

GPC



Gen² 2026



ICLR

An expressive and tractable deep generative model
for genetic variation data

Prateek Anand

Computer Science PhD Student, UCLA

ICLR 2026 Workshop on Generative AI in Genomics

April 27, 2026



What is genetic variation data?

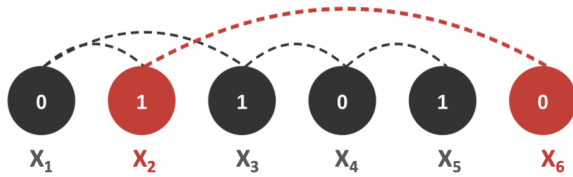
Individuals

Variants (SNPs)

| | X_1 | X_2 | X_3 | X_4 | X_5 | X_6 | X_7 | X_8 |
|-----------------|-------|-------|-------|-------|-------|-------|-------|-------|
| ID ₁ | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 |
| ID ₂ | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 |
| ID ₃ | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| ID ₄ | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 |
| ID ₅ | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 1 |
| ID ₆ | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |

Three key challenges

Generating synthetic data that captures complex correlation



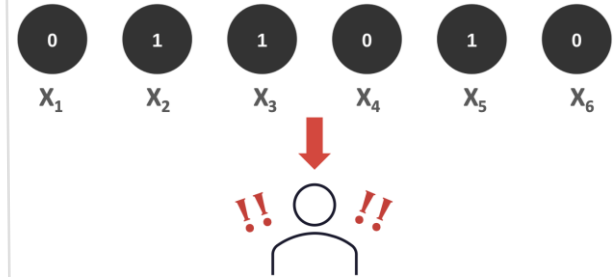
SNPs are correlated across different length scales

Imputing missing data



Many SNPs are unobserved; need to infer missing genotypes

Generating synthetic data that preserves privacy

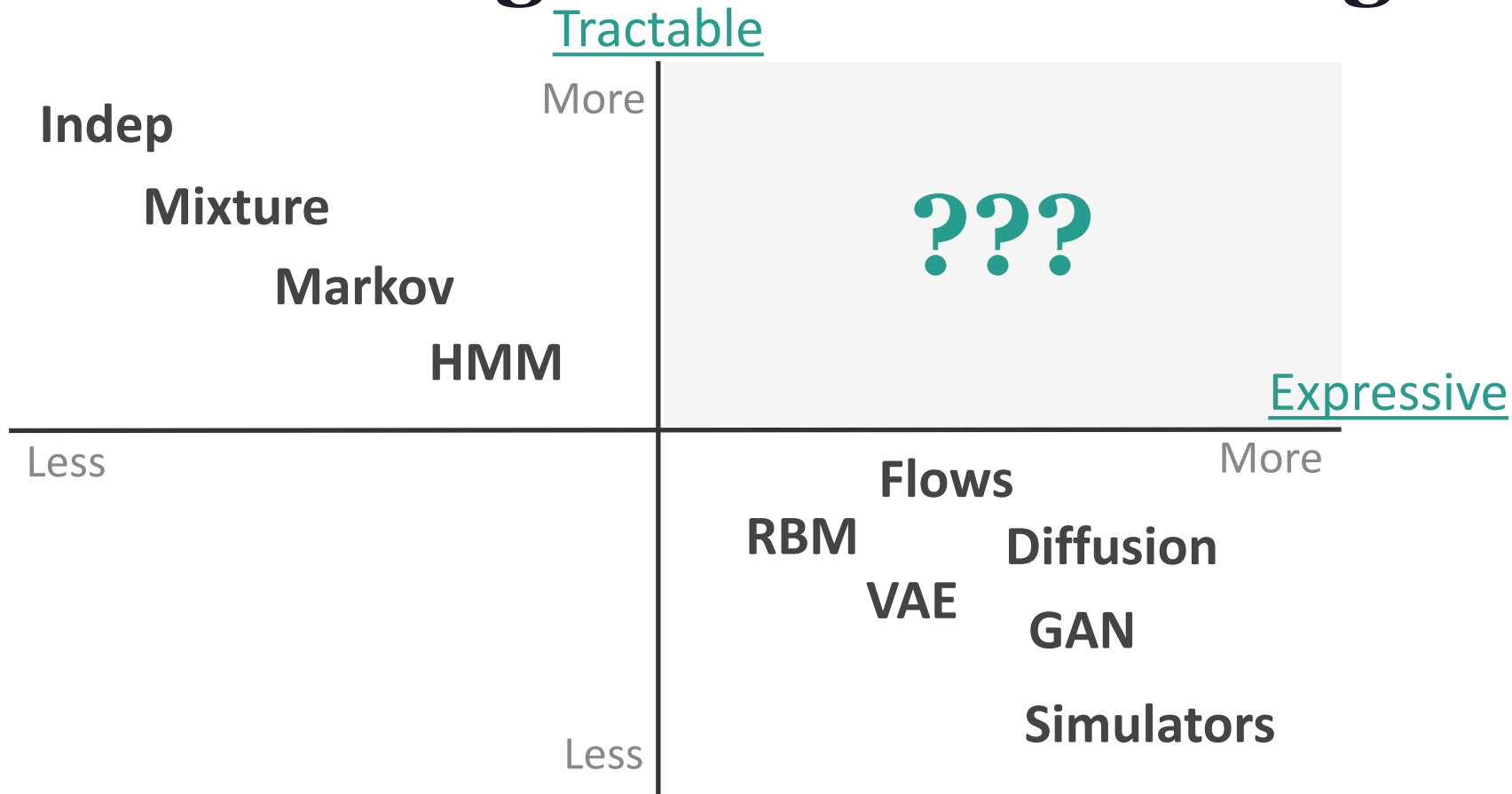


Real genomes can be linked back to individuals

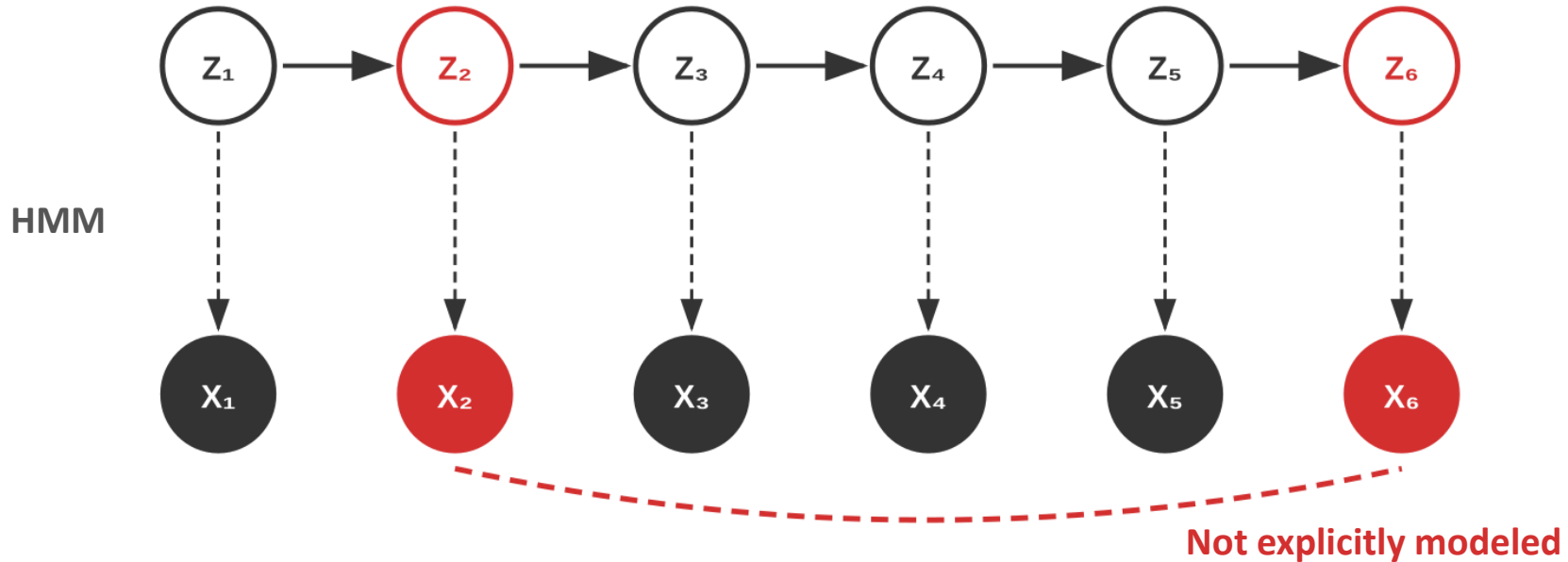
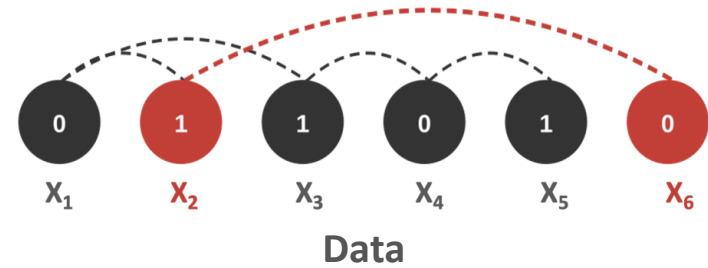
A good model should capture correlation structure, enable imputation, and preserve privacy

Current models struggle to address these while also being tractable

Tradeoffs in generative modeling

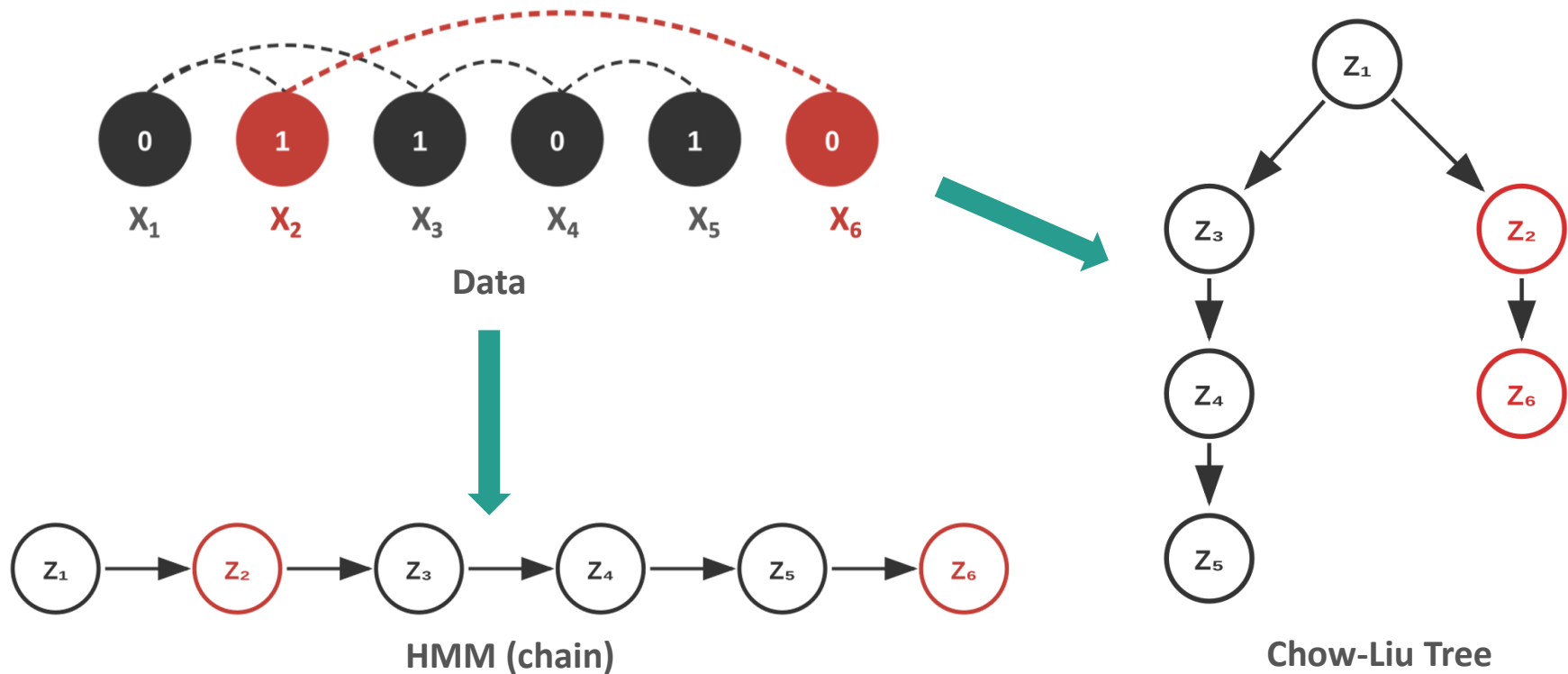


Hidden markov models



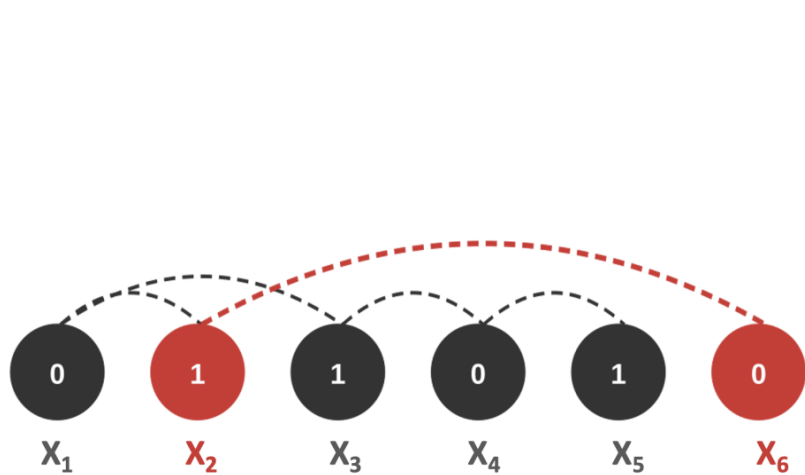
Chain structure forces information between distant SNPs through all intermediates

From chains to trees

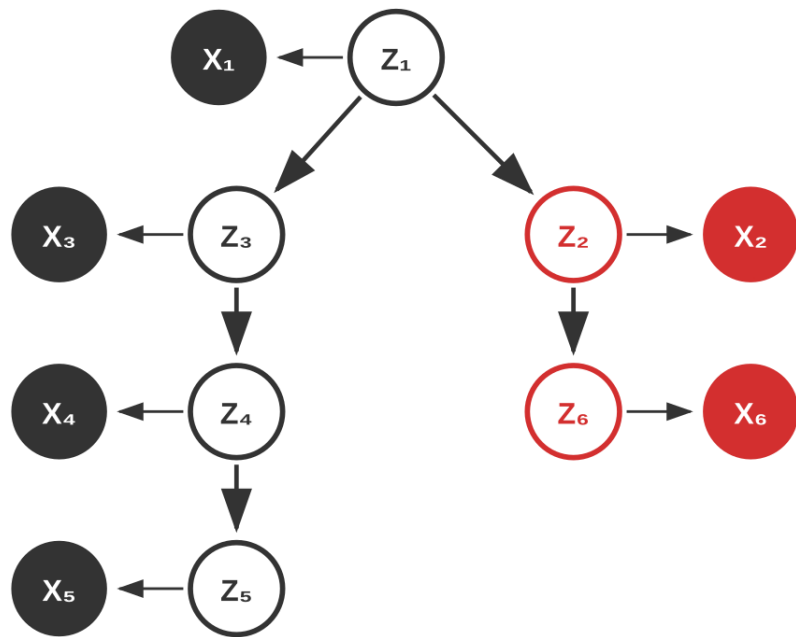


Tree structure places correlated SNPs close together regardless of genomic position

Genetic Probabilistic Circuits (GPC)



Data

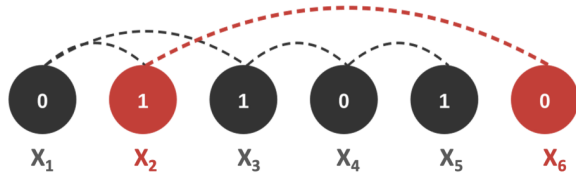


Hidden
Chow-Liu Tree

Compile HCLT \rightarrow **probabilistic circuit** \rightarrow exact likelihoods, efficient training and inference

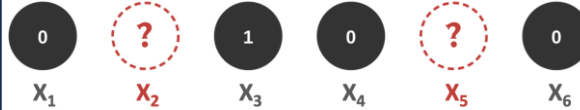
Three key challenges

Generating synthetic data that captures complex correlation



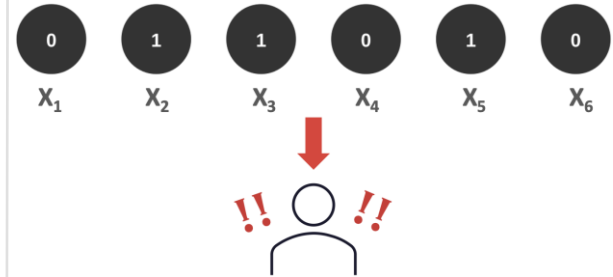
SNPs are correlated across different length scales

Imputing missing data



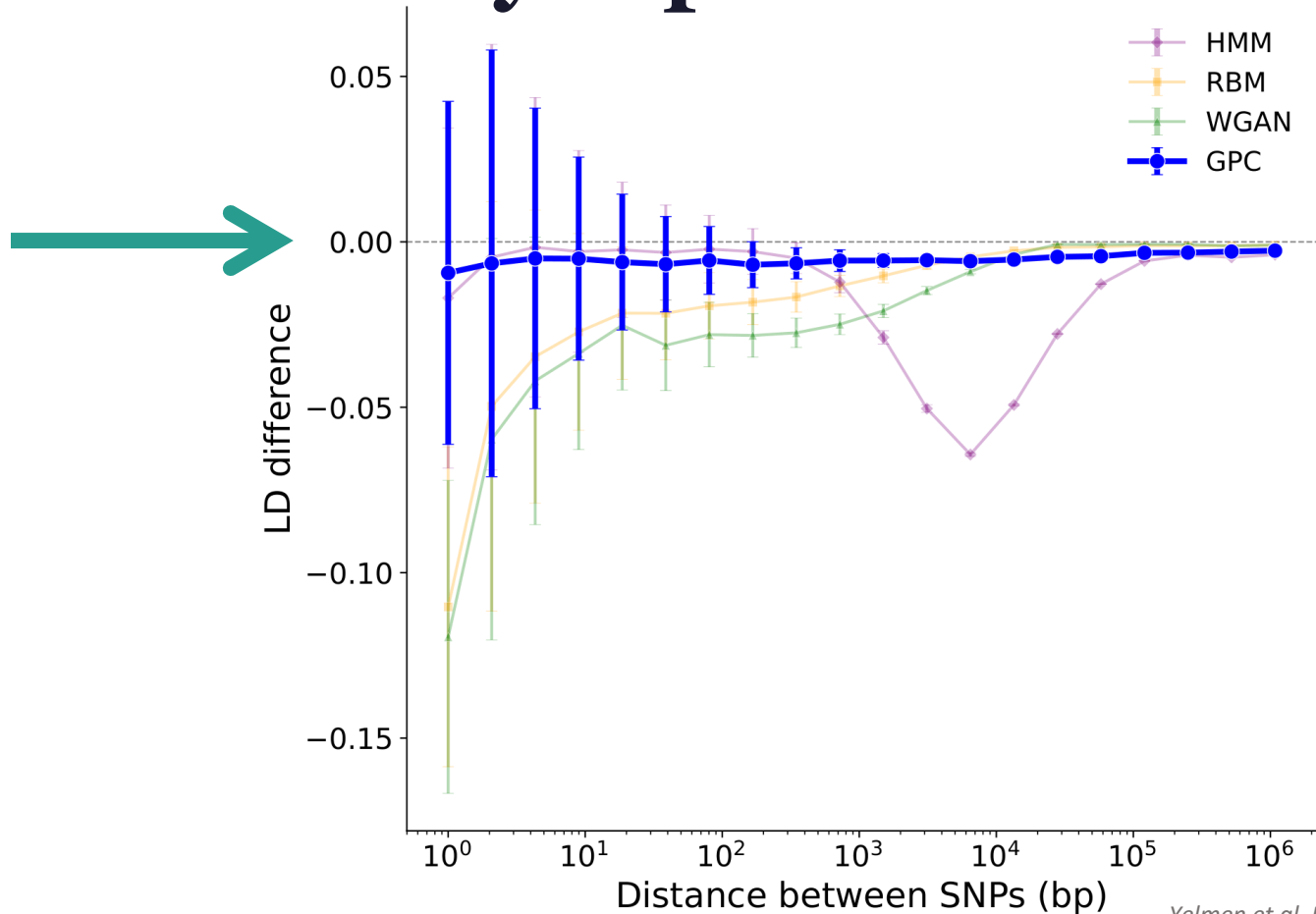
Many SNPs are unobserved; need to infer missing genotypes

Generating synthetic data that preserves privacy

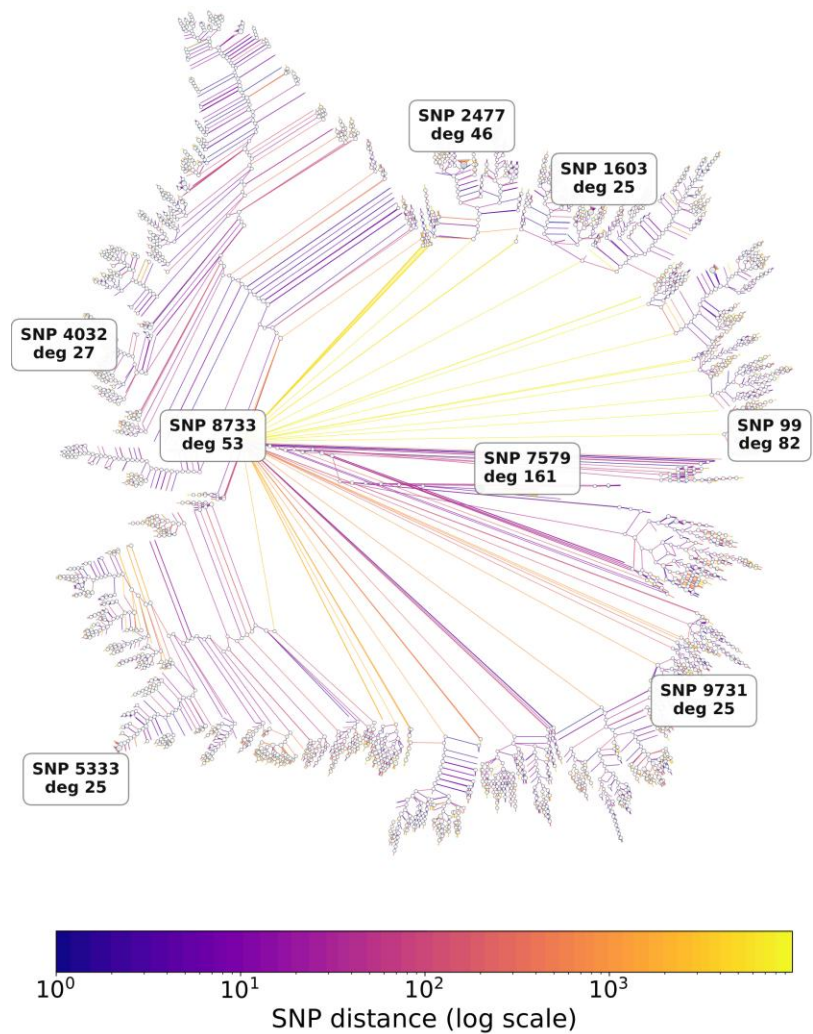
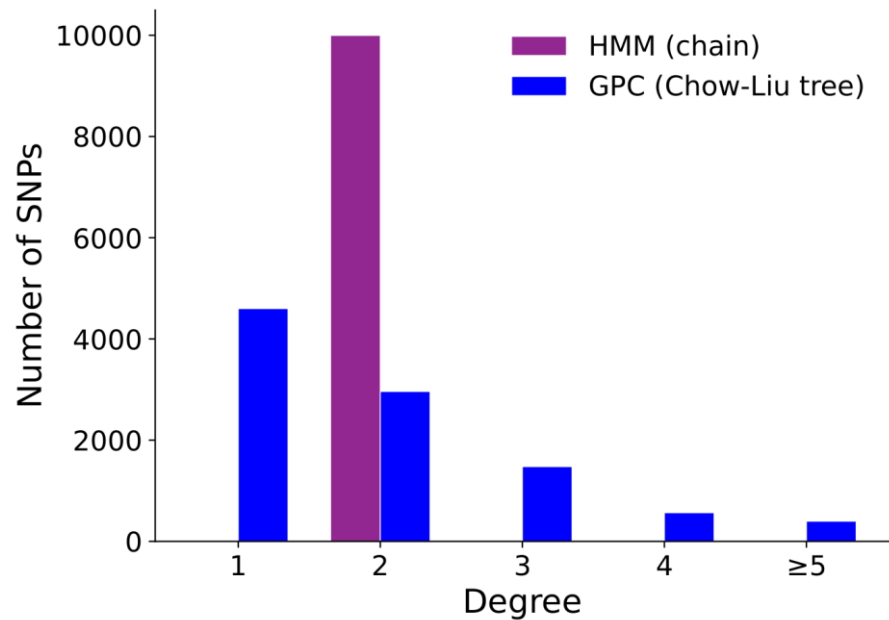


Real genomes can be linked back to individuals

GPC closely reproduces correlation

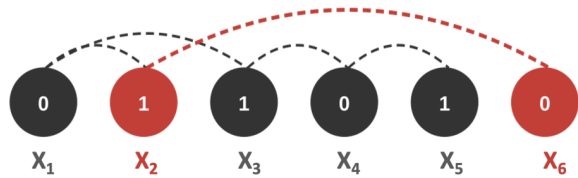


GPC tree structure reflects complex correlation



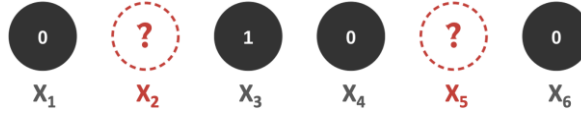
Three key challenges

Generating synthetic data that captures complex correlation



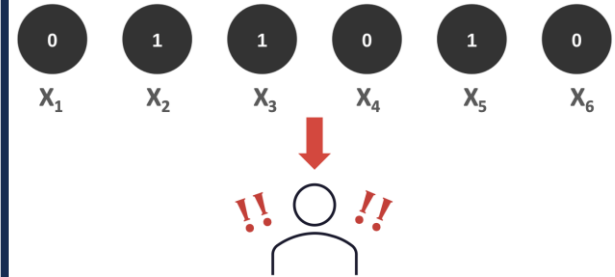
SNPs are correlated across different length scales

Imputing missing data



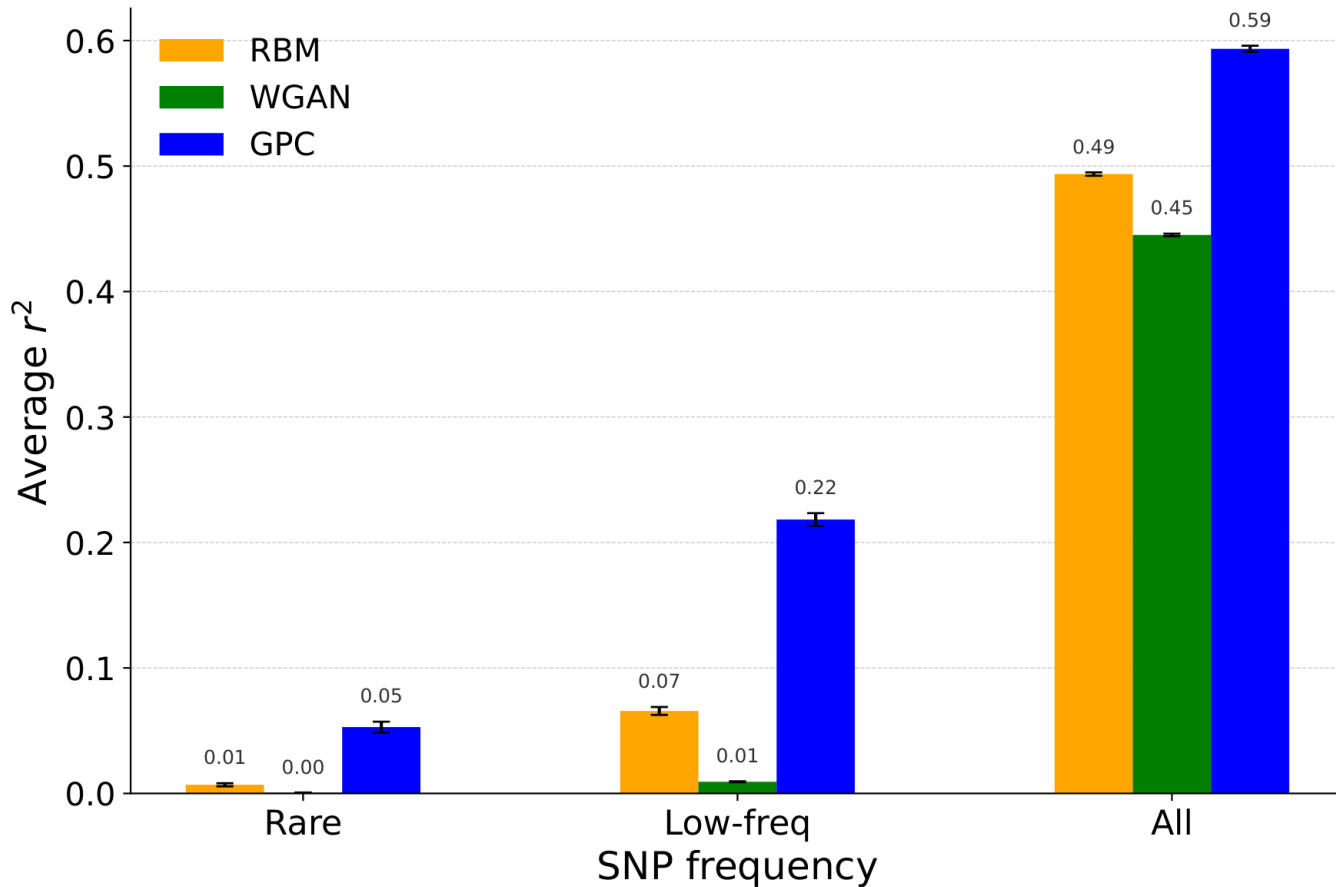
Many SNPs are unobserved; need to infer missing genotypes

Generating synthetic data that preserves privacy



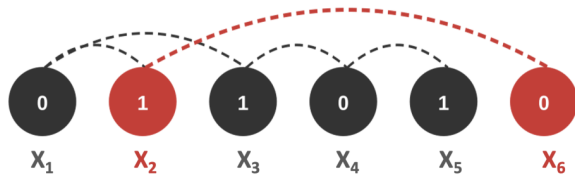
Real genomes can be linked back to individuals

GPC achieves best imputation among generative baselines



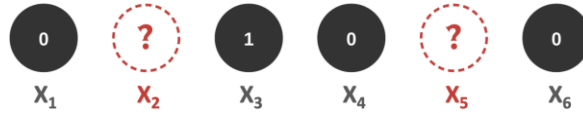
Three key challenges

Generating synthetic data that captures complex correlation



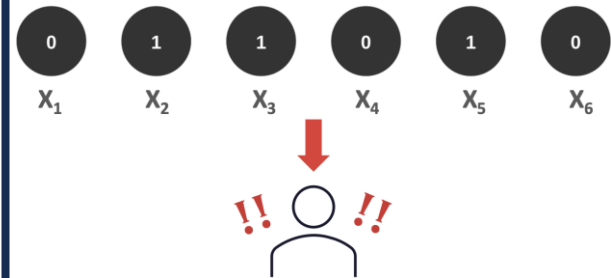
SNPs are correlated across different length scales

Imputing missing data



Many SNPs are unobserved; need to infer missing genotypes

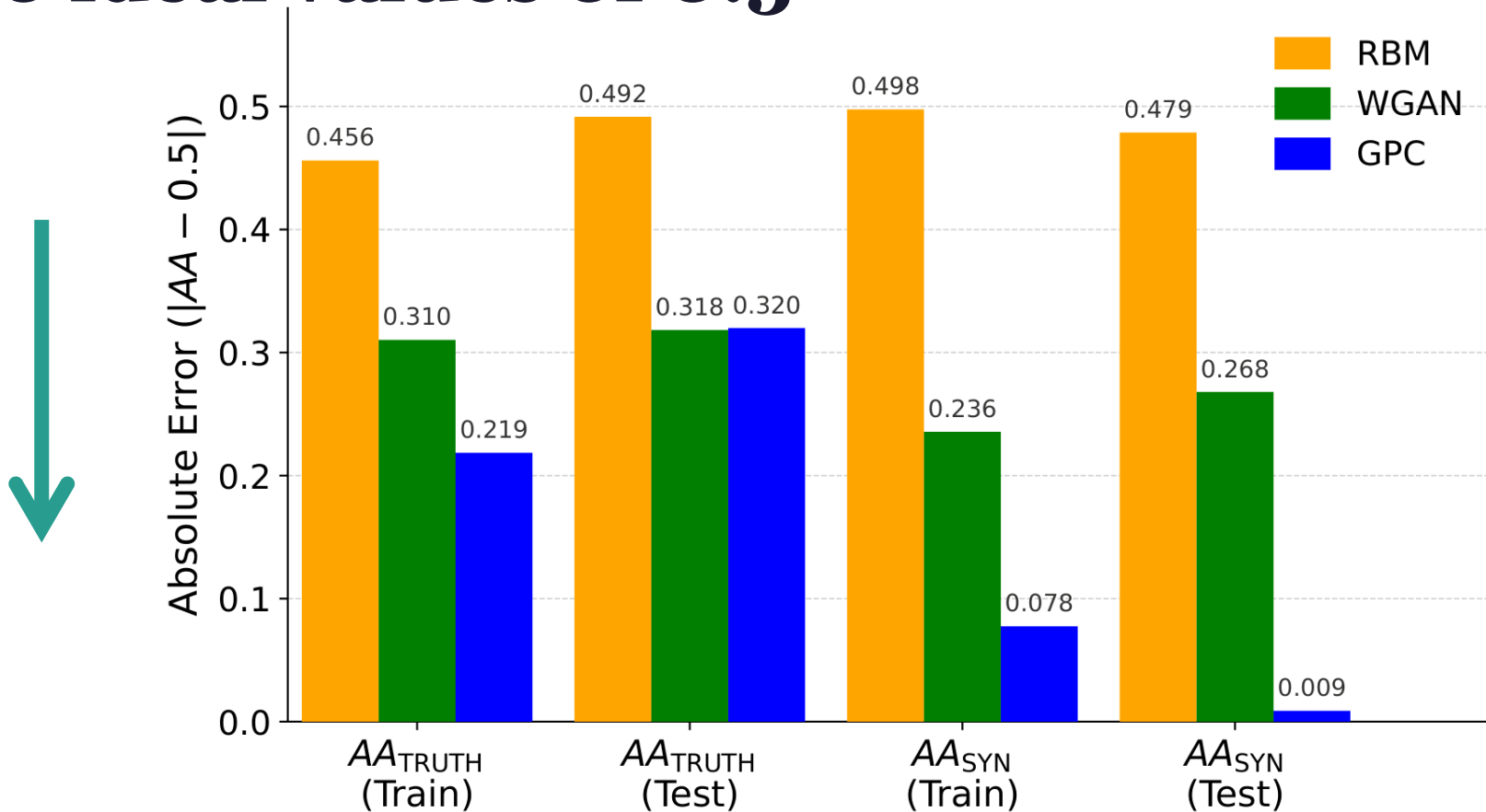
Generating synthetic data that preserves privacy



Real genomes can be linked back to individuals

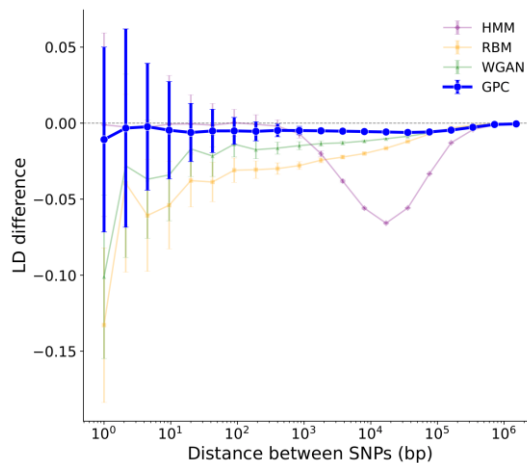
Adversarial Accuracy (AA)
Based on nearest neighbors

GPC balances privacy and utility: closest to ideal values of 0.5



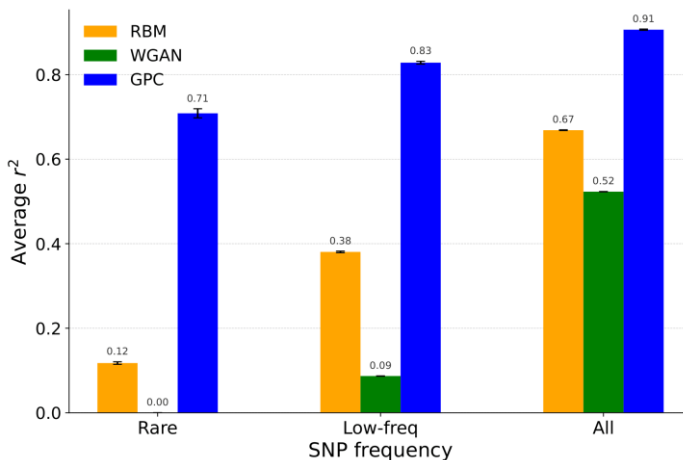
All findings replicate in UK Biobank

Generating synthetic data that captures complex correlation

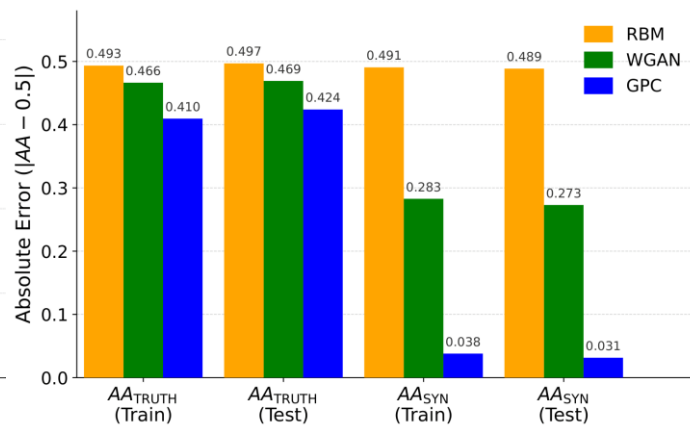


Close to 0

Imputing missing data



Generating synthetic data that preserves privacy



GPC takeaways

1

Expressive, tractable, and efficient: tree structure captures LD at all scales with exact inference

2

Best imputation accuracy among generative models

3

Better privacy-utility tradeoff for synthetic data than existing methods

Acknowledgments



Anji Liu

NUS



Meihua Dang

Stanford



Boyang Fu

Harvard



Xinzhu Wei

Cornell



Guy Van den Broeck*

UCLA



Sriram Sankararaman*

UCLA



Sriram Lab

UCLA

* Joint supervision

prateek@cs.ucla.edu | <https://prateekanand2.github.io/> | github.com/sriramlab/GPC

Thank you!



Preprint



GitHub

Come find me at Poster Session 1 Board 26!