# Automatic evolutionary inference using Generative Adversarial Networks

## Sara Mathieson

Assistant Professor

HAVERFORD COLLEGE

Department of Computer Science

NIH, NHGRI

*Machine Learning in Genomics Workshop*

April 13, 2021

# Central question in population genetics: data -> quantify evolution
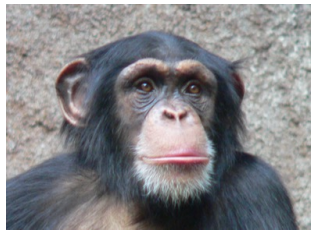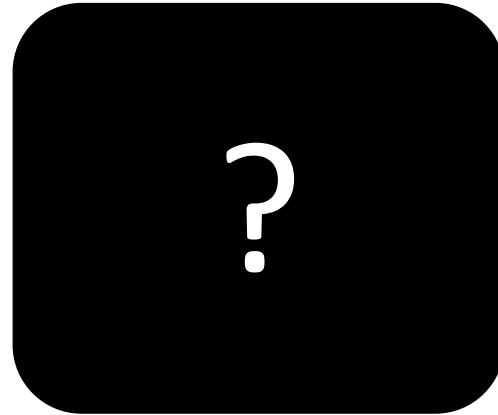
INPUT

OUTPUT

**Sites or SNPs**

samples/haplotypes

```
1000010001000001011000010100100101110000010
1100001110110011110100001011110001100100011
0100100101000100000000000001111011100001001
1101100000110010000000000100010001100011100 11
0010110010000000000011010111101001000111000
0010001010000000111000110000010101110010000
0000000010101001101001010110010010000001 0000
0000000010110001000000100000001000110110000
1010110100101001001100000001111001100000011
0000100011101000000000011101001010010100111
0000000111001000011100000011111000000110101
0100011000011000000000101000001011110000101
0100000010101000100001000010100000001001011
0000010000010101011000010110110100000100100
```
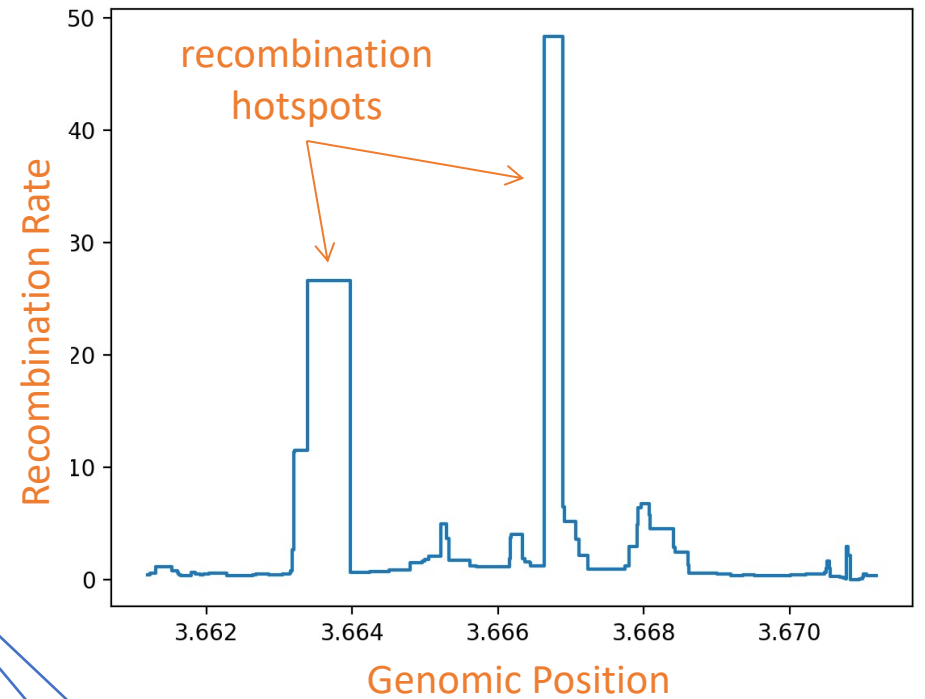
?


recombination hotspots

Recombination Rate

Genomic Position

- Population size changes
- Natural selection
- Mutation rate variation
-   Migration, admixture, introgression
- Heritable traits and diseases

Images: wikipedia

# Central question in population genetics: data -> quantify evolution
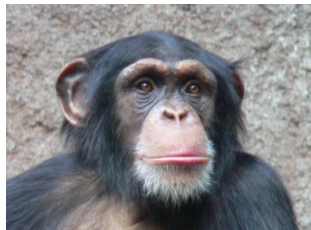
**Sites or SNPs**

**samples/haplotypes**

```
10000100010000010110000101001001011110000010
11000011101100111101000010111100011001100011
01001001010001000000000000001111011100001001
11011000001100100000000010001000110001110011
00101100100000000000110101111010010010111000
00100010100000000111000110000010101110010000
00000001010100110101001010110010010000010000
00000000101100010000001000000010001101100000
10101101001010010011000000011110011000000011
00001000111010000000001110100101001010100111
00000001110010000011100000011111100000110101
01000110000110000000101000001011110000101
01000000101010001000010001010000001001011
00000100000101010110000101101101000001000100
```
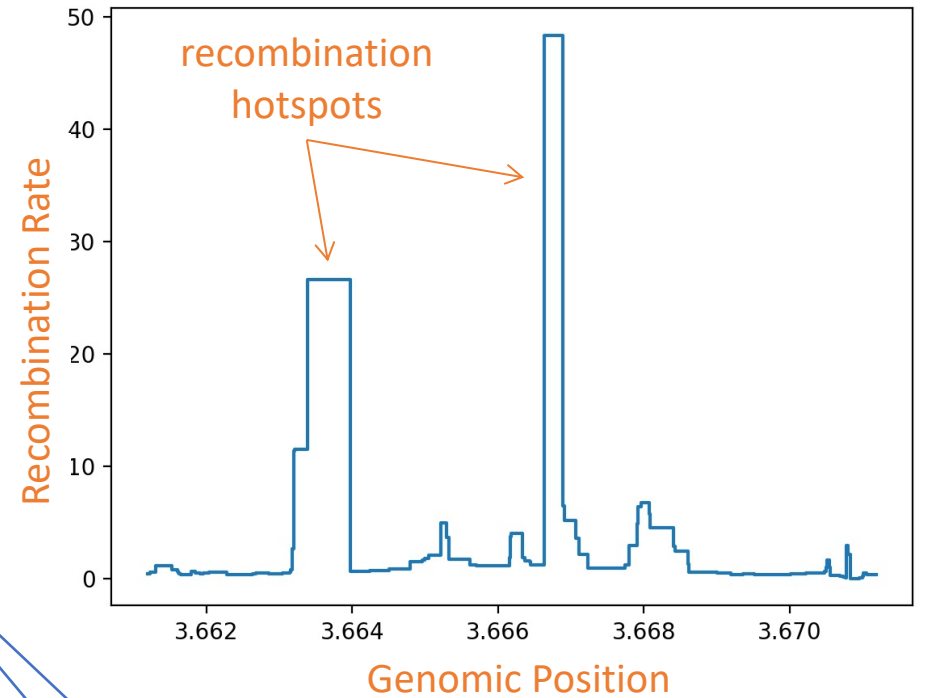
Fast?
Flexible?
Machine
learning?



recombination hotspots

Recombination Rate

50

40

30

20

10

0

3.662    3.664    3.666    3.668    3.670

Genomic Position

- Population size changes

- Natural selection

- Mutation rate variation

- Migration, admixture, introgression

- Heritable traits and diseases

Images: wikipedia

# Outline

- Shift to machine learning in population genetics

- Shift away from summary statistics to "raw" data

- GANs and adversarial training
  - pg-gan algorithm for creating realistic simulated data

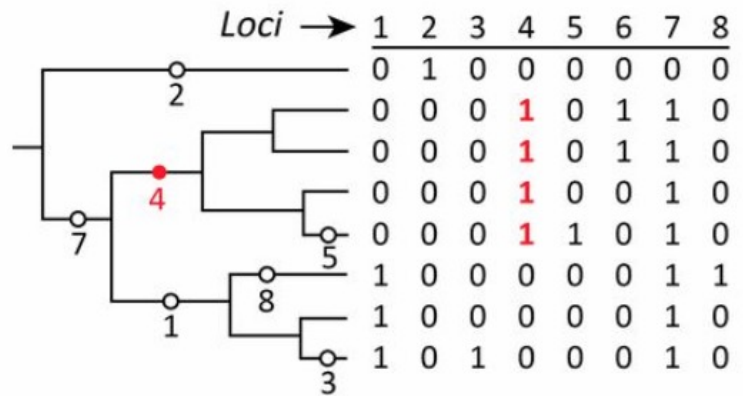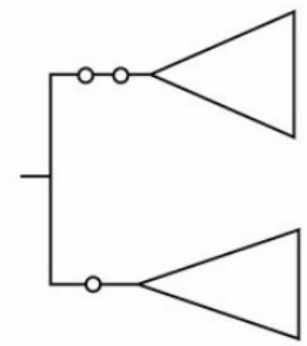- Results on human data from Africa, Europe, and East Asia

# Outline

- Shift to machine learning in population genetics

- Shift away from summary statistics to "raw" data

- GANs and adversarial training
  - pg-gan algorithm for creating realistic simulated data

- Results on human data from Africa, Europe, and East Asia

# 2013: Using machine learning to infer selection

**A**  Genealogy & SNP Matrix

**B**  Neutral Evolution

**C**  Positive Selection

# 2013: Using machine learning to infer selection

**Method**:
support vector
machines (SVM)



region of real data

Regions under
selection (simulated)

Neutral regions
(simulated)

Optimal hyperplane

Maximum margin
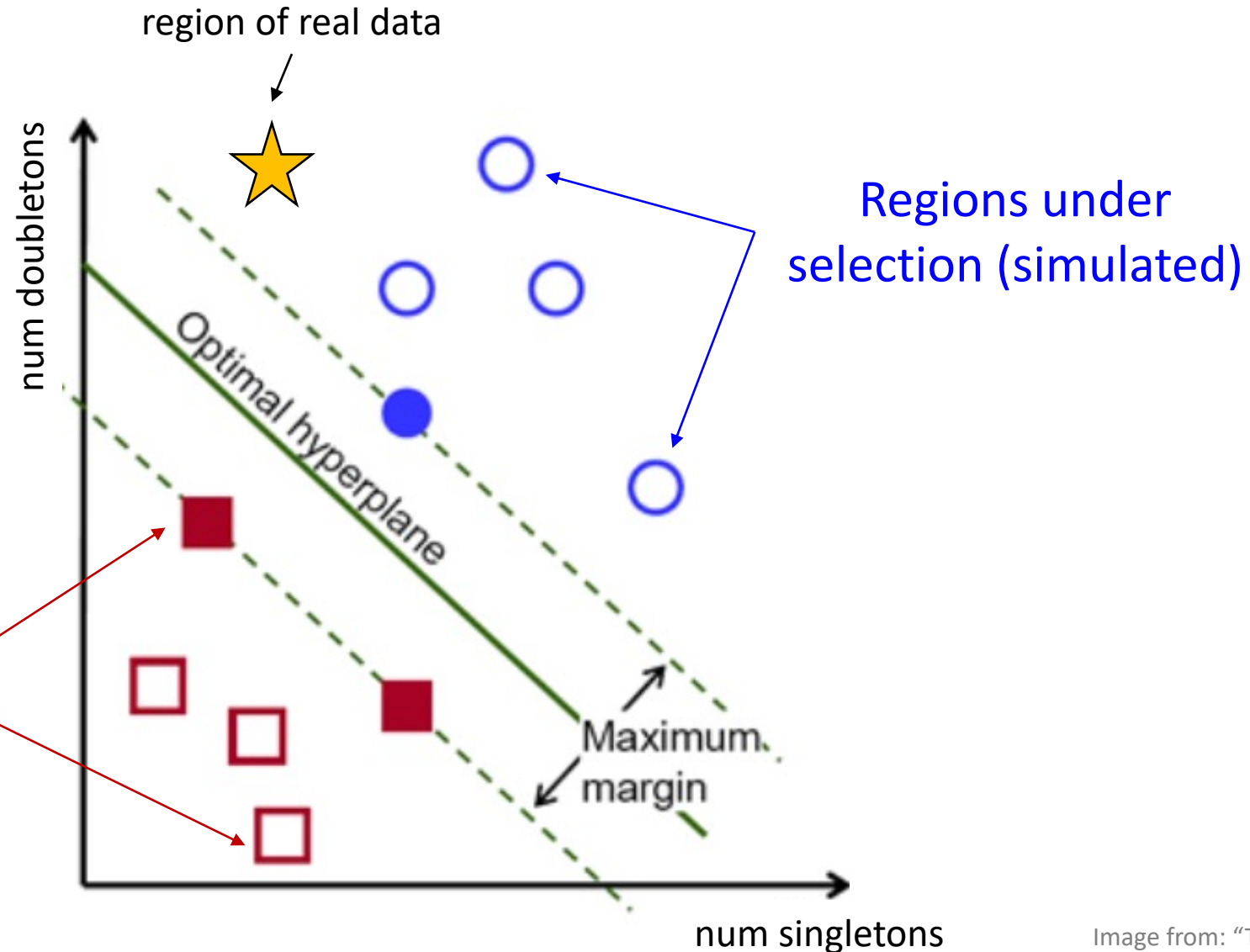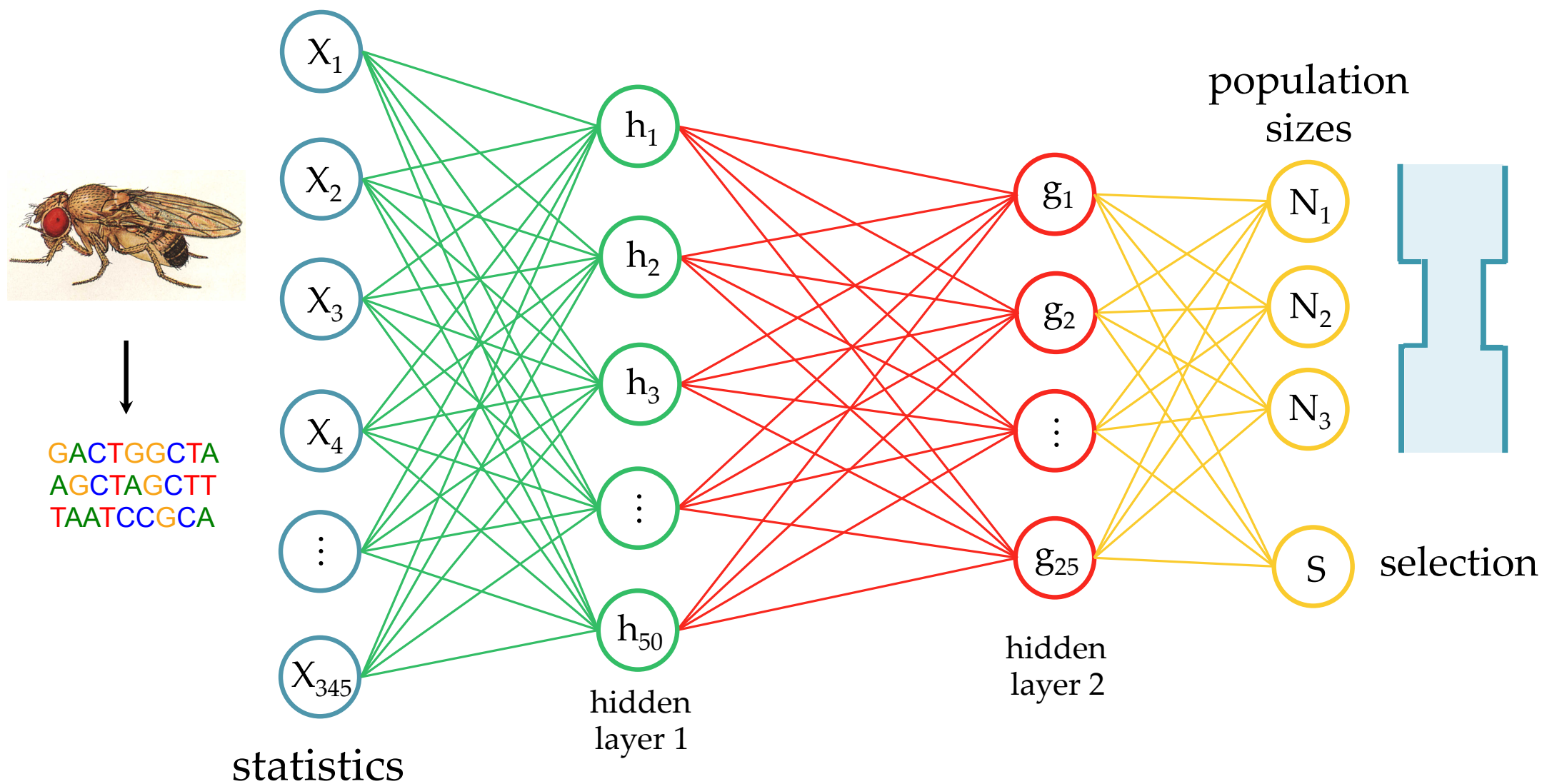
num doubletons

num singletons

Image from: "Towards Data Science"

# Which summary statistics to use?

▶ Number of segregating sites     **3 stats**

▶ Tajima's $D$     **3 stats**

▶ Folded site frequency spectrum (SFS)     **150 stats**

▶ Length distribution between segregating sites     **48 stats**    **345 total**

▶ Identity-by-state (IBS) tract length distribution     **90 stats**

▶ Linkage disequilibrium (LD) distributions     **48 stats**

▶ Haplotype frequency statistics     **3 stats**

Example summary statistics from "Deep learning for population genetic inference", Sheehan and Song, 2016
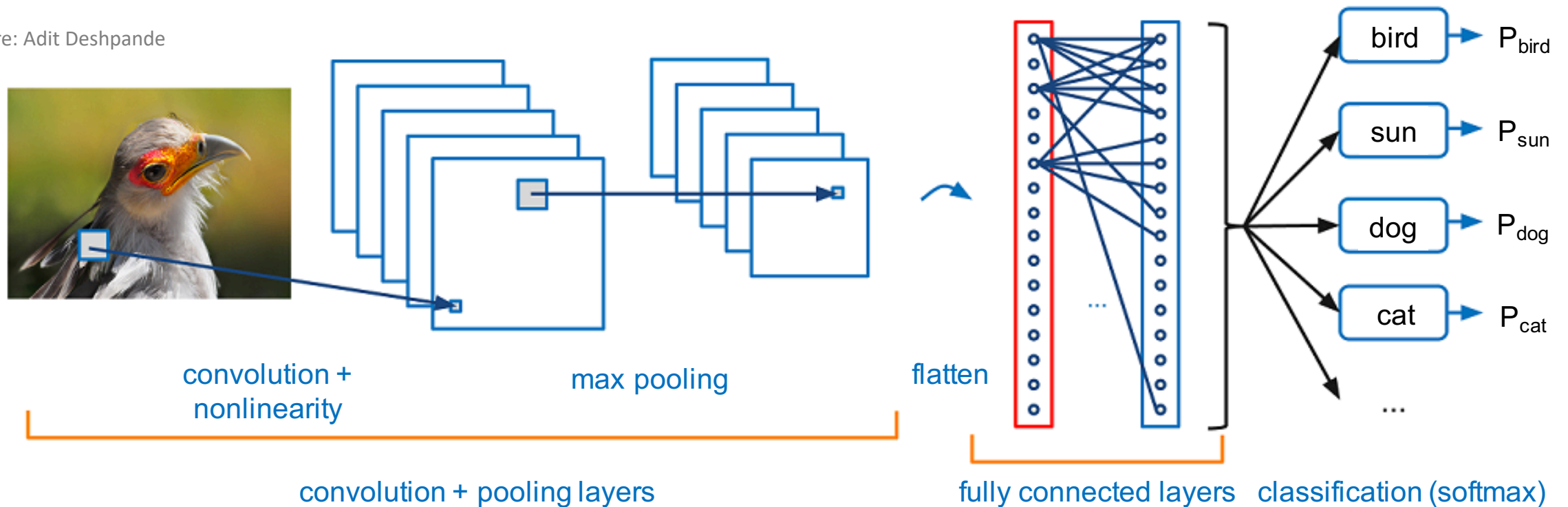
# 2016: deep learning with summary statistics



"Deep learning for population genetic inference", Sheehan and Song, *PLOS Comp Bio*, 2016

Sara Mathieson

# Outline

- Shift to machine learning in population genetics

- **Shift away from summary statistics to "raw" data**

- GANs and adversarial training
  - pg-gan algorithm for creating realistic simulated data

- Results on human data from Africa, Europe, and East Asia

Figure: Adit Deshpande



convolution + nonlinearity

max pooling

flatten

convolution + pooling layers

fully connected layers   classification (softmax)

bird   $P_{bird}$

sun   $P_{sun}$

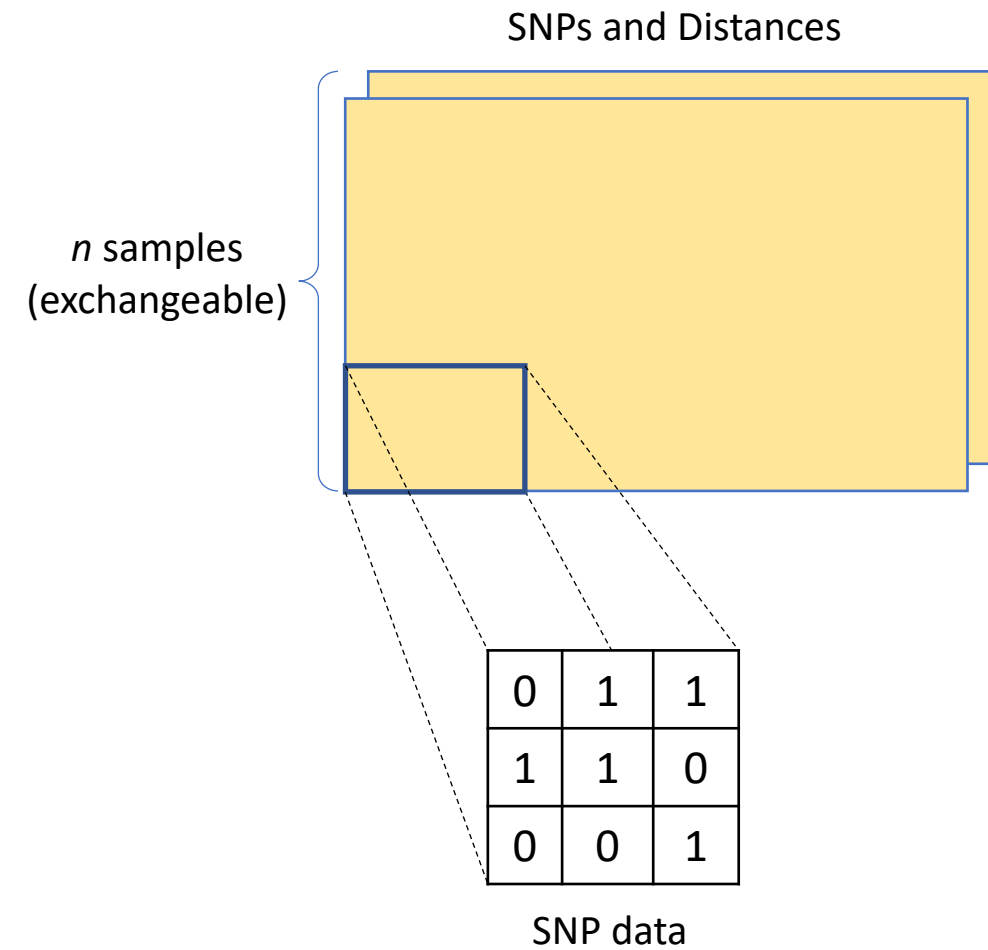dog   $P_{dog}$

cat   $P_{cat}$

...

## Issues

1. Image CNNs are optimized for different local features

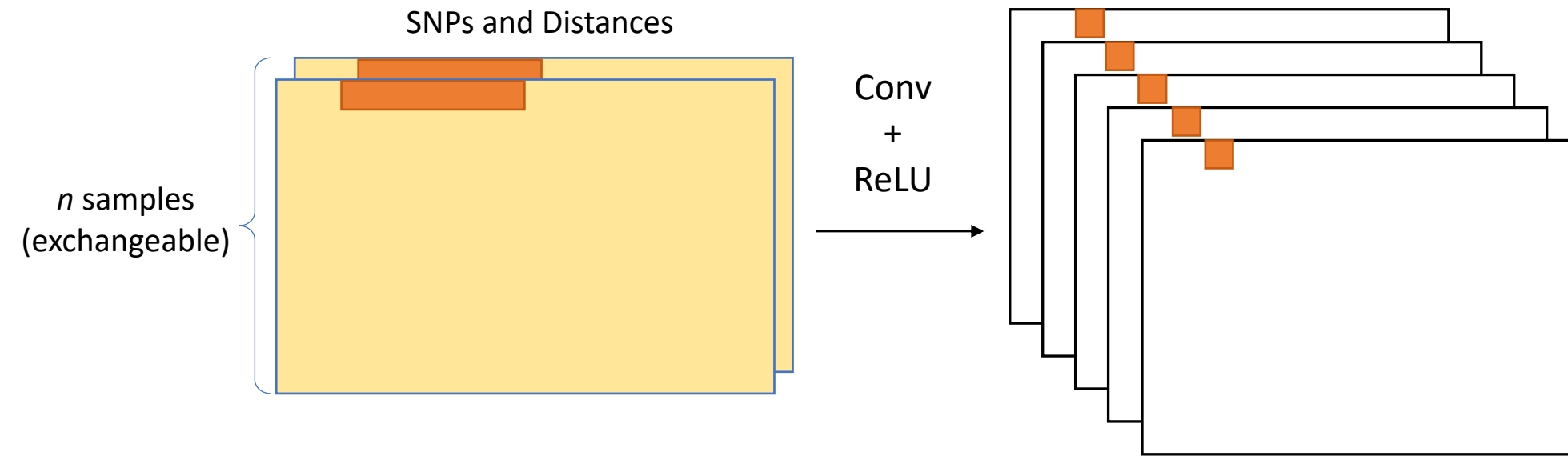2. For unstructured populations, sample (row) order doesn't matter

Flagel, Brandvain, Schrider. "The unreasonable effectiveness of convolutional neural networks in population genetic inference."
*Molecular biology and evolution,* 2018

Chan, Perrone, Spence, Jenkins, Mathieson, Song. "A Likelihood-Free Inference Framework for Population Genetic Data using Exchangeable Neural Networks"
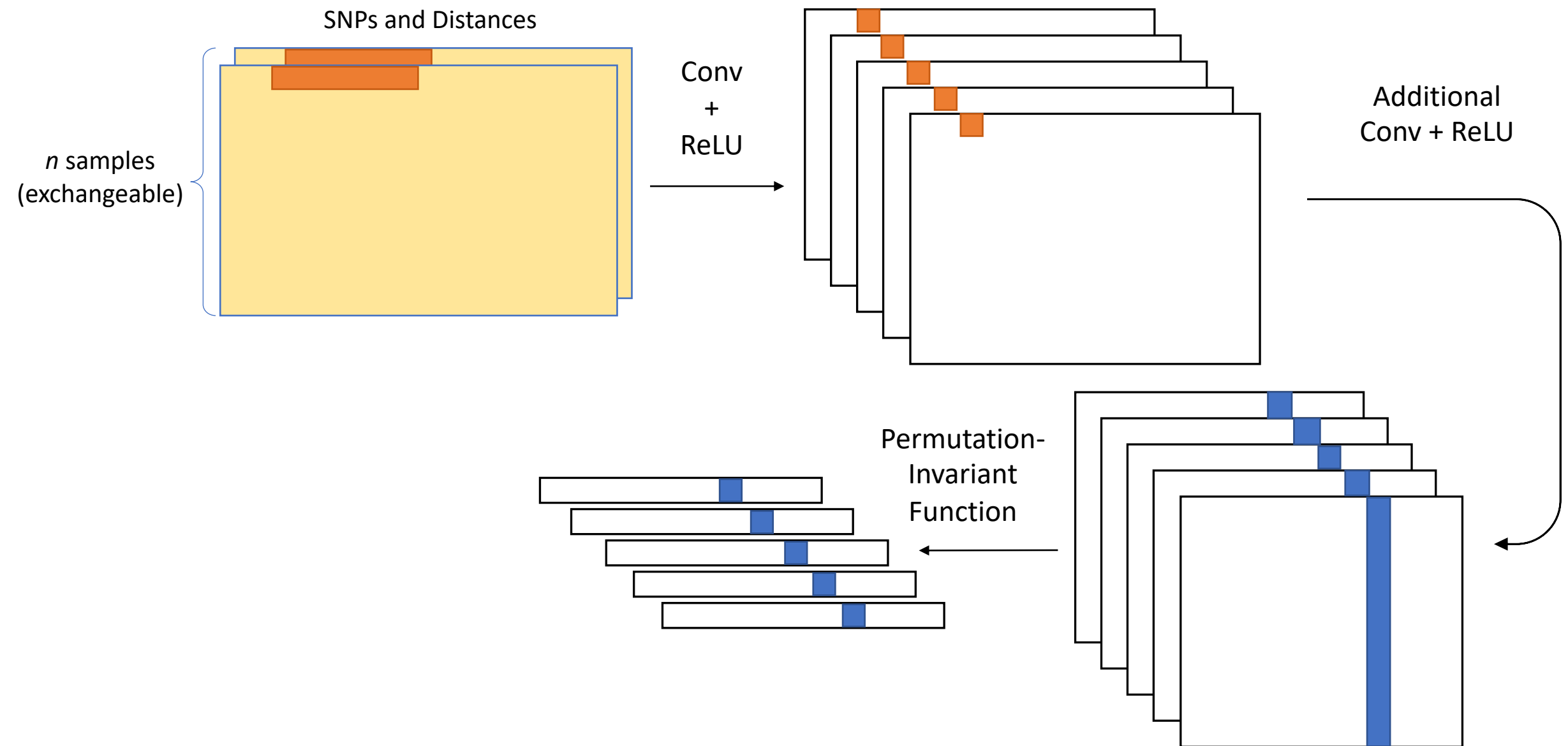*NeurIPS*, 2018, https://github.com/popgenmethods/defiNETti

SNPs and Distances

*n* samples
(exchangeable)

| 0 | 1 | 1 |
|---|---|---|
| 1 | 1 | 0 |
| 0 | 0 | 1 |

SNP data

SNPs and Distances

*n* samples
(exchangeable)

Conv
+
ReLU

# 2018: CNN for "raw" population genetic data

SNPs and Distances

*n* samples
(exchangeable)

Conv
+
ReLU

Additional
Conv + ReLU

Permutation-
Invariant
Function



Sara Mathieson

SNPs and Distances

*n* samples (exchangeable)

Conv + ReLU

Additional Conv + ReLU

Permutation-Invariant Function

Flatten + Fully Connected

Additional FC layers

Output: evolutionary parameter of interest

# Impact of permutation-invariant architecture (recombination hotspots)
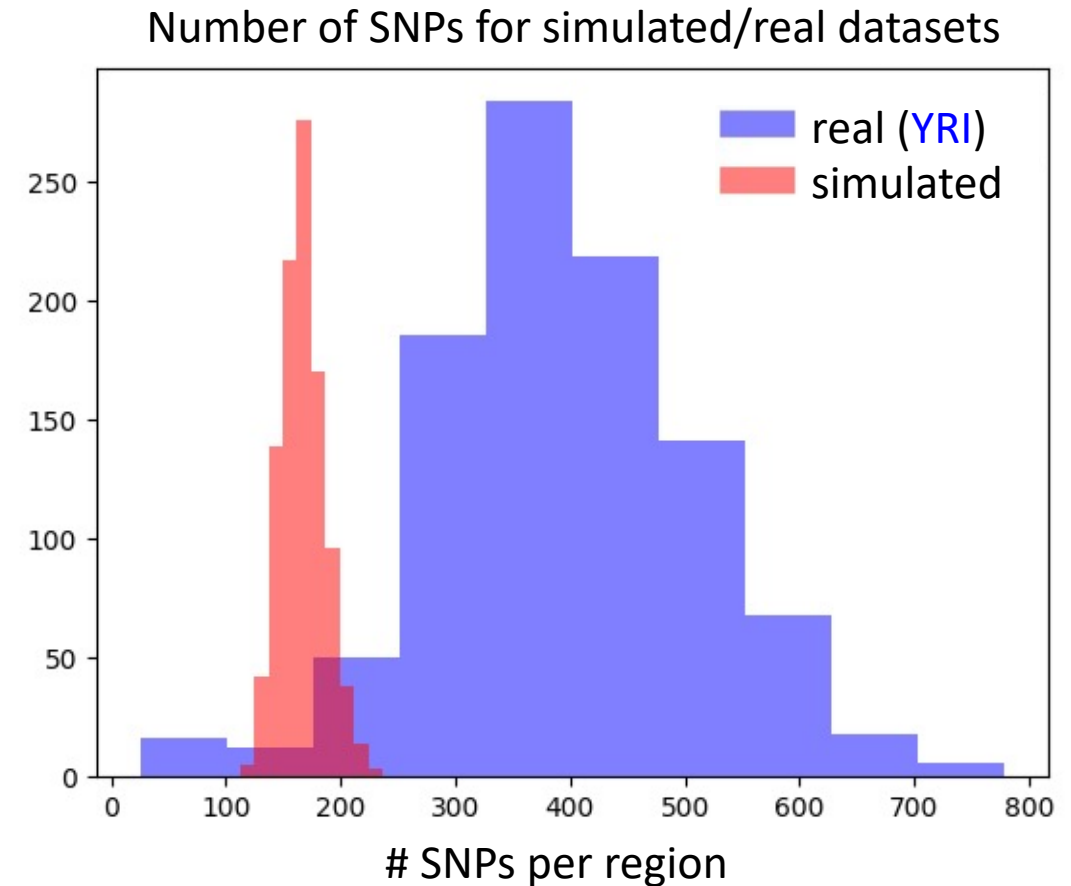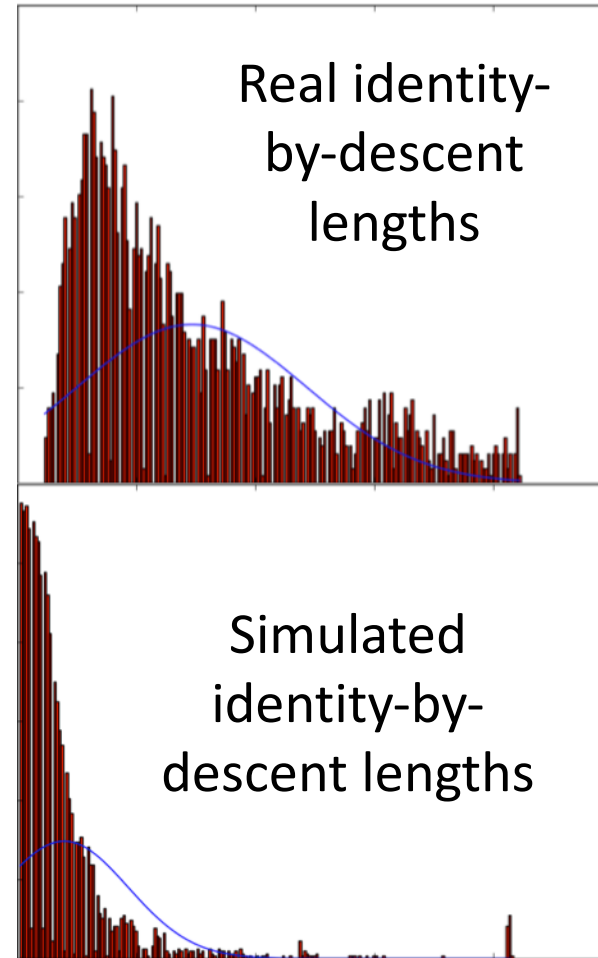


Testing Accuracy for Human Recombination Maps

# Outline

- Shift to machine learning in population genetics

- Shift away from summary statistics to "raw" data

- GANs and adversarial training
  - pg-gan algorithm for creating realistic simulated data

- Results on human data from Africa, Europe, and East Asia

# Even using good simulation programs, it is difficult to match real data

**High-quality simulated data is crucial!**

- Develop intuition

- Validate methods

- Provide training data for machine learning methods

- Popular simulators: SLiM, msprime

Real identity-by-descent lengths

Simulated identity-by-descent lengths

Number of SNPs for simulated/real datasets

real (YRI)
simulated

# SNPs per region

YRI: Yoruba in Ibadan, Nigeria

# Idea behind GANs (Generative Adversarial Networks)



Which is "real" and
which is "fake"?

*Centre de Estudios Borjanos/AFP/Getty Images*

# Idea behind GANs (Generative Adversarial Networks)
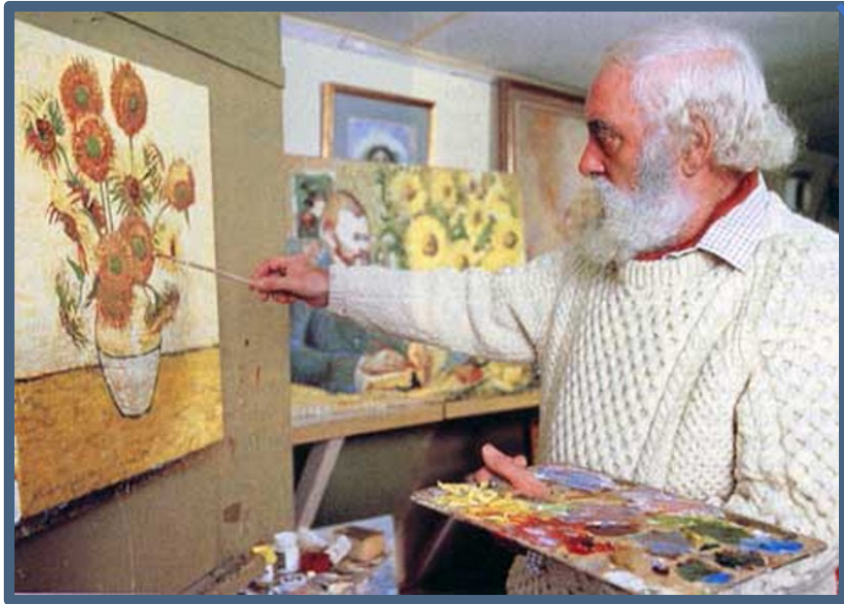


Which is "real" and which is "fake"?

https://webartacademy.com/fake-picasso

# Idea behind GANs (Generative Adversarial Networks)



Which is "real" and which is "fake"?



Photo: Courtesy International Foundation for Art Research (IFAR).

Sara Mathieson

# Idea behind GANs (Generative Adversarial Networks)



Latent random variable

Generator

Real world images

Sample

Sample

Discriminator

Generator ("forger") tries to create realistic artwork

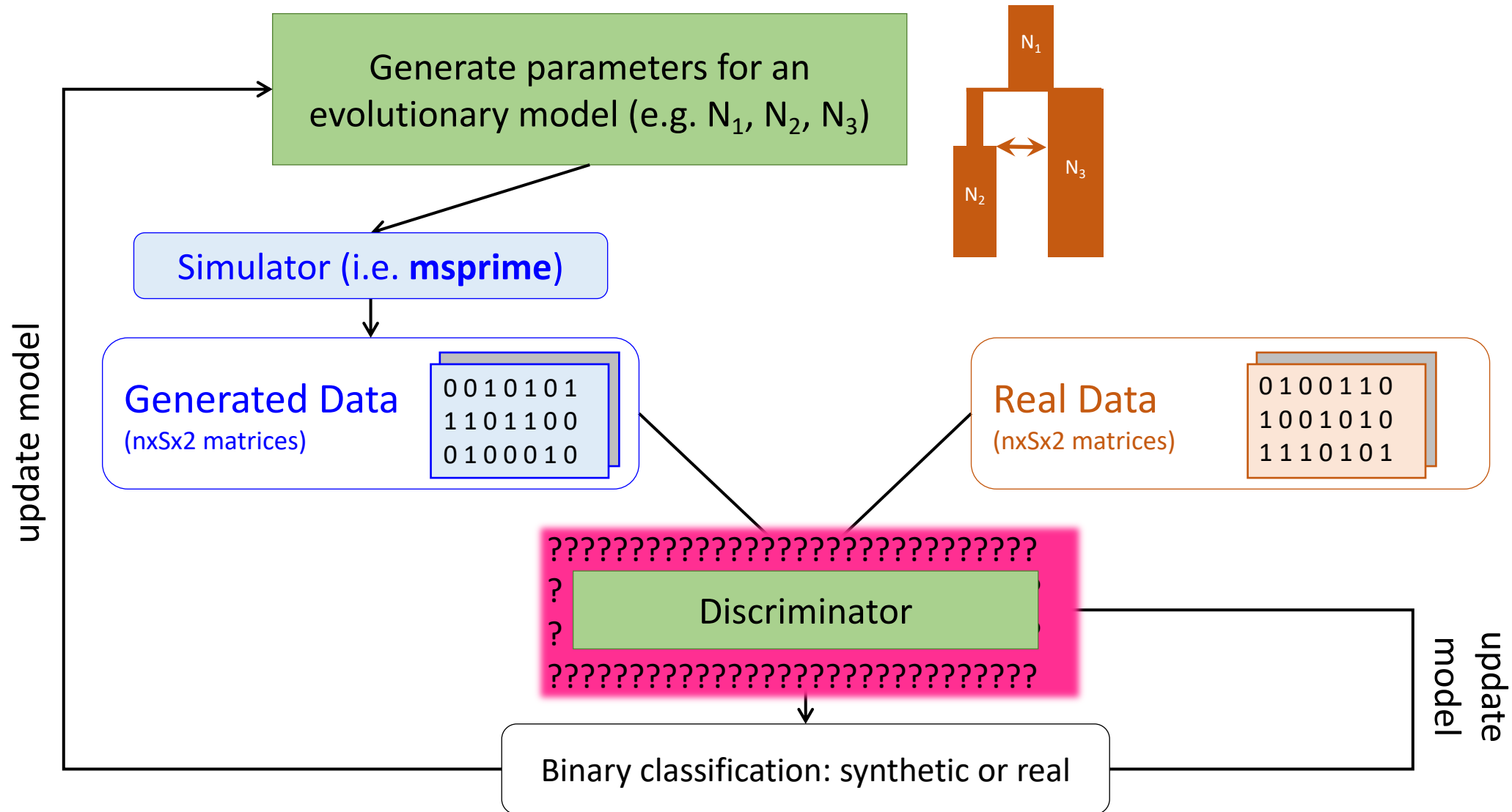Discriminator ("art critic") tries to identify real vs. fake

Fake 🔴🟢 Real

Loss

People Images / Getty Images

GAN diagram: Adapted from Kevin McGuinness

(Source: Lost in the Louvre)

Sara Mathieson

# pg-gan algorithm overview

# pg–gan algorithm overview

????????????????????????????????????????????
Generate parameters for an evolutionary model (e.g. $N_1$, $N_2$, $N_3$)
????????????????????????????????????????????

$N_1$

$N_3$

$N_2$

Simulated annealing algorithm

Temperature cools linearly

Change one parameter each iteration

Simulator (i.e. **msprime**)

update model

Generated Data
(nxSx2 matrices)

```
0 0 1 0 1 0 1
1 1 0 1 1 0 0
0 1 0 0 0 1 0
```

Real Data
(nxSx2 matrices)

```
0 1 0 0 1 1 0
1 0 0 1 0 1 0
1 1 1 0 1 0 1
```

Discriminator

update model

Binary classification: synthetic or real

# pg-gan algorithm overview

Generate parameters for an evolutionary model (e.g. $N_1$, $N_2$, $N_3$)

$N_1$
$N_2$
$N_3$

Simulator (i.e. **msprime**)

Generated Data
(nxSx2 matrices)

0010101
1101100
0100010

Real Data
(nxSx2 matrices)

0100110
1001010
1110101

update model

???????????????????????????????
?
?
????????????????????????????????

Discriminator

update model

Binary classification: synthetic or real

# pg-gan discriminator architecture: extend to multiple populations



YRI: Yoruba in Ibadan, Nigeria

CEU: Utah residents with Northern and Western European ancestry

# Example of failed GAN training



Discriminator classifies
everything as real

Generator cannot
learn and reduce loss

# Example of successful GAN training



Generator not fooling discriminator

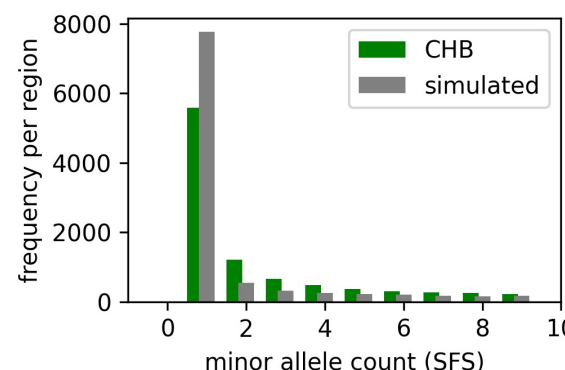Generator and discriminator are balanced
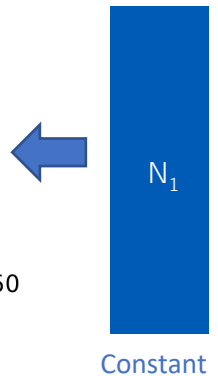
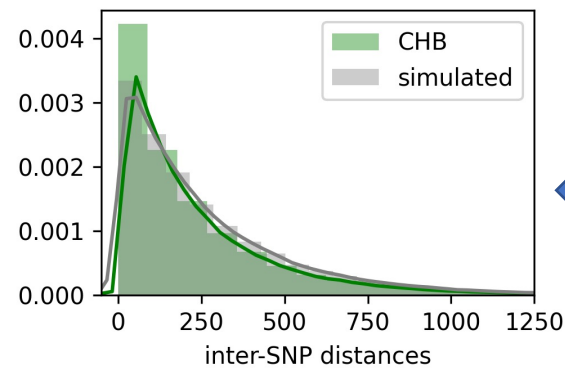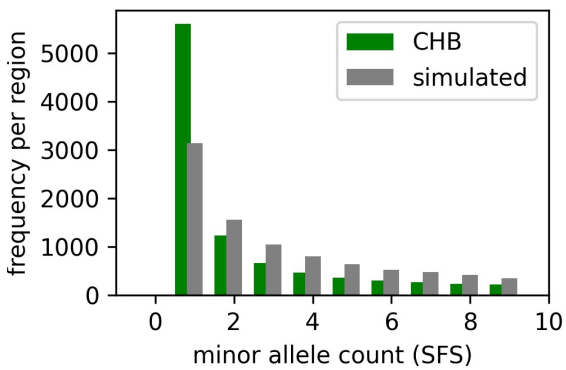Discriminator easily able to tell training from simulated
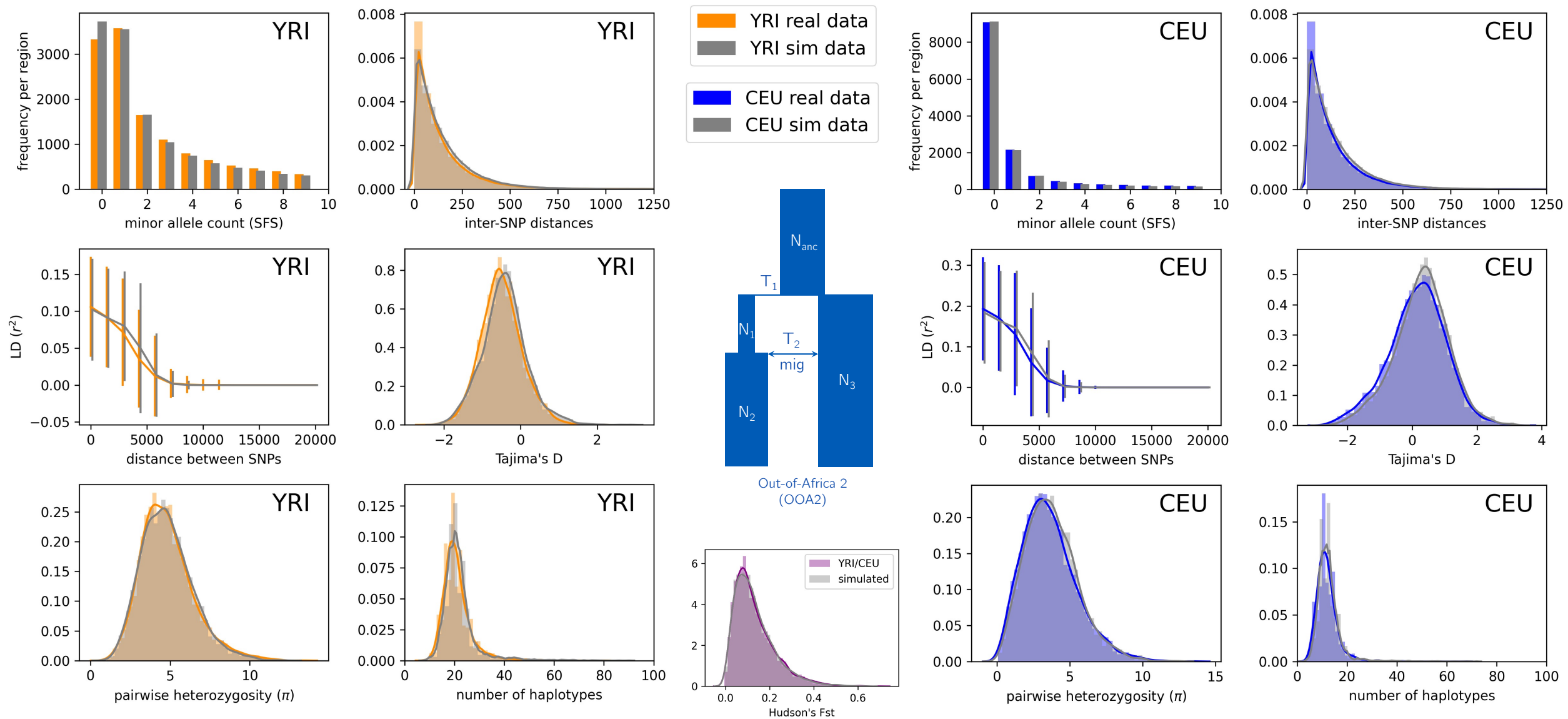
Discriminator is often confused

# Outline

- Shift to machine learning in population genetics

- Shift away from summary statistics to "raw" data

- GANs and adversarial training
  - pg-gan algorithm for creating realistic simulated data

- **Results on human data from Africa, Europe, and East Asia**

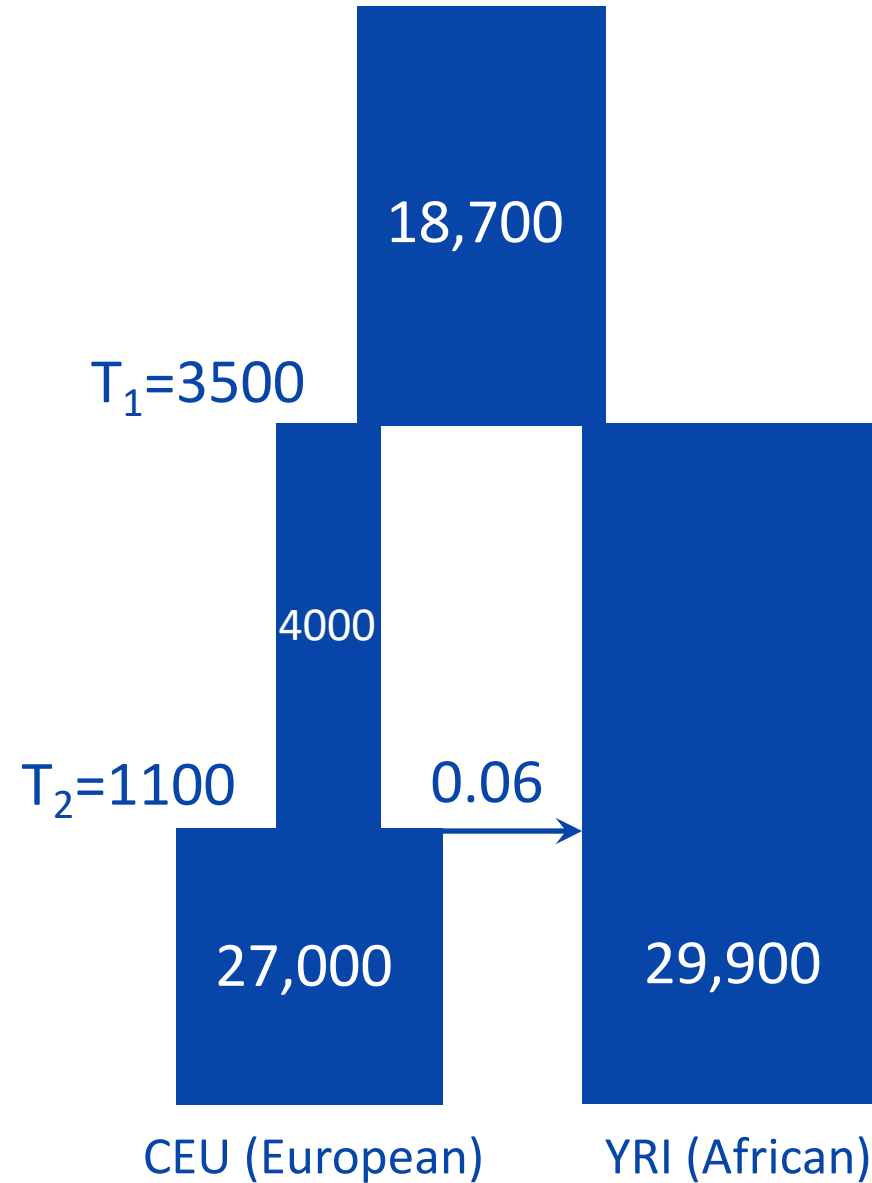# CHB: 1-param model

# CHB: 5:param model

# Simulated data under our GAN-inferred model matches real data

# YRI/CEU split inference

- Time measured in generations

- Out-of-African bottleneck apparent



$T_1=3500$

18,700

4000

$T_2=1100$    0.06

27,000    29,900

CEU (European)    YRI (African)

# Conclusion for Machine Learning in Population Genetics
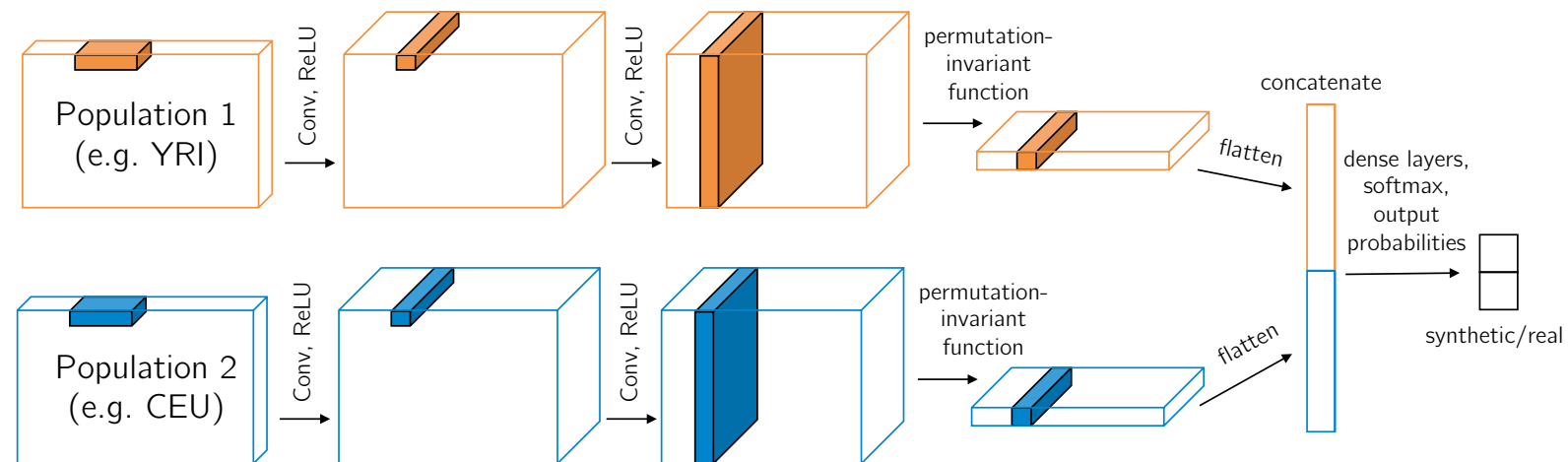
**Future directions for** pg-gan

- Apply to understudied populations

- Overcome data imbalance

**Where are we going?**

- Keep the data in mind

- ML methods need to be more interpretable

- Combine ML with evolutionary modeling

- Unsupervised learning

# Thank you!

- Jeffrey Chan
- Nhung Hoang
- Paul Jenkins
- Michael Kourakos
- Hunter Lee
- Iain Mathieson
- Valerio Perrone
- Yun S. Song
- Jeffrey Spence
- Zhanpeng Wang
- Jiaping Wang

# END