

Deep generative model of genetic variation data improves imputation accuracy in private populations

Prateek Anand¹, Anji Liu², Meihua Dang³, Boyang Fu⁴, Xinzhu Wei⁵, Guy Van den Broeck^{1,*}, Sriram Sankararaman^{1,6,7,*}

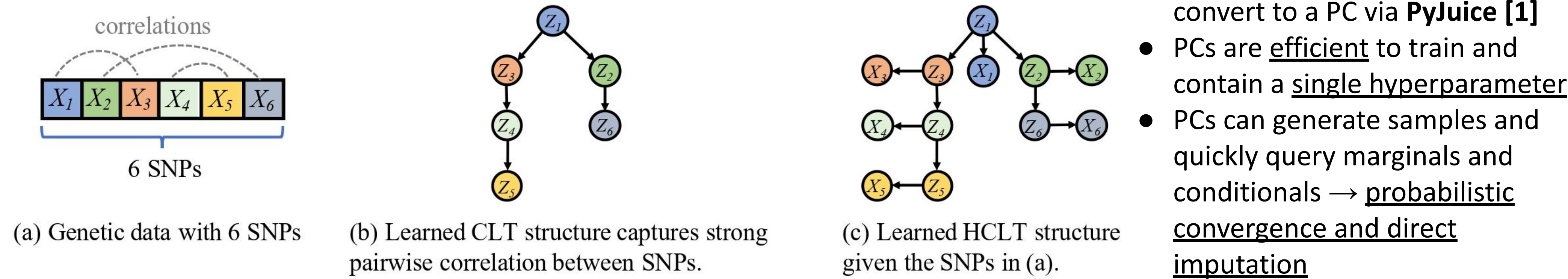
*Equal Contribution ¹Department of Computer Science, UCLA ²Department of Computer Science, National University of Singapore ³Department of Computer Science, Stanford University

⁴Department of Biomedical Informatics, Harvard Medical School ⁵Department of Computational Biology, Cornell University ⁶Department of Human Genetics, UCLA ⁷Department of Computational Medicine, UCLA

Introduction

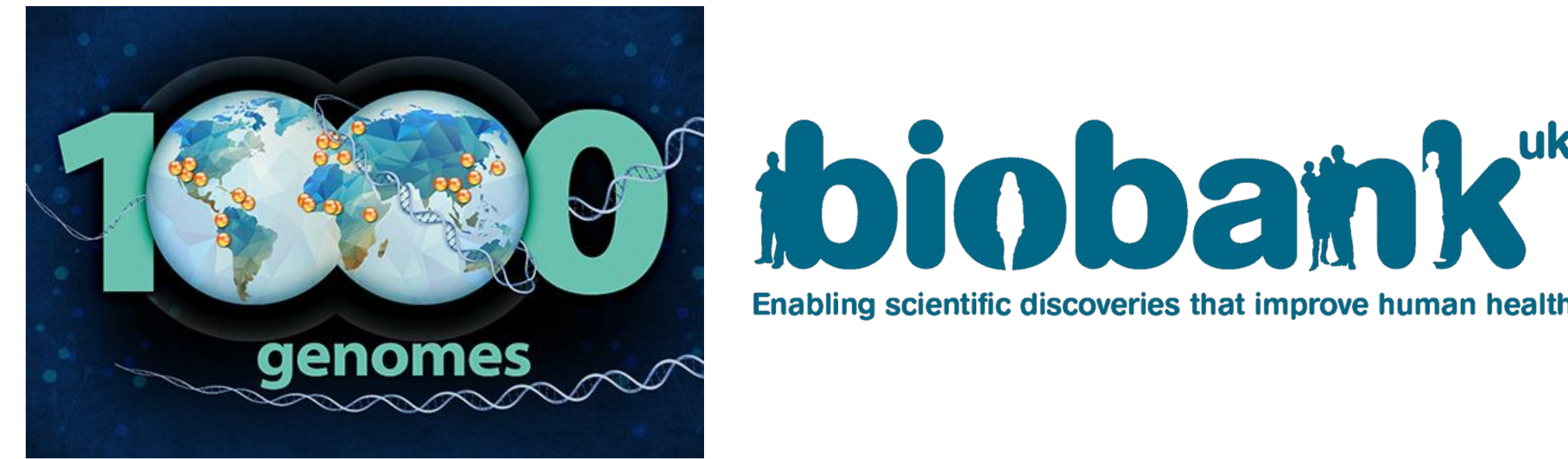
- Artificial Genomes (AGs) enable sharing and benchmarking without exposing sensitive genetic data, supporting reproducibility and equity in population genomics
- Coalescent and PAC/HMM approaches are interpretable and tractable but often computationally intensive, and may not fully capture fine-scale genomic dependencies
- GANs, VAEs, and RBMs offer expressiveness but lack tractable inference, reliable likelihood computation, and are challenging to train and tune
- We propose a novel **tractable and expressive deep generative model** of genetic variation based on **Hidden Chow–Liu Trees (HCLTs)** represented as **Probabilistic Circuits (PCs)**, enabling efficient learning and inference

Methods



Data

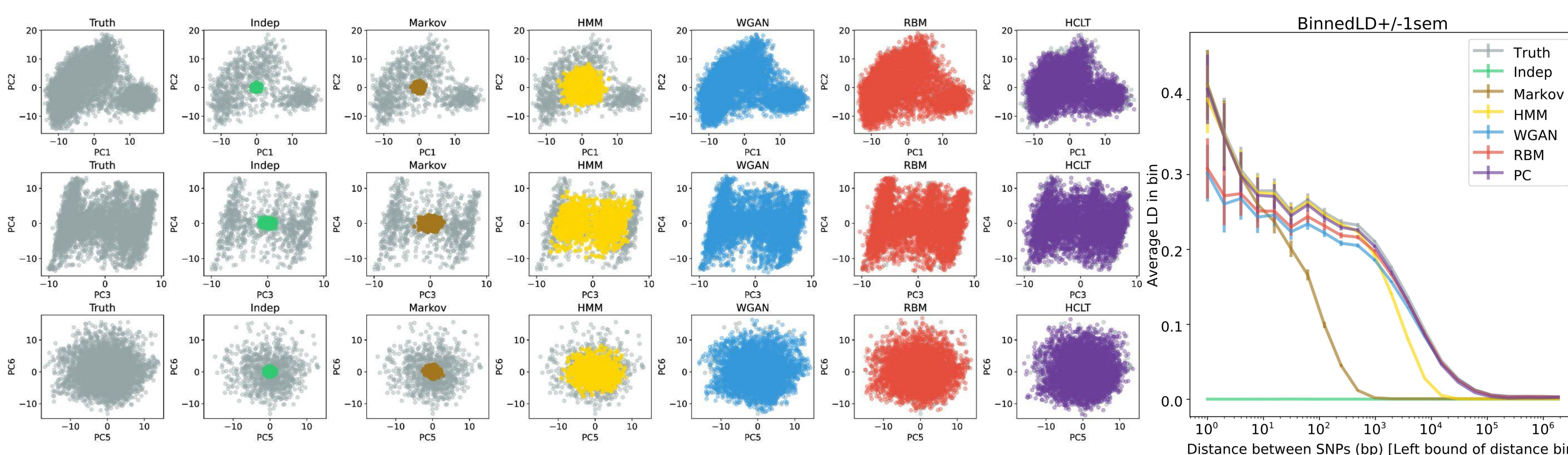
- Phased haplotype data (0/1)
- 1000 Genomes Project Phase 3 (low-coverage)
 - 15:27379578 - 15:29625035 (5008 x 10000)
- UK Biobank
 - 22:29456546 - 22:32665772 (26924 x 9820)
- 1000 Genomes Project Phase 3 (high-coverage)
 - 15:27134431 - 15:29332831 (5008 x 14670)



Reconstructing Global and Local Population Structure

Dataset	Category	INDEP	MARKOV	HMM	PC
1KG	train LL	-2386.81	-1806.33	-591.08	-202.51
	test LL	-2404.51	-1819.96	-599.88	-265.06
	#params	10.00k	39.99k	163774.98k	88473.73k
UKBB	train LL	-1642.62	-1360.86	-554.88	-120.10
	test LL	-1648.03	-1362.16	-554.38	-127.75
	#params	9.82k	39.27k	160825.86k	88850.56k

- Top:** Log-likelihood of models that support tractable likelihood inference
- Bottom:** Left shows top 6 principal components of ground truth vs. samples, while right shows binned LD as a function of SNP distance



Preserving Privacy of the Training Data

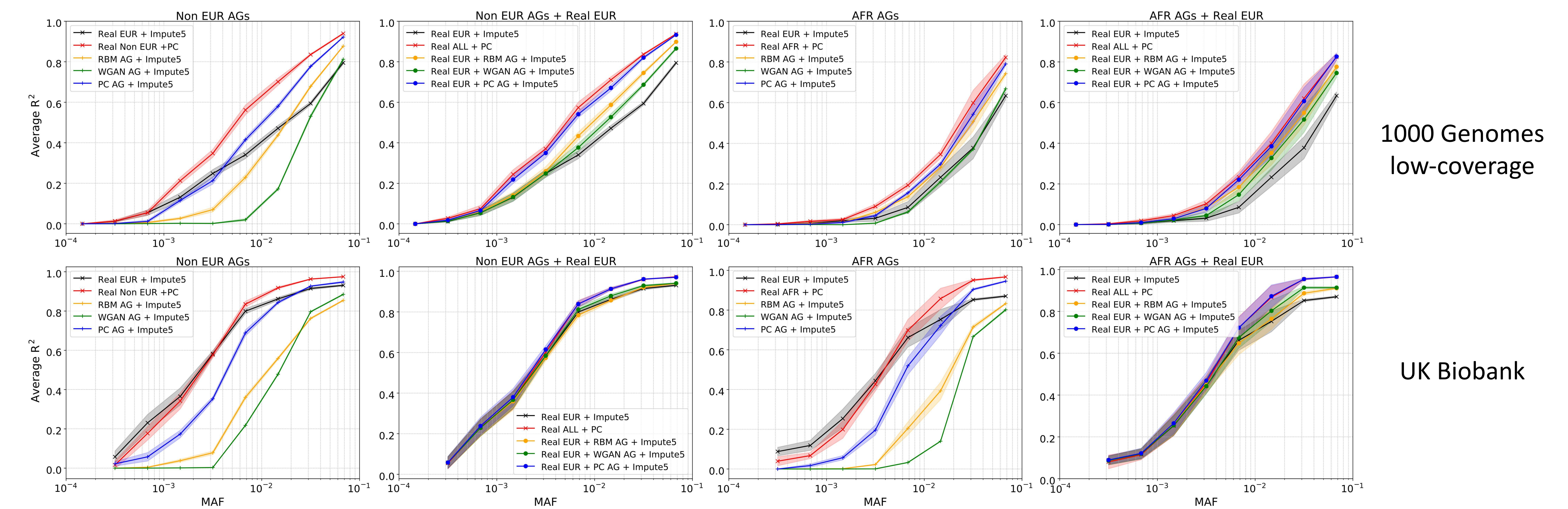
$$AATS = \frac{1}{2} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{1}(d_{TS}(i) > d_{TT}(i)) + \frac{1}{n} \sum_{i=1}^n \mathbf{1}(d_{ST}(i) > d_{SS}(i)) \right)$$

Dataset	Metric	RBM	WGAN	PC
1KG	AA _{TRUTH} (Train)	0.9561	0.8103	0.7185
	AA _{TRUTH} (Test)	0.9928	0.7764	0.7680
	AA _{SYN} (Train)	0.0024	0.7356	0.4225
	AA _{SYN} (Test)	0.0276	0.7847	0.5304
UKBB	AA _{TRUTH} (Train)	0.9954	0.9674	0.9204
	AA _{TRUTH} (Test)	0.9962	0.9688	0.9198
	AA _{SYN} (Train)	0.0064	0.7768	0.5324
	AA _{SYN} (Test)	0.0160	0.7582	0.4630

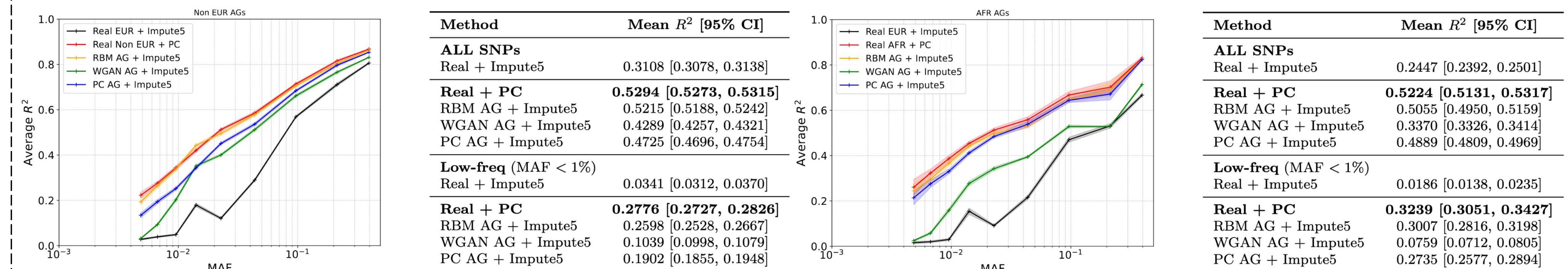
- Nearest Neighbor Adversarial Accuracy [2]**
- T = Truth, S = Synthetic
- AA_{TRUTH}: Low → synthetic mimics real too closely (overfit, low privacy). High → synthetic is far from real (high privacy, lower utility).
- AA_{SYN}: Low → synthetic copies real samples (low privacy). High → synthetic has its own structure (high privacy, may reduce utility if too different).
- Ideal:** Both around **0.5** → realistic synthetic data that balances **privacy and utility**
- Only report components since average can hide poor performance

Improving Imputation Accuracy in Private Populations

- Imputation of one SNP at a time using AGs as reference panel in Impute5 pipeline (or combined with real EUR data)
- Black line represents Impute5 only using real EUR data as reference



- Imputation of all SNPs not present on InfiniumOmni2-5-8v1-4 A1 array using 1000 Genomes high-coverage as training
- Out of 14670 SNPs, 1302 are present on the array, and remaining are imputed simultaneously (91% missingness)



Conclusion

- We present a novel **deep generative model** for genetics based on **PCs**
- Fast & tractable:** Easy to train and use at scale
- Realistic:** Matches key genetic patterns, with more accurate LD
- Privacy-preserving:** Balances similarity and separation effectively
- Boosts imputation:** Improves accuracy for underrepresented groups
- Direct inference:** Enables imputation and likelihood evaluation
- Enables access:** Supports open, equitable genomic research

References

- Anji Liu, Kareem Ahmed, and Guy Van den Broeck. Scaling tractable probabilistic circuits: A systems perspective. In Proceedings of the 41th International Conference on Machine Learning (ICML), jul 2024.
- Andrew Yale, Saloni Dash, Ritik Dutta, Isabelle Guyon, Adrien Pavao, and Kristin P. Bennett. Generation and evaluation of privacy preserving synthetic health data. Neurocomputing, 416:244–255, 2020.