



Putting the HAL in Haldane: leveraging deep learning for population genetics

Andrew Kern
Institute of Ecology and Evolution
University of Oregon

Leo Breiman's Two Cultures

the logic of data analysis



Leo Breiman's Two Cultures

Data Modeling Culture



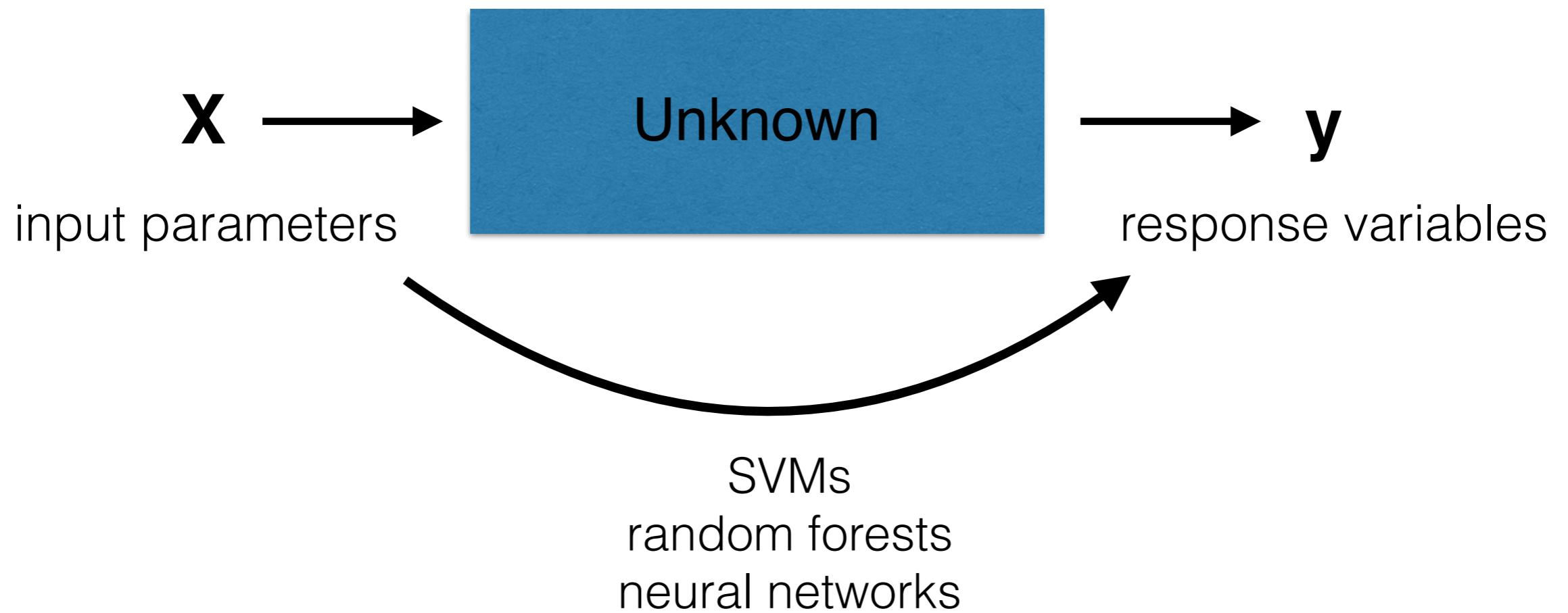
e.g. linear regression

Focus on stochastic model to explain
how $f(x) \rightarrow y$

98% of Statistics

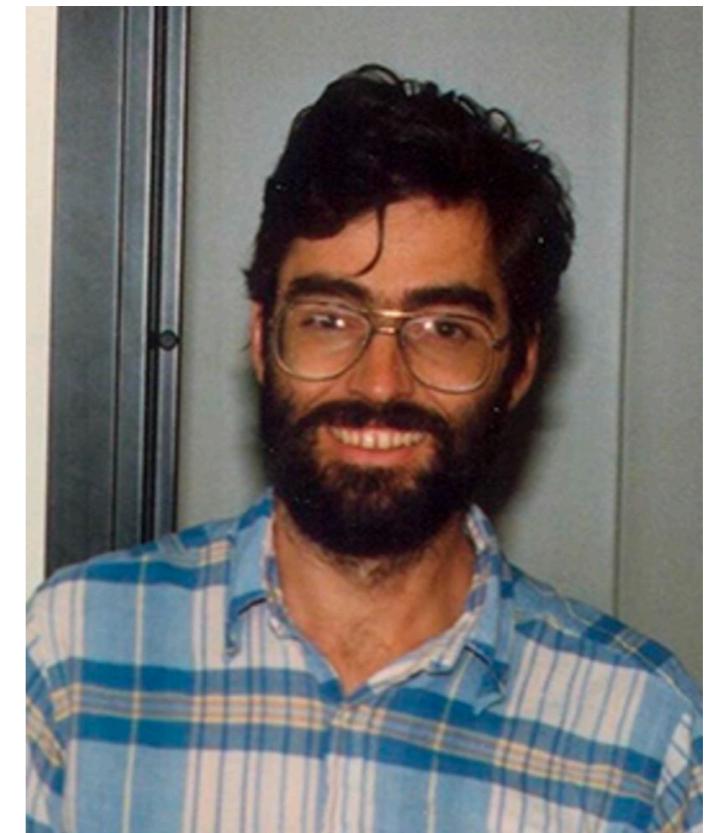
Leo Breiman's Two Cultures

Algorithmic Modeling Culture
(machine learning)



Ignore probabilistic generative model $f(x) \rightarrow y$

Analysis in Population Genetics? One culture!



Major focus on stochastic models, param estimation from those models

Analysis in Population Genetics? One culture!

A Genomic Map of the Effects of Linked Selection in Drosophila

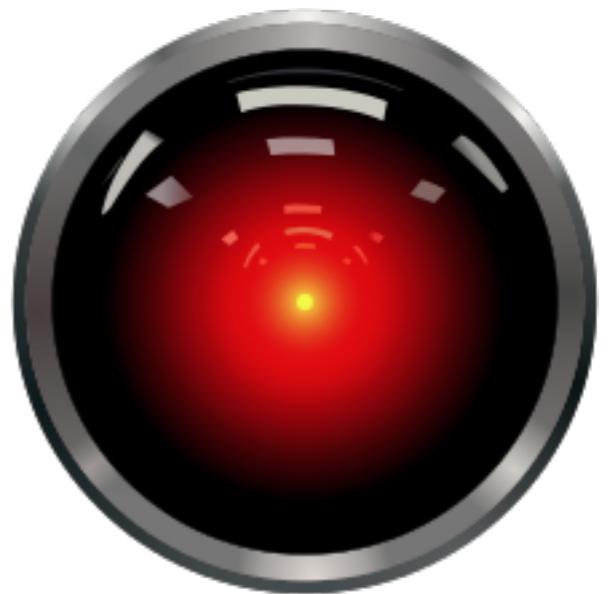
Eyal Elyashiv , Shmuel Sattath, Tina T. Hu, Alon Strutsovsky, Graham McVicker, Peter Andolfatto, Graham Coop, Guy Sella 

$$\pi(x) = \frac{\pi_0 \cdot (u(x)/\bar{u})}{\pi_0 \cdot (u(x)/\bar{u}) + 1/B(x) + S(x; \bar{N}_e, T)},$$

Model	Background selection and classic sweeps with added "missing substitutions"	Background selection and classic sweeps without "missing substitutions"
ΔCL	3.9×10^{-4}	3.8×10^{-4}
R^2		
1 Mb	0.71	0.70
100 kb	0.44	0.44
10 kb	0.26	0.25
1 kb	0.20	0.20

Deep learning has a lot to offer to popgen

- variable level of model dependance
- Powerful simulation-based inference
- automated feature extraction



Outline

- RNNs for estimating recombination model-based
- DNNs for estimating geographic origin model-free
- CNNs for estimating spatial params model-based

Jeff Adrion



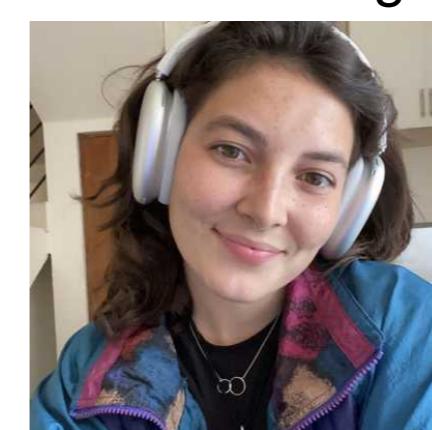
CJ Battey



Clara Rehmann



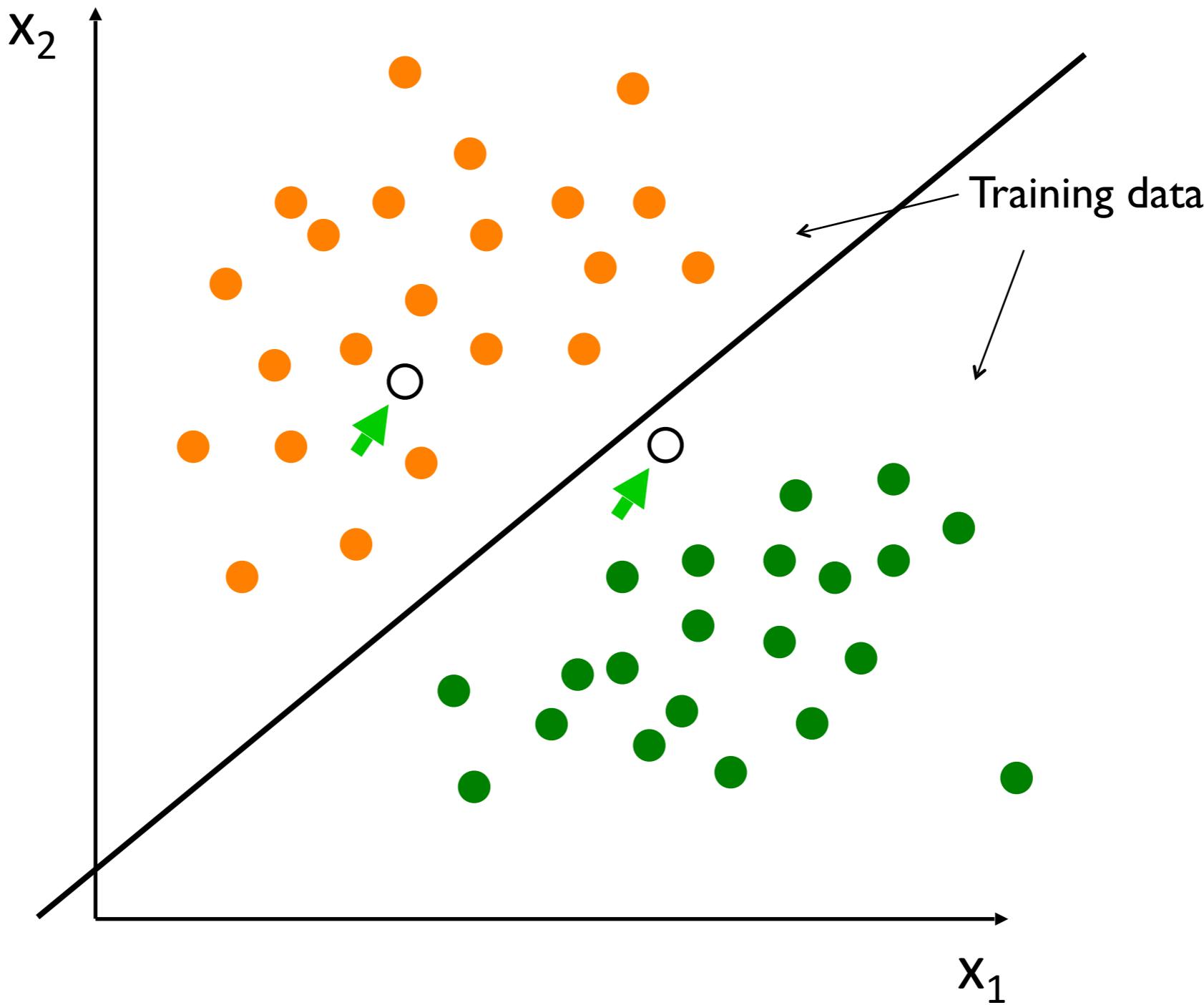
Jordan Rodriguez



Chris Smith

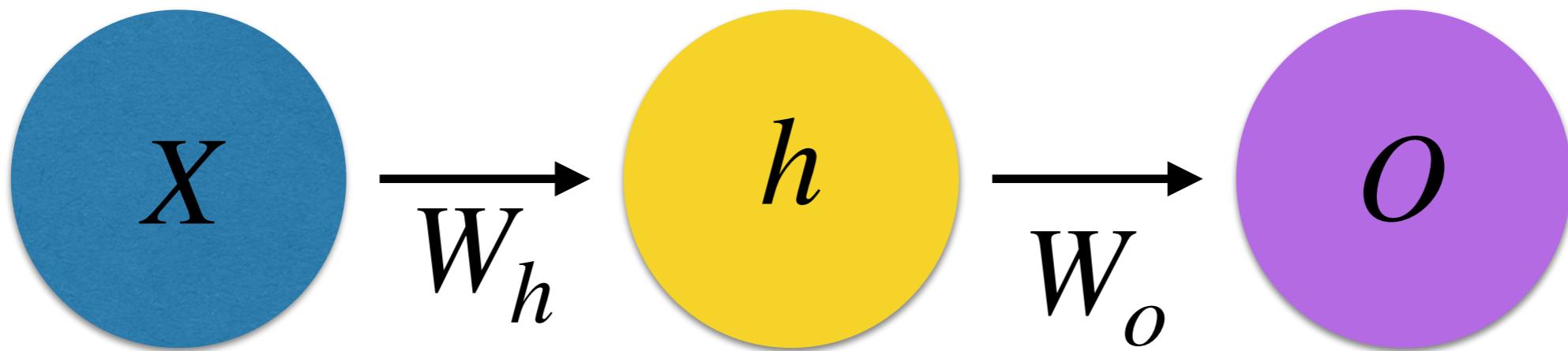


Supervised Machine Learning



toy neural network

Input Layer Hidden Layer Output Layer

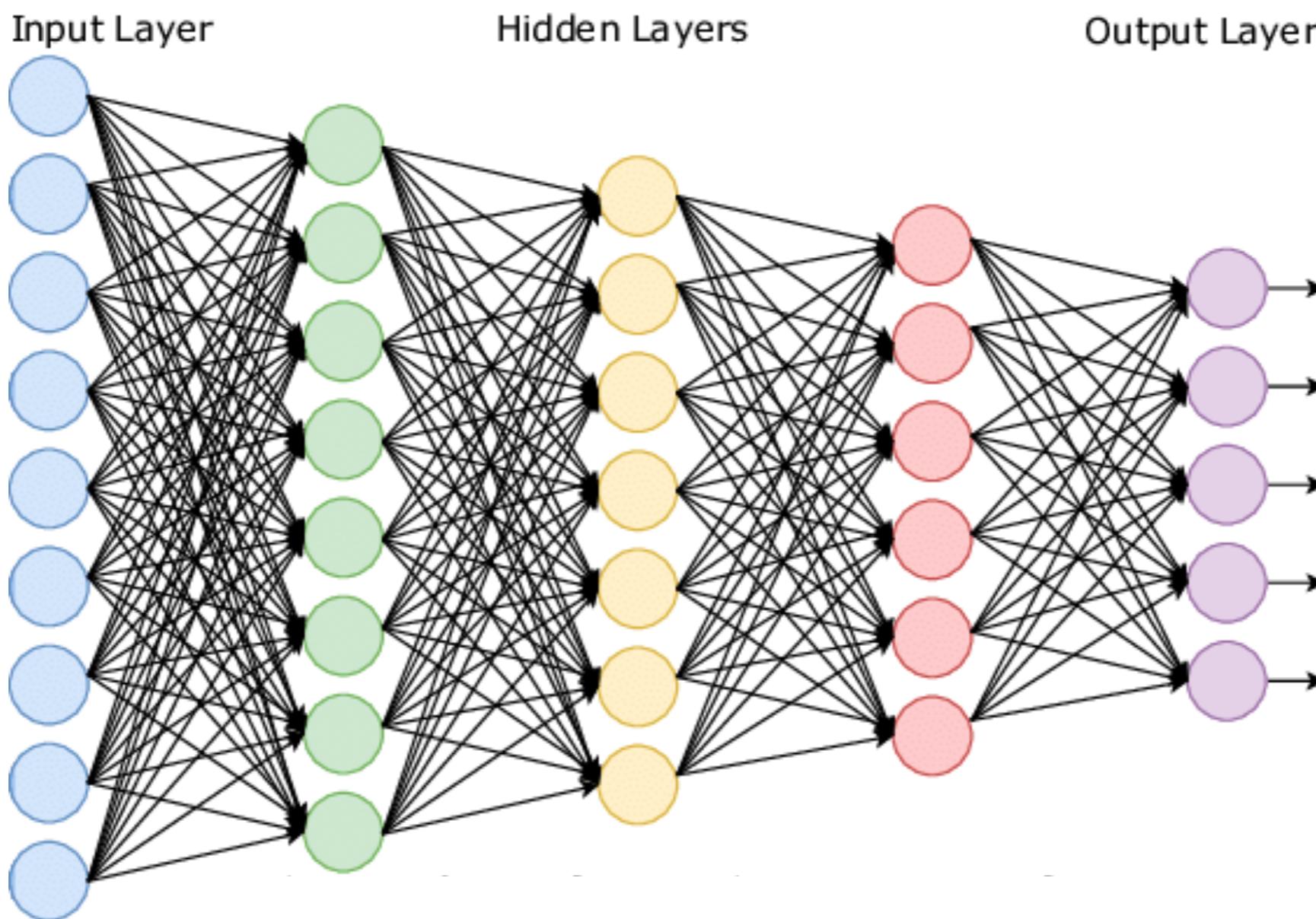


Feed-forward

$$O = f(f(X \cdot W_h + b_h) \cdot W_o + b_o)$$

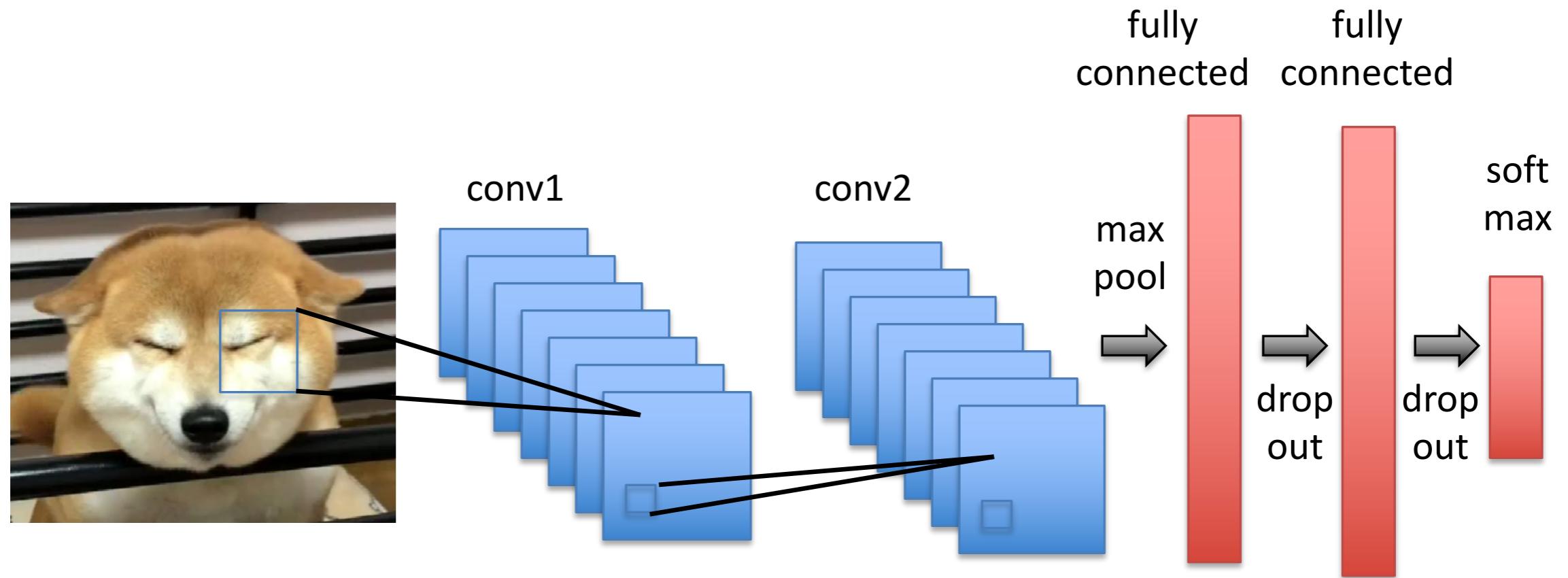
Linear algebra, no magic

Deep Neural Networks



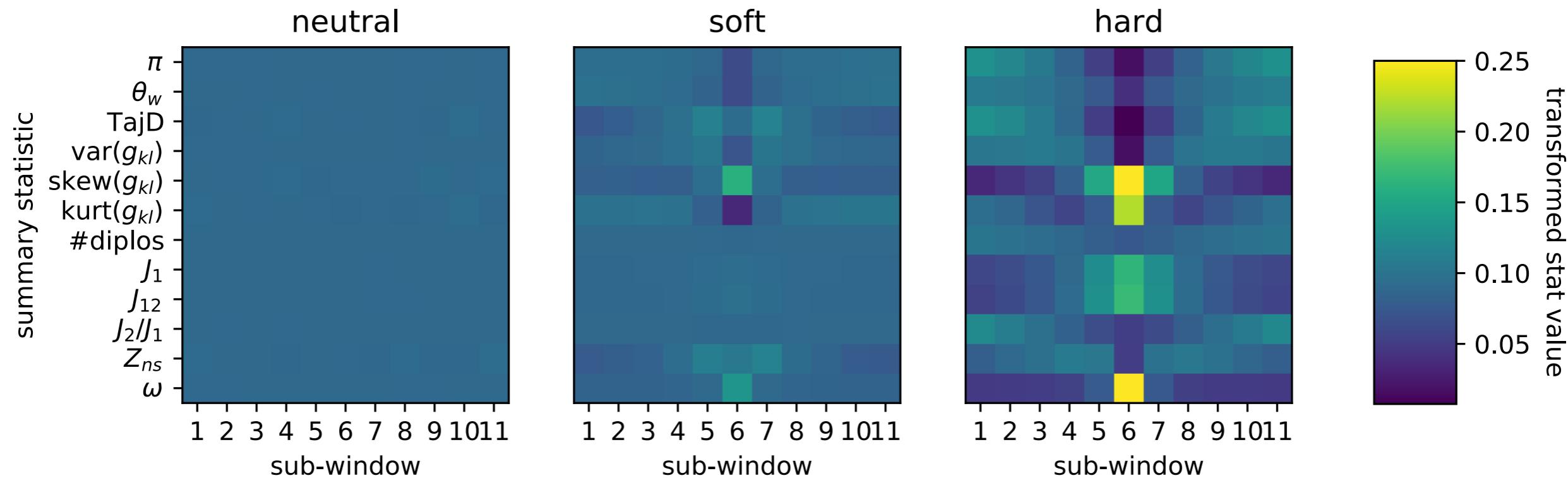
“multilayer perceptron”

Convolutional Neural Network



Slide “receptive filter” along image, take dot product, repeat

popgen as image recognition



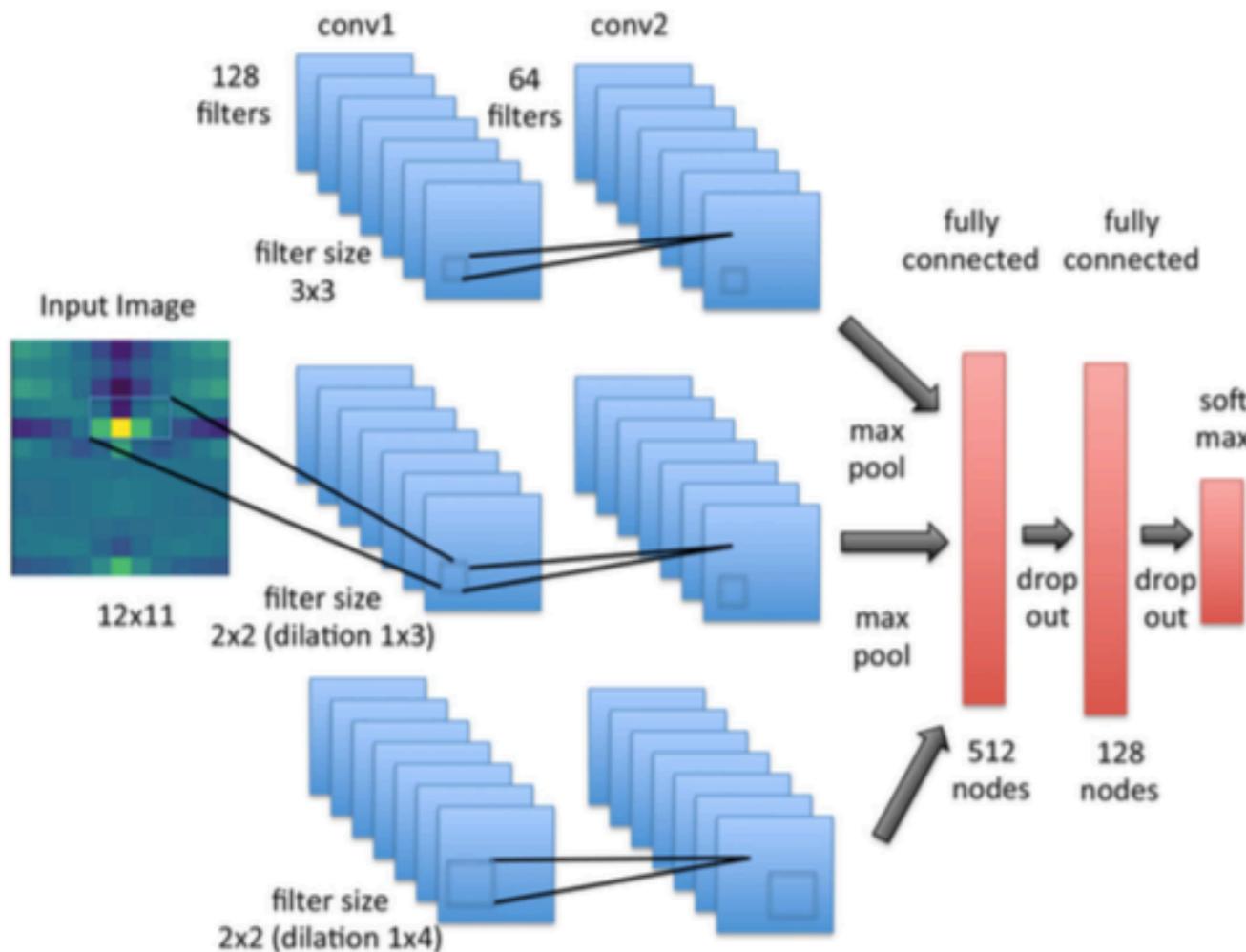
diploS/HIC: An Updated Approach to Classifying Selective Sweeps

Andrew D. Kern¹ and Daniel R. Schrider

Department of Genetics, Rutgers University, Piscataway, NJ 08854

ORCID IDs: 0000-0003-4381-4680 (A.D.K.); 0000-0001-5249-4151 (D.R.S.)

diploS/HIC



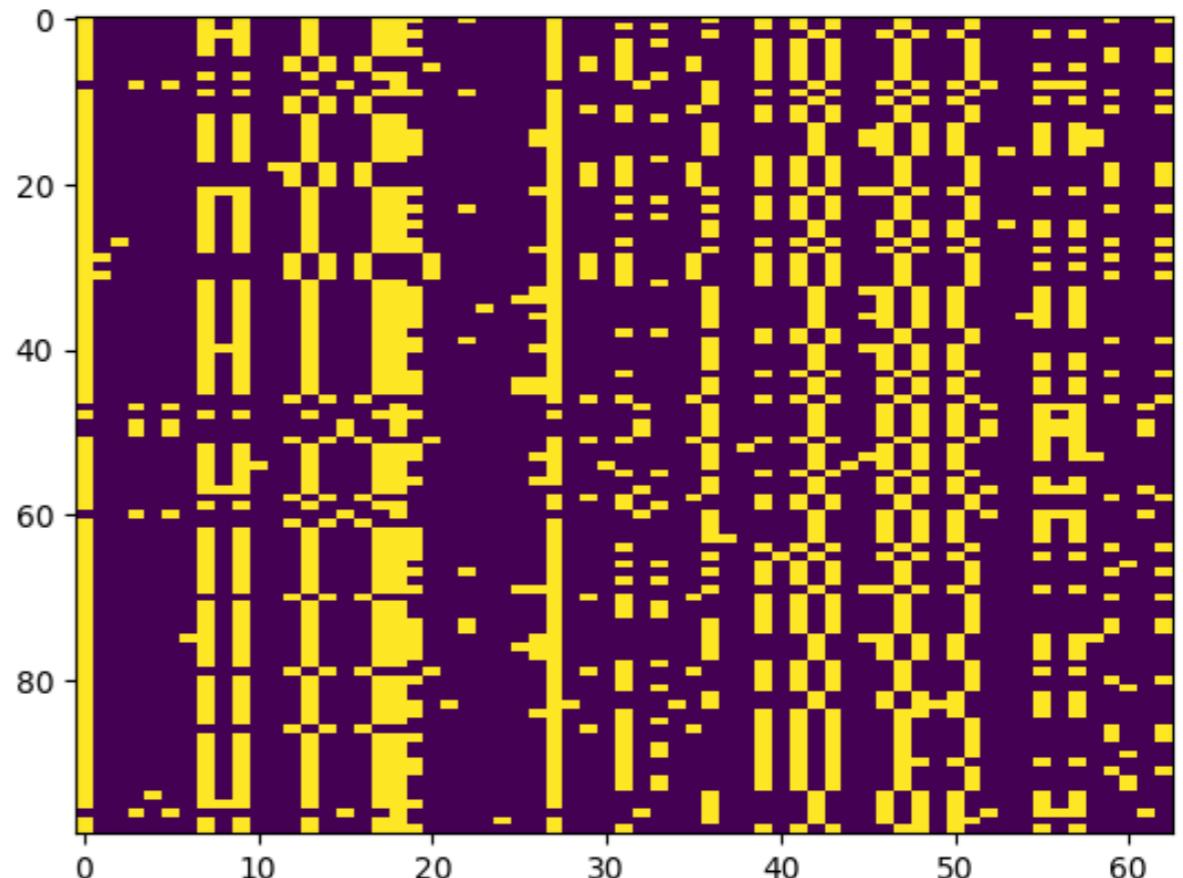
diploS/HIC: An Updated Approach to Classifying Selective Sweeps

Andrew D. Kern¹ and Daniel R. Schrider
Department of Genetics, Rutgers University, Piscataway, NJ 08854
ORCID IDs: 0000-0003-4381-4680 (A.D.K.); 0000-0001-5249-4151 (D.R.S.)

using CNNs provides added accuracy & robustness

Alignment “images” for Deep Learning

- Automated feature extraction
- Using convolutional neural nets (deep learning) quickly have gained new traction on many problems



Dan Schrider (now UNC)



A Likelihood-Free Inference Framework for Population Genetic Data using Exchangeable Neural Networks

Jeffrey Chan, Valerio Perrone, Jeffrey P. Spence, Paul A. Jenkins, Sara Mathieson, Yun S. Song

(Submitted on 16 Feb 2018)

The Unreasonable Effectiveness of Convolutional Neural Networks in Population Genetic Inference

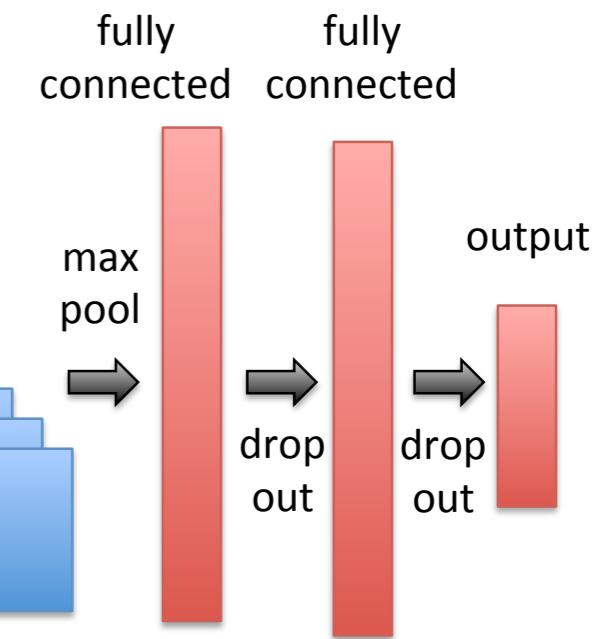
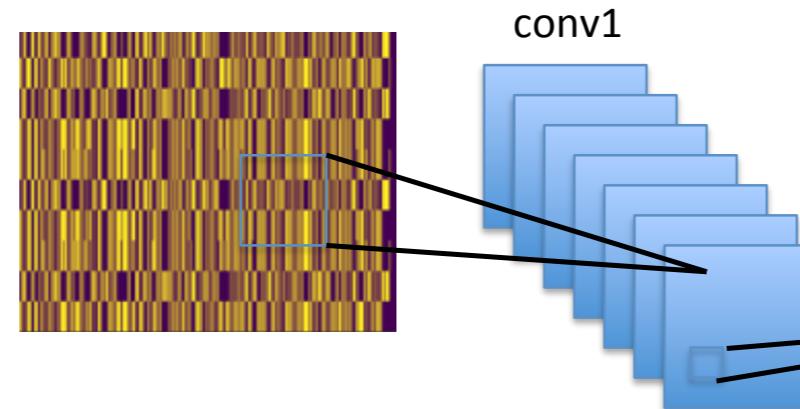
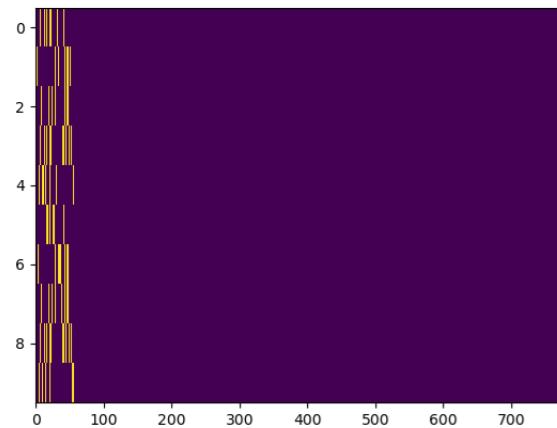
Lex Flagel, Yaniv J Brandvain, Daniel R Schrider

doi: <https://doi.org/10.1101/336073>

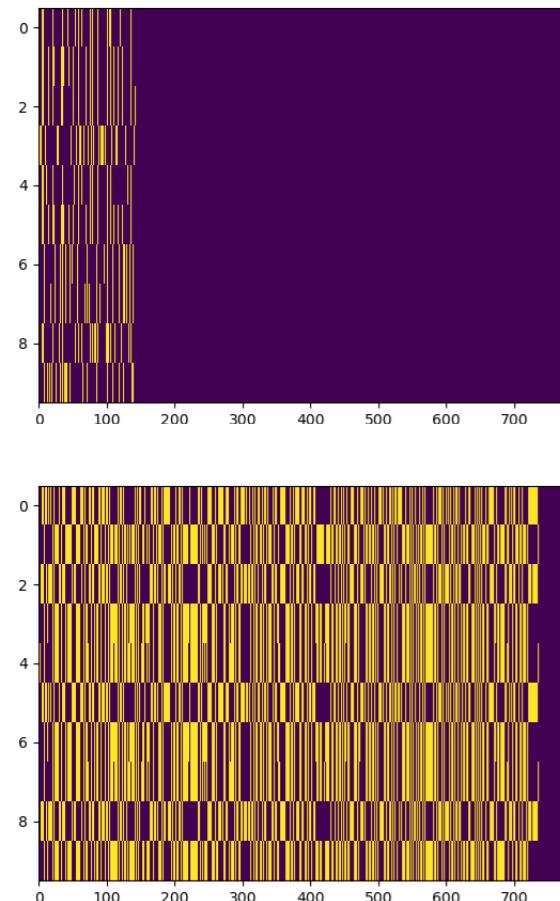
Convolutional Neural Nets

Population genetic parameter estimation: Θ

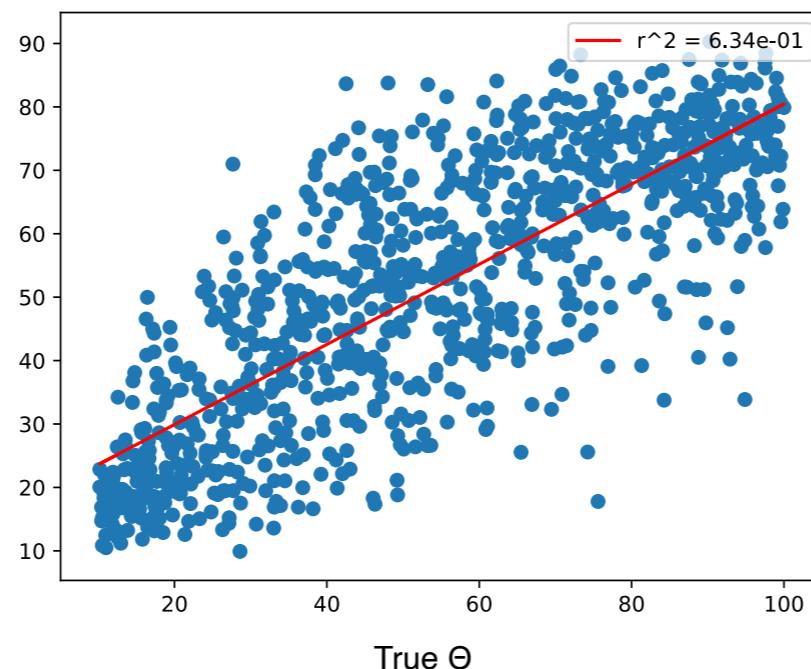
Alignment images



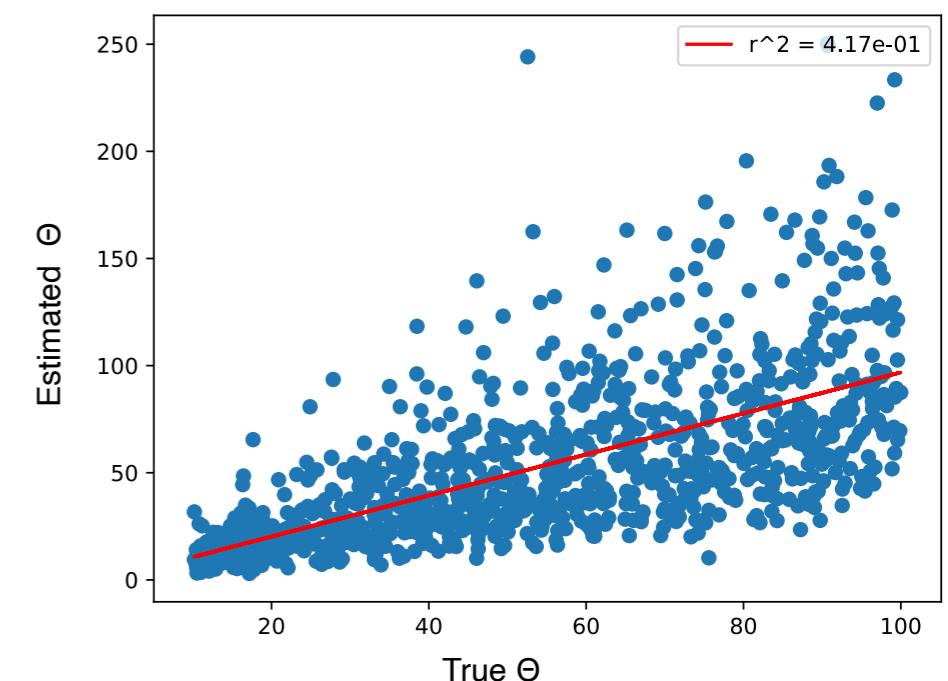
CNN estimator



Estimated Θ



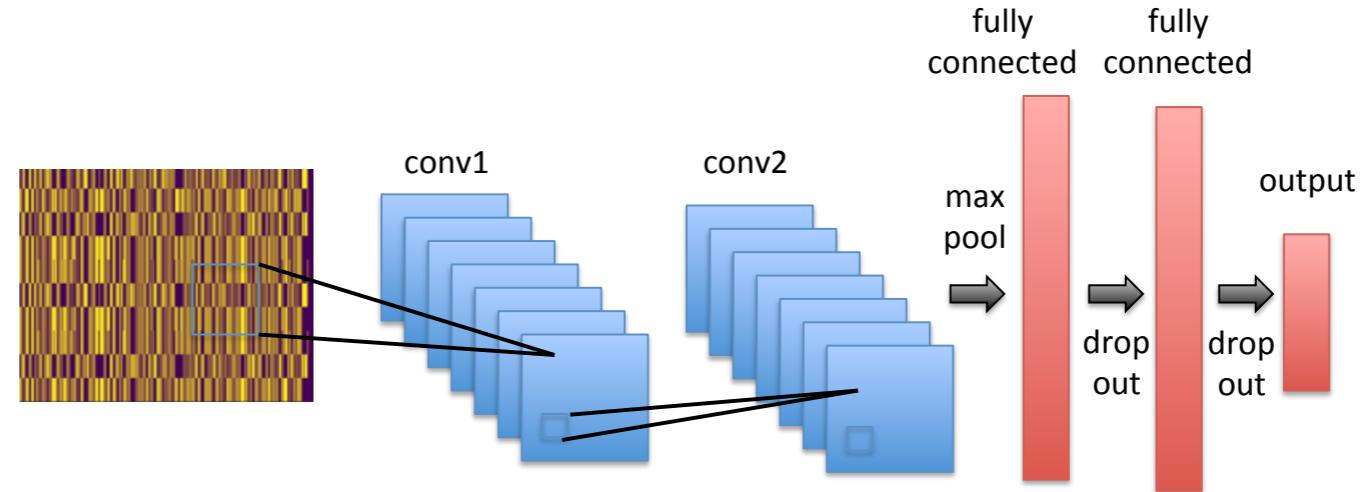
Tajima's π estimator



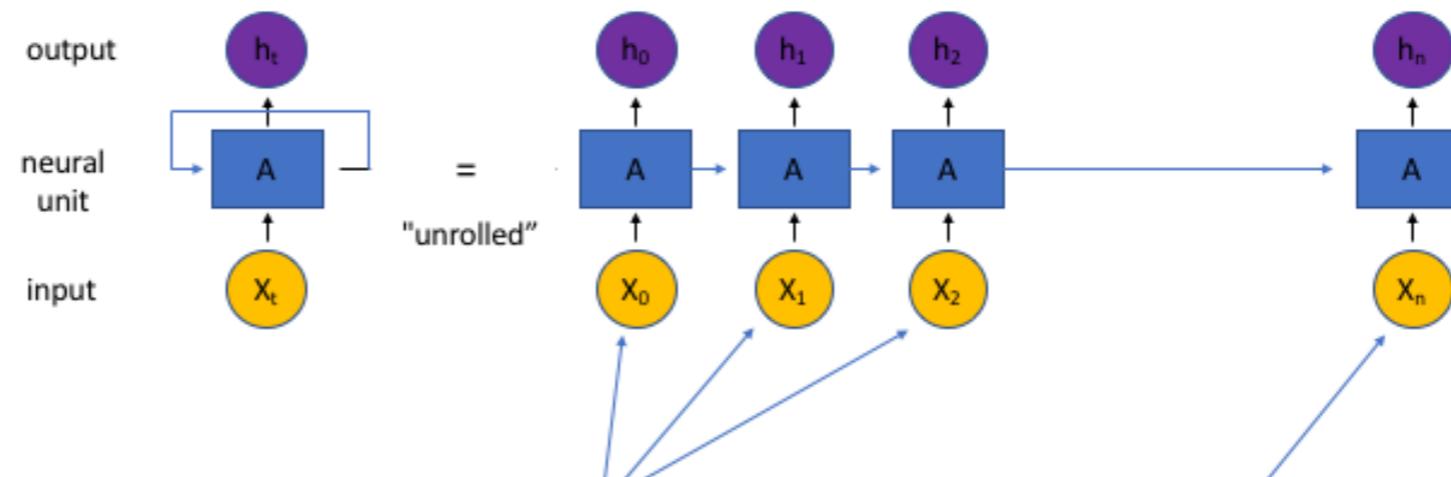
Deep Neural Nets

What is the best architecture to use?

CNNs?



RNNs?



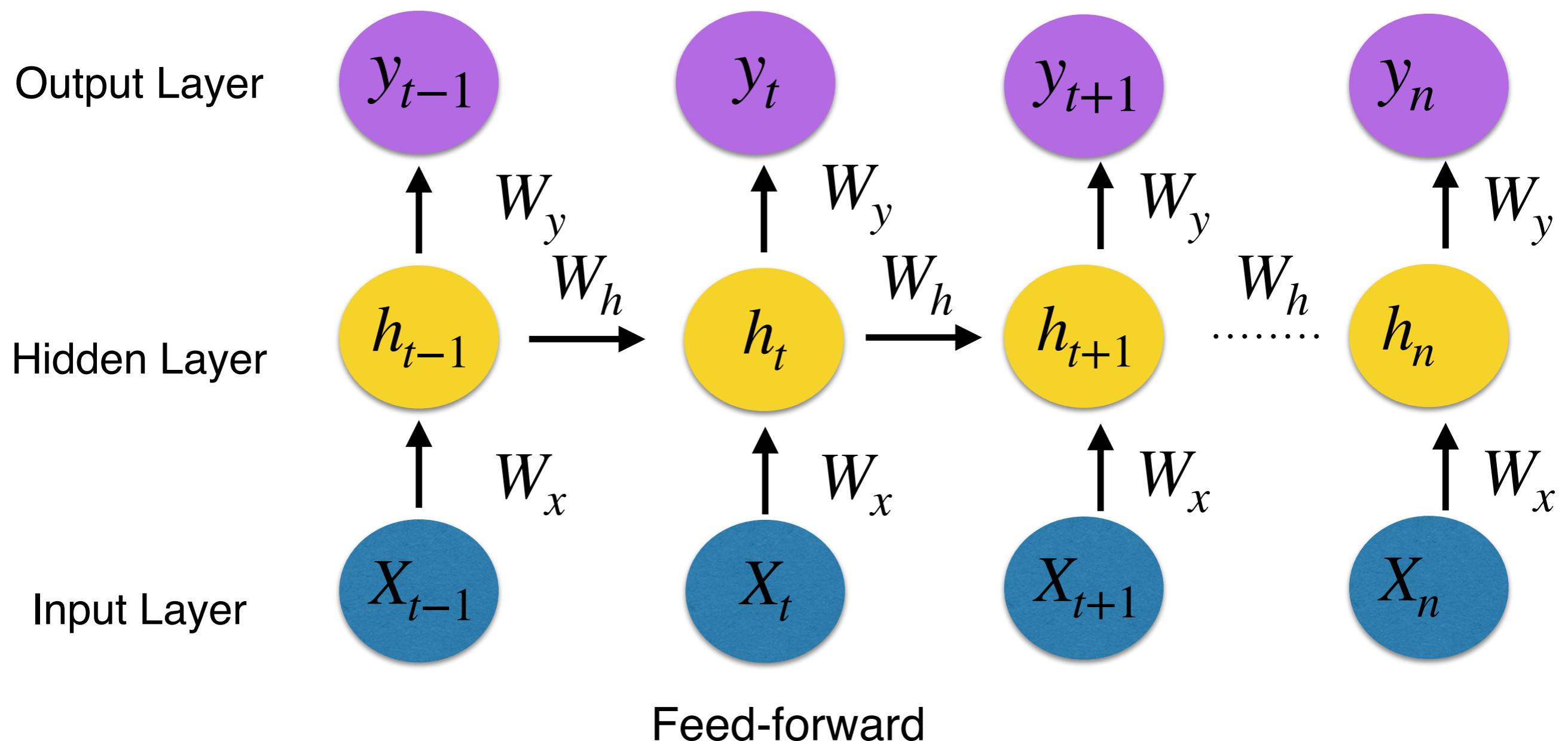
Sequence alignment

```
1001100111111100010010010001001000001010  
1001100111111100010010010001001000001010  
10001001011110100000010000011001100101010  
11001001011110100100010100001011010001011  
00000110000000000000001101000000000000010100  
00100110000000001100110100000100000010100  
0000011000000000000000101001100000001010100  
0000011000000000000000101000000001010100  
0000011000000000000000101000000000000001010100
```

Jeff Adriion



recurrent neural network



$$h_t = f(W_h \cdot h_{t-1} + W_x \cdot x_t)$$

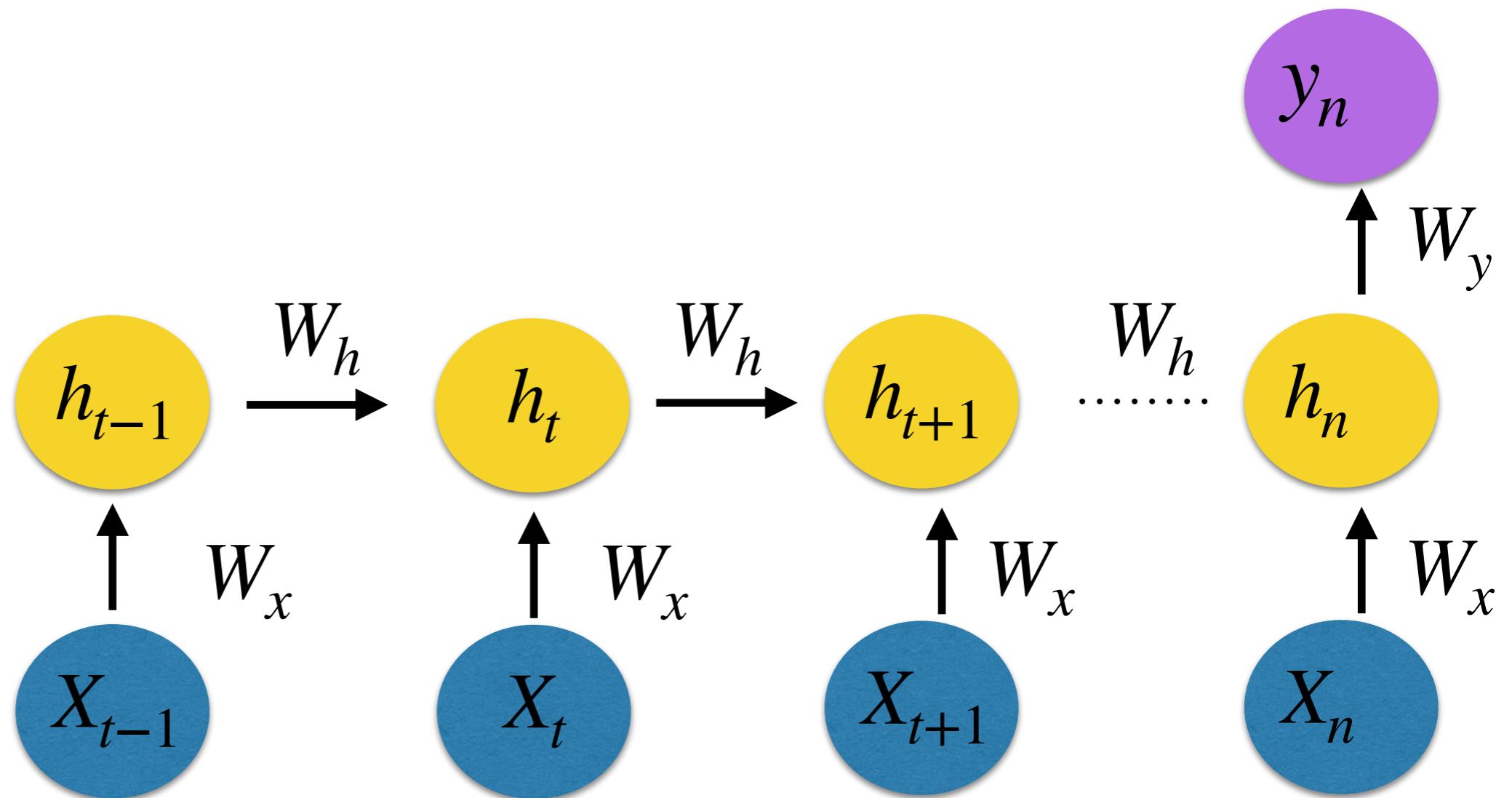
$$\hat{y}_t = \text{softmax}(W_y \cdot h_t)$$

recurrent neural network

Output Layer

Hidden Layer

Input Layer



Feed-forward

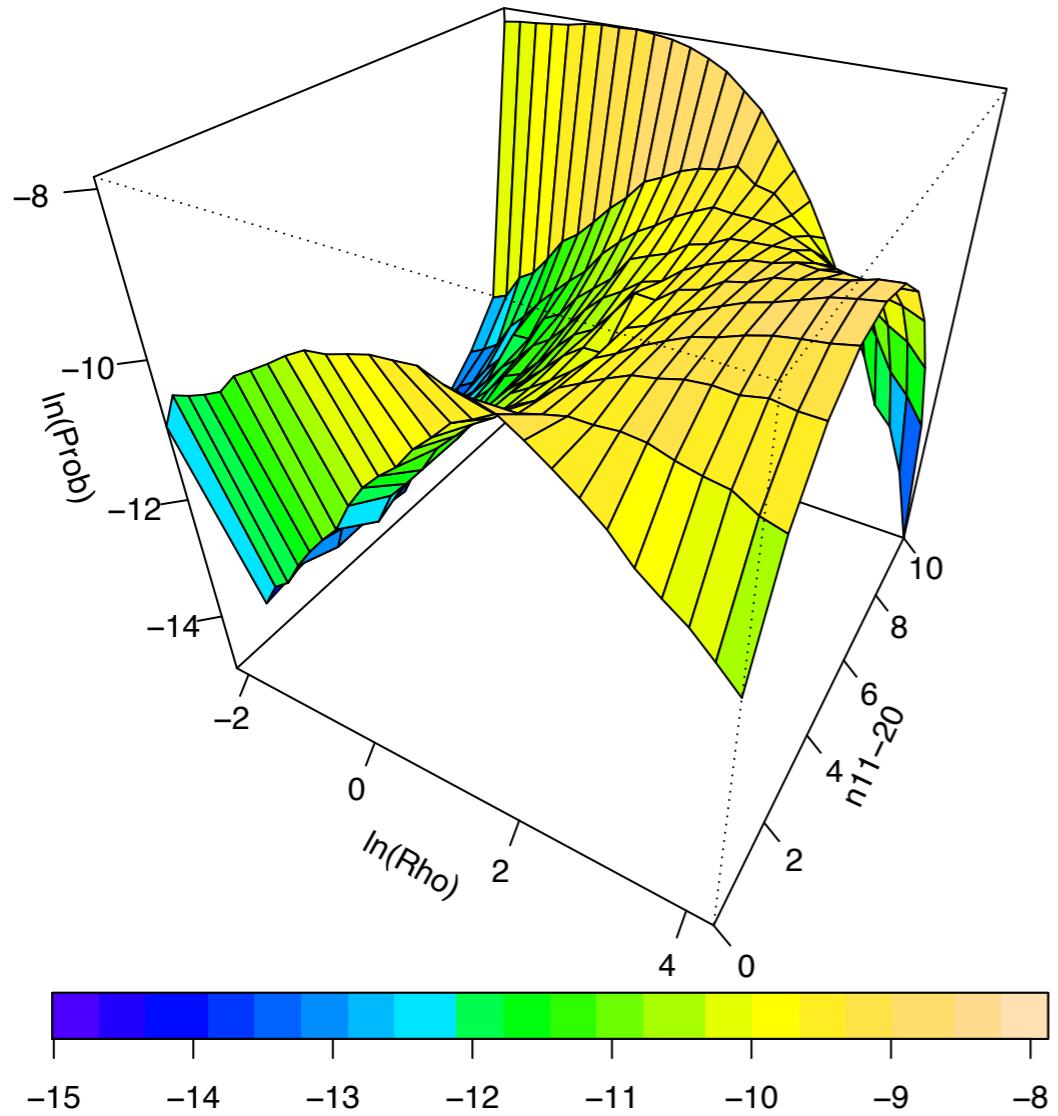
$$h_t = f(W_h \cdot h_{t-1} + W_x \cdot x_t)$$

$$\hat{y}_t = \text{softmax}(W_y \cdot h_t)$$

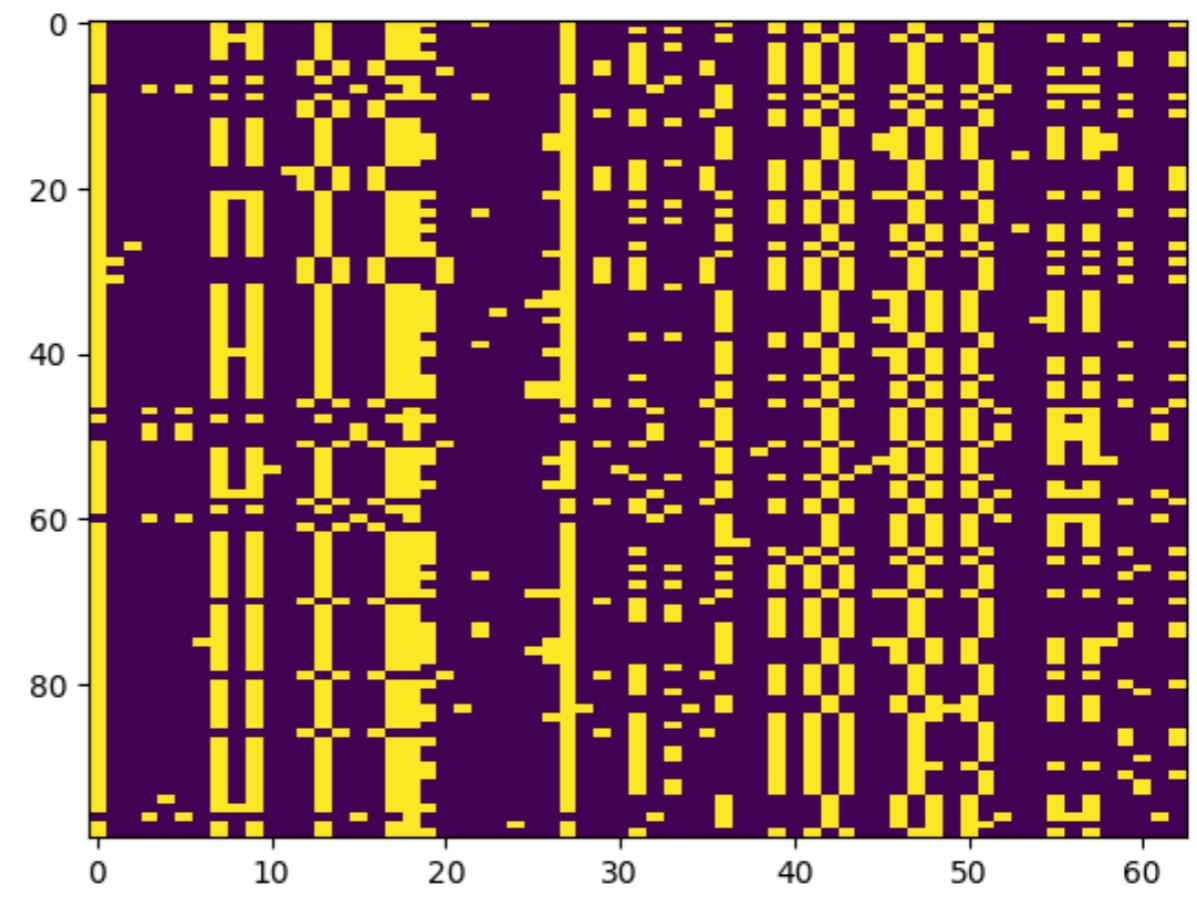
Estimating Recombination Rate with Deep Learning

Population genetic parameter estimation: $4Nr$

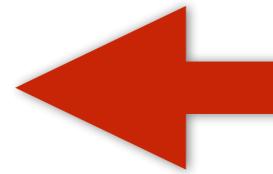
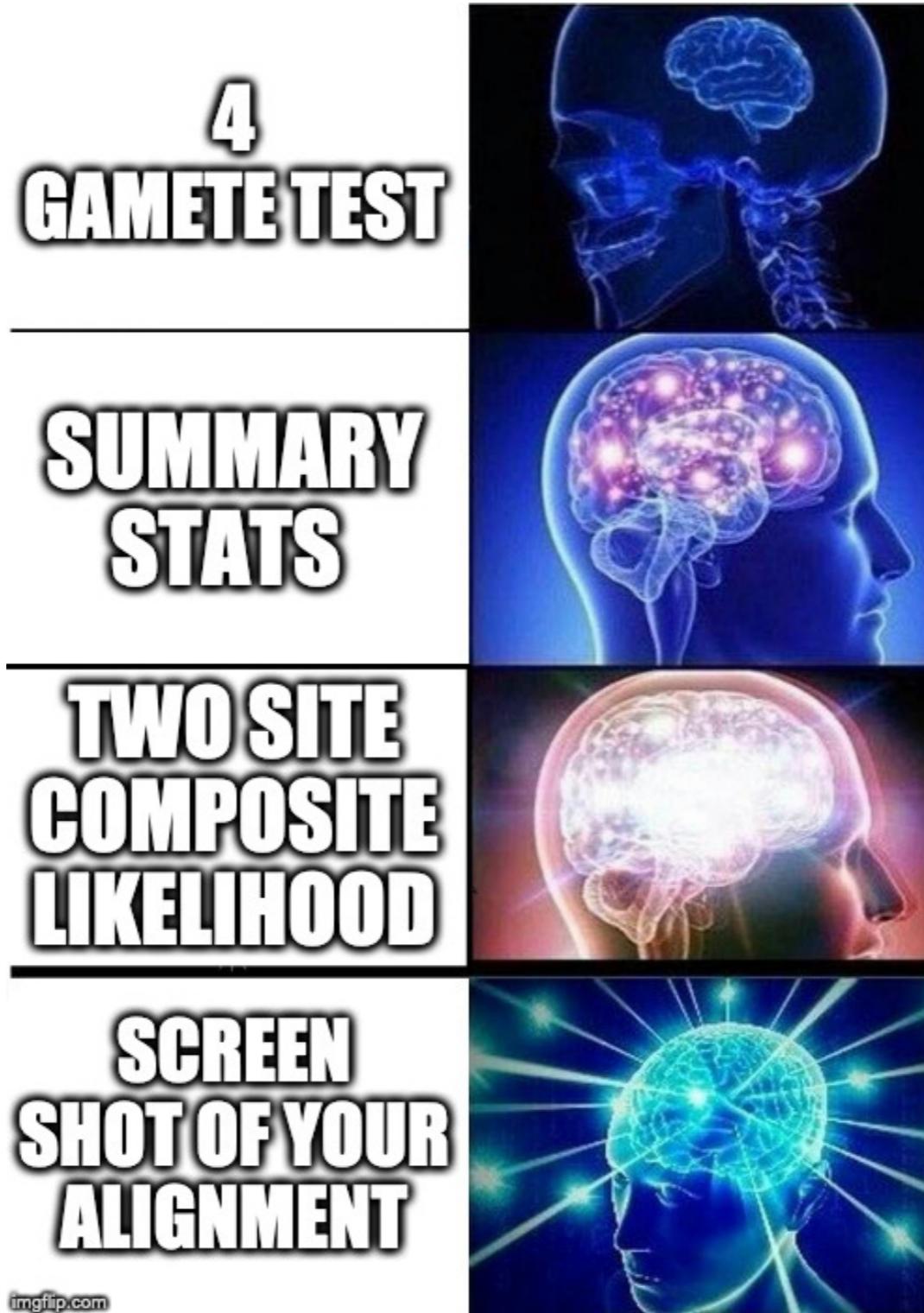
Composite Likelihood— LDhat



Deep Neural Nets

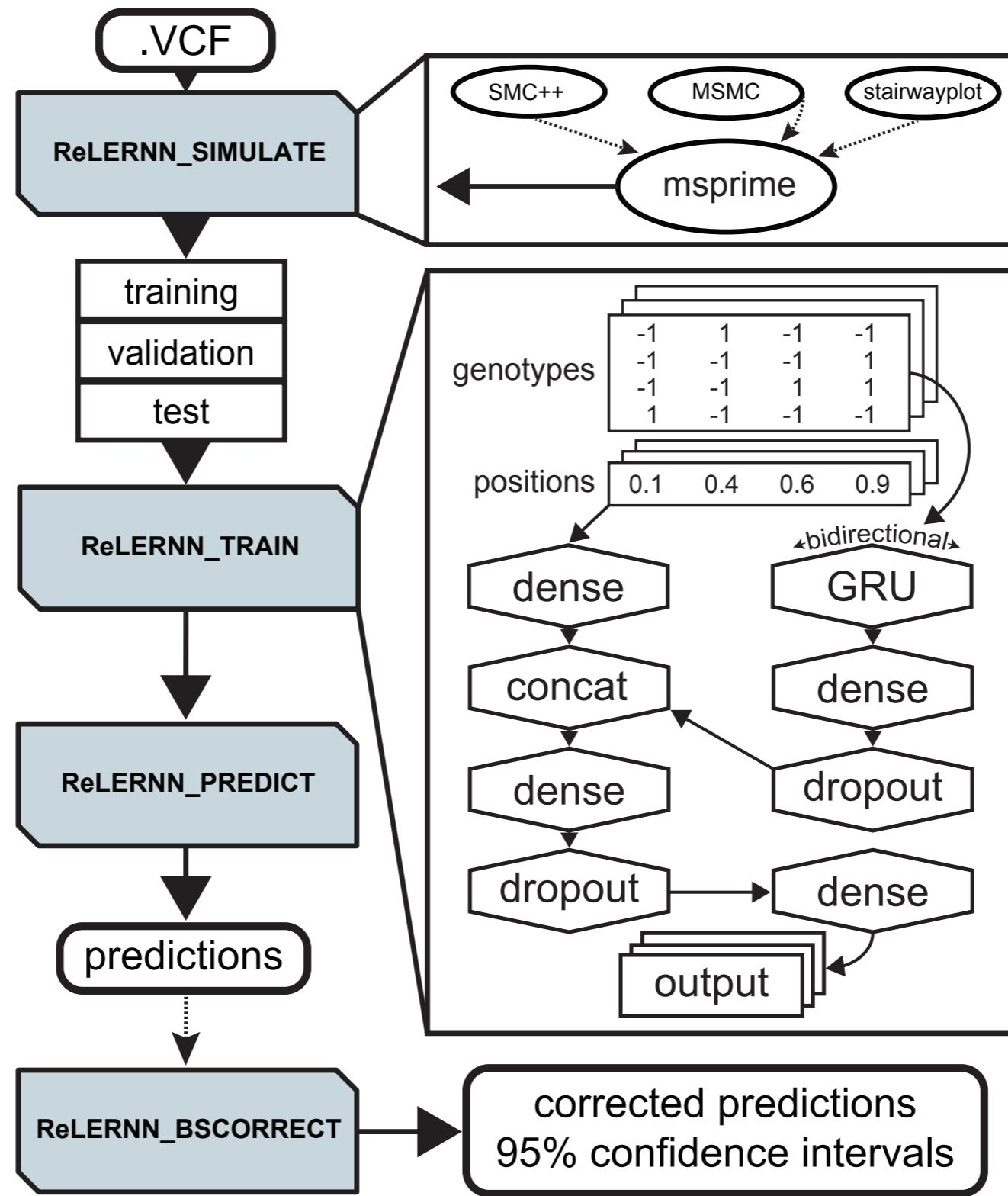


Meme history of estimating recombination rate



Where we
are now

ReLERNN

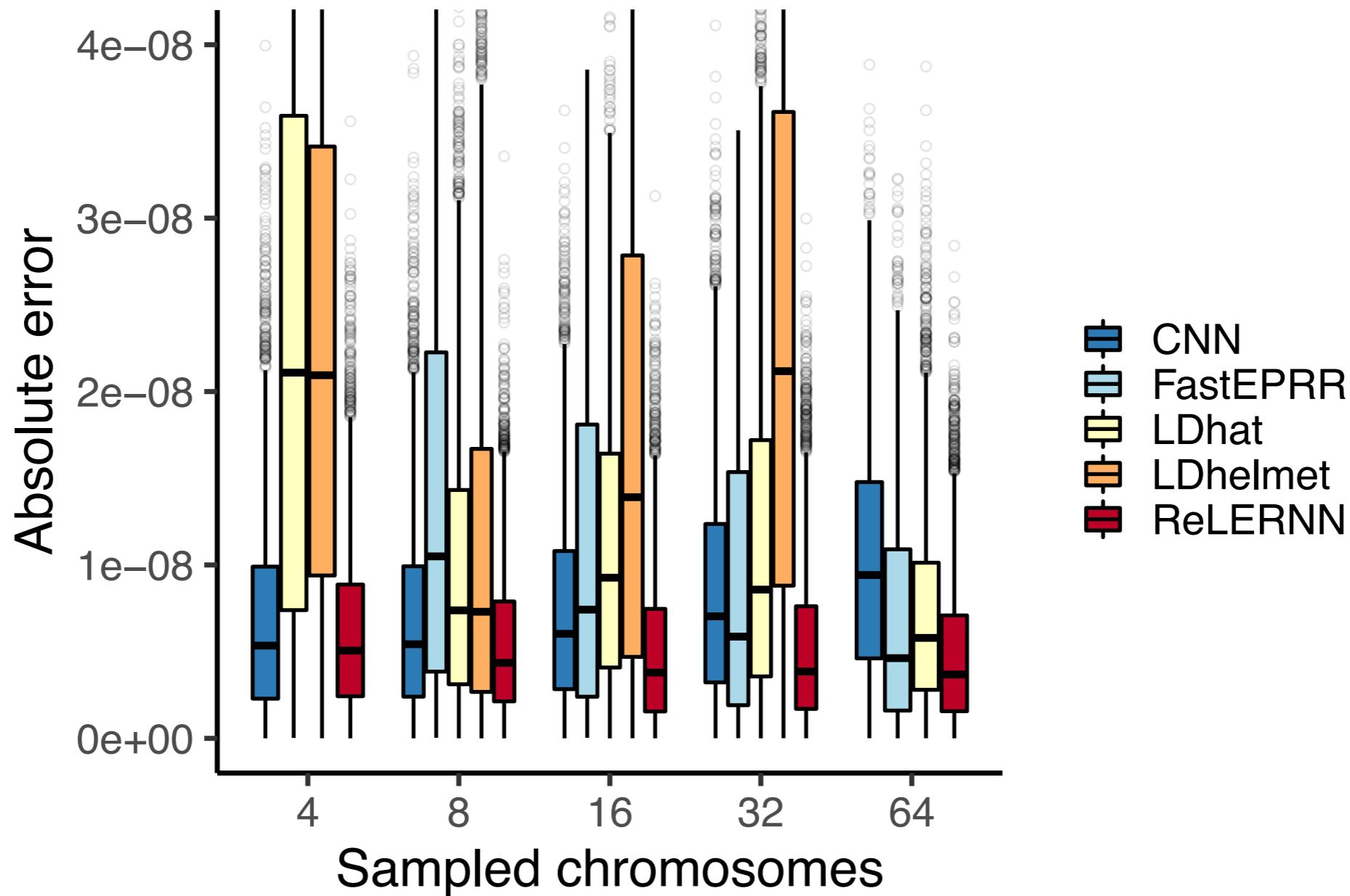


Jeff Adriion



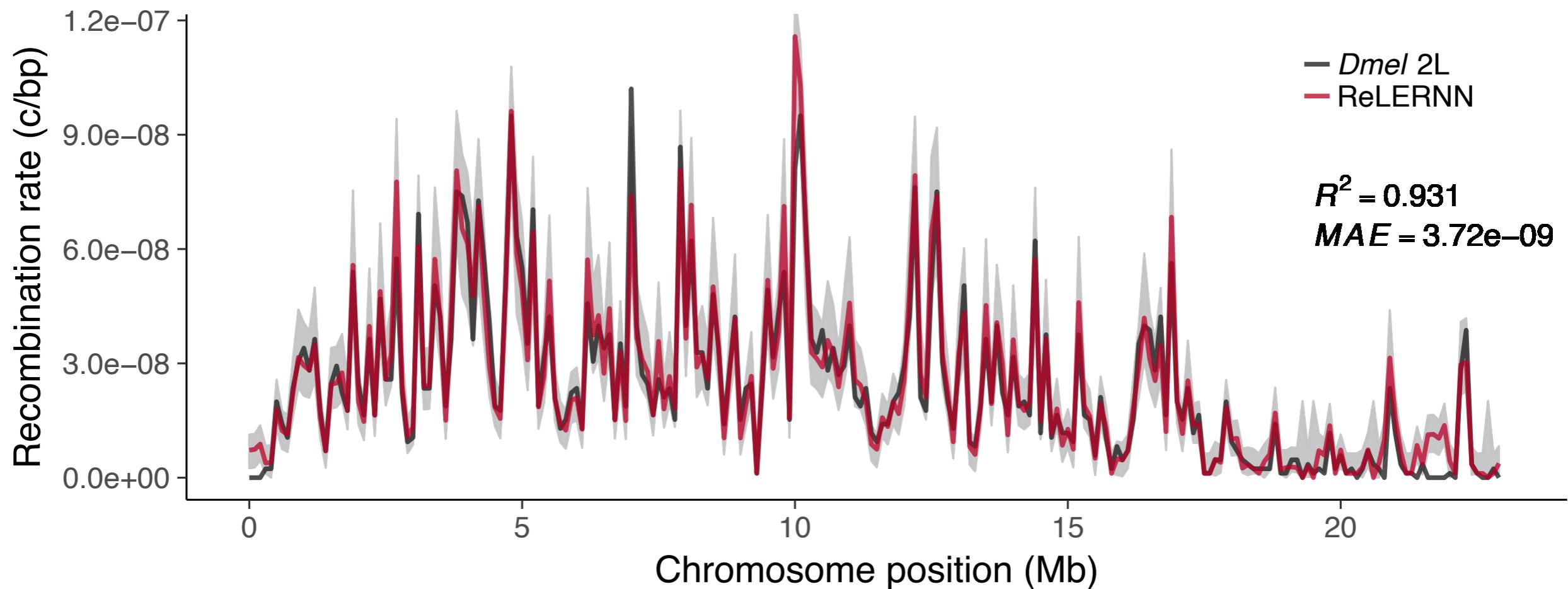
Estimating Recombination Rate with Deep Learning

Population genetic parameter estimation: $4Nr$



Estimating Recombination Rate with Deep Learning

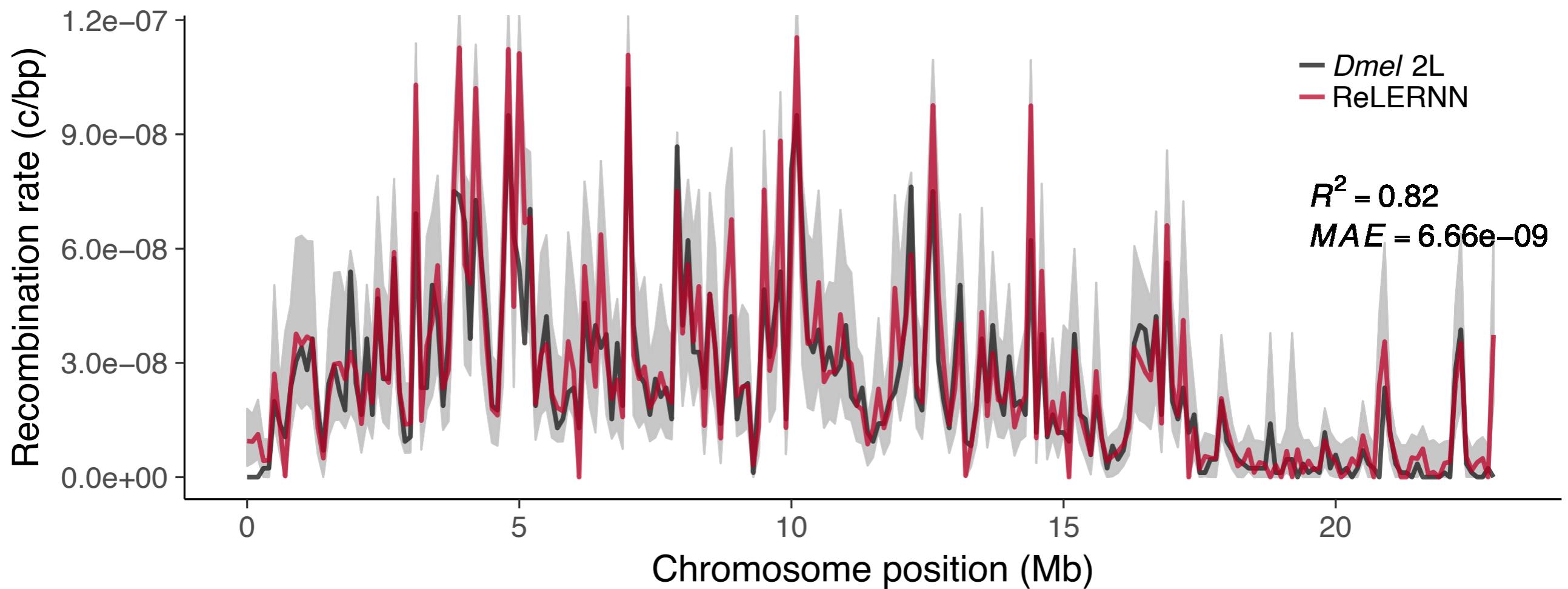
Simulated Landscape of Recombination



Comeron et al. (2012) genetic map from Drosophila (n=20)

Estimating Recombination Rate with Deep Learning

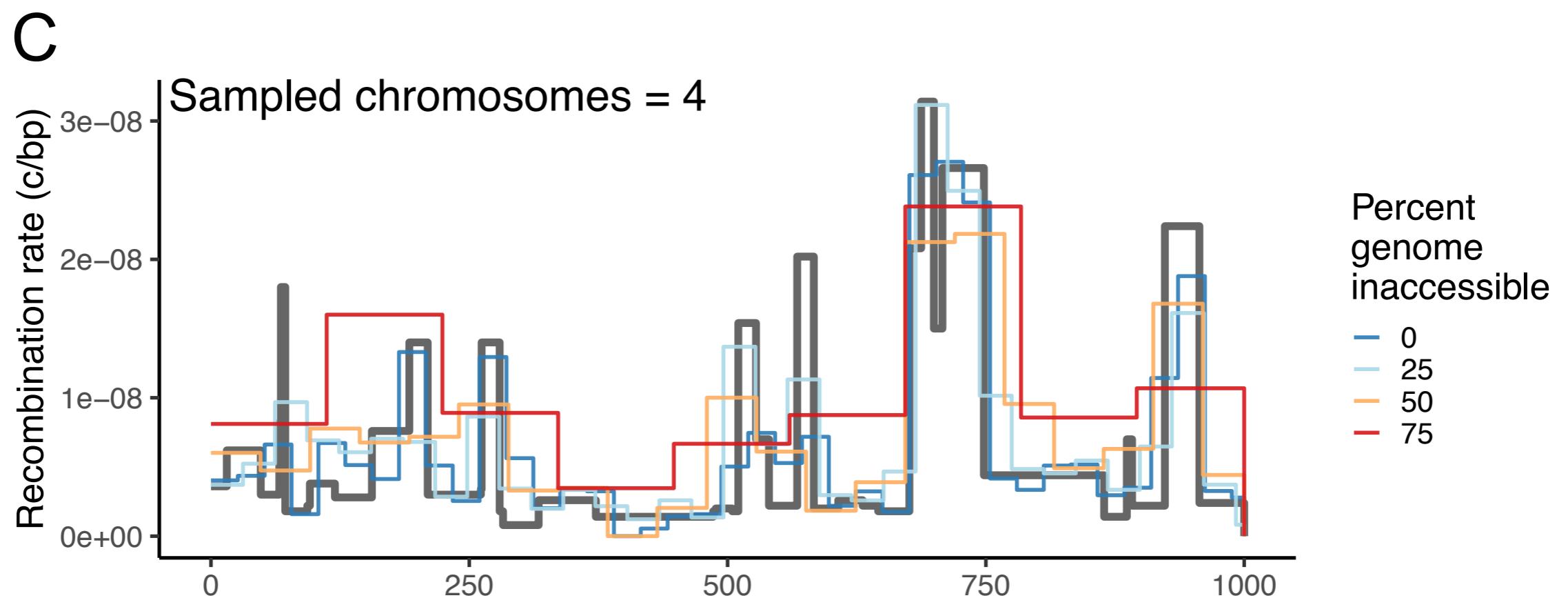
Simulated Landscape of Recombination



Comeron et al. (2012) genetic map from Drosophila (n=4)

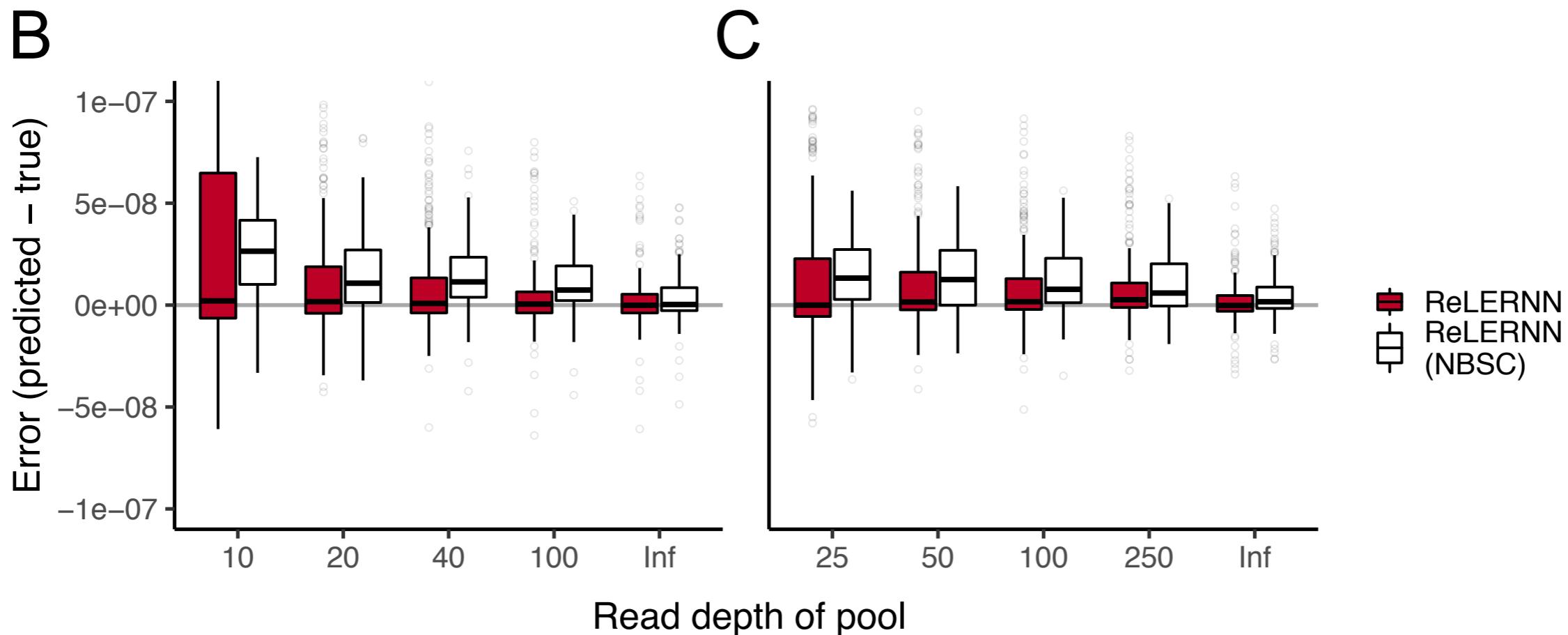
Estimating Recombination Rate with Deep Learning

ReLERNN does well with ‘bad data’



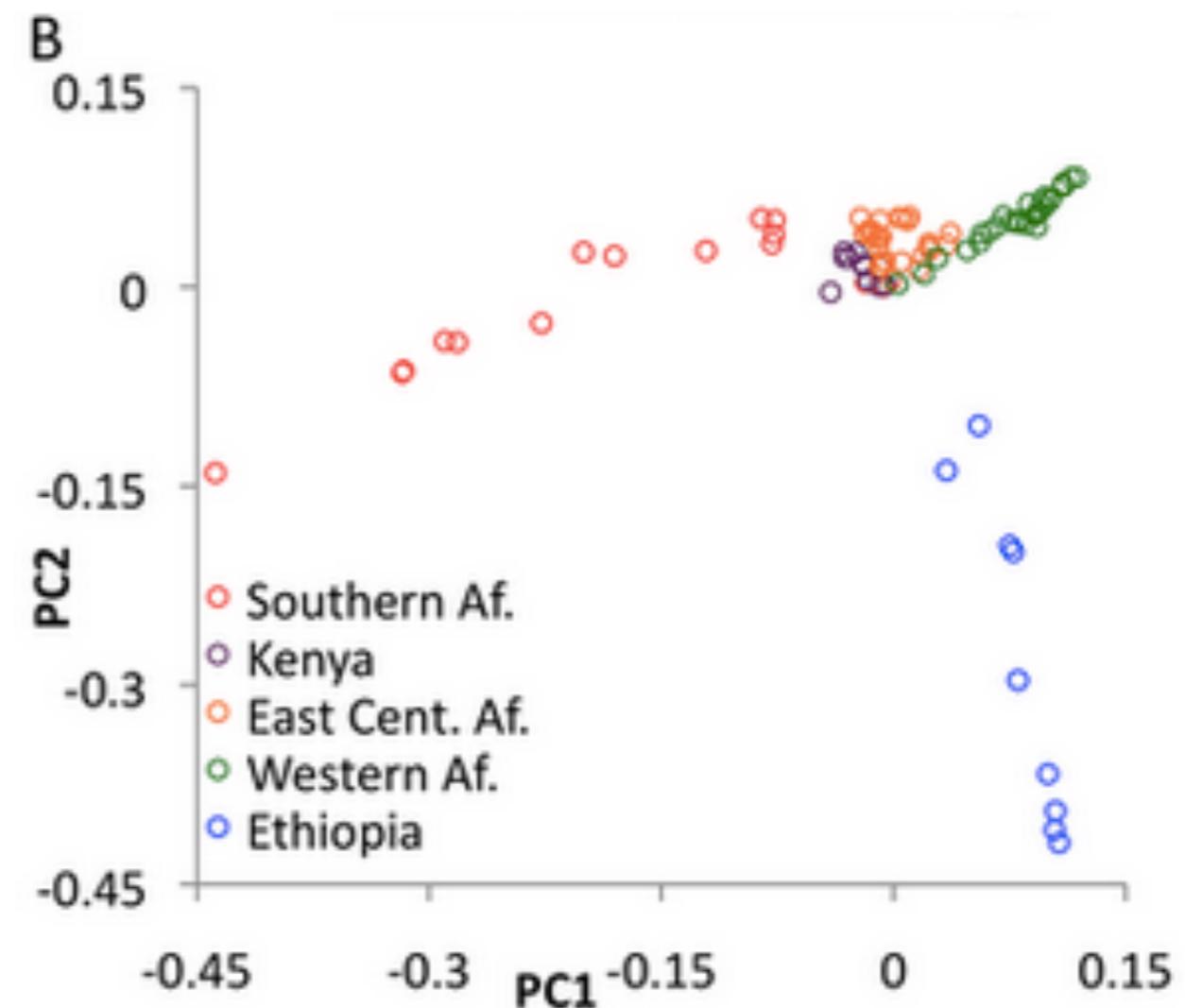
Estimating Recombination Rate with Deep Learning

ReLERNN does well with ‘bad data’



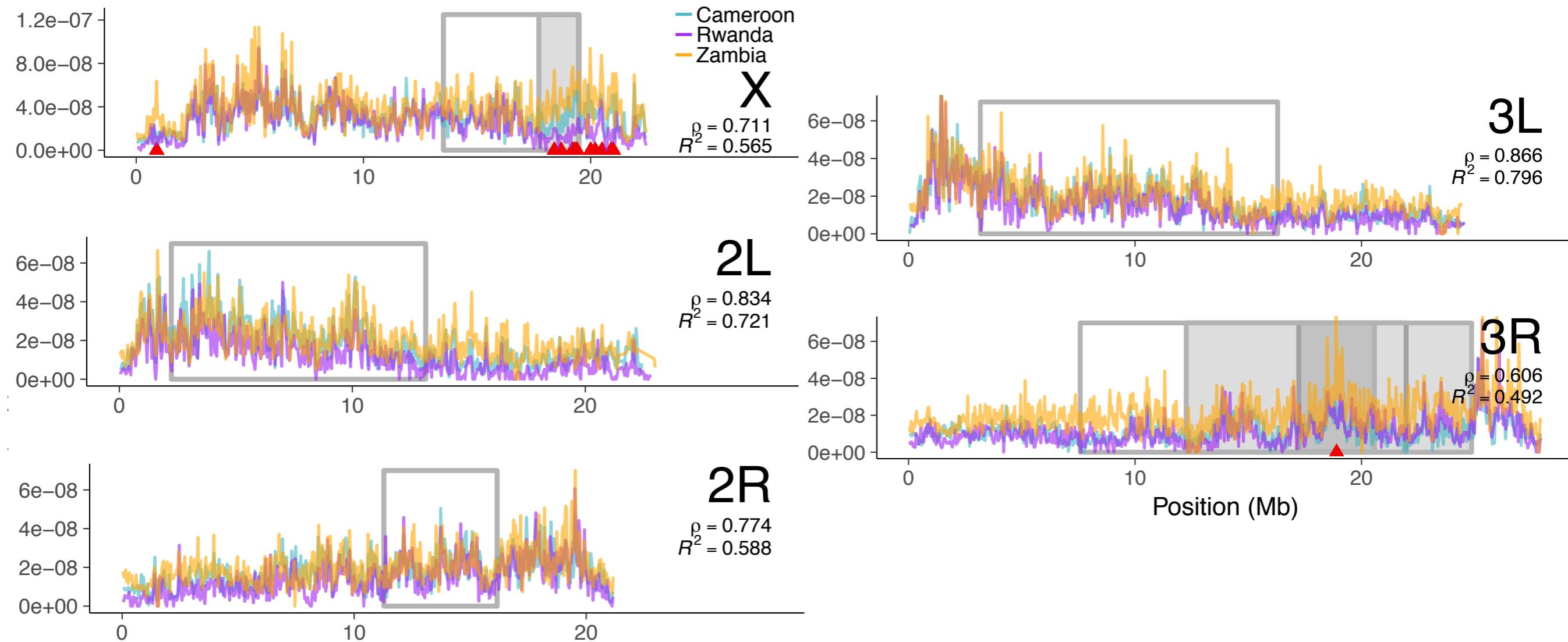
Can even use ReLERNN for pool-seq data!

Application to African Drosophila Samples



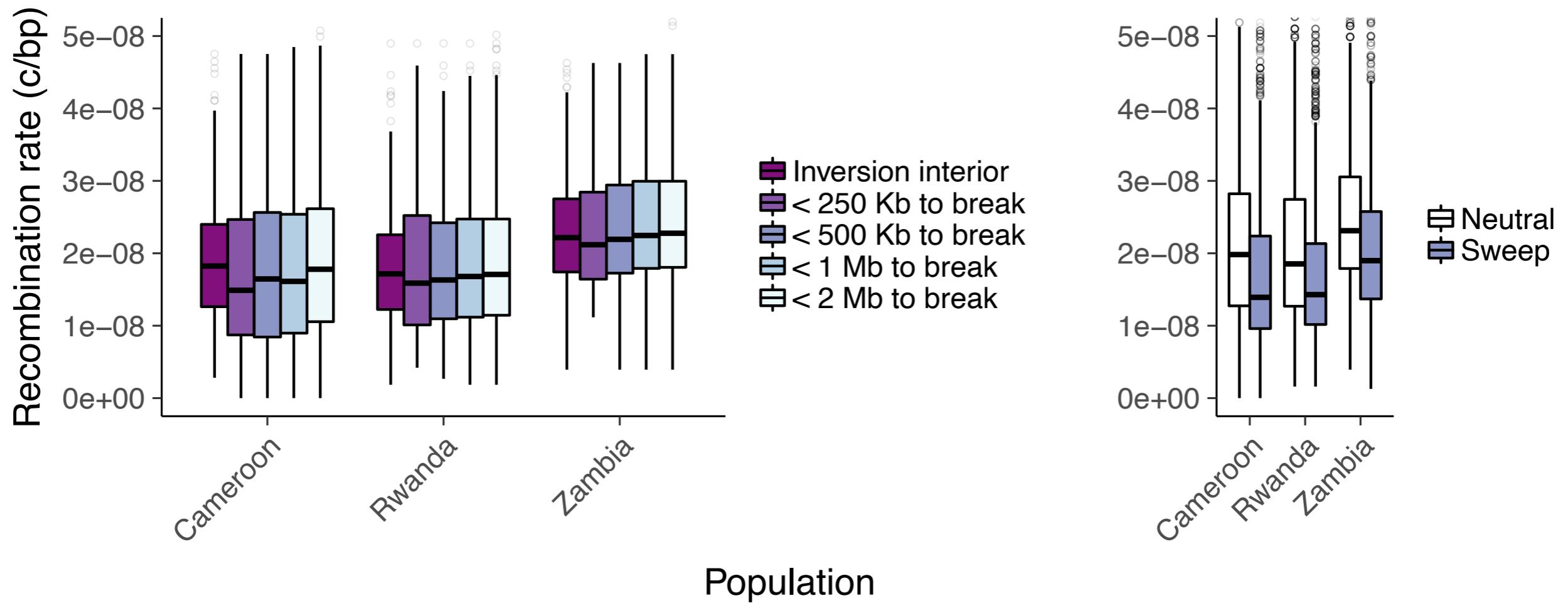
Data from 3 populations sampled in Pool et al (2012)

Application to African Drosophila Samples



Major Outliers Associated with Segregating Inversions

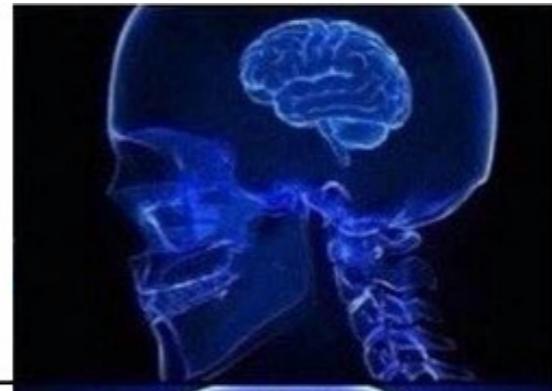
Both inversions and sweeps affect inferred recombination



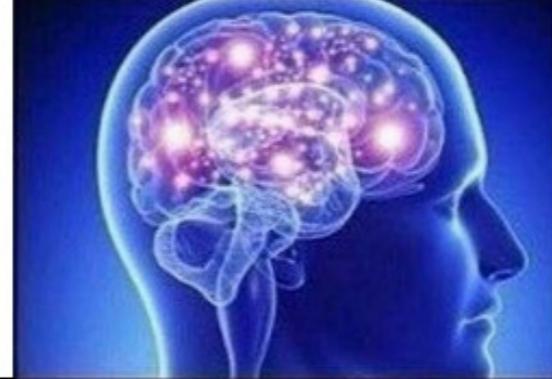
Effect of cosmopolitan inversion $\ln(2L)t$

ReLERN Summary

4
GAMETE TEST



**SUMMARY
STATS**



**TWO SITE
COMPOSITE
LIKELIHOOD**



**SCREEN
SHOT OF YOUR
ALIGNMENT**



imgflip.com

- Early days for automated feature extraction for popgen
- To do: which architectures? which tasks? uncertainty?

Predicting the landscape of recombination using deep learning

Jeffrey R Adrion , Jared G Galloway, Andrew D Kern Author Notes

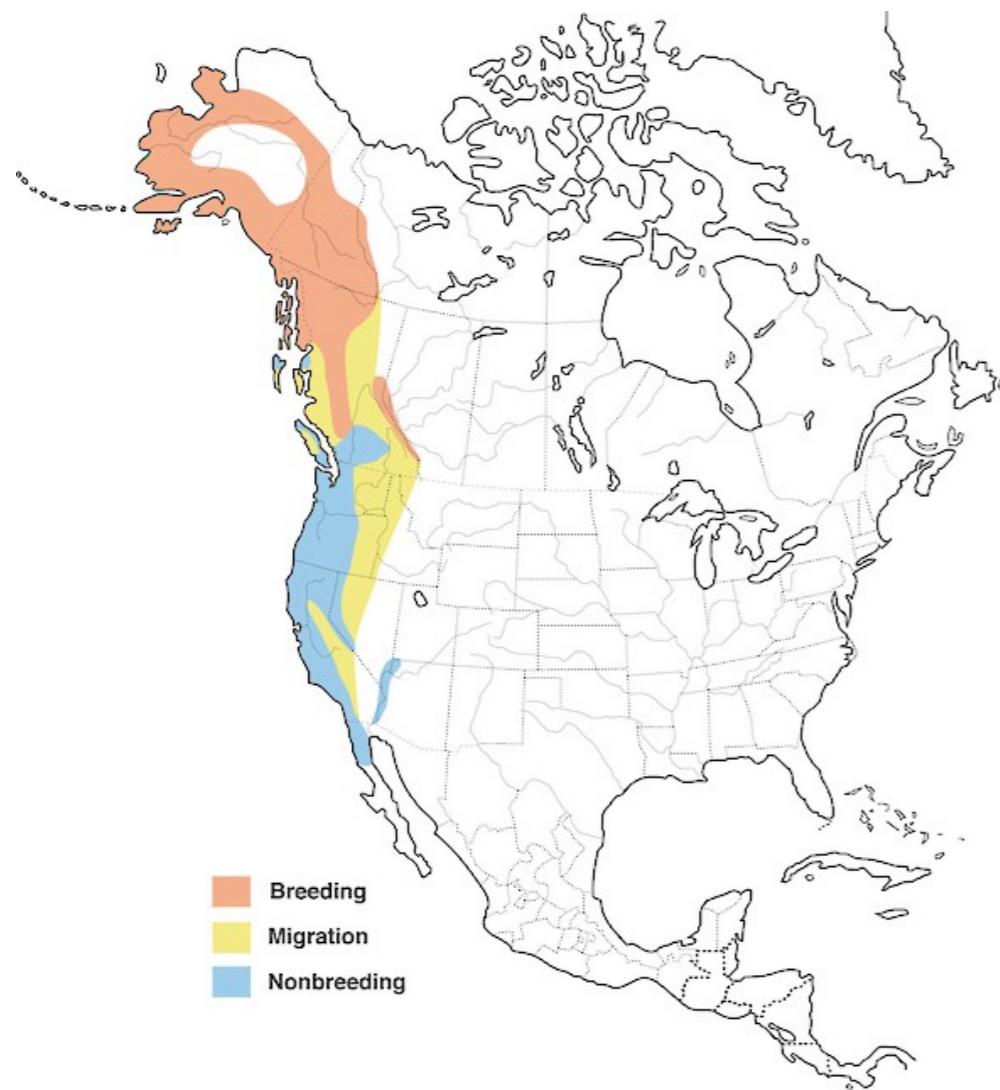
Molecular Biology and Evolution, msaa038, <https://doi.org/10.1093/molbev/msaa038>

Published: 20 February 2020 Article history ▾

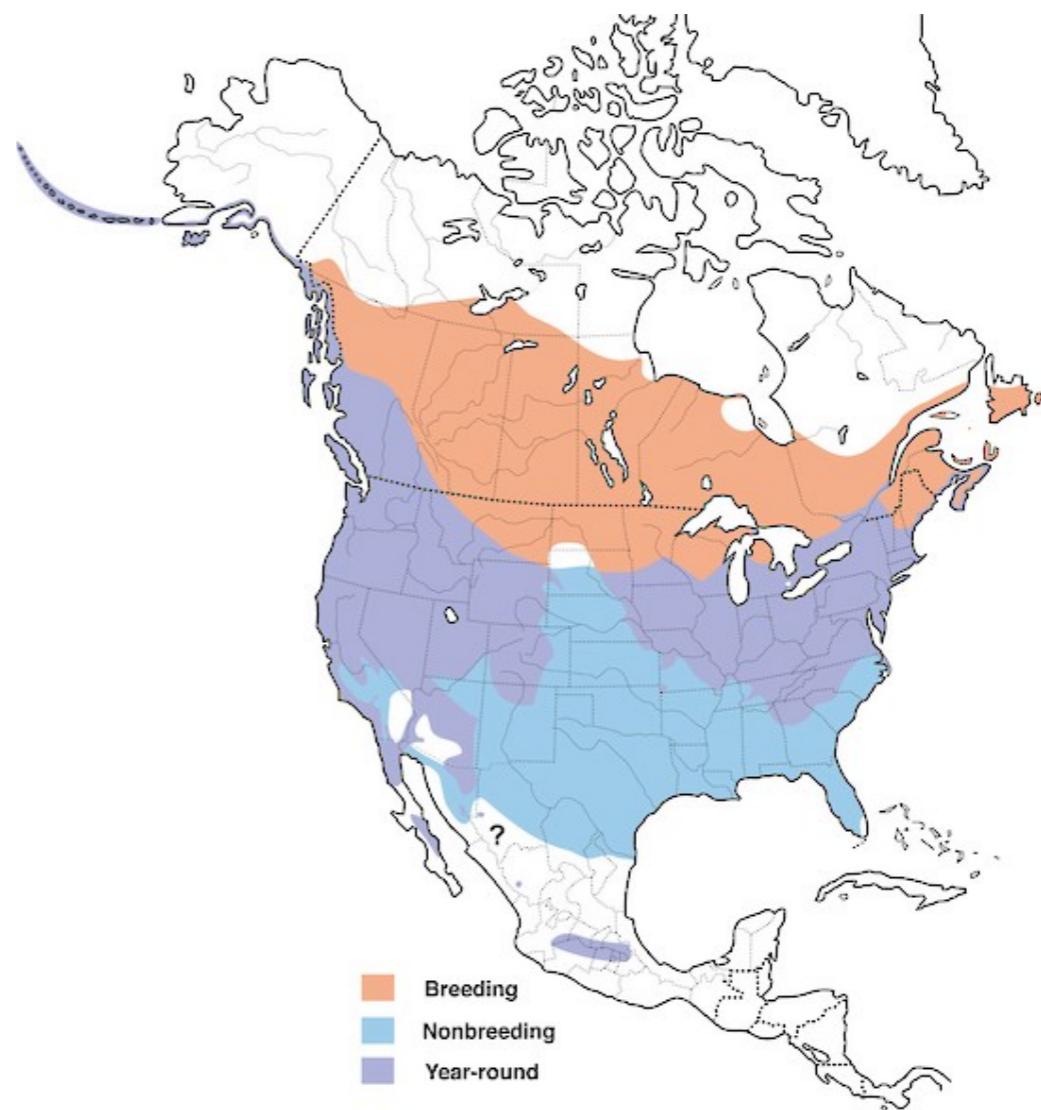


CJ Battey

Organisms live in space



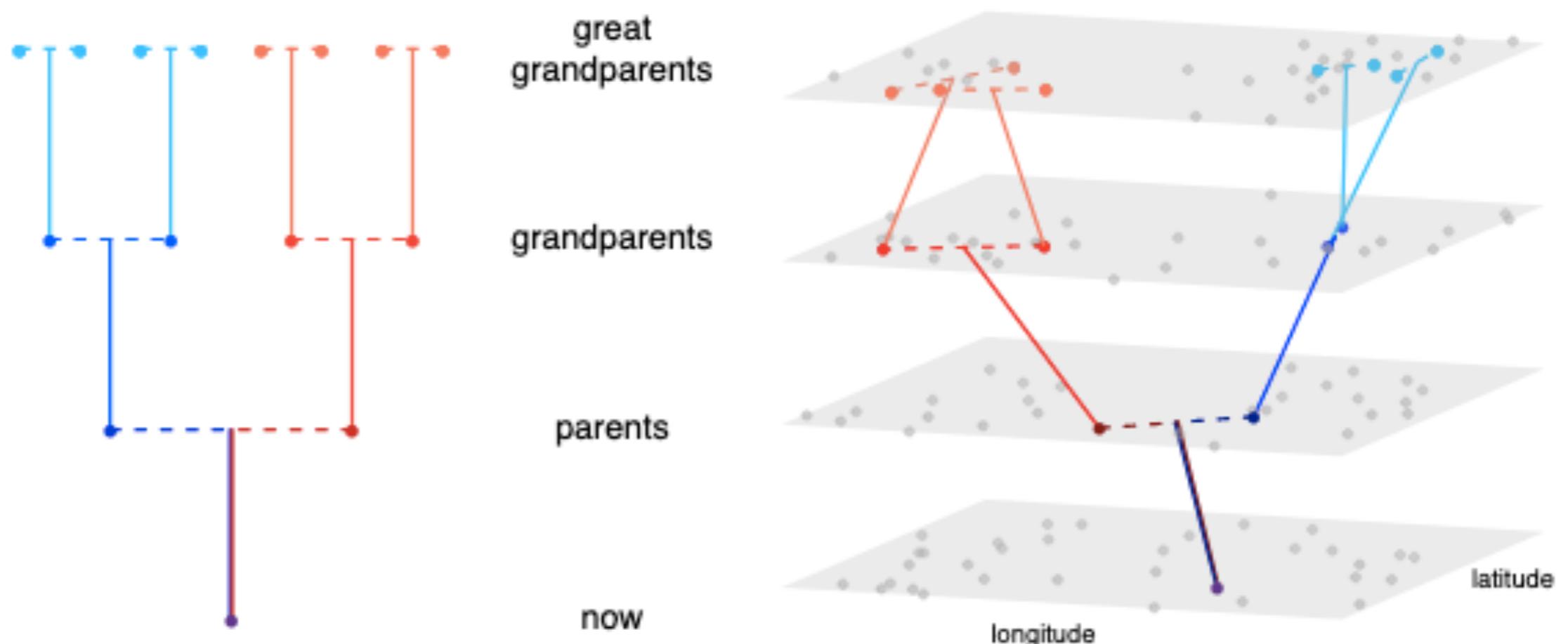
Golden crowned sparrow



Song sparrow



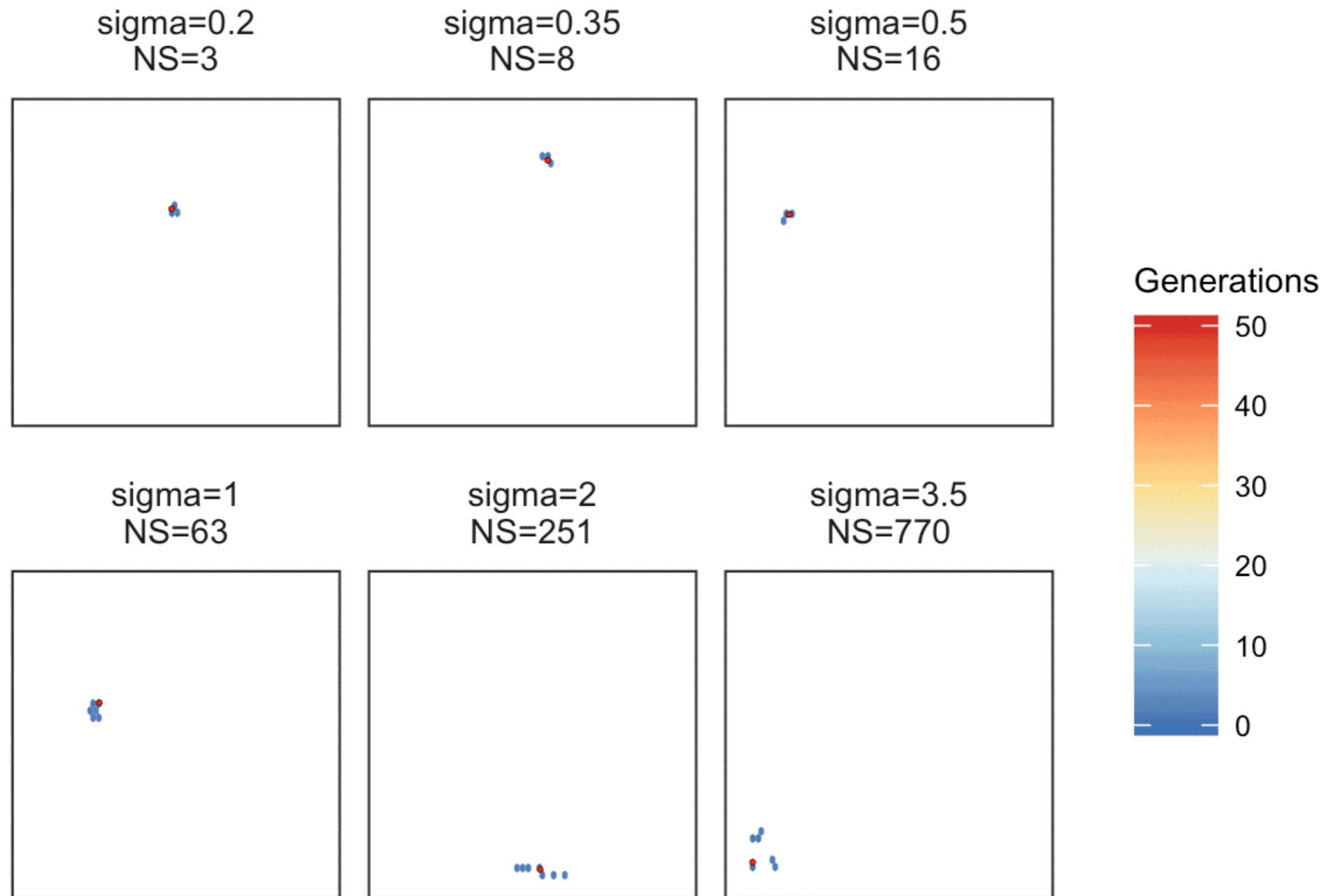
Space is the Place



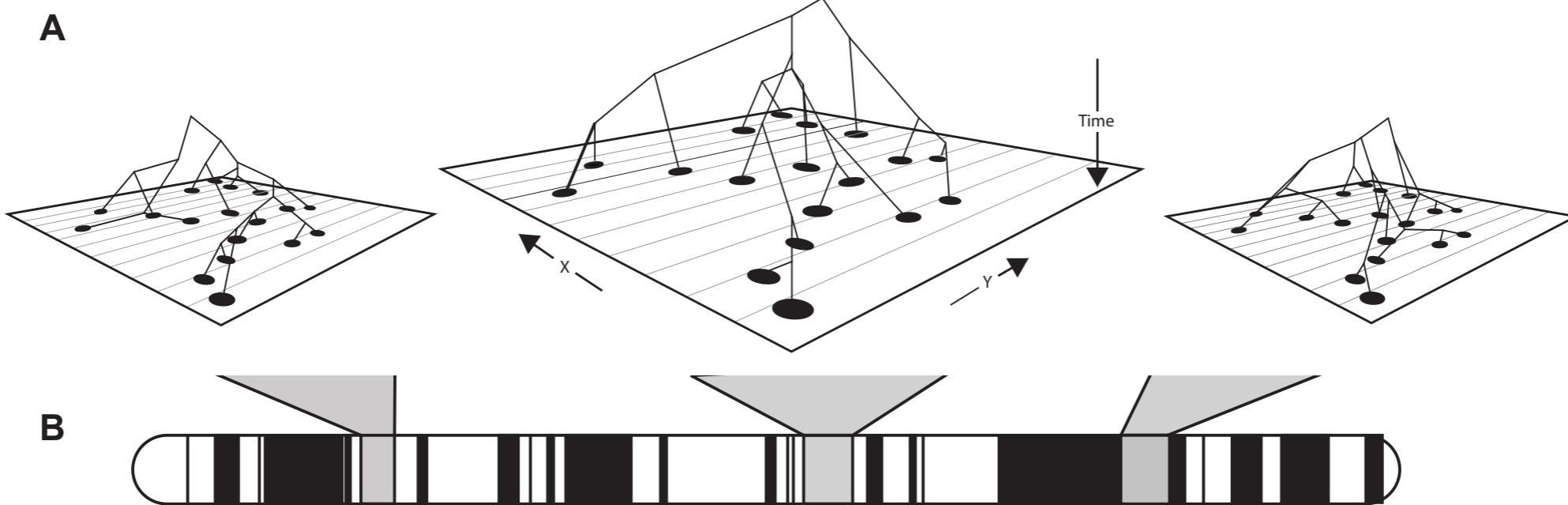
From Bradburd and Ralph (2019)

Space is the Place

Genealogical Ancestors by Neighborhood Size

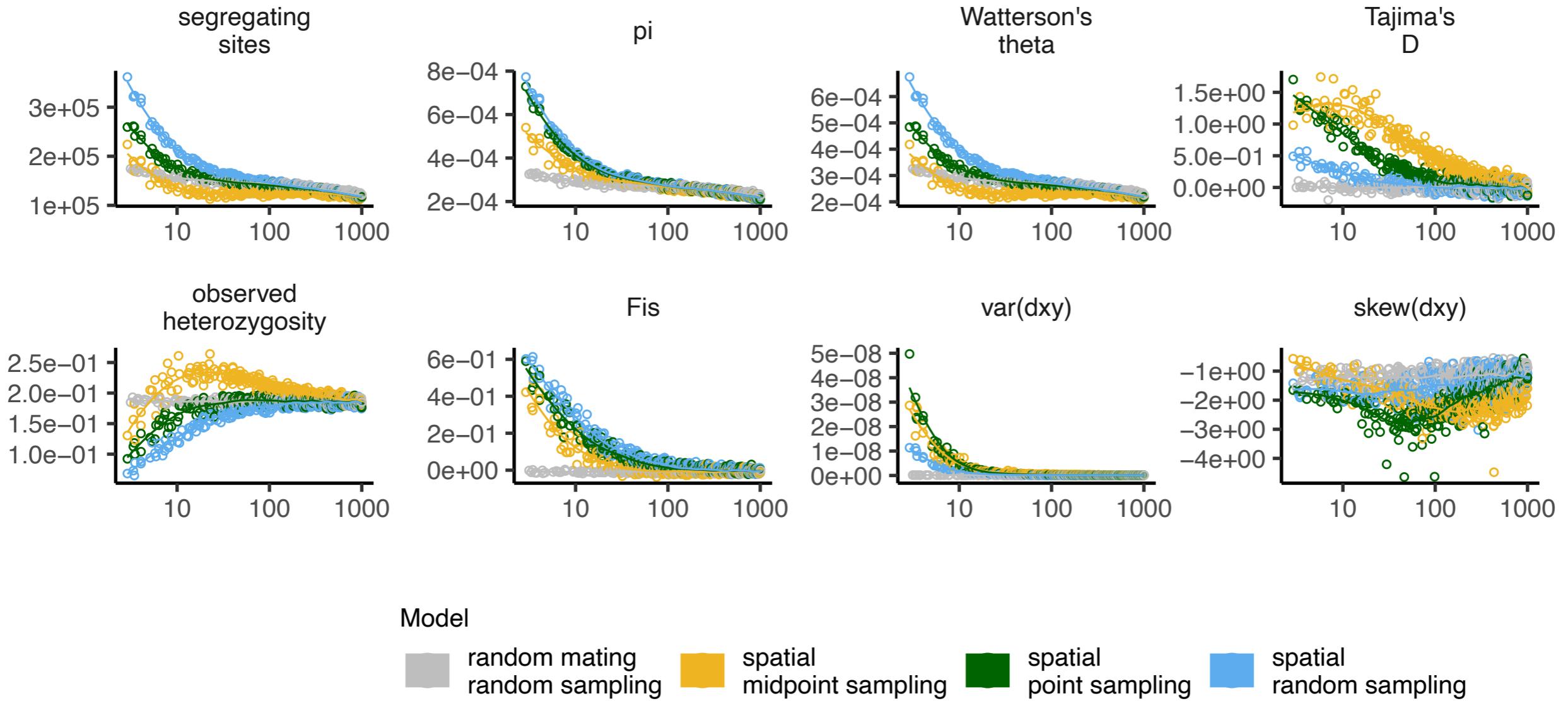


Space is the Place



Recombination means different spatial ancestry
at different points along chromosomes

Space is the Place



Space matters for genetic variation

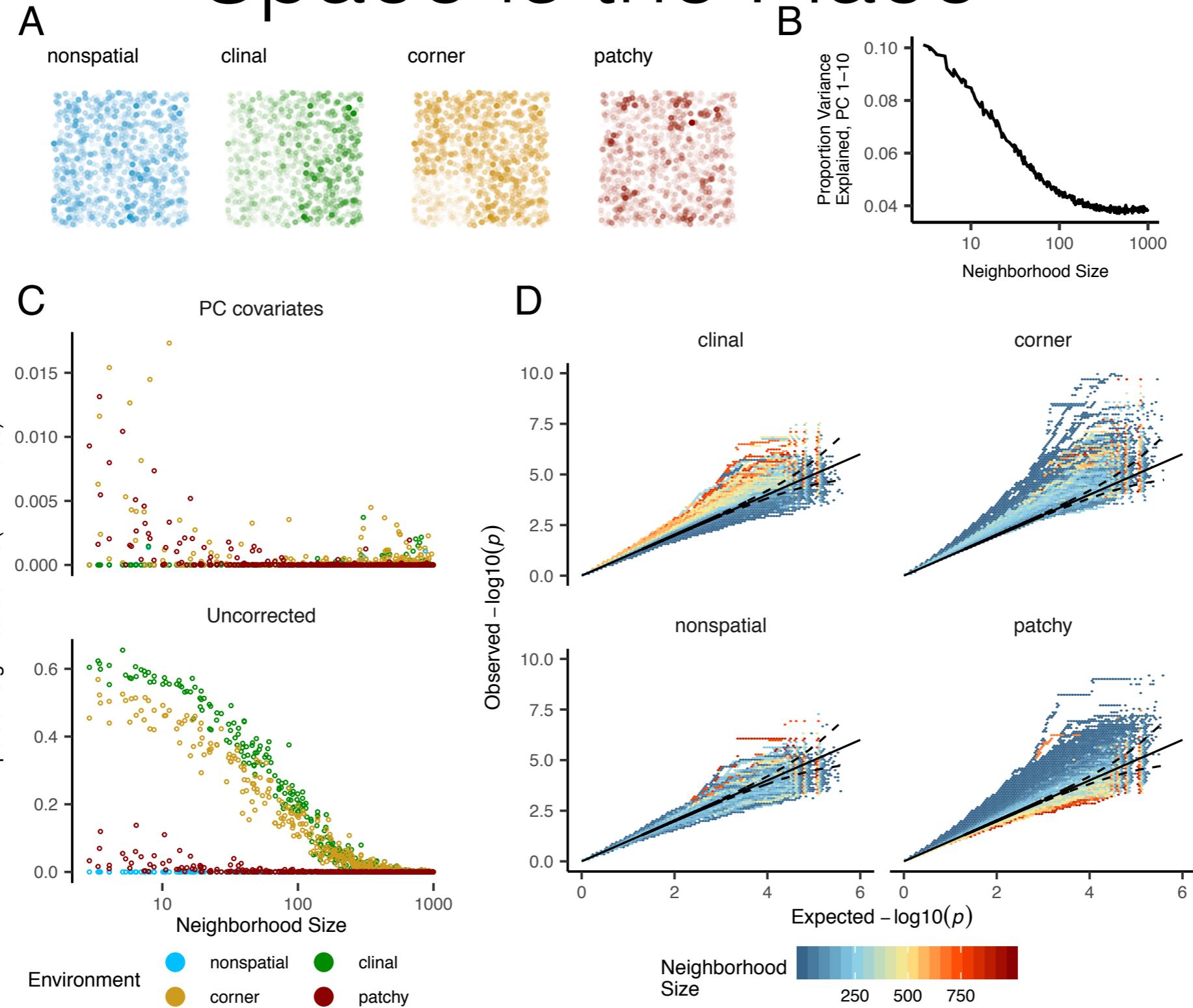
JOURNAL ARTICLE

**Space is the Place: Effects of Continuous Spatial Structure
on Analysis of Population Genetic Data** FREE

C J Battey ✉, Peter L Ralph, Andrew D Kern

Genetics, Volume 215, Issue 1, 1 May 2020, Pages 193–214,

Space is the Place



from Battey et al 2020

Geography in our genes

POLITICS AUGUST 7, 2012

Gibson Guitars and Feds Settle in Illegal Wood Case

KATE SHEPPARD

Bio | Follow



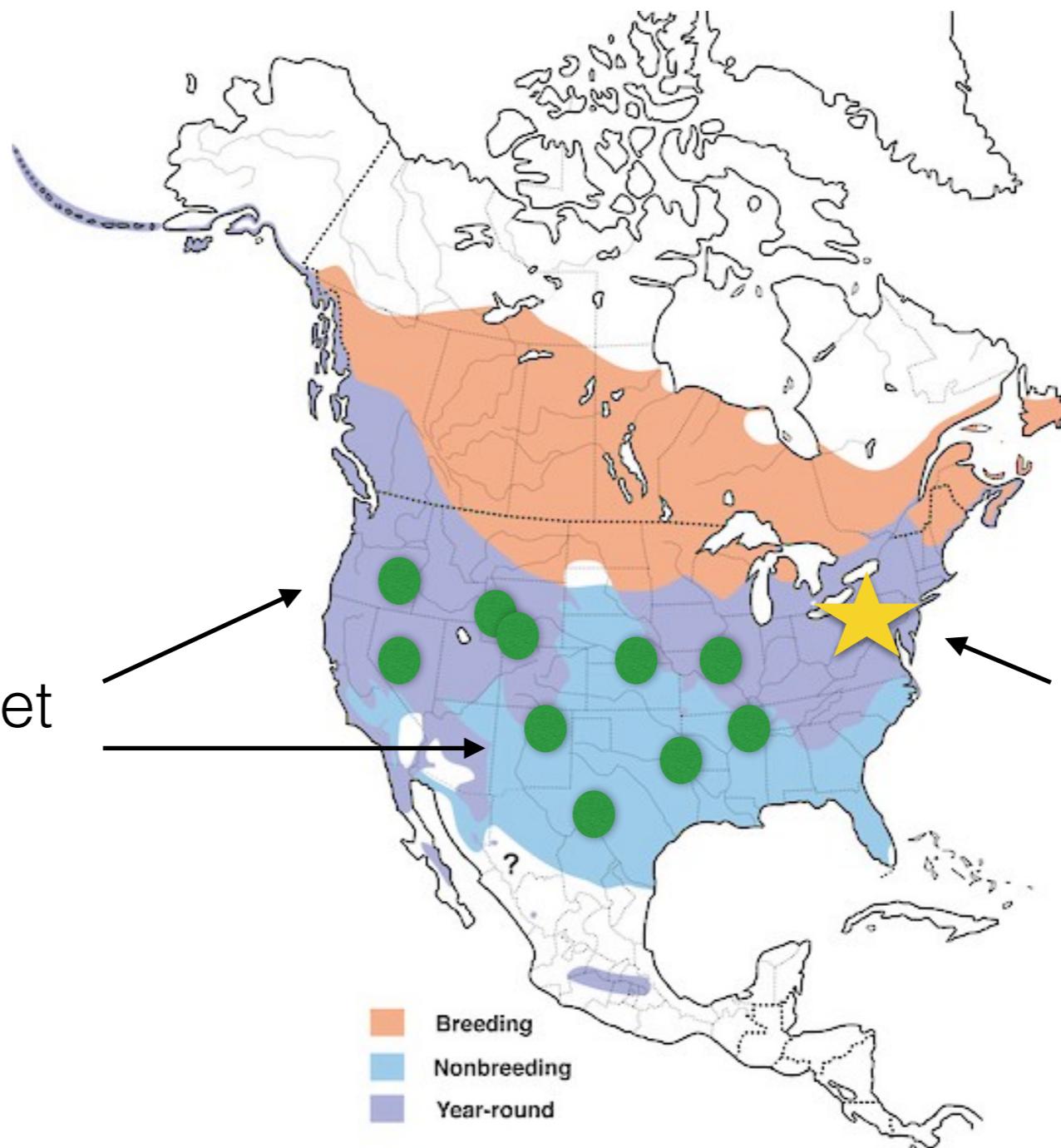
Can we exploit this?

Predicting geography

Song sparrow

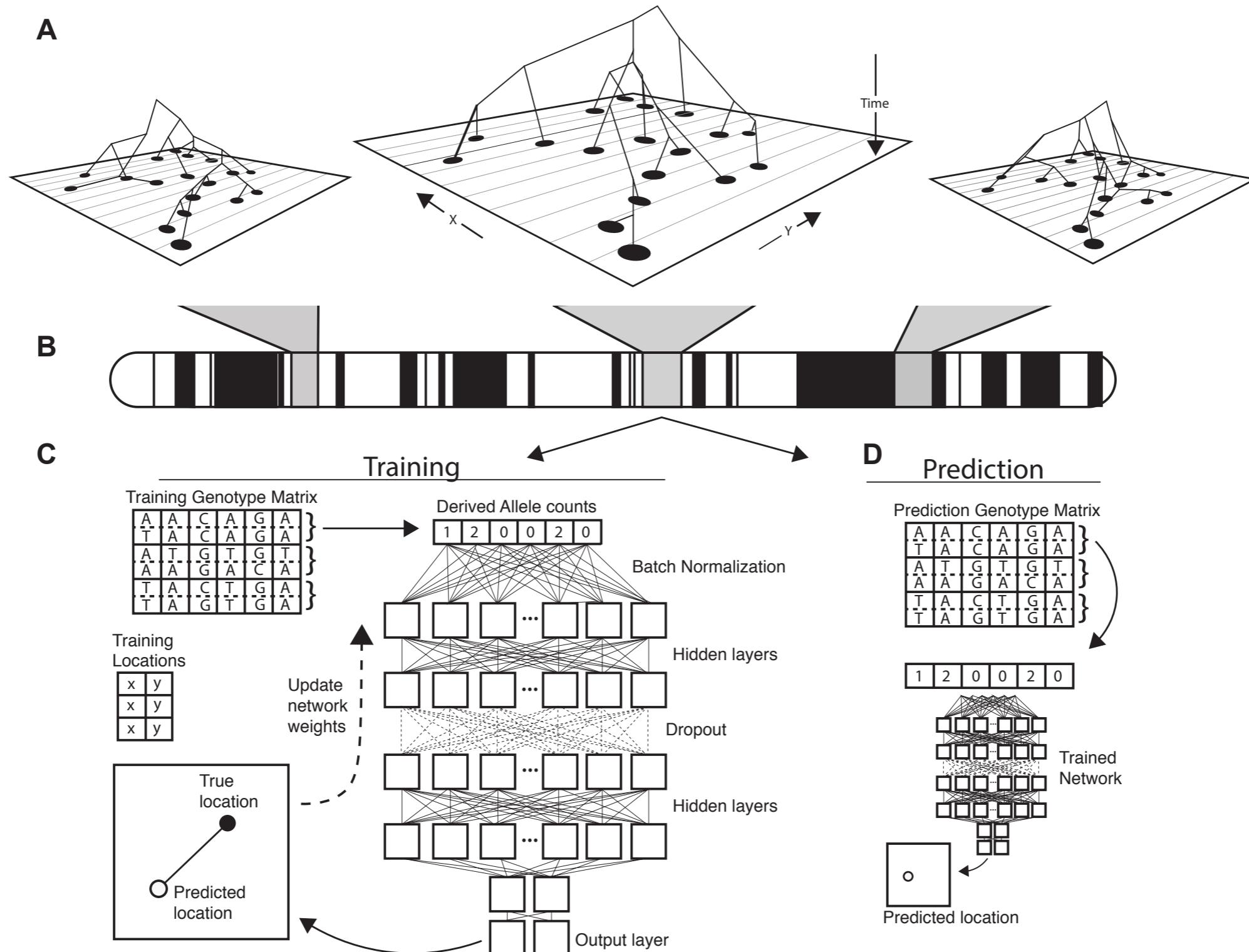


Training set



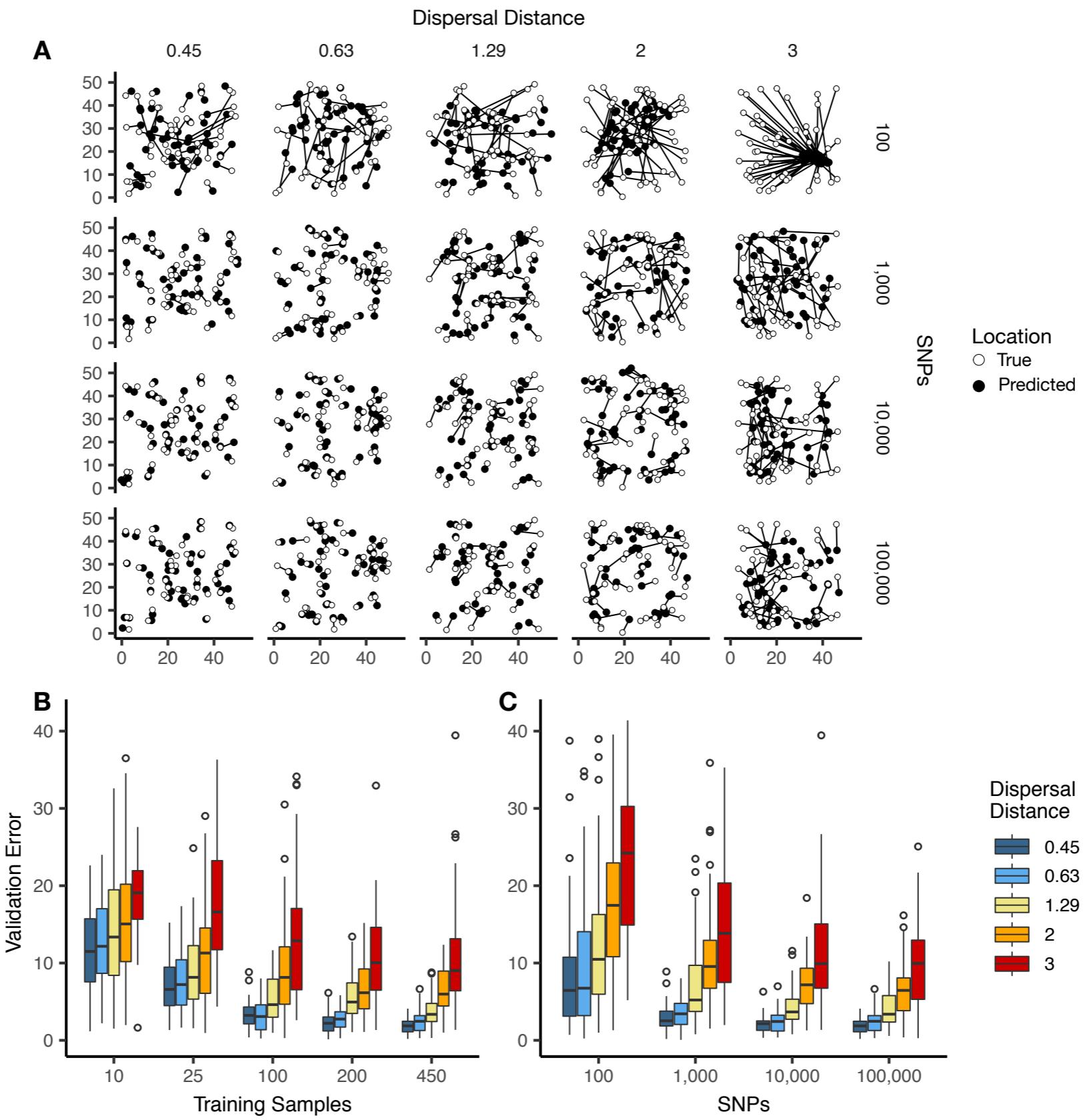
MODEL FREE

Locator— (deep) learning space

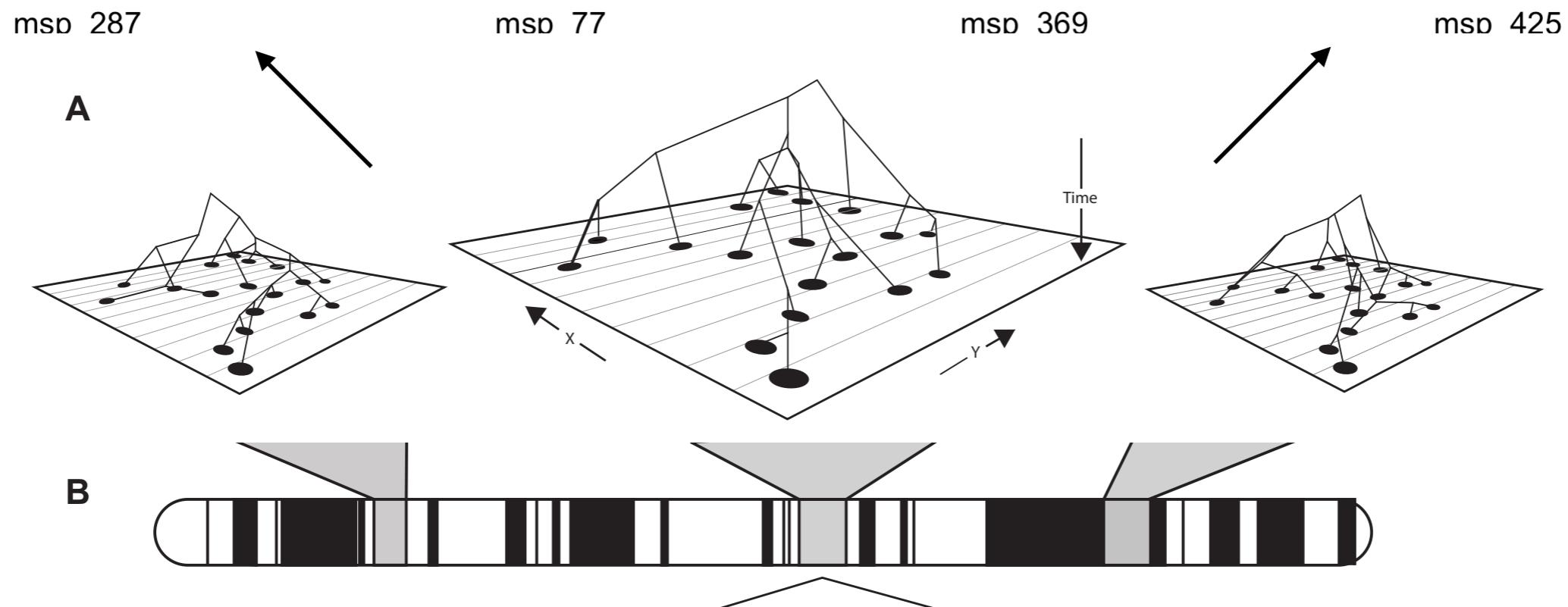
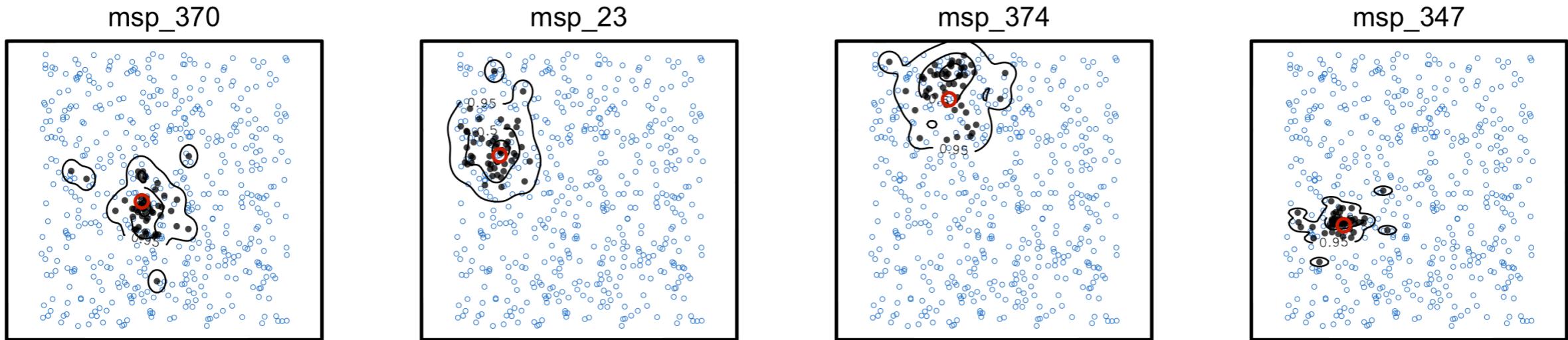


MODEL FREE

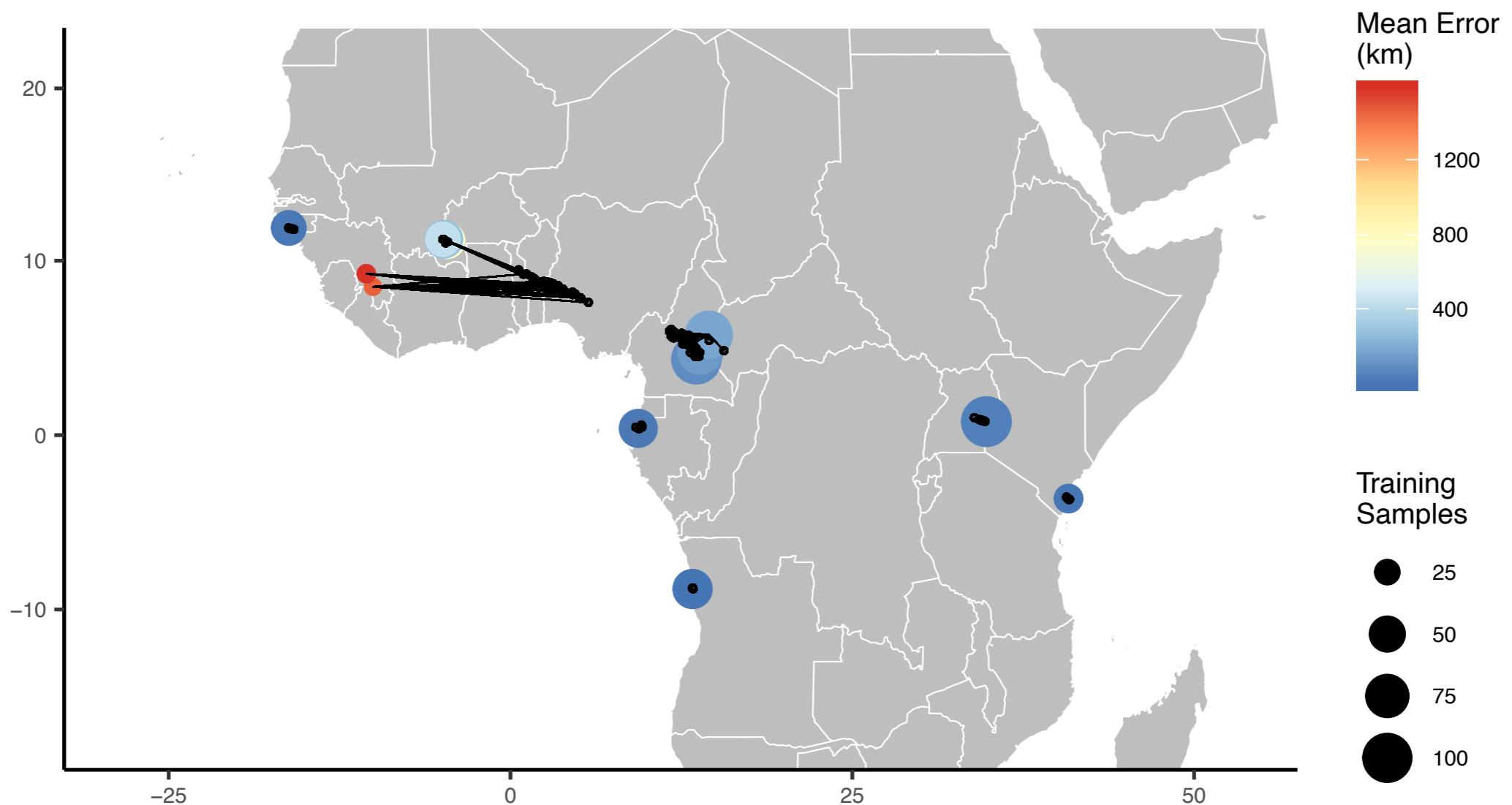
Locator – (deep) learning space



Locator – (deep) learning space

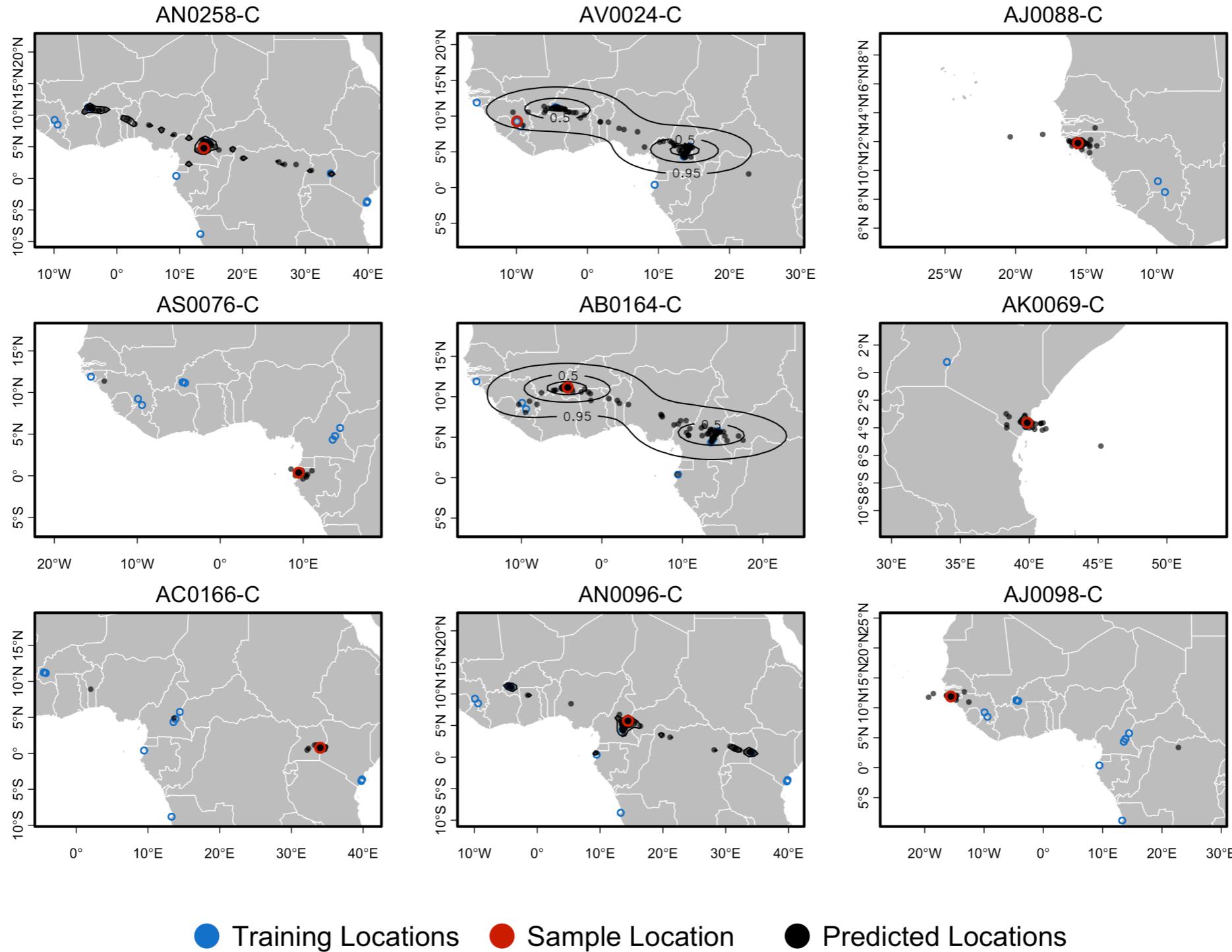


Locator – (deep) learning space



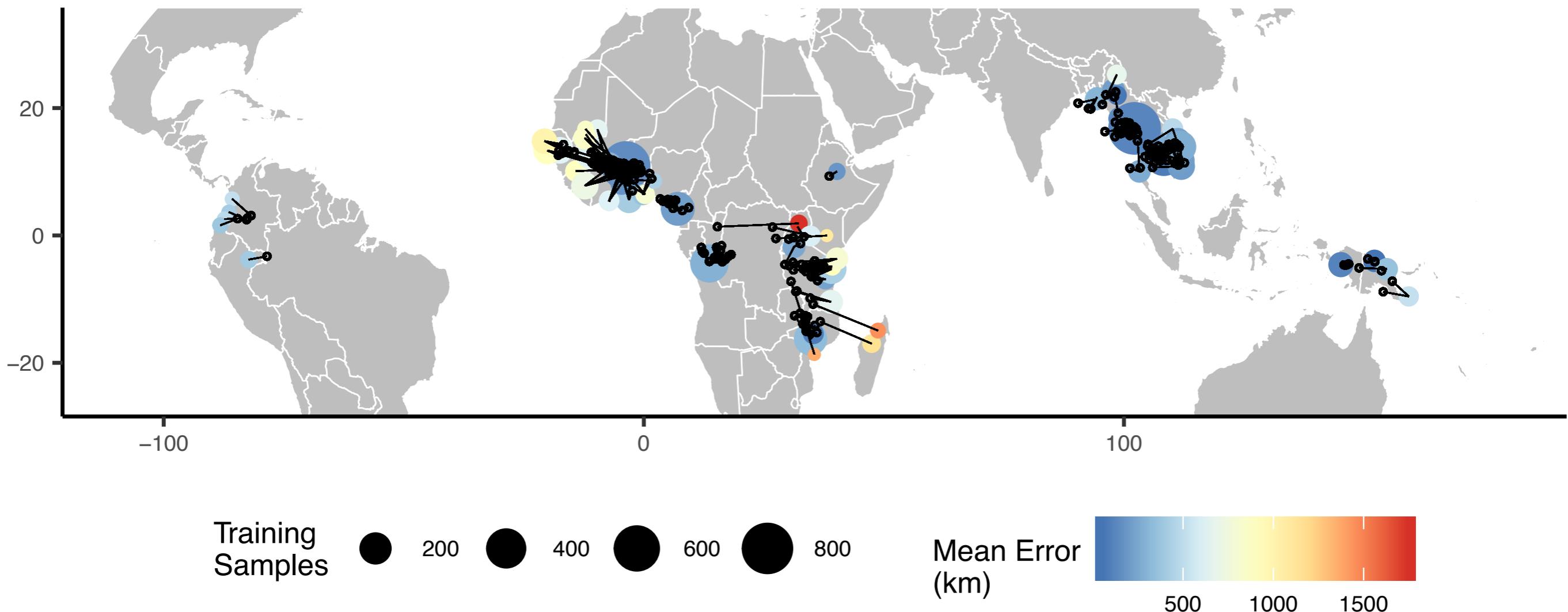
Anopheles gambiae - ag1000g data
median error = 5.7km

Locator – (deep) learning space



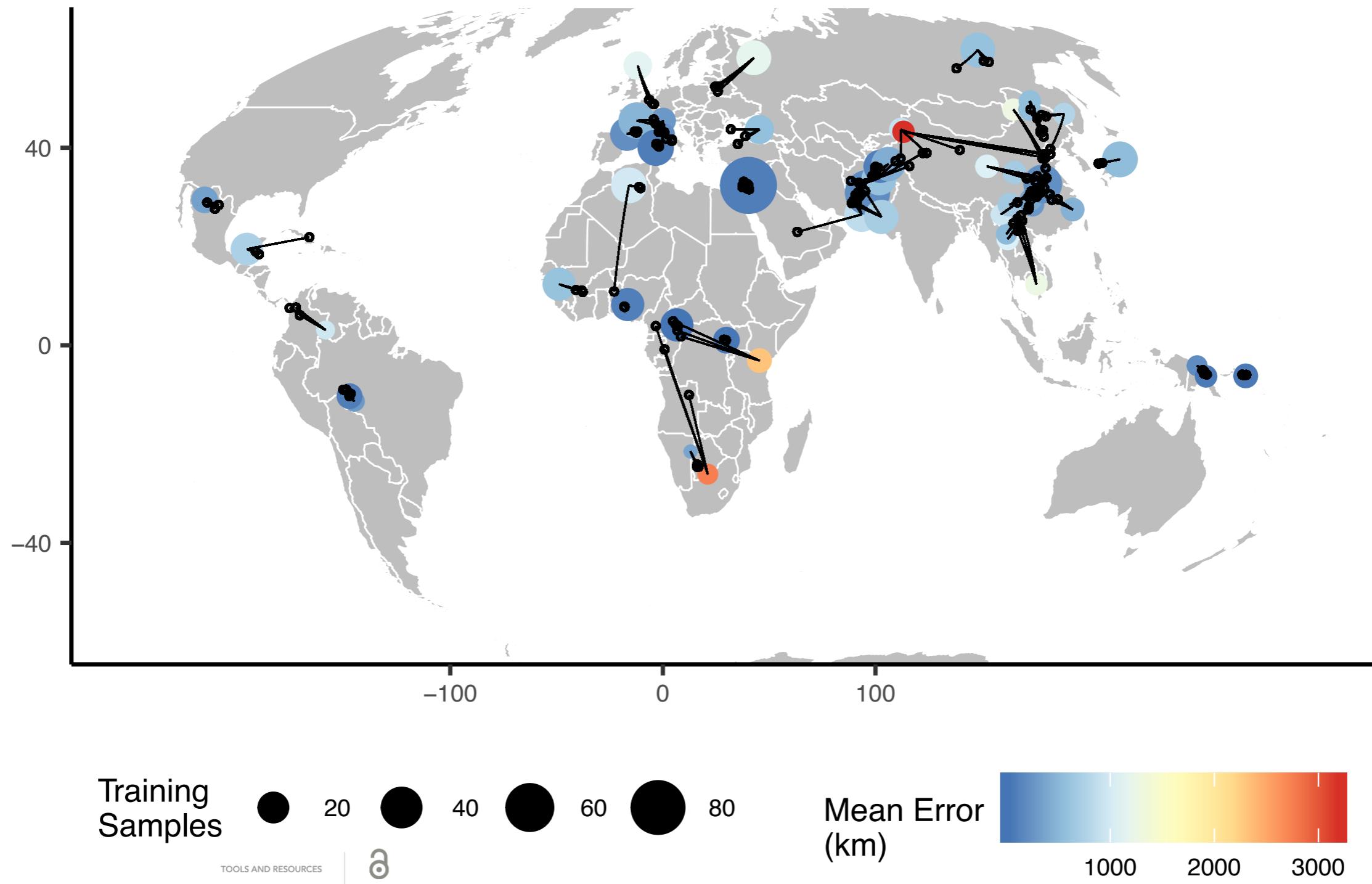
Anopheles gambiae - ag1000g data

Locator— (deep) learning space

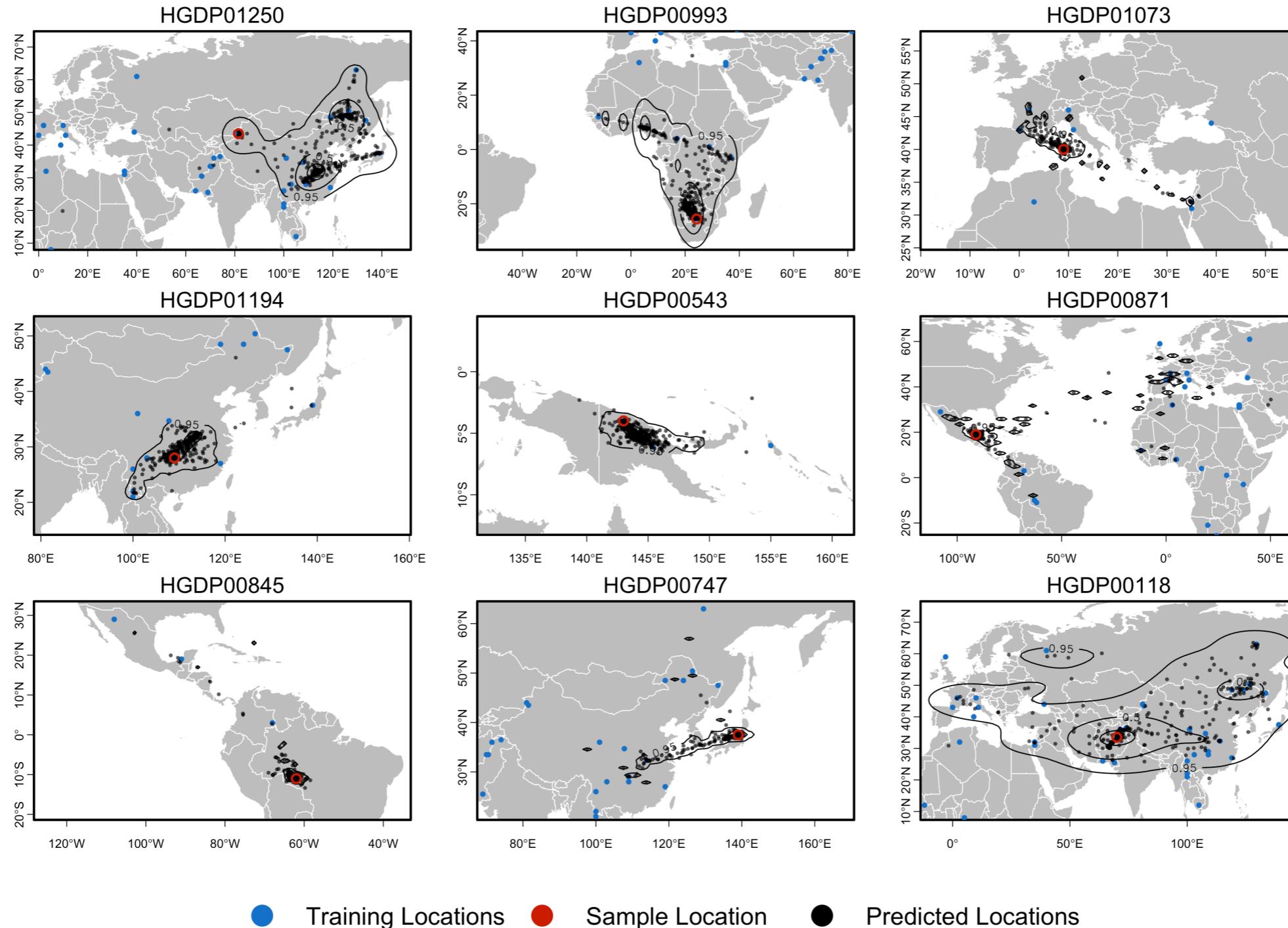


Plasmodium falciparum- Pf7K dataset
median error = 16.9 km

Locator – (deep) learning space

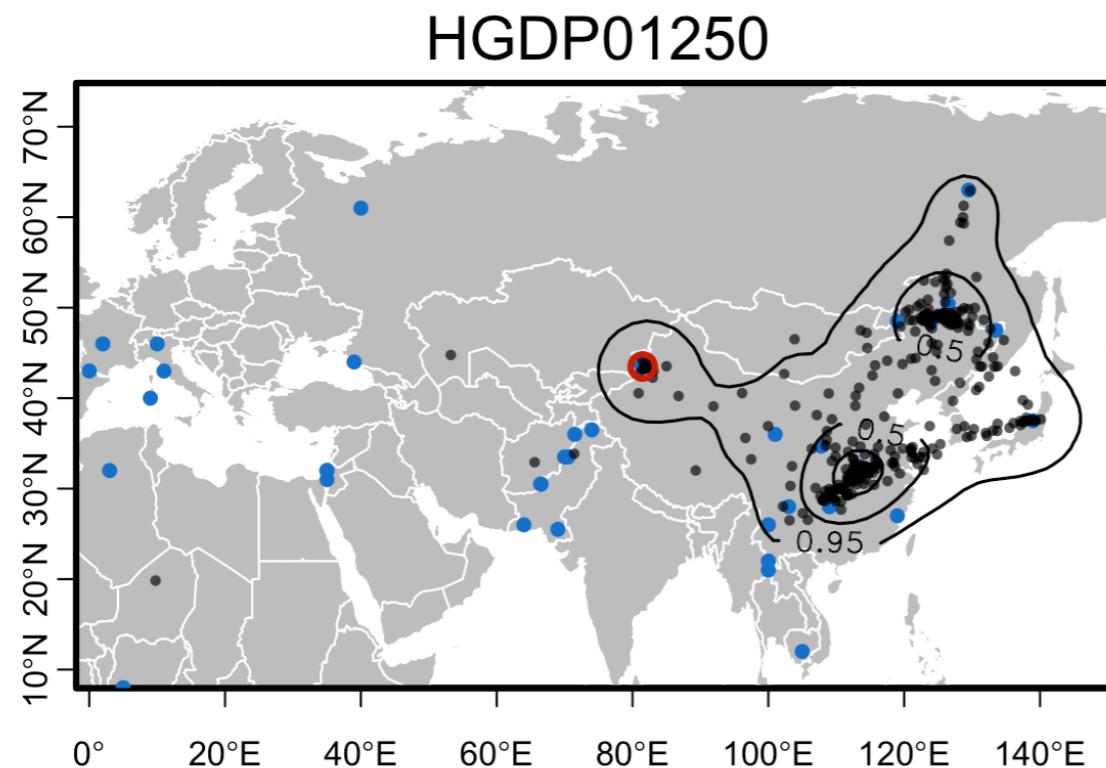


Locator – (deep) learning space

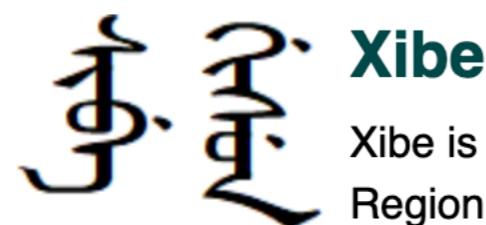


Humans - HGDP

Locator – (deep) learning space



Xibe individual

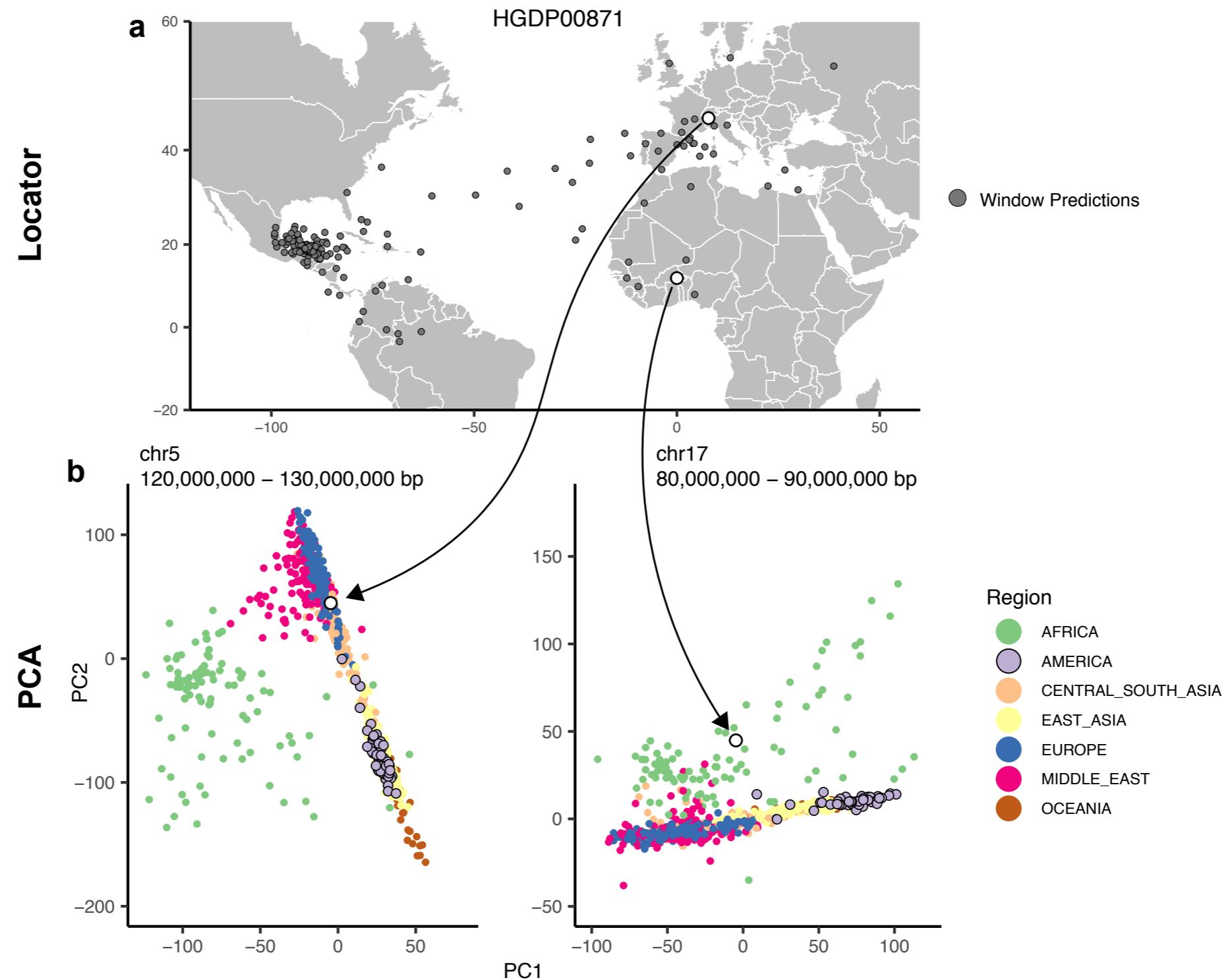


Xibe

Xibe is a Tungusic language spoken in Xinjiang Uyghur Autonomous Region in north west China by about 30,000 people. It is closely related to Manchu, though the Xibe people consider themselves a separate ethnic group. The Xibe were moved to the region in 1764 by the Ch'ing emperor Qianlong. The language is also known as Sibe, Xibo or Sibo.

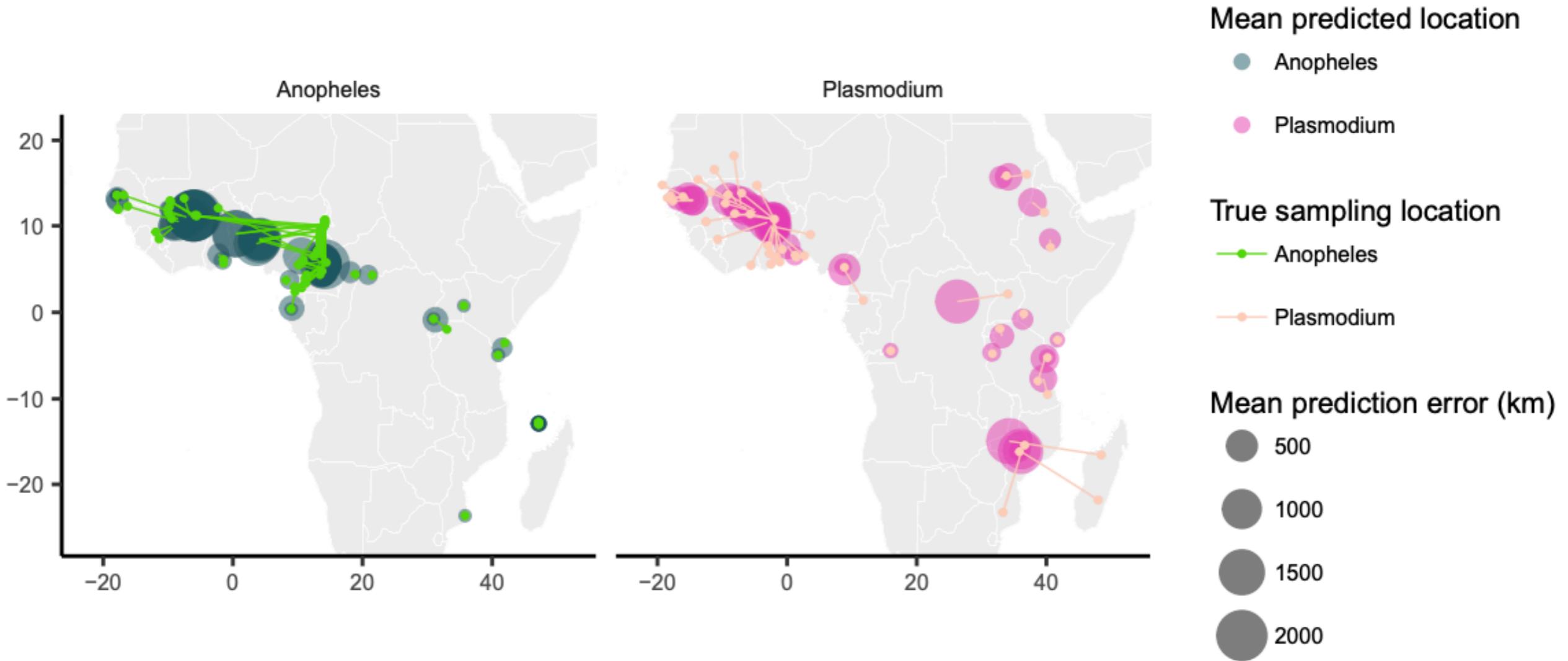
Humans - HGDP

Locator— interpretation



Mayan outlier windows?

Co-geography in a host-parasite system

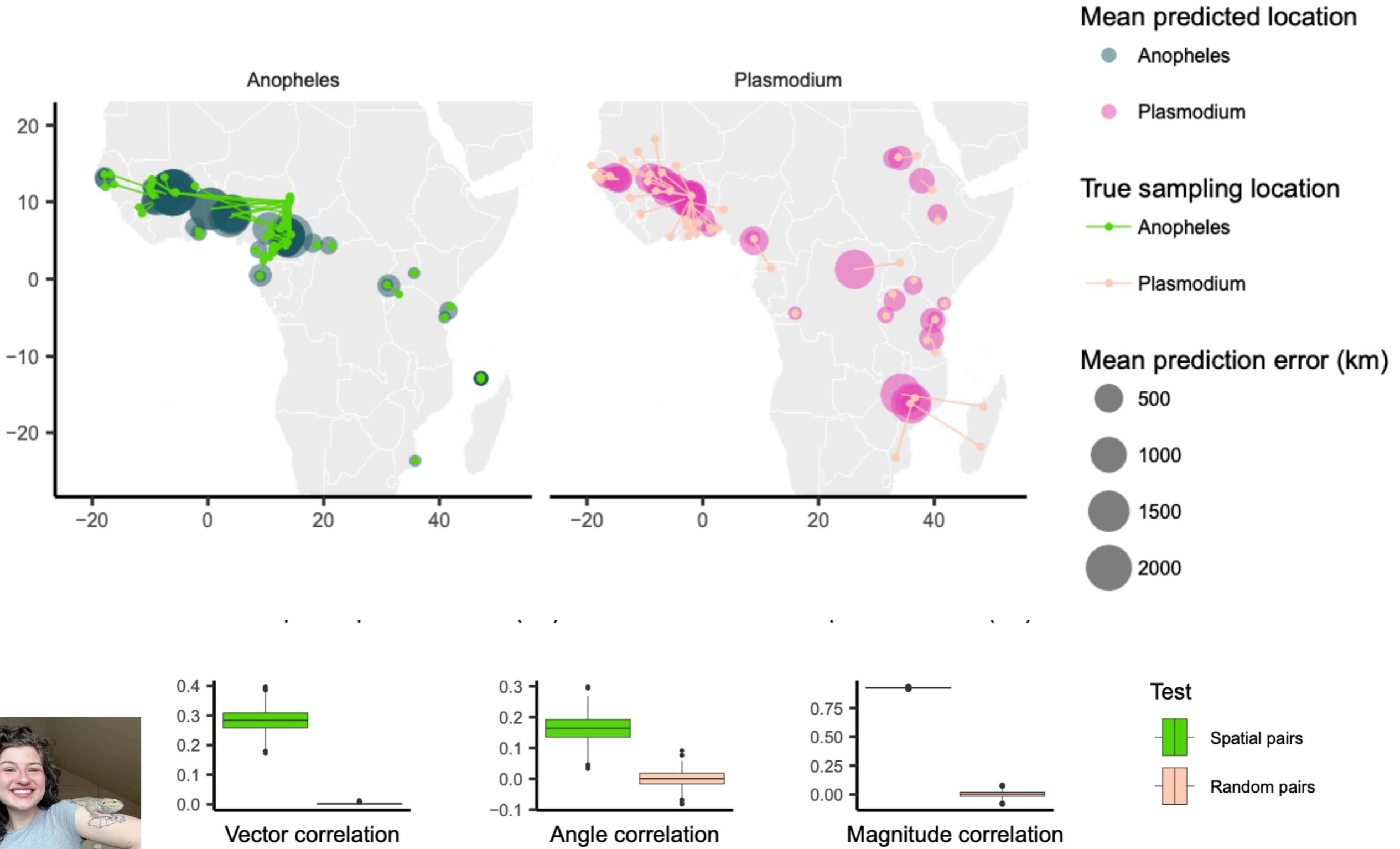


are geographies of Plasmodium
& Anopheles linked?



Clara Rehmann

Interpreting Locator residuals



Interpreting Locator residuals

Facial recognition is accurate if you're a white dude

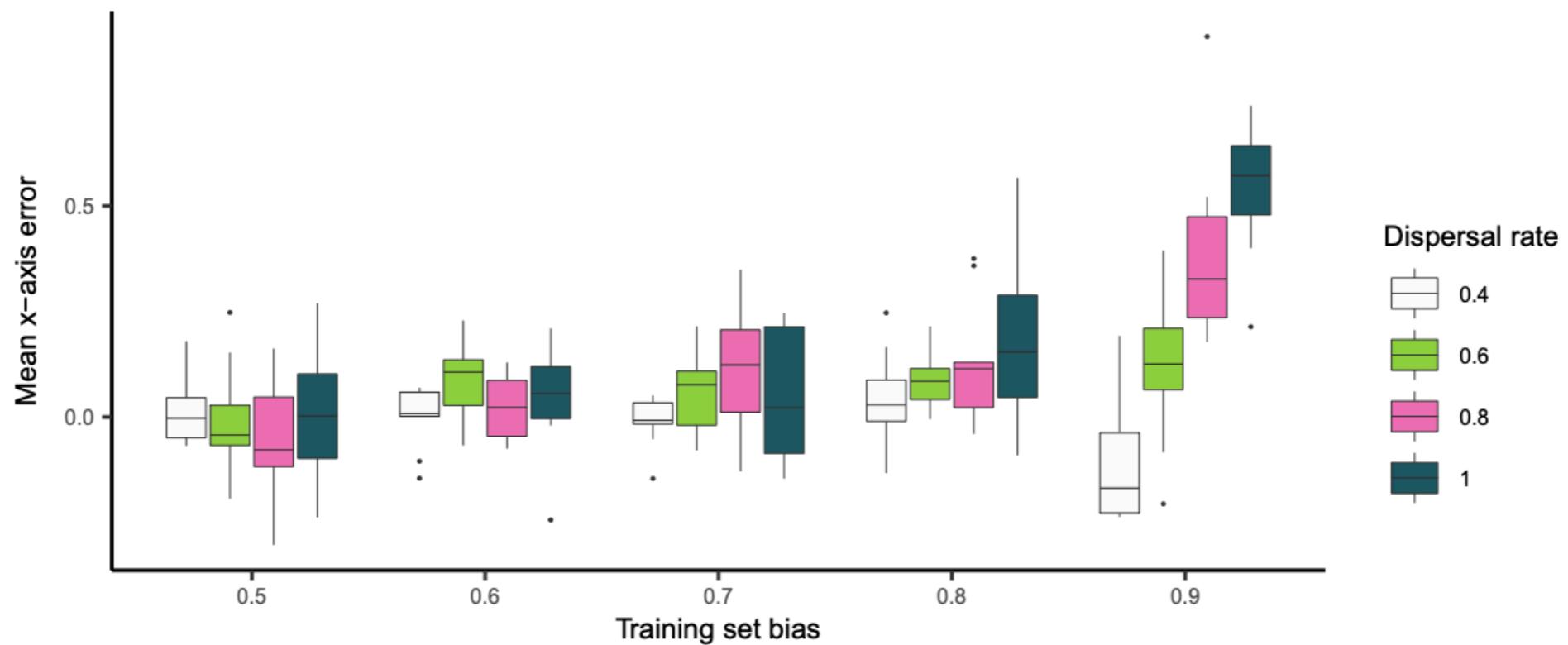
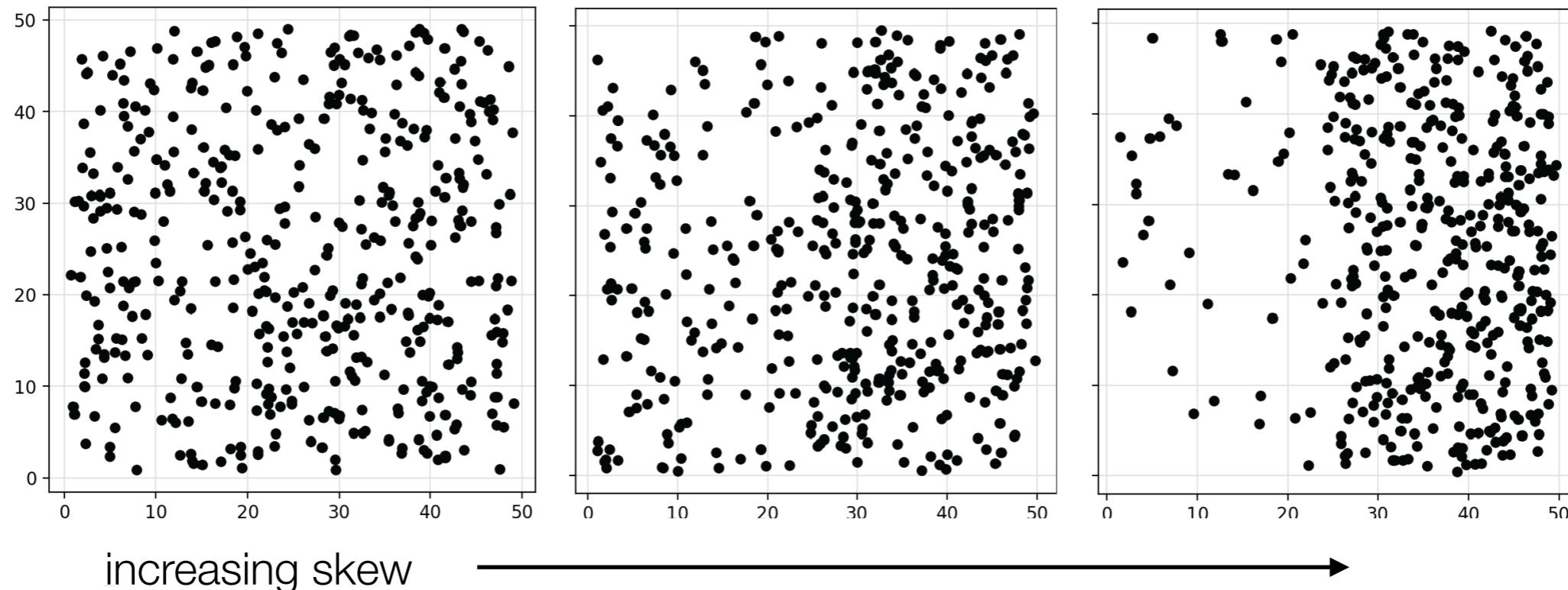


Gender was misidentified in **up to 1 percent of lighter-skinned males** in a set of 385 photos.

Gender was misidentified in **35 percent of darker-skinned females** in a set of 271 photos.

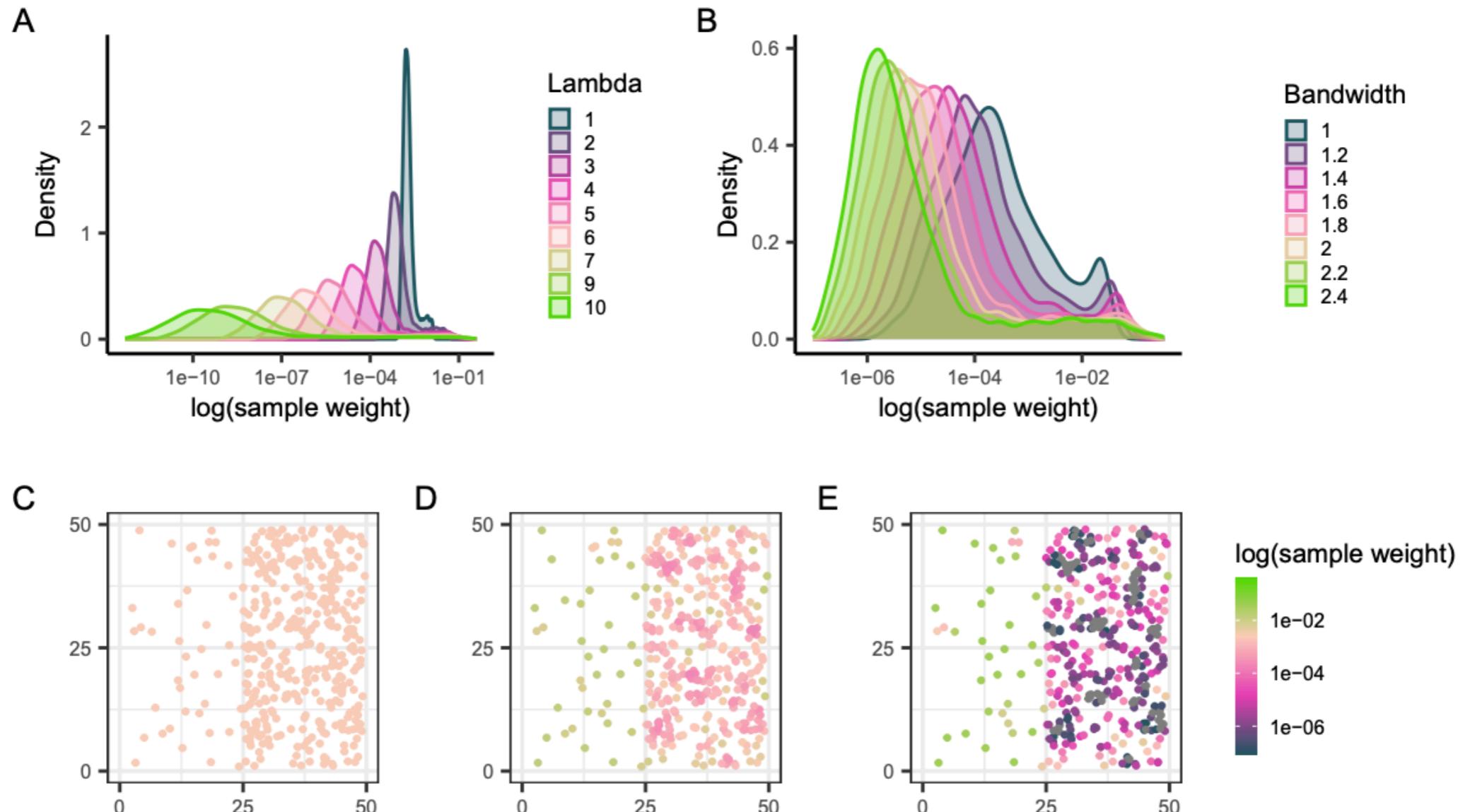
biased training set leads to biased predictions

Interpreting Locator residuals



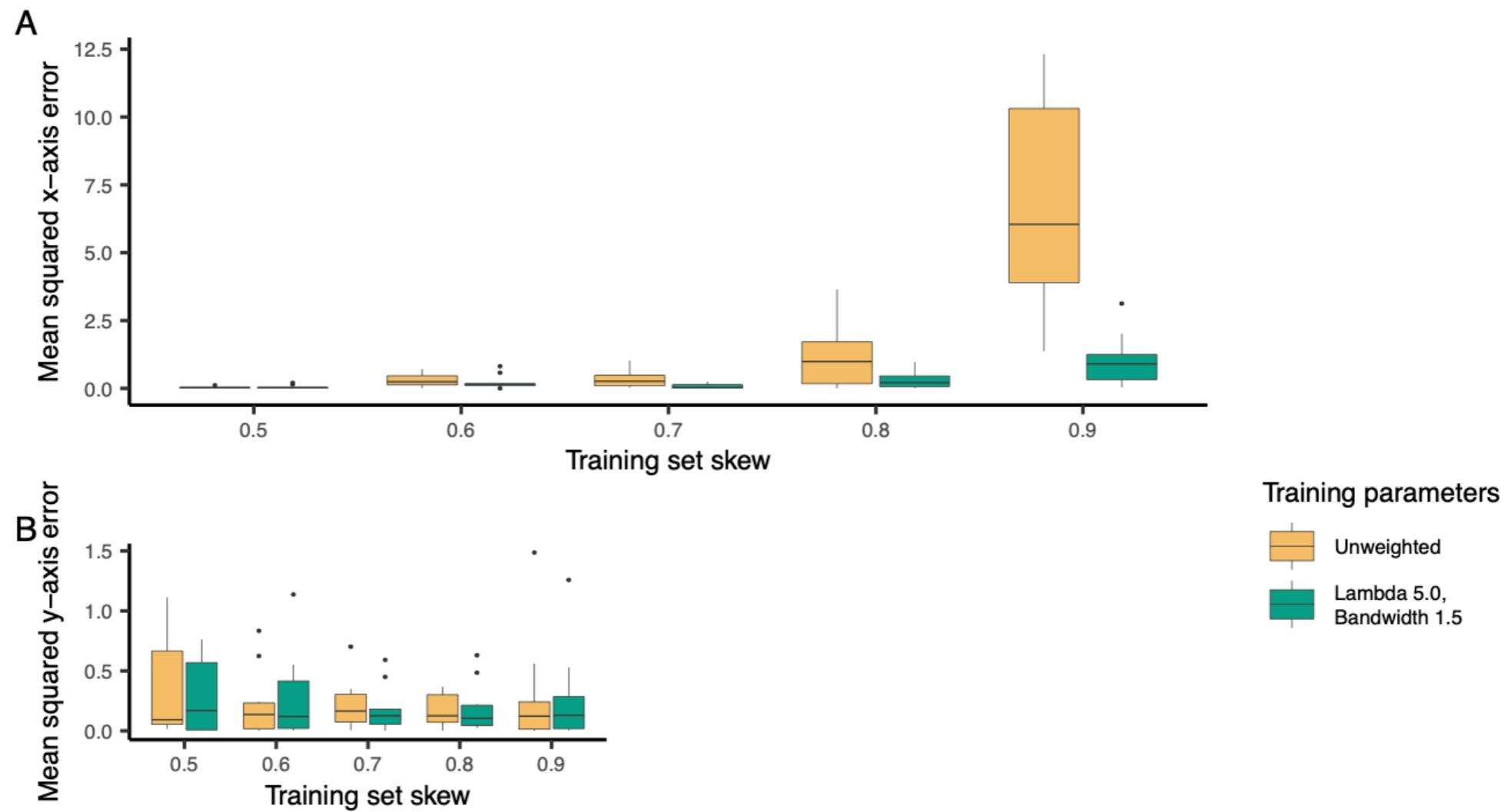
biased density of samples in training set leads to biased predictions

Interpreting Locator residuals



Can correct bias by weighting training loss function

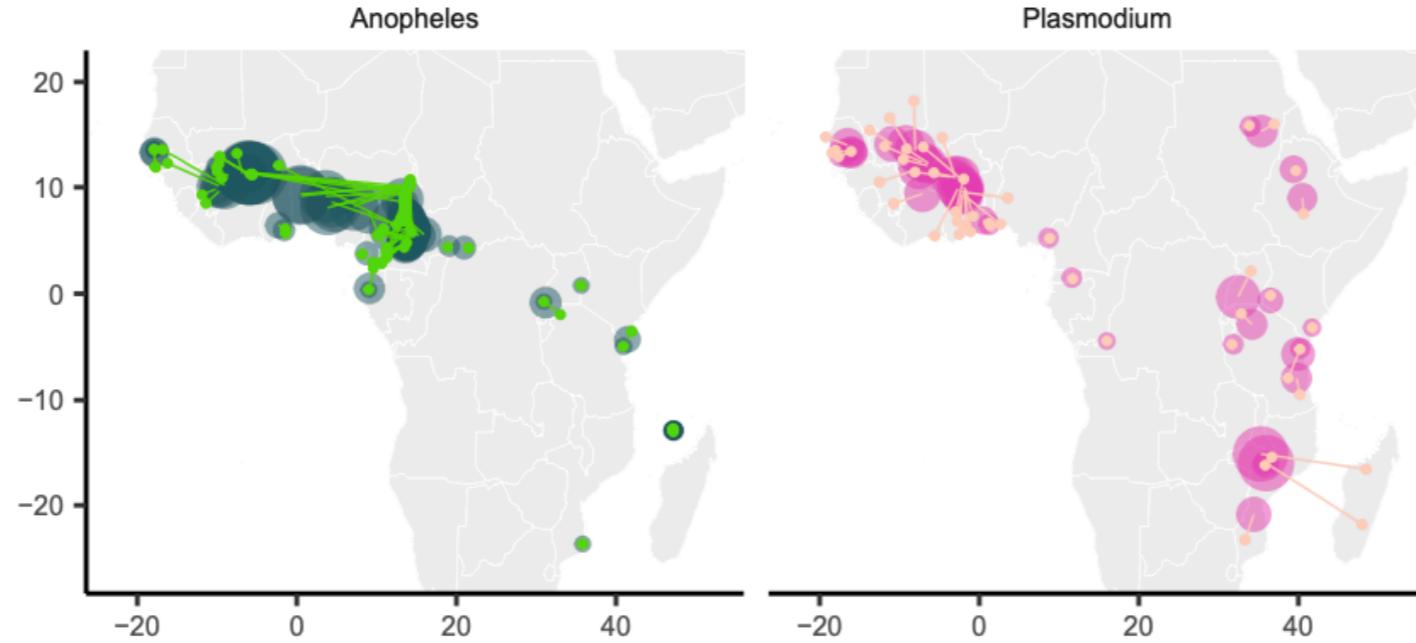
Interpreting Locator residuals



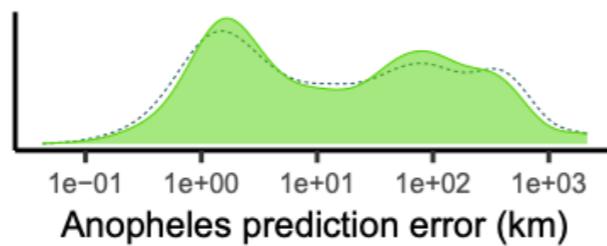
Can correct bias by weighting training loss function

Co-geography in a host-parasite system

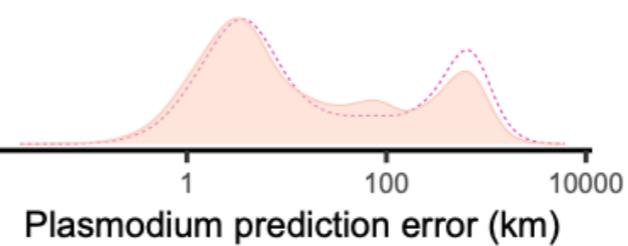
A



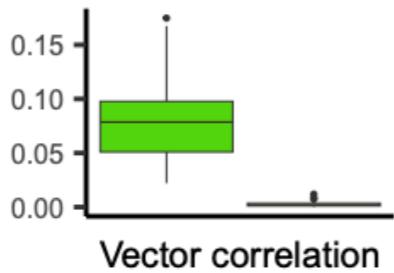
B



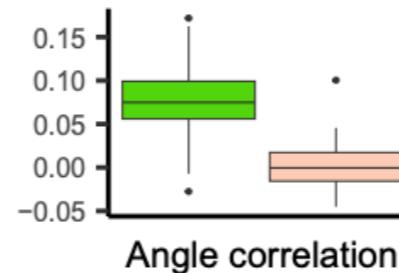
C



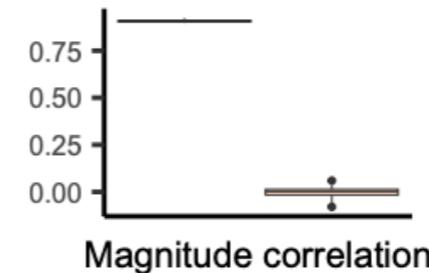
D



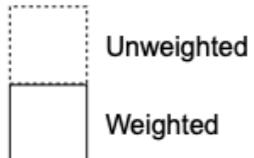
E



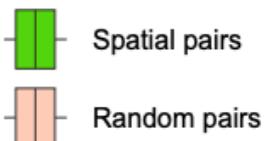
F



Training



Test



After correction, correlation remains!

Clara Rehmann

Co-geography in a host-parasite system

Evaluating evidence for co-geography in the *Anopheles–Plasmodium* host–parasite system

3

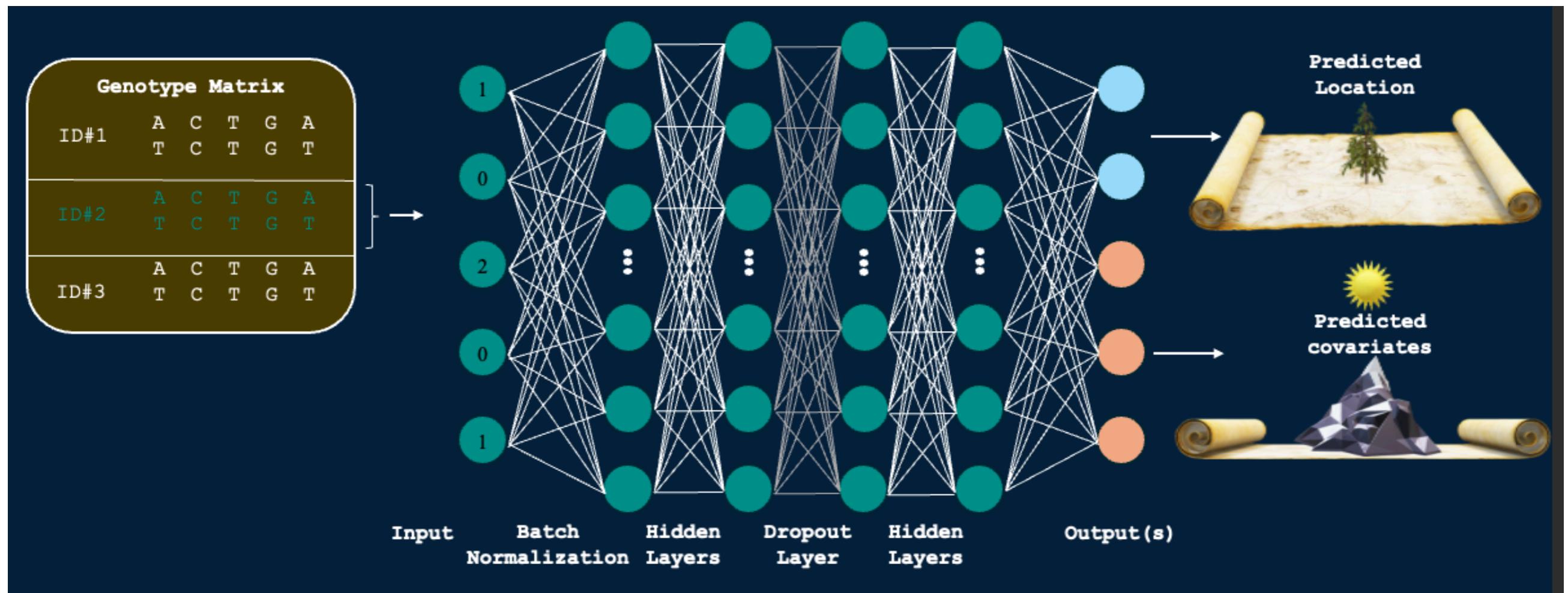
Clara T Rehmann , Peter L Ralph, Andrew D Kern Author Notes

G3 Genes|Genomes|Genetics, Volume 14, Issue 3, March 2024, jkae008,
<https://doi.org/10.1093/g3journal/jkae008>



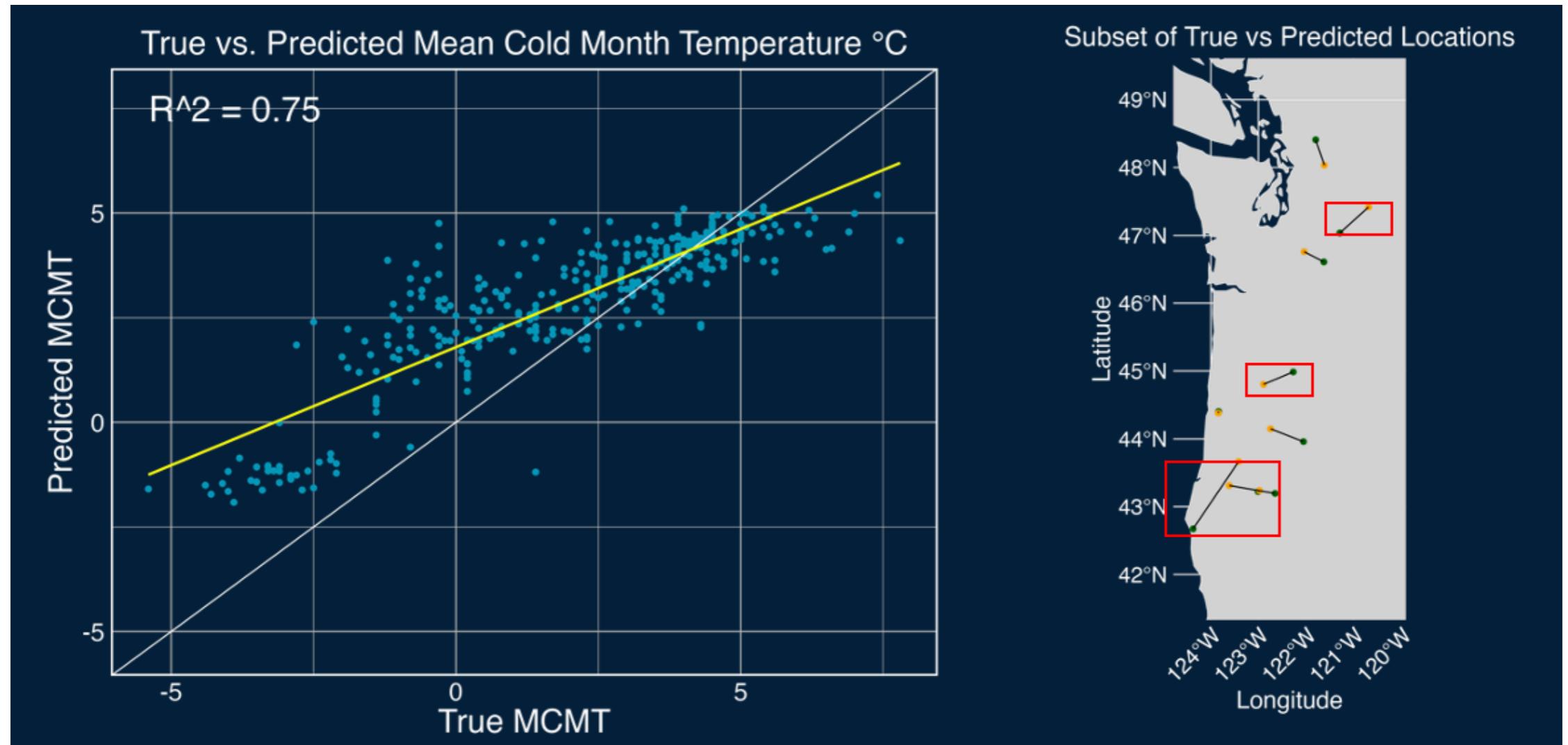
Clara Rehmann

Ecolocator



Extending the locator model to local adaptation

Ecolocator

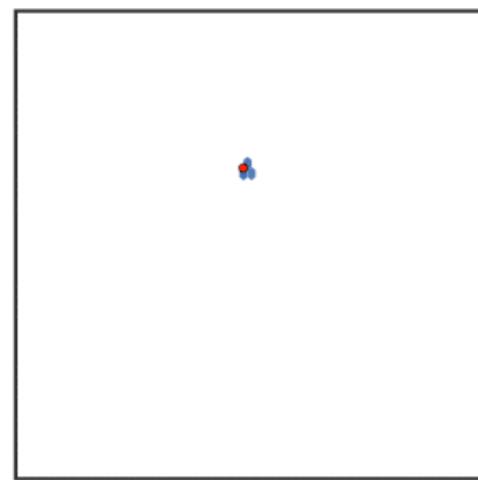


Extending the locator model to local adaptation

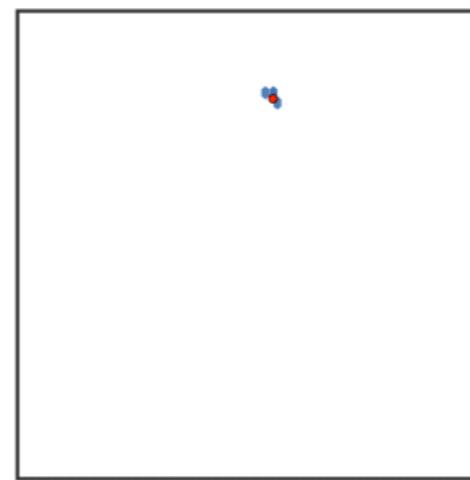
Estimating σ

Genealogical Ancestors by Neighborhood Size

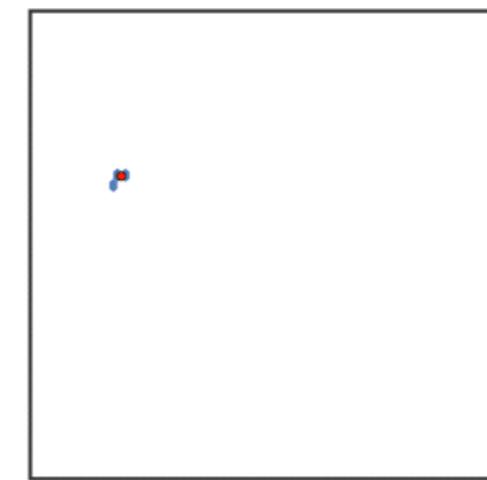
sigma=0.2
NS=3



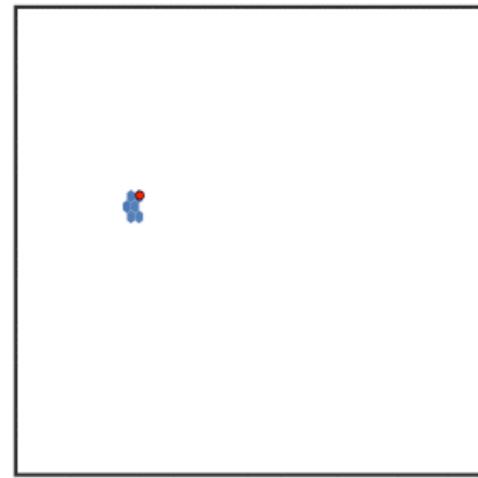
sigma=0.35
NS=8



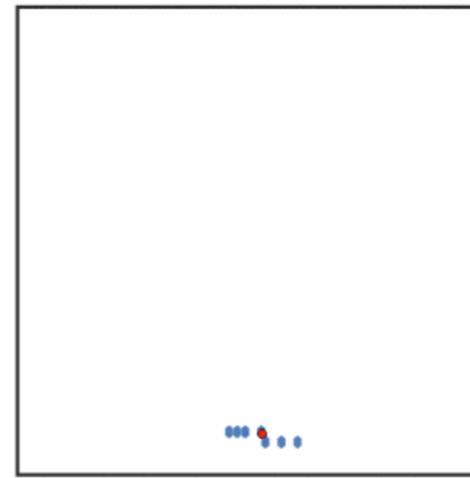
sigma=0.5
NS=16



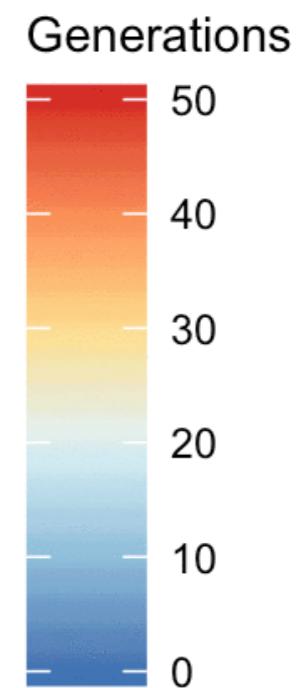
sigma=1
NS=63



sigma=2
NS=251



sigma=3.5
NS=770

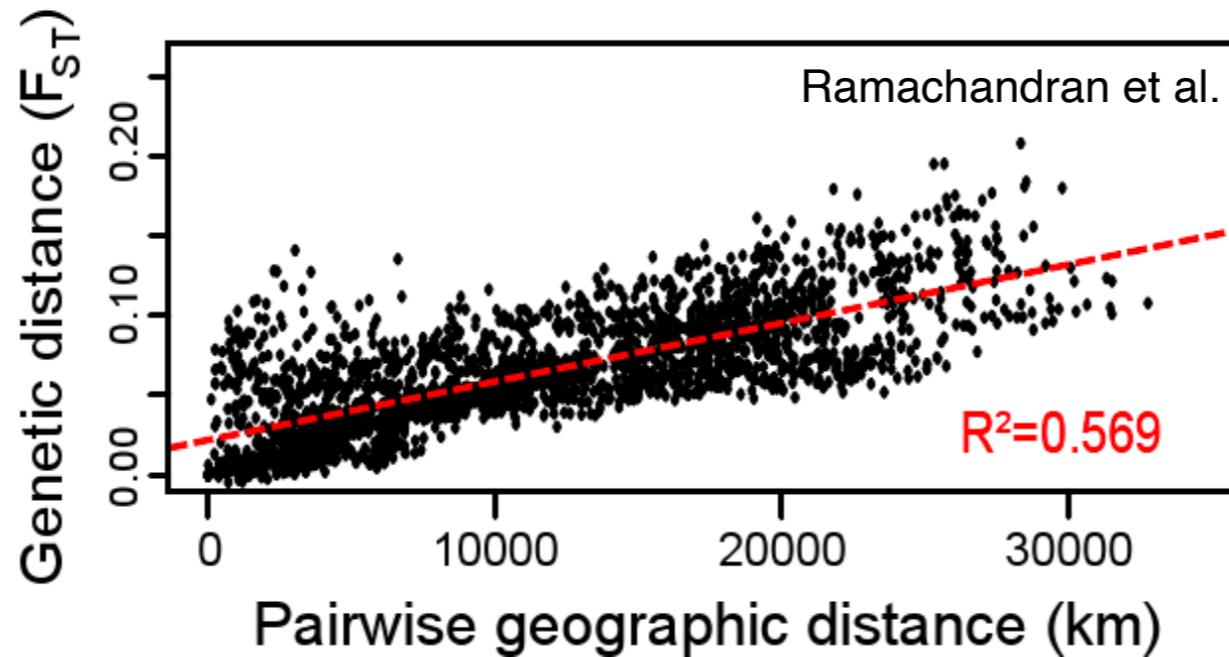


Chris Smith



Estimating σ

Isolation by distance



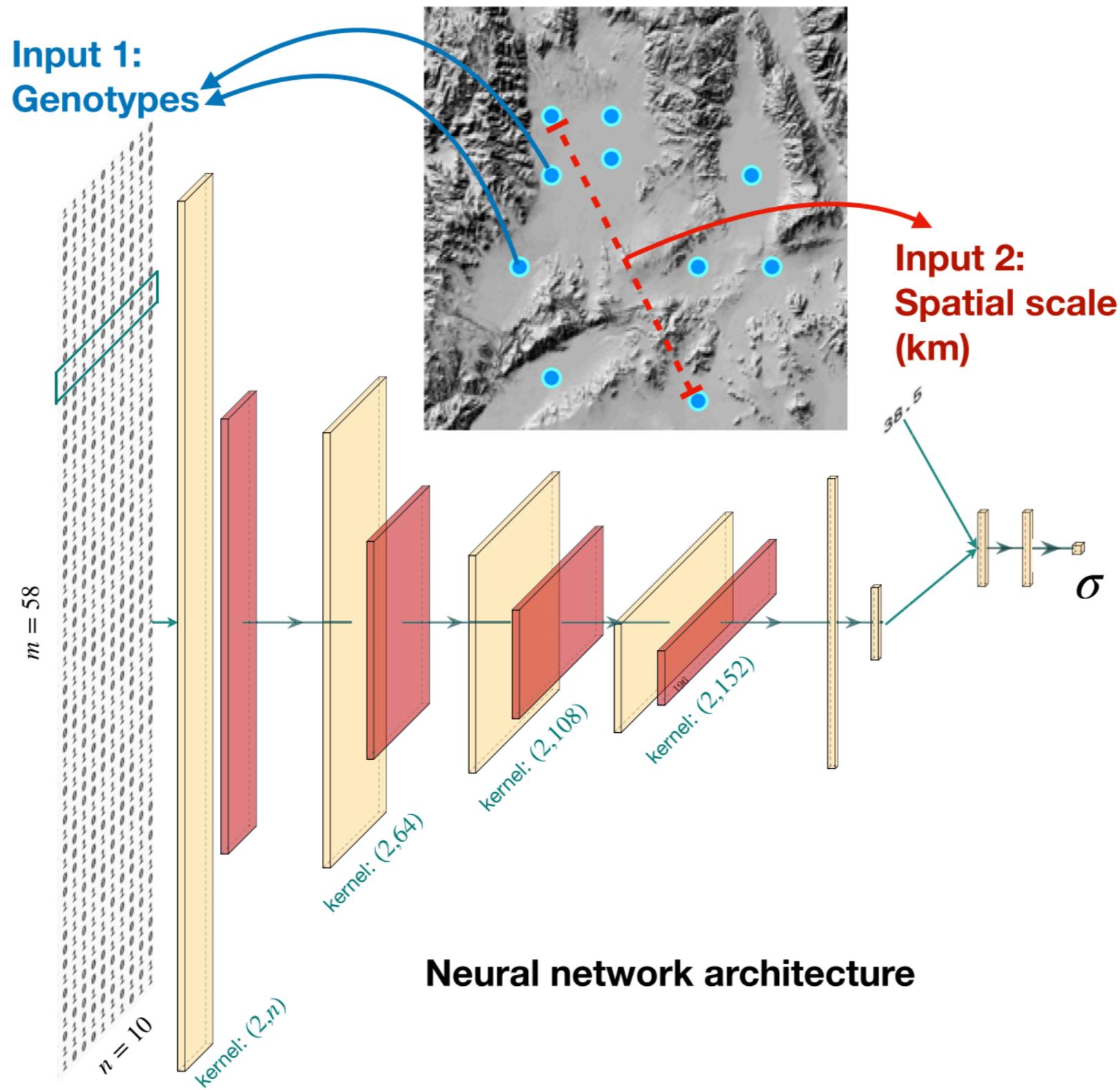
Classic result

Rousset's method – fit this regression, slope is approx $\frac{1}{4N\pi\sigma^2}$

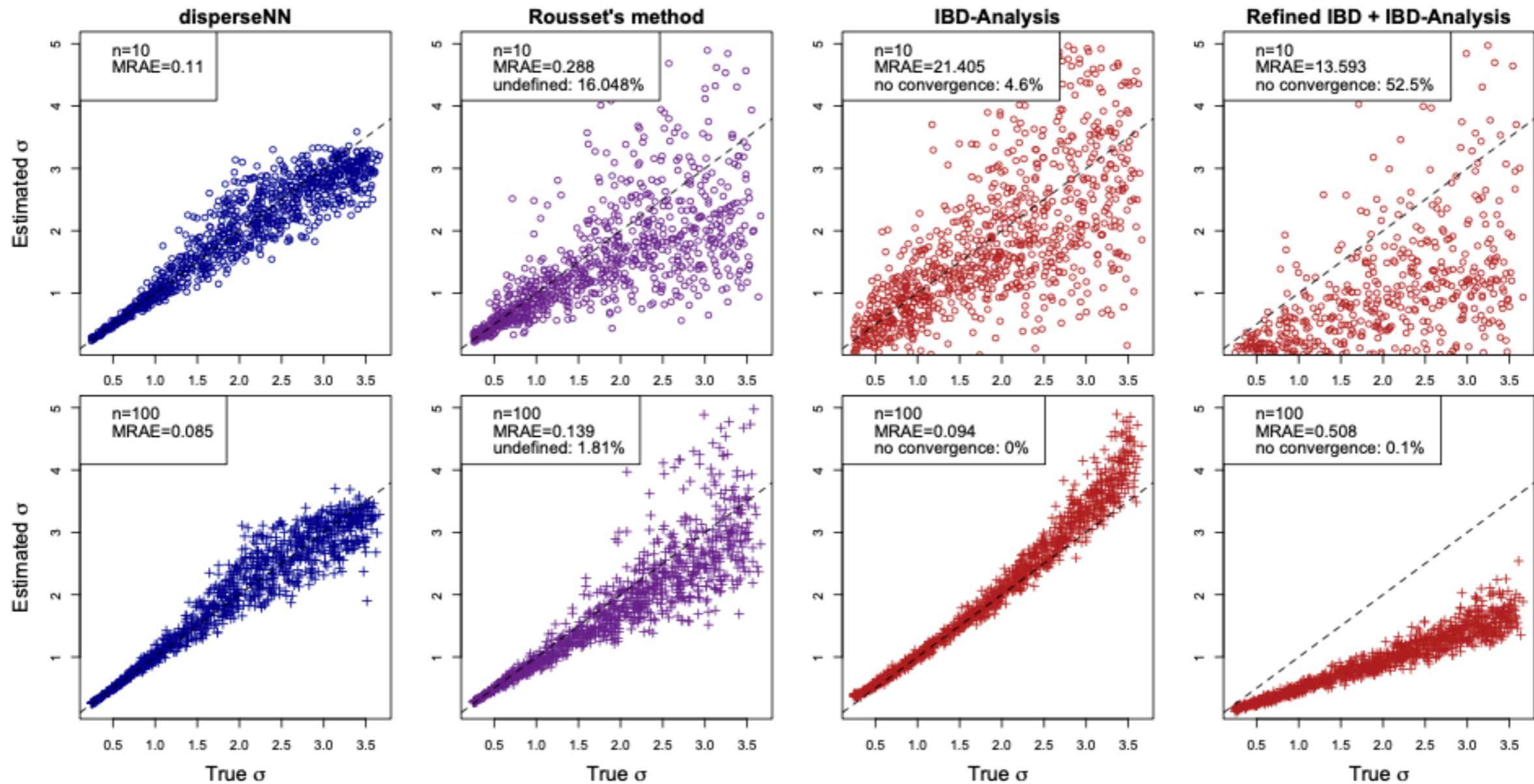
BUT

need to know local N!

disperseNN

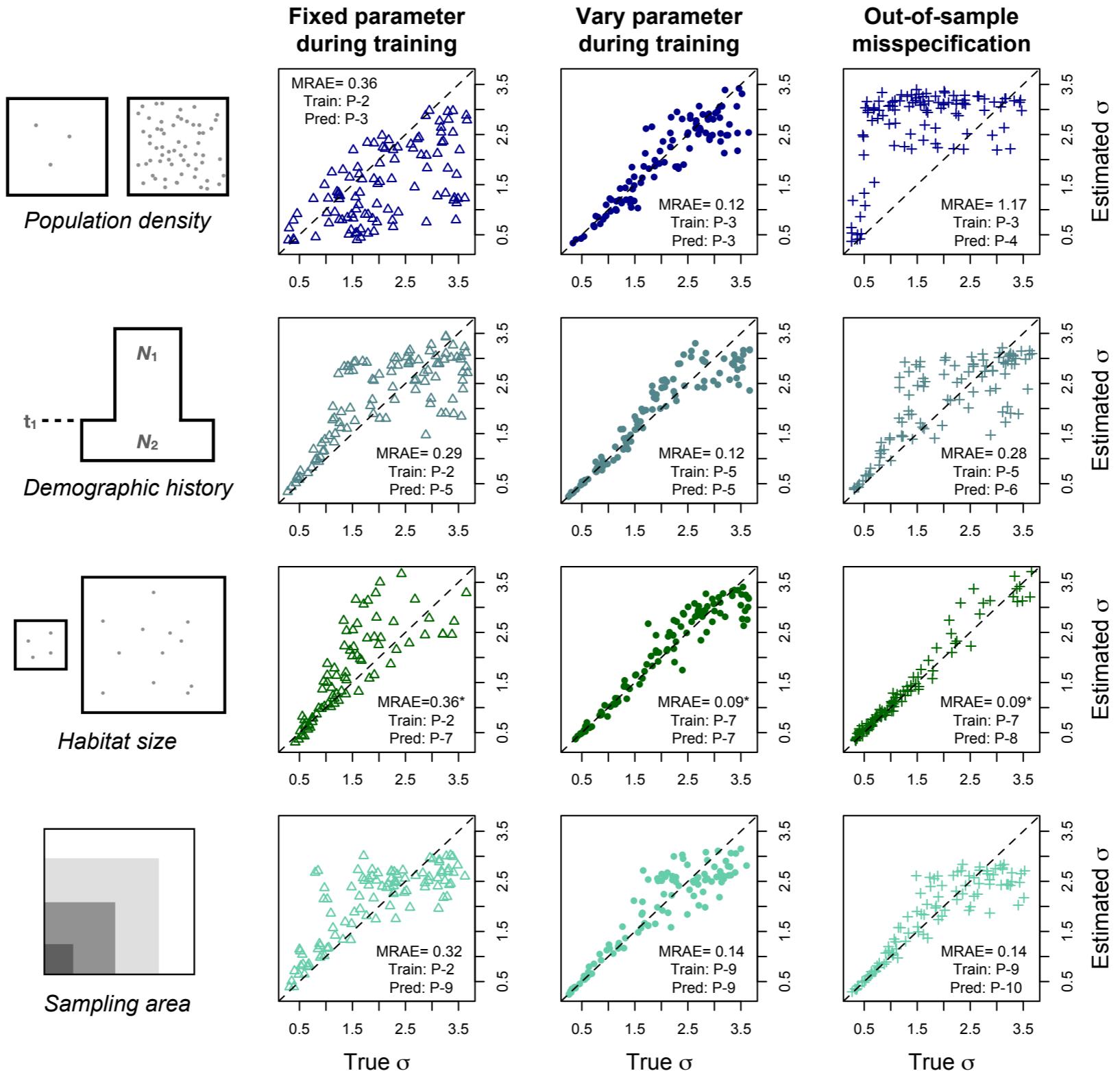


disperseNN



disperseNN works really well, particularly at small sample size
(assume perfect knowledge of N or perfect IBD tract for other me)

disperseNN



disperseNN sensitive to misspecification but can train our way out of it (mostly)

disperseNN

Species	Common name	Region	σ (km)	95% CI (km)	Previous (km)	N_{loc}	n	S (km)	M. dist.
<i>Zosterops borbonicus</i>	Réunion grey white-eye	Réunion	4.97	(1.76, 13.83)	NA	295	41	62	4.59
<i>Peromyscus leucopus</i>	white-footed mouse	New York	0.77	(0.32, 1.67)	0.03-0.11	-231	12	38	8.15
<i>Anopheles gambiae</i>	African malaria mosquito	Cameroon	10.29	(2.00, 48.03)	0.04-0.5	52	29	278	9.62
<i>Bombus bifarius</i>	two-form bumble bee	Washington	14.75	(5.60, 37.28)	1.2-5	1,147	14	273	10.47
<i>Bombus vosnesenskii</i>	yellow-faced bumble bee	California	7.70	(1.21, 38.11)	1.2-5	3,944	18	169	11.83
<i>Hippoglossus hippoglossus</i>	Atlantic halibut	Canada	4.29	(0.71, 33.85)	NA	-5,546	11	193	14.59
<i>Crassostrea virginica</i>	eastern oyster	Canada	1.52	(0.72, 4.31)	21.9	1,435	13	187	19.69
<i>Canis lupus</i>	grey wolf	N. America	15.68	(2.36, 107.3)	98-147	35	13	721	25.42
<i>Helianthus petiolaris</i>	prairie sunflower	Kansas	1.00	(0.39, 3.52)	0.156	9	11	204	45.28
<i>Zosterops olivaceus</i>	Réunion olive white-eye	Réunion	1.05	(0.27, 4.36)	NA	2,392	10	50	45.97
<i>Helianthus argophyllus</i>	silverleaf sunflower	Texas	1.04	(0.38, 4.08)	0.156	57	30	307	86.49
<i>Arabidopsis thaliana</i>	thale cress	Spain	1.36	(0.28, 5.05)	0.001	35	35	80	198.25
<i>Arabidopsis thaliana</i>	thale cress	Sweden	0.44	(0.20, 0.93)	0.001	84	84	325	428.17

Empirical estimates from diverse set of organisms

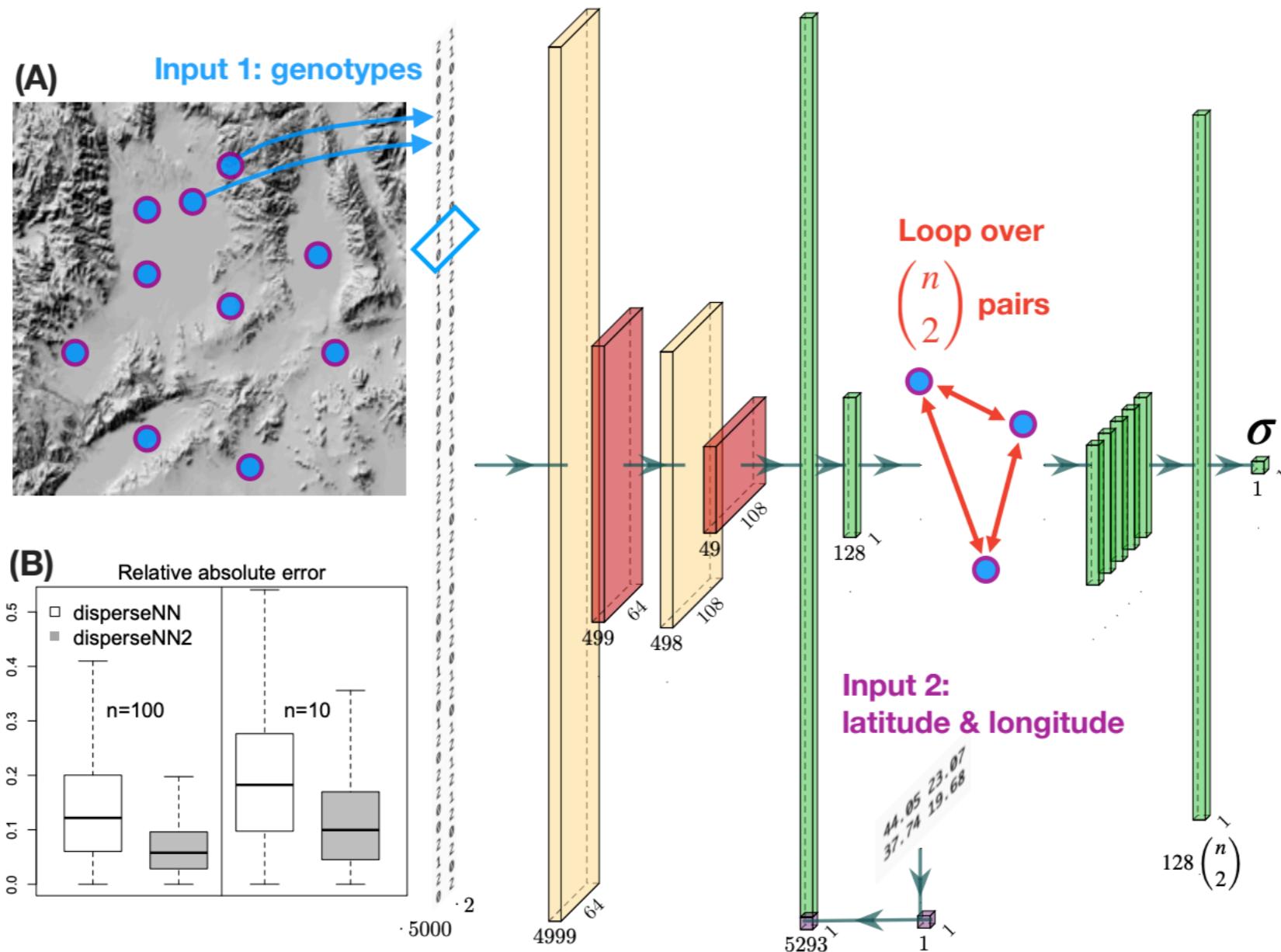
Dispersal inference from population genetic variation using a convolutional neural network

Chris C. R. Smith , * Silas Tittes , Peter L. Ralph , Andrew D. Kern 

Institute of Ecology and Evolution, University of Oregon, Eugene, OR 97403, USA

*Corresponding author: Institute of Ecology and Evolution, University of Oregon, Eugene, OR 97403, USA. Email: chriscs@uoregon.edu

disperseNN2: electric boogaloo



**disperseNN2: a neural network for estimating
dispersal distance from georeferenced polymorphism
data**

Chris C. R. Smith & Andrew D. Kern

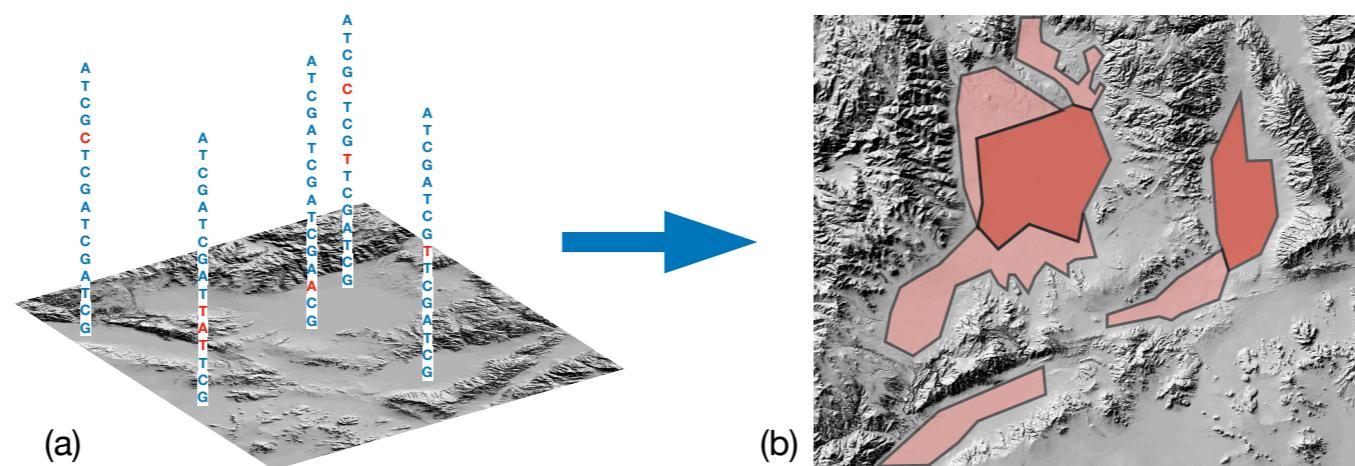
BMC Bioinformatics 24, Article number: 385 (2023) | Cite this article

mapNN



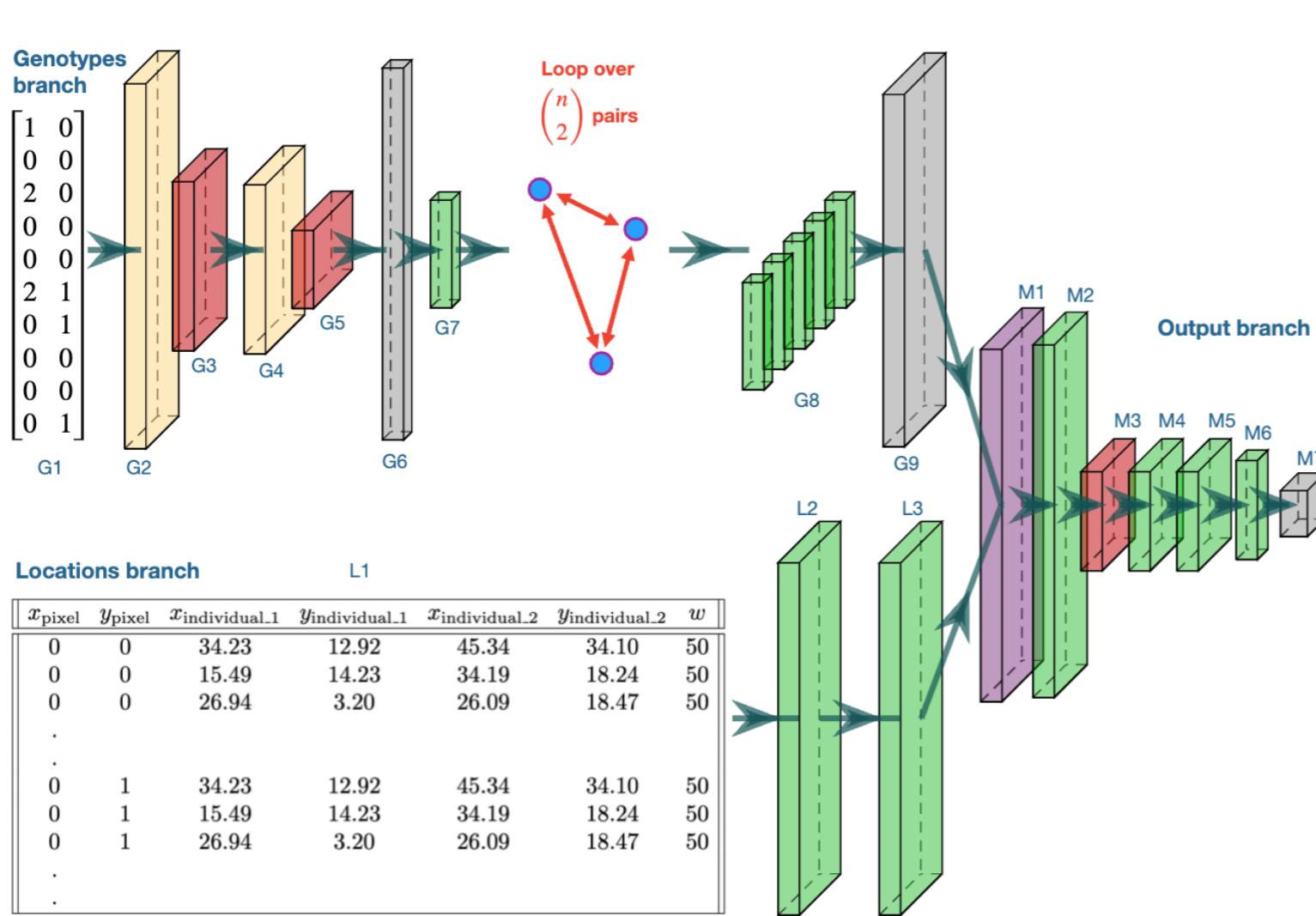
But dispersal need not be homogenous across space!

mapNN

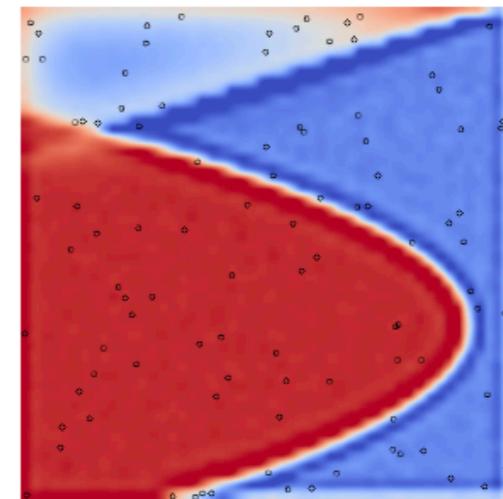
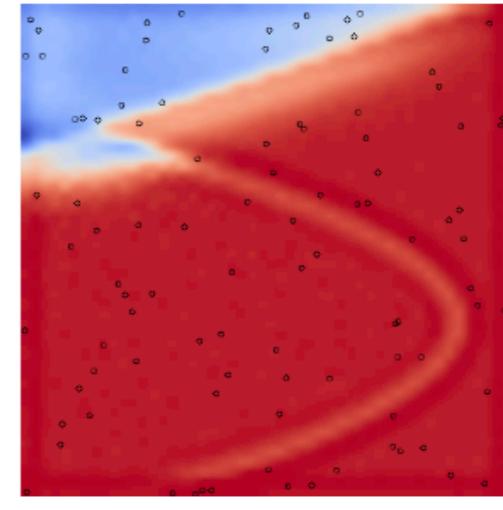


Predicting maps of dispersal and density with a segmentation network

mapNN



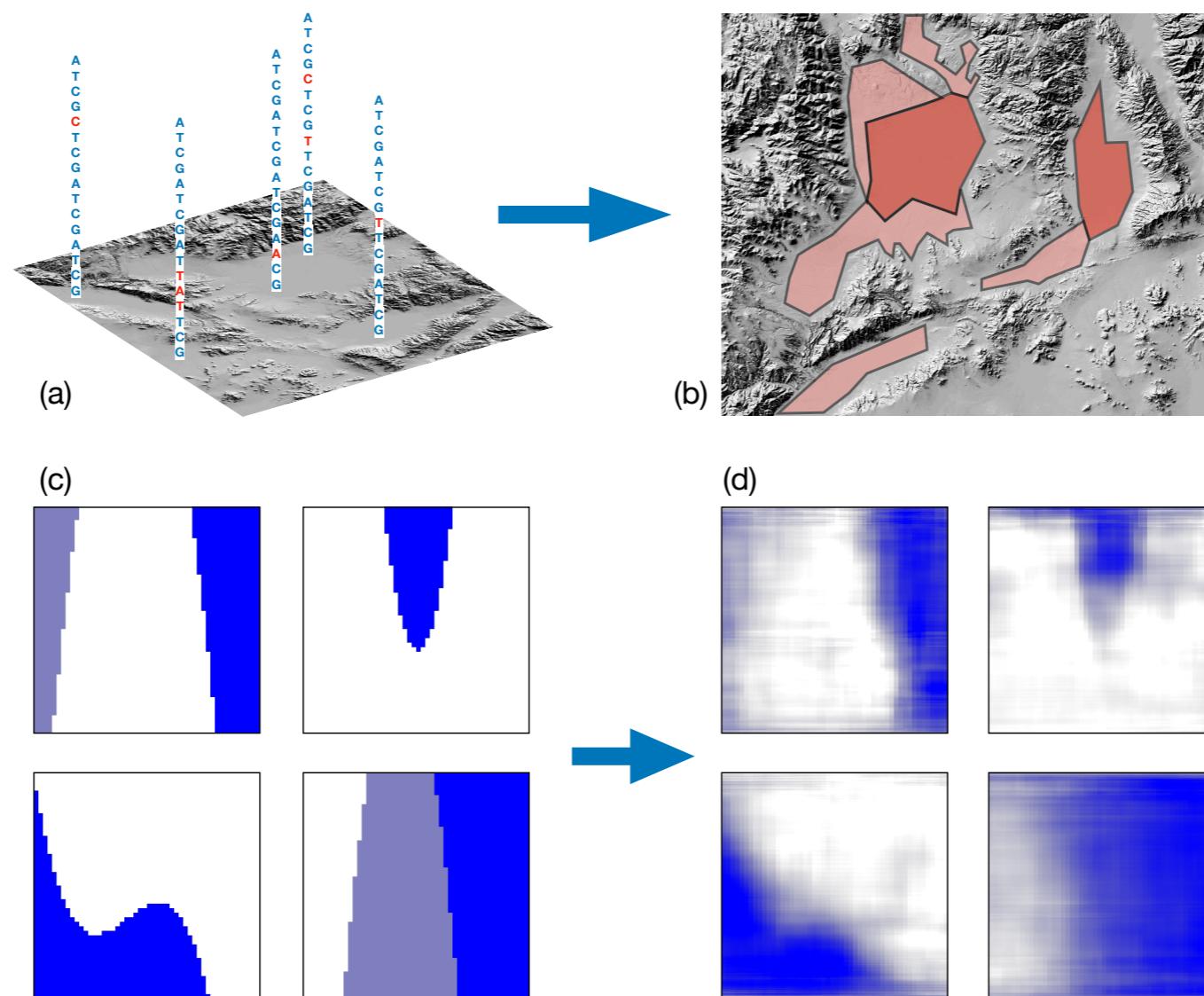
Dispersal



Density

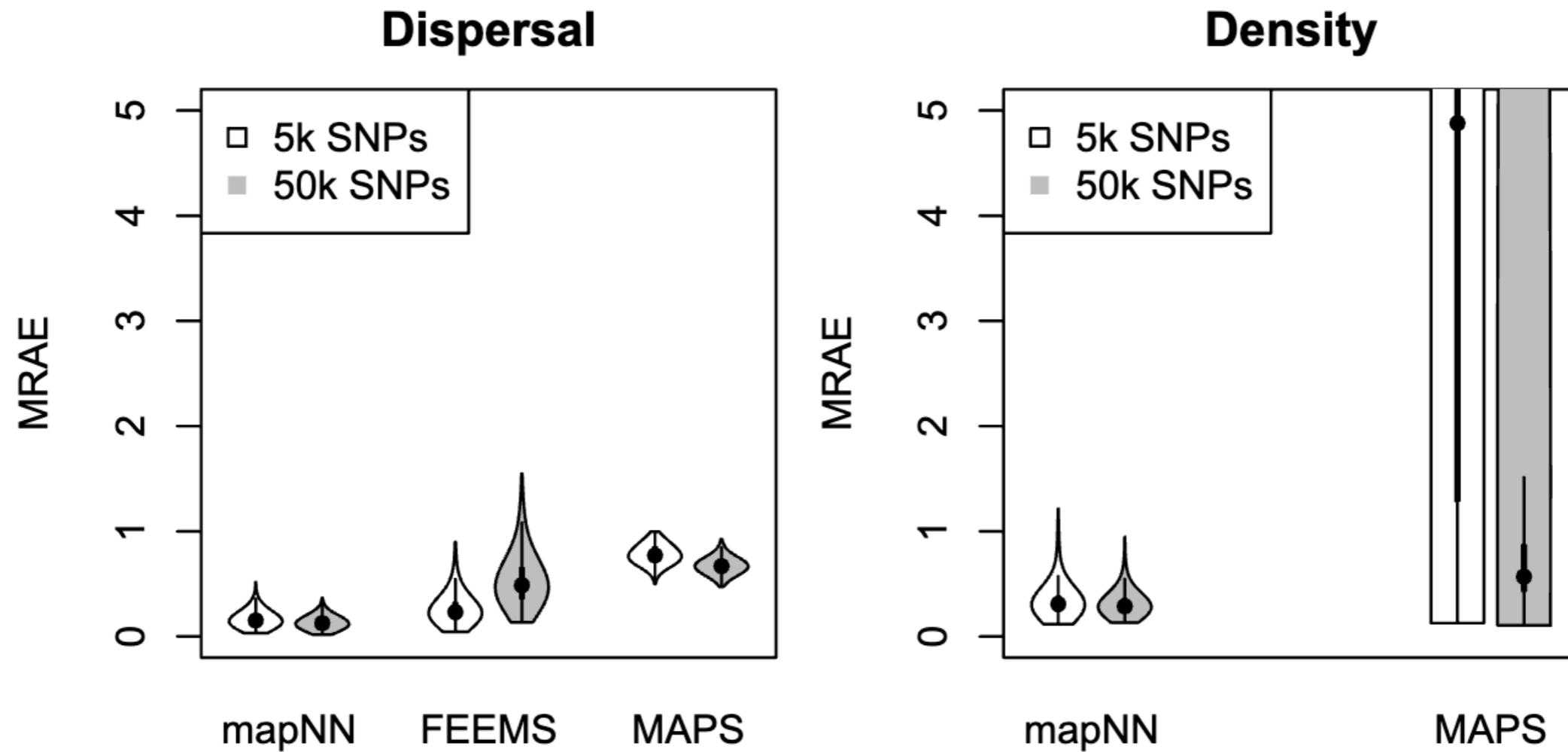
Predicting maps of dispersal and density with a segmentation network

mapNN



Predicting maps of dispersal with an segmentation network

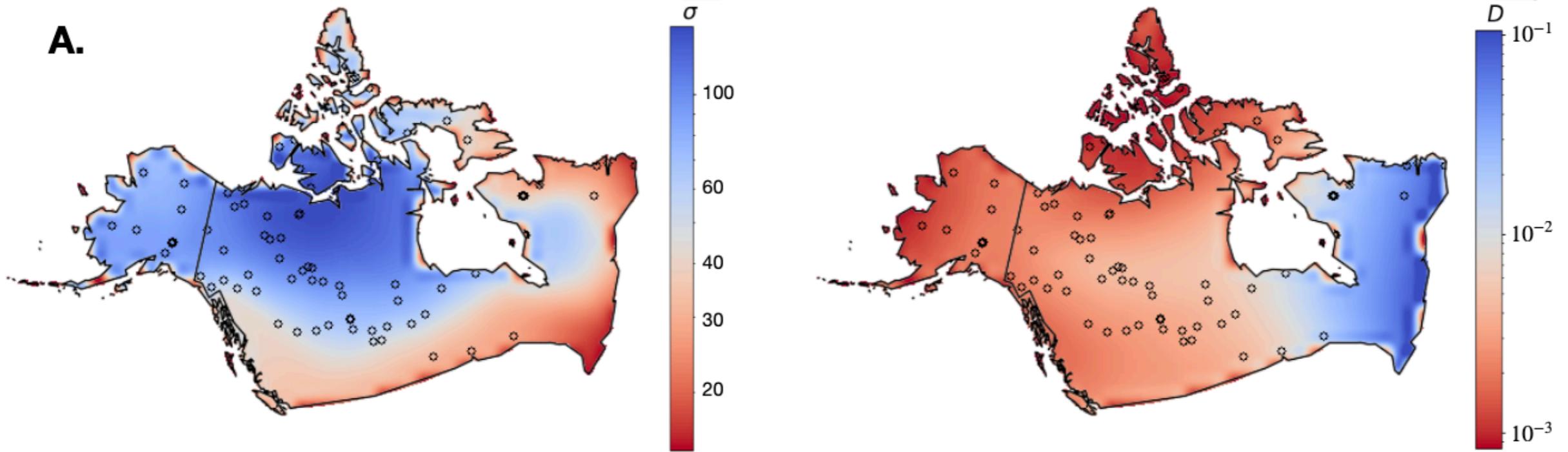
mapNN



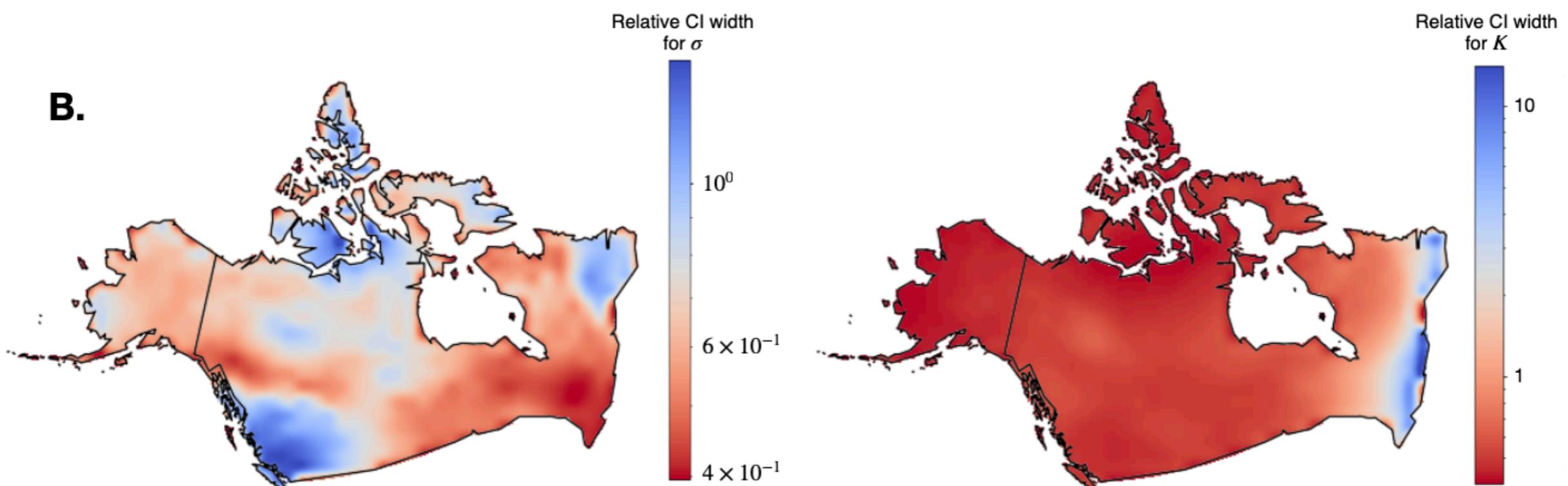
Outcompetes only other existing methods that are based on likelihood

mapNN

A.



B.



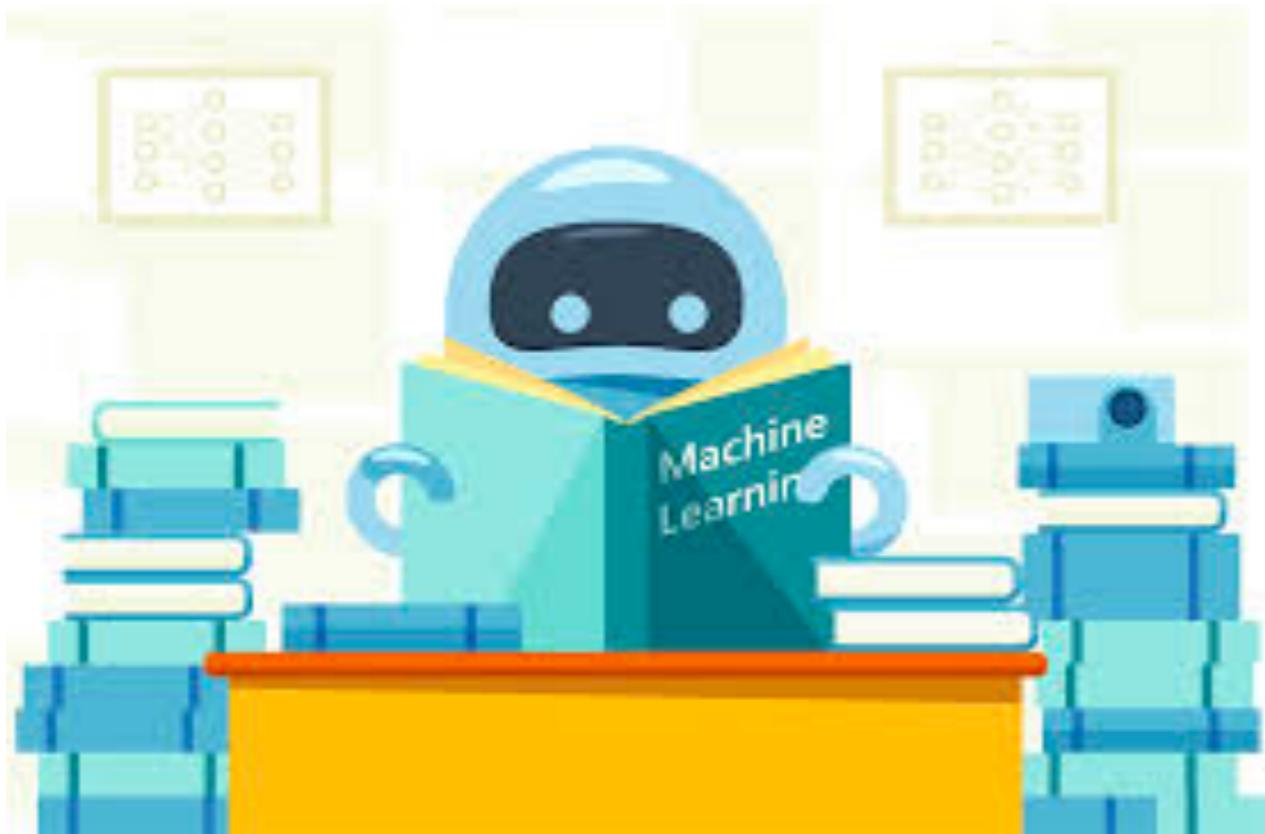
Applied to 97 wolf genotypes collected from North America

Space summary



- Now can perform very realistic, continuous space simulations (Battye et al. 2020)
- Space Matters! Summary stats, GWAS
- Can use DNN with georeferenced genomes for model free inference (Locator / Ecolocator)
- biased training an issue!
- Can estimate dispersal independent of density using a CNNs (disperseNN / disperseNN2)
- mapNN provides SOTA estimates of heterogeneous maps of dispersal/density

take aways



- lots of ways that machine learning can help evolutionary biology / genetics
- huge growth potential— few people working in the area
- many open questions still:
 - best architectures? prob depends
 - way to deal with uncertainty?
 - ways to build DL models that are robust?

Thanks!

Jeff Adrion
CJ Battey
Gabby Coffing
Clara Rehmann
Chris Smith
Jared Galloway

Dan Schrider
Chuck Langley
Peter Ralph

all members of the Kern-Ralph lab

