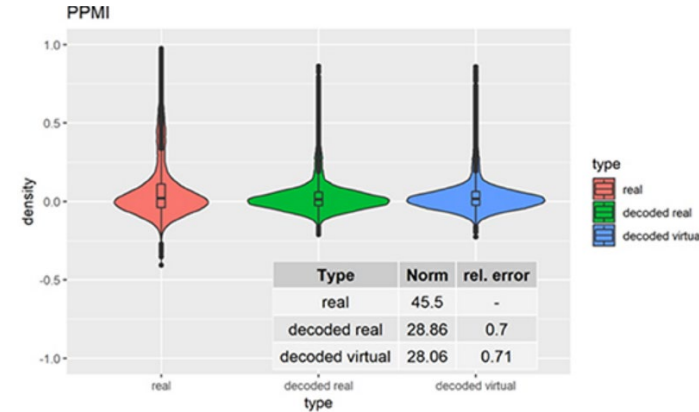
a) SP513



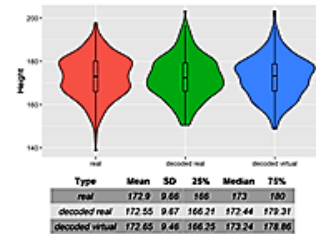
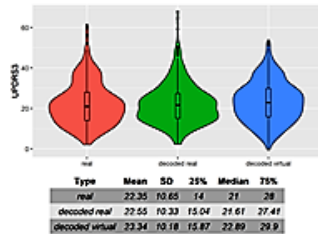
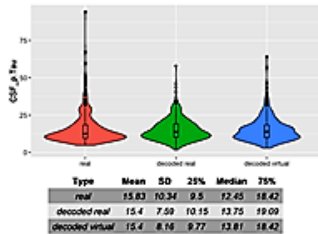
Marginal distributions



Pearson correlations

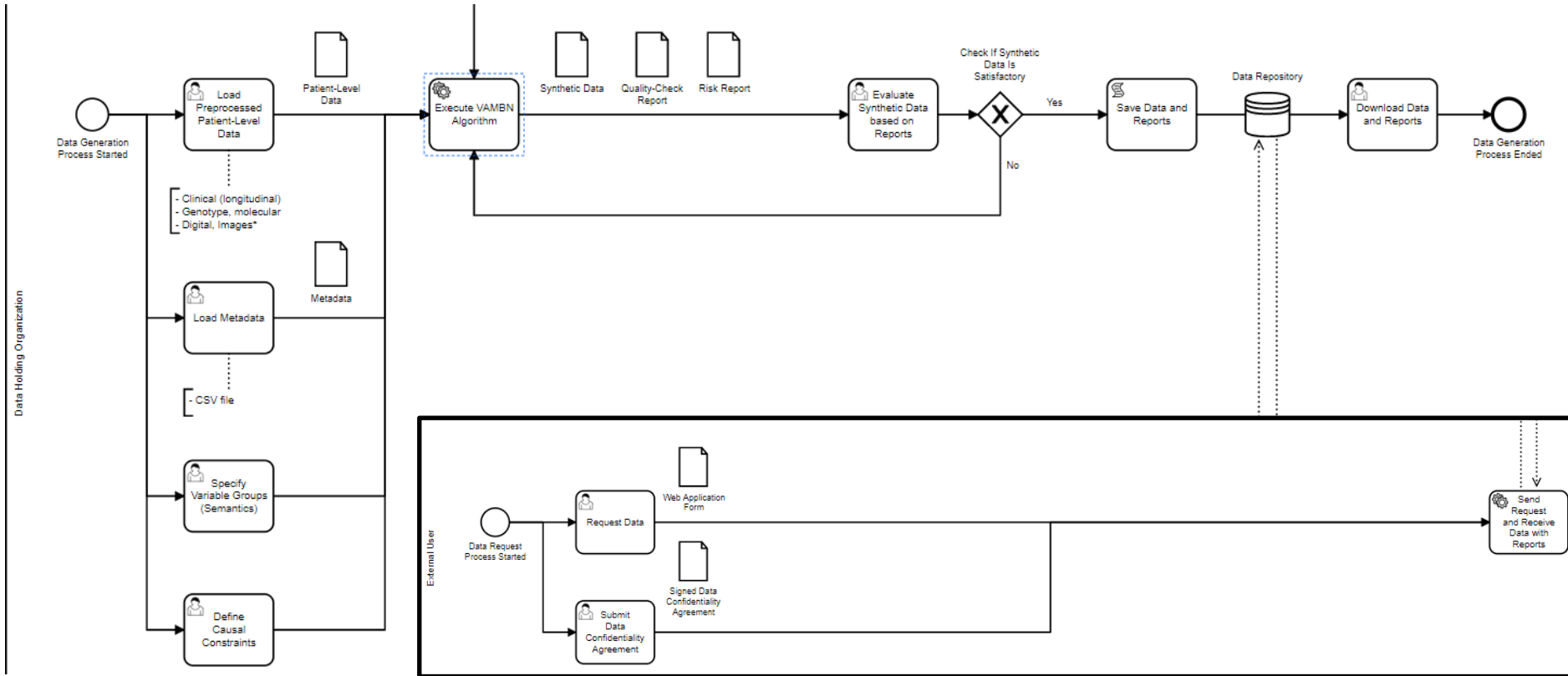


b) PPMI



Guarantee on the probability to compromise an individual's privacy

# Envisioned Use of Synthetic Data Generator in NFDI4Health



## Conclusion

Access to pseudonymized patient-level data is one of **THE** obstacles for Biomedical Data Science, specifically in Germany

- Comparison to US: market places for data, e.g. IBM

Data silos hinder progress in medical research

- Data biases
- Data within each organization often too small

Sufficiently realistic synthetic data might be a way to facilitate data sharing and thus help data driven biomedical research

