# Ensuring security and privacy of large-scale datasets: a synthetic data perspective

Jaideep Vaidya

Management Science and Information Systems
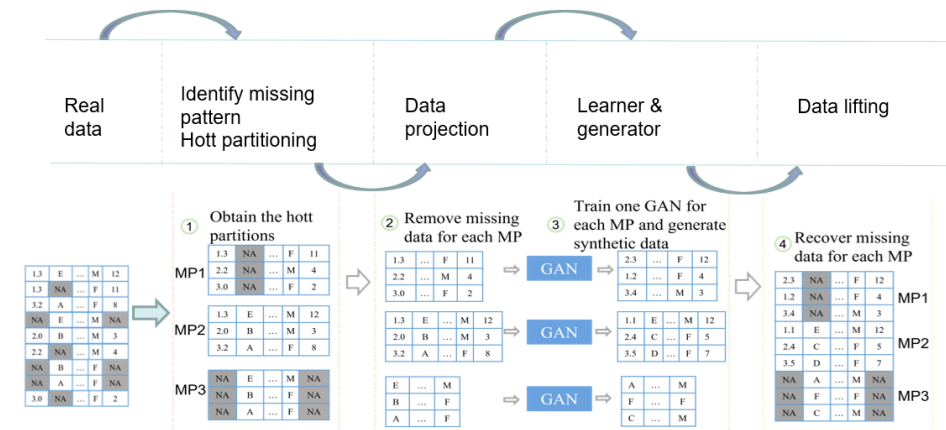
jsvaidya@rutgers.edu

**RUTGERS–NEW BRUNSWICK**
**Artificial Intelligence and Data Science Collaboratory**

- What can you learn from data?

- Lots and lots!
  - Identifying personal genomes by surname inference [GMGHE'13]
  - Identifying participation in complex DNA mixture [H et al.'08]
  - EEG and MEG data can leak financial and identity related information [IHE'18]
  - Modeling 3D facial shape from DNA [C et al.'14]
  - Identifying participation from a summary statistic [DSSUV'15]
  - Identifying participants in the Personal Genome Project by Name [SAW'13]
  - Identifying users in Netflix prize data [NS'08]
  - Identifying "complete patient's info." from aggregate and anonymized queries [VSJO'13]
- The list goes on…

- What can you do about this?
  - Access control
  - Differential privacy
  - Secure Multiparty Computation
  - Synthetic Data
  - …

Recent work: Generate synthetic datasets **conditioned on the missing patterns** that preserves both data and missing data distribution

**Approach 1: HottGAN →** $P(Xo|M)P(M)$



More details: X. Wang, H. Asif, and J. Vaidya. 2023. Preserving Missing Data Distribution in Synthetic Data. WWW '23 https://doi.org/10.1145/3543507.3583297

Genomic Data: Deep generative models (DNADiffusion, CTGAN, ...)

More details: Wang, X., Min, S. and Vaidya, J., 2024. Descriptor: Synthetic Genomic Dataset with Diverse Ancestry (SynGen6). *IEEE Data Descriptions*.