

# Overcoming Data Limitations in Animal Genetics with GAN-based Synthetic Genotypes

Sihan Xie<sup>1</sup>, Thierry Tribout<sup>1</sup>, Blaise Hanczar<sup>2</sup>, Julien Chiquet<sup>3</sup>, Eric Barrey<sup>1</sup>

<sup>1</sup> Université Paris-Saclay, AgroParisTech, INRAE, GABI, 78350 Jouy-en-Josas, France <sup>2</sup> IBISC, University Paris-Saclay (Univ. Evry) <sup>3</sup> UMR MIA-Paris, Université Paris-Saclay, AgroParisTech, INRAE

## Context

Real genotype data is hard to access due to privacy, cost, and regulatory constraints, posing challenges for animal genetics research. We overcome these issues by employing deep learning models to generate synthetic genotype conditioned on phenotype.

## Data

We train our model using a Holstein cow dataset that contains genotype and milk production traits from 93484 cows, featuring 50161 single-nucleotide polymorphism (SNPs) distributed across 29 chromosomes.



## Model

We developed a Wasserstein Generative Adversarial Network with Gradient Penalty (WGAN-GP) that leverages a Gumbel-Softmax mechanism to generate discrete synthetic genotype data conditioned on milk production phenotype:

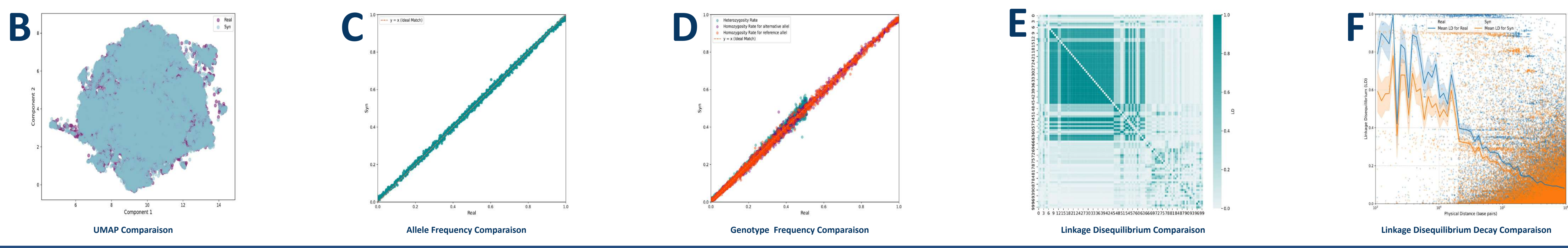
$$\min_G \max_D \underbrace{\mathbb{E}_{x \sim P_r}[D(x|c)] - \mathbb{E}_{z \sim P_z}[D(G(z|c))]}_{\text{Wasserstein distance between real and synthetic distribution}} - \underbrace{\lambda \mathbb{E}_{\hat{x} \sim P_{\hat{x}}}(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2}_{\text{enforce Lipschitz constraint}}$$

where G is a Generator Network and D is a Discriminator Network.

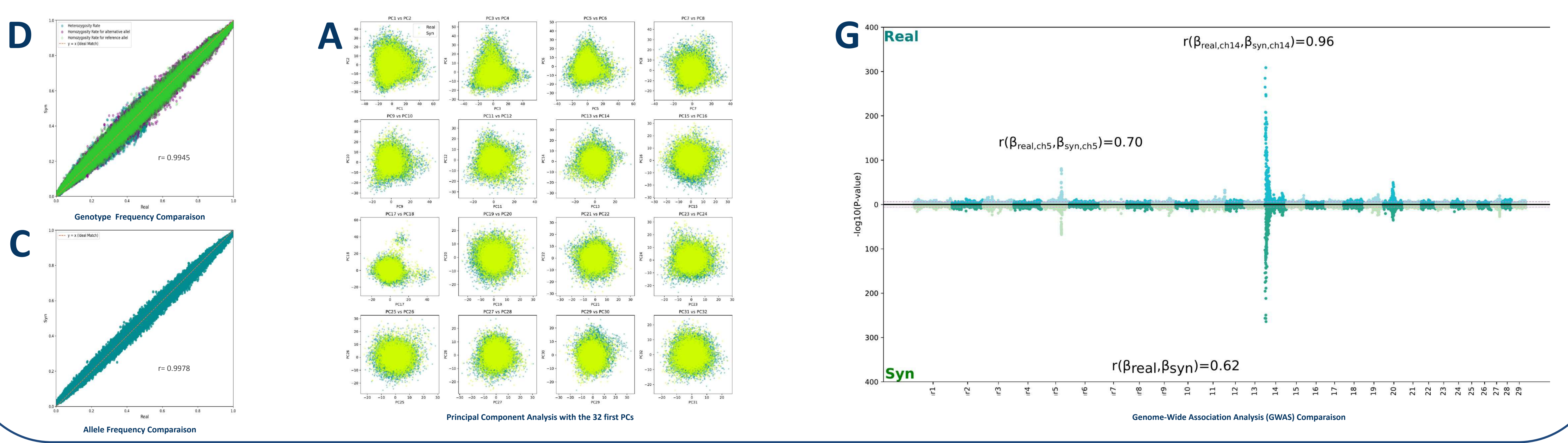
## Quantitative Metrics from Deep Learning Literature

Results for Chromosome 14 (1771 SNPs)					Results for All Chromosomes (50161 SNPs)				
	Precision (%) ↑	Recall (%) ↑	F1 (%) ↑	Correlation Score (%) ↑		Precision (%) ↑	Recall (%) ↑	F1 (%) ↑	Correlation Score (%) ↑
Vanilla GAN	80.88 ± 0.23	57.97 ± 0.45	67.53 ± 0.24	72.60 ± 0.13	Vanilla GAN	100 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.52 ± 0.01
WGAN-GP	99.64 ± 0.04	99.88 ± 0.02	99.76 ± 0.02	98.65 ± 0.02	WGAN-GP	92.00 ± 0.16	99.93 ± 0.01	95.80 ± 0.09	83.32 ± 0.06

## Results for Chromosome 14 (1771 SNPs)



## Results for All Chromosomes (50161 SNPs)



## Conclusion

By comparing results for real and synthetic data across several quantitative metrics, PCA (Fig. A), UMAP (Fig. B), allele frequency (Fig. C), genotype frequency (Fig. D), linkage disequilibrium and its decay with physical distance (Figs. E–F), and GWAS (Fig. G), we demonstrate that our **WGAN-GP combined with Gumbel Softmax model generated synthetic genotype data accurately capture key genetic structures while preserving critical genotype–phenotype associations.**

We thank the INRAE GABI GBOS Team for providing and preparing the dataset used in this study.