

Use of external data to accelerate evidence generation

Digital twins and external controls

Prof. Dr. Holger Fröhlich

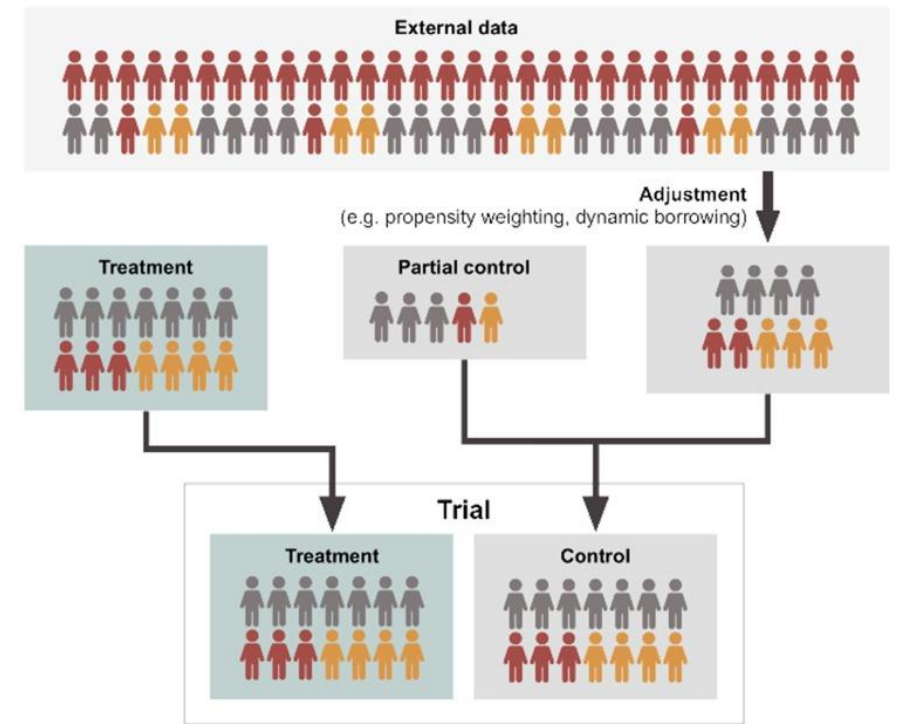
Head of AI & Data Science Group, Deputy Head of Department of Bioinformatics
Fraunhofer Institute for Algorithms and Scientific Computing (SCAI)

What are external controls and digital twins?

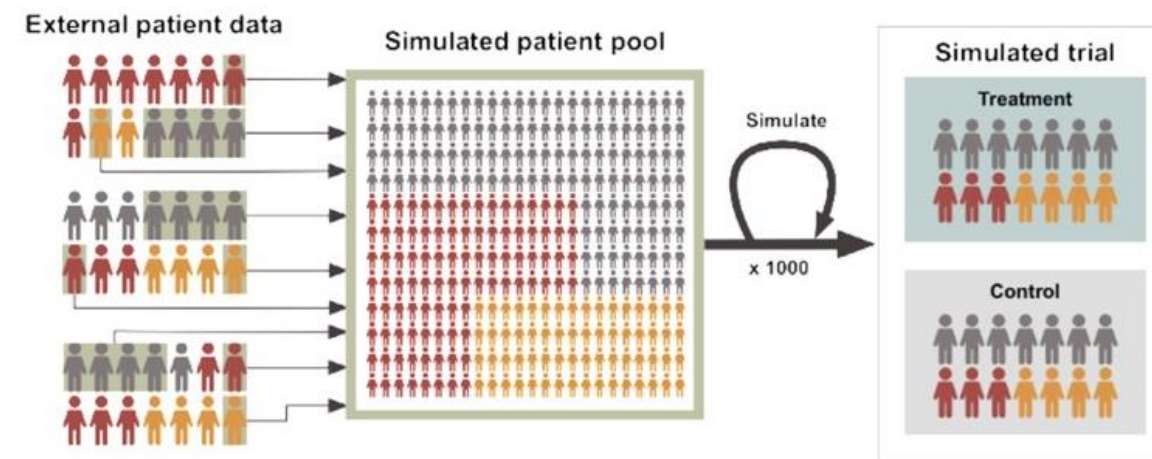
- Real-world data can be used to construct **external controls** for clinical trials

„A **digital twin** is a digital model of an intended or actual real-world physical product, system, or process (a *physical twin*) that serves as the effectively indistinguishable digital counterpart of it for practical purposes, such as simulation, integration, testing, monitoring, and maintenance.“
(Wikipedia)

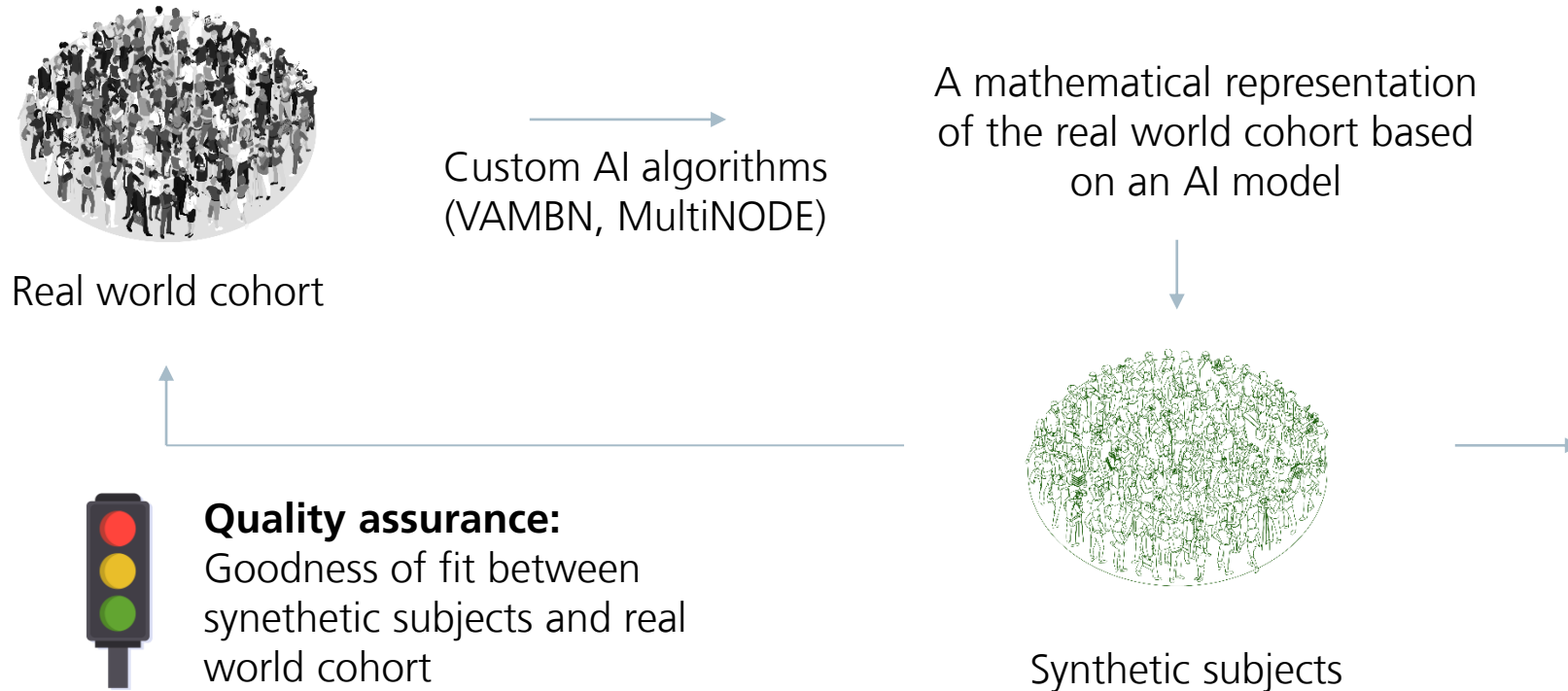
- Key concept is **Simulation**
- How to simulate a patient?**



Digital twin simulation



AI for Simulation of Synthetic Patient Trajectories



Data sharing

Facilitate sharing of sensitive data across organizations

Trial design

Simulation

- synthetic control arm
- "what, if" scenarios

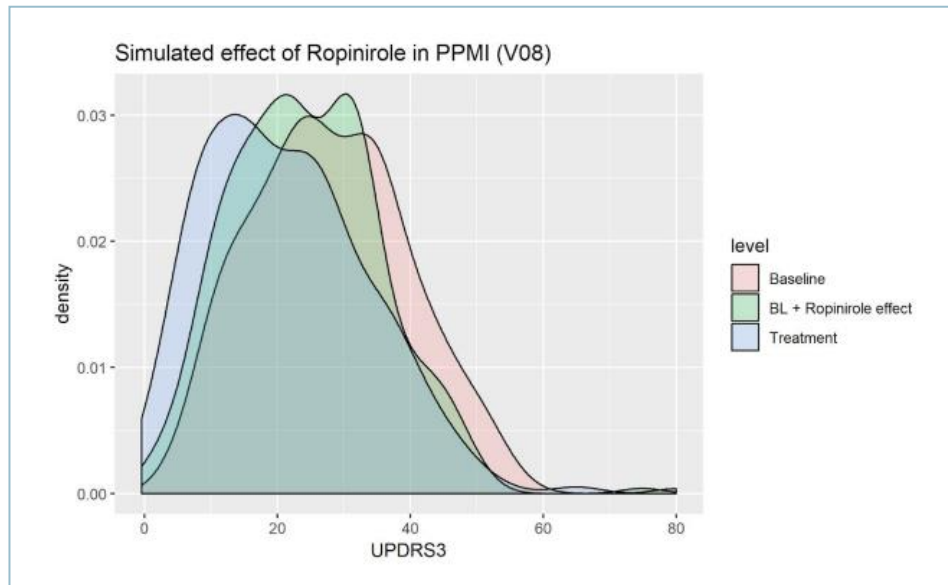
Exploration & Extrapolation

- inclusion / exclusion criteria
- Interpolation / extrapolation
- Statistical power

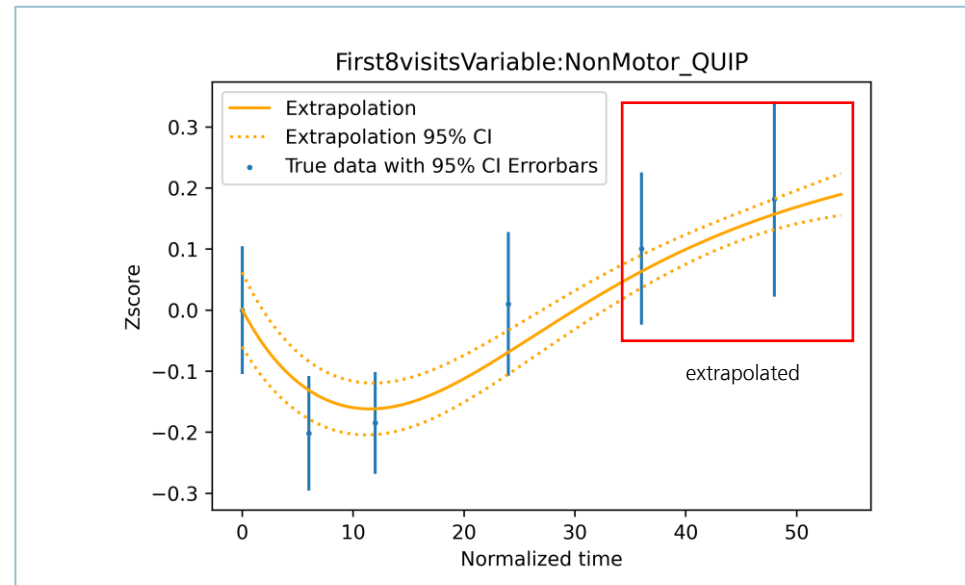


Utility of AI Generated Synthetic Data

- Trained AI models allow
 - Simulating counterfactual scenarios, e.g. treatment with a given drug
 - Statistical power calculations
 - Interpolation / extrapolation between visits, also on continuous time scale

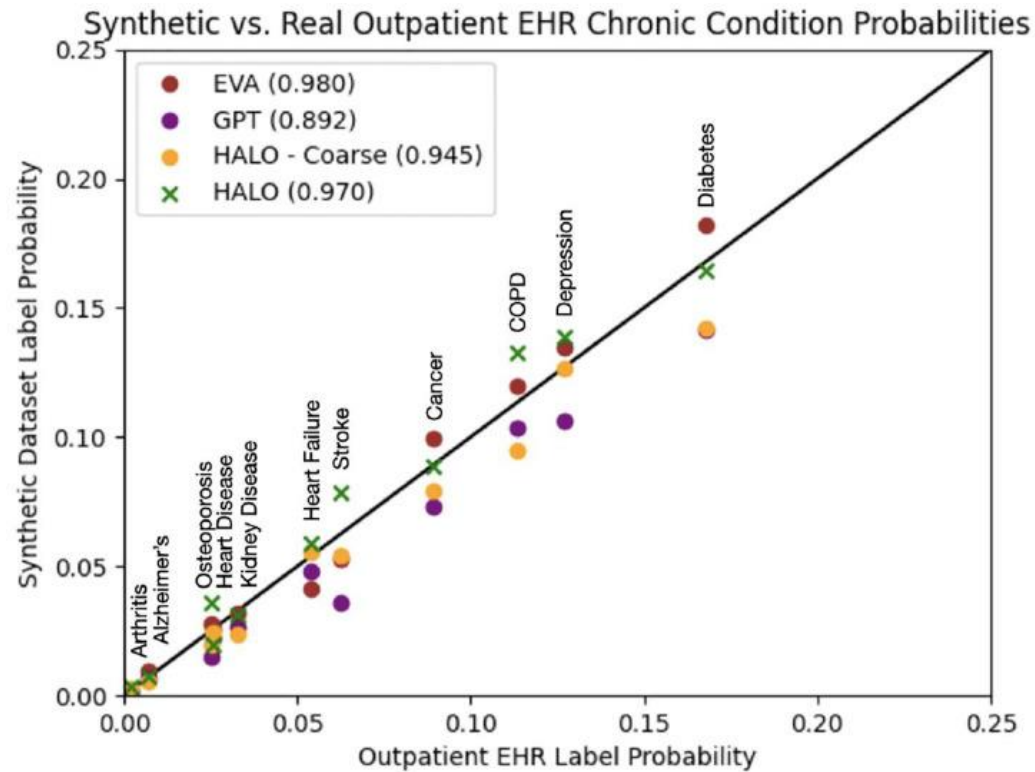


Gootjes-Dreesbach, ... & Fröhlich, *Frontiers in Big Data: Medicine and Public Health*, 2020



Wendland, ... & Fröhlich, *npj digital medicine*, 2022

AI Generated Synthetic EHRs



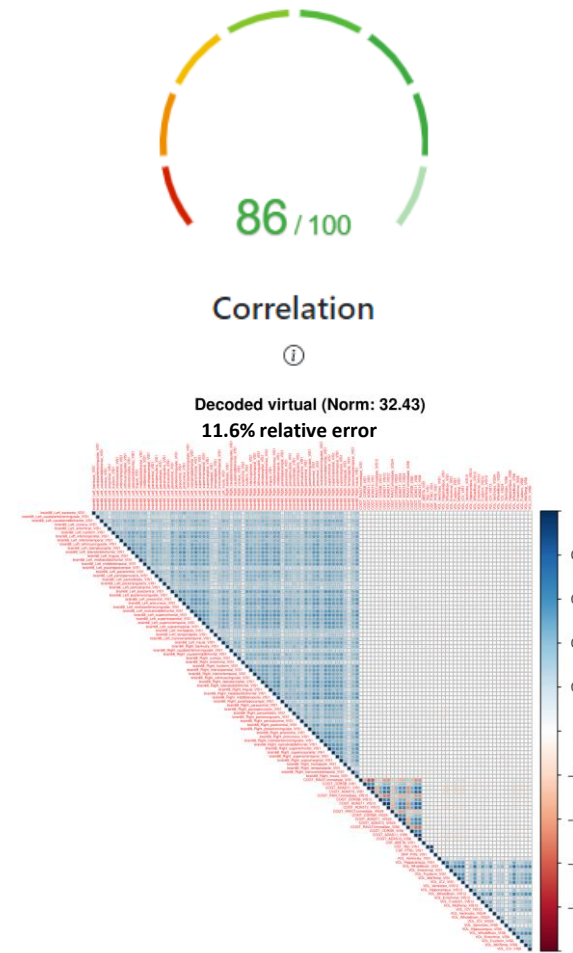
Synthesize *Extremely High-dimensional Longitudinal Electronic Health Records* via Hierarchical Autoregressive Language Model

Brandon Theodorou¹, Cao Xiao², Jimeng Sun^{1*}
University of Illinois Urbana-Champaign.¹
Relativity Inc.²

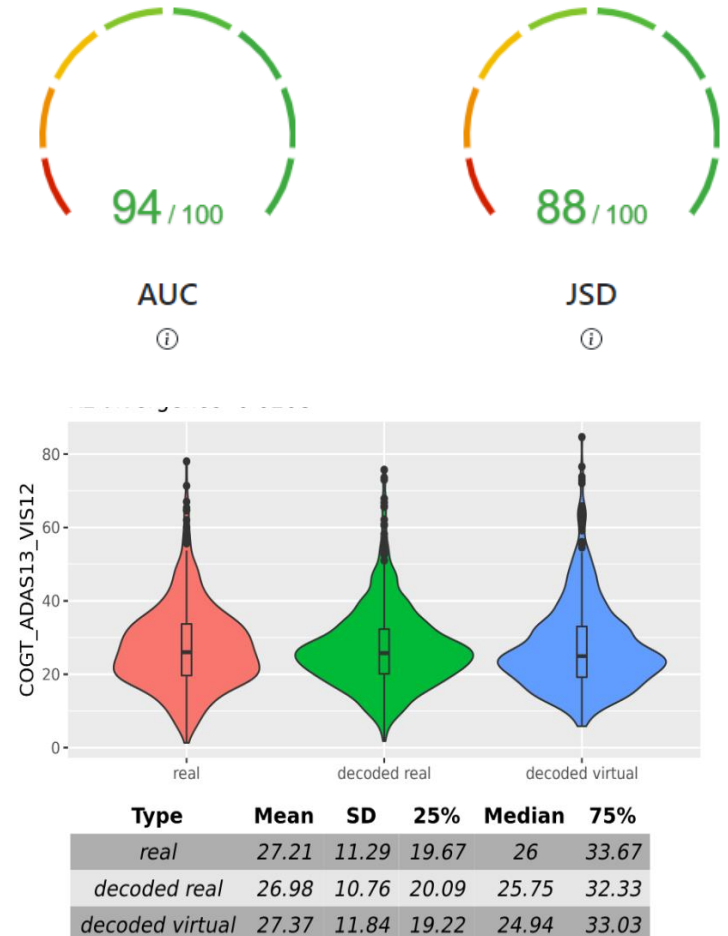
* To whom correspondence should be addressed: jimeng@illinois.edu

How realistic is synthetic data?

- Syndat: A web-based tool for quality assessment
 - Statistical distributions of individual variables should be similar to real data
 - ML classifier (Random Forest) should misclassify majority of synthetic as real patients
 - Correlation structure of synthetic data should be similar to real data
- Users can upload their own synthetic data



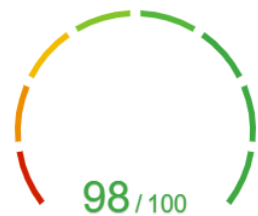
Output evaluation results:



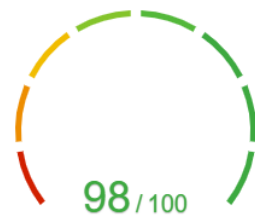
How privacy preserving is synthetic data?

- **Singling out risk:** Can we single out a real individual based on a rare combination of attributes in the synthetic data?
- **Linkability risk:** Can we link together two or more records (either in the same dataset or in different ones) belonging to the same individual or group of individuals?
- **Inference risk:** Can an attacker confidently guess (infer) the value of an unknown attribute of a real data record?

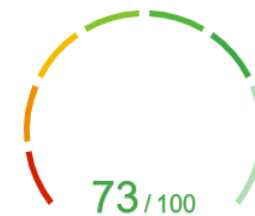
Giomi et al., A Unified Framework for Quantifying Privacy Risk in Synthetic Data, arXiv, 2022



Singling Out Risk Score

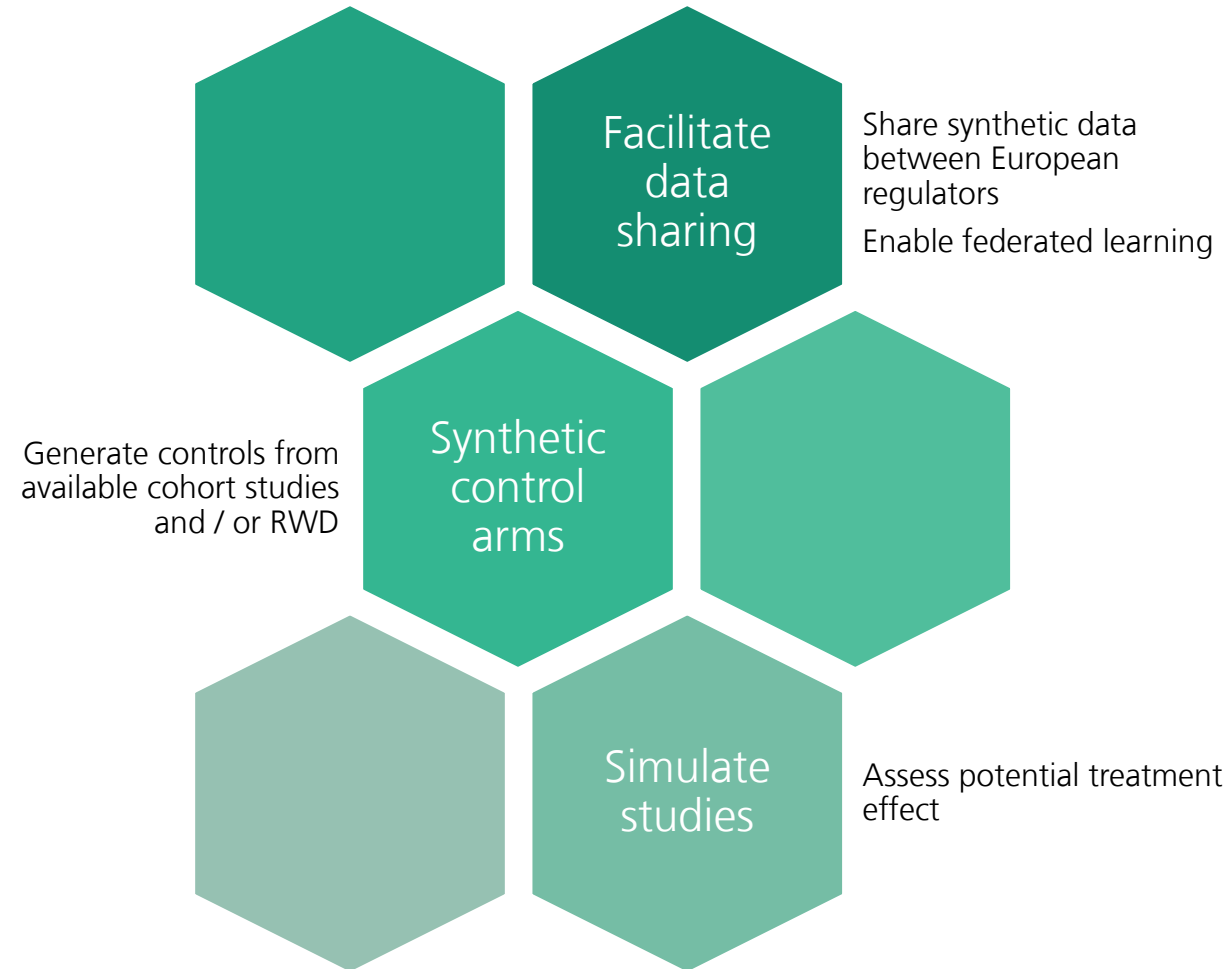


Linkability Risk Score



Inference Risk Score

Summary: Where is the potential benefit for regulators?



Contact

Prof. Dr. Holger Fröhlich

Head of AI & Data Science Group

Deputy Head of Department of Bioinformatics

Fraunhofer Institute for Algorithms and Scientific Computing (SCAI)

Tel. +49 2241 14-4206

holger.froehlich@scai.fraunhofer.de