# Representing long-range genetic similarity on a background of spatially heterogeneous IBD

Vivaswat Shastry, John Novembre
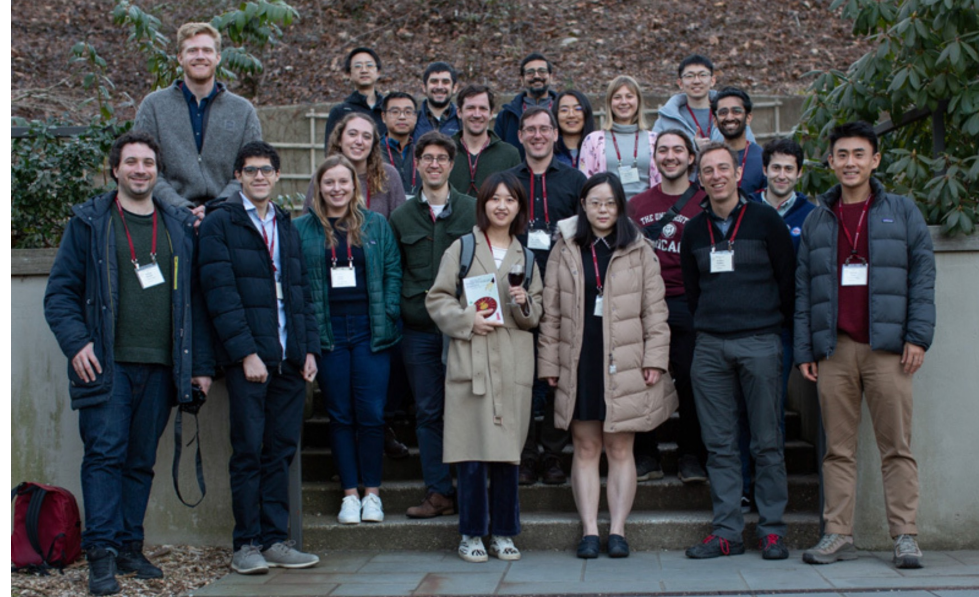
SMBE 2024
PUERTO VALLARTA

# Acknowledgements

Labs:
- Berg
- Novembre
- Steinrücken

Community:

Program in Computational Biology

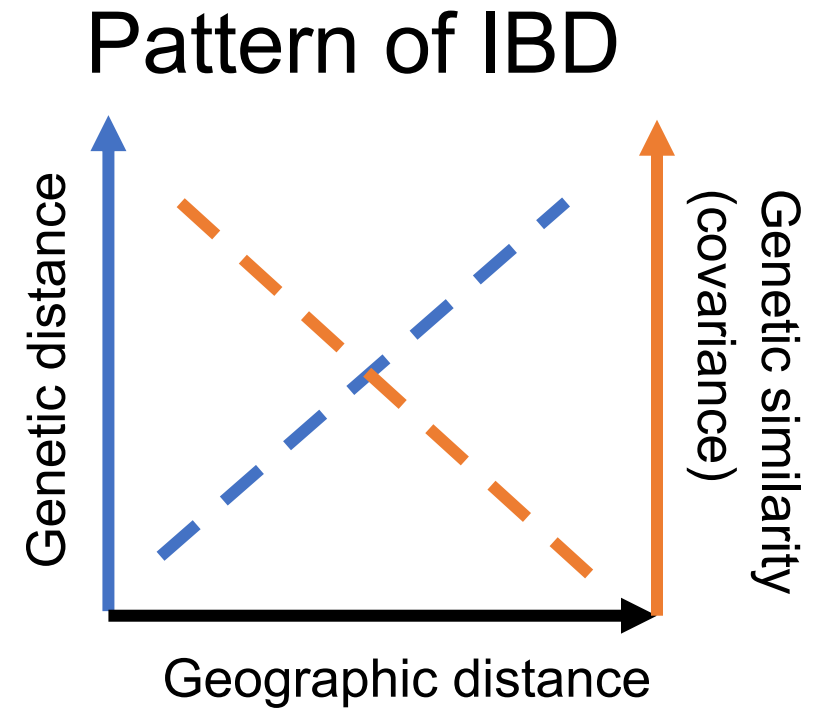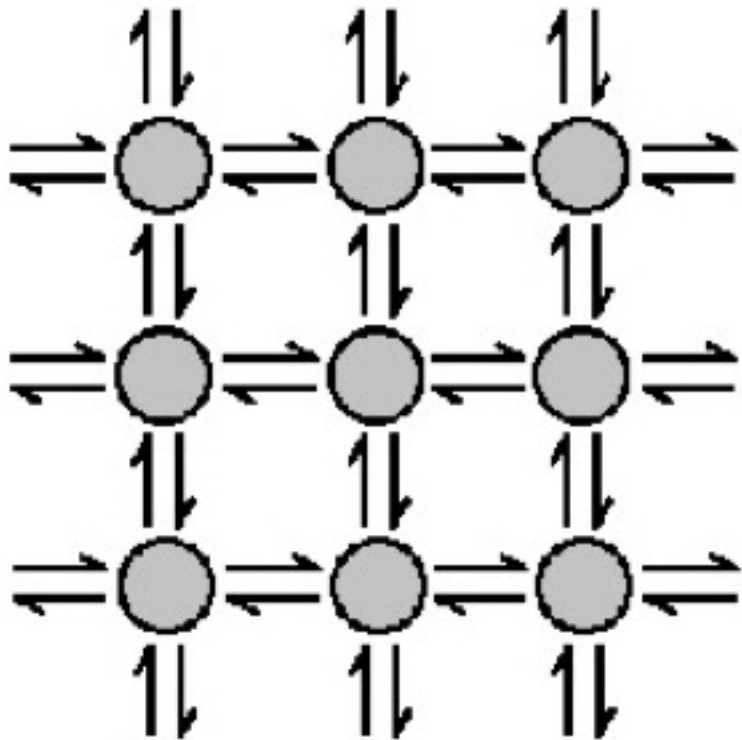Biological Sciences Division

# Outline

- Introduction + Motivation

- Spatially heterogeneous models of isolation-by-distance (IBD)

- Our model for long-range gene flow events (`FEEMSmix`)

- Results:
  1. Simulations
  2. North American grey wolves
  3. Afro-Eurasian panel of humans

# Outline

- **Introduction + Motivation**

- Spatially heterogeneous models of isolation-by-distance (IBD)

- Our model for long-range gene flow events (`FEEMSmix`)

- Results:
    1. Simulations
    2. North American grey wolves
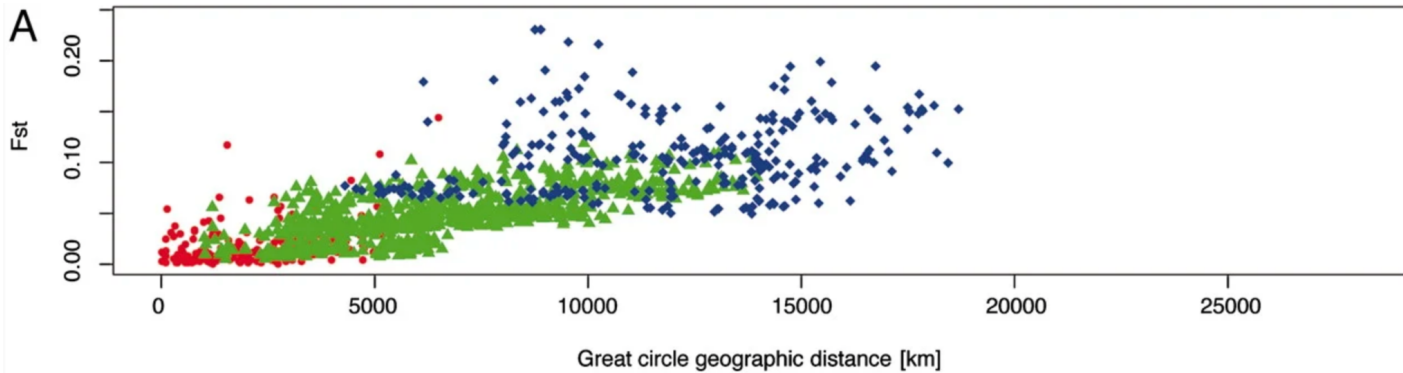    3. Afro-Eurasian panel of humans

# Major mode of gene flow in most species: IBD

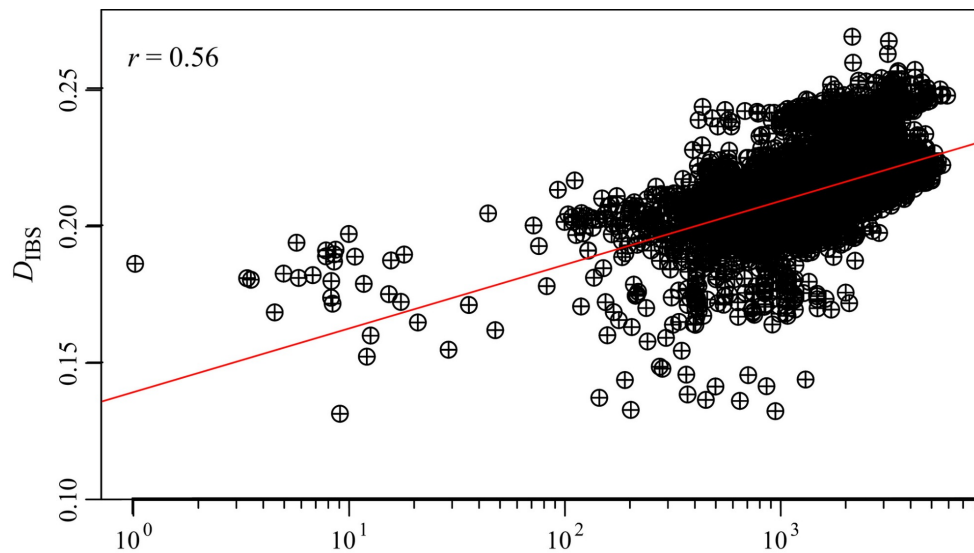"…accumulation of local genetic differences under geographically restricted dispersal."
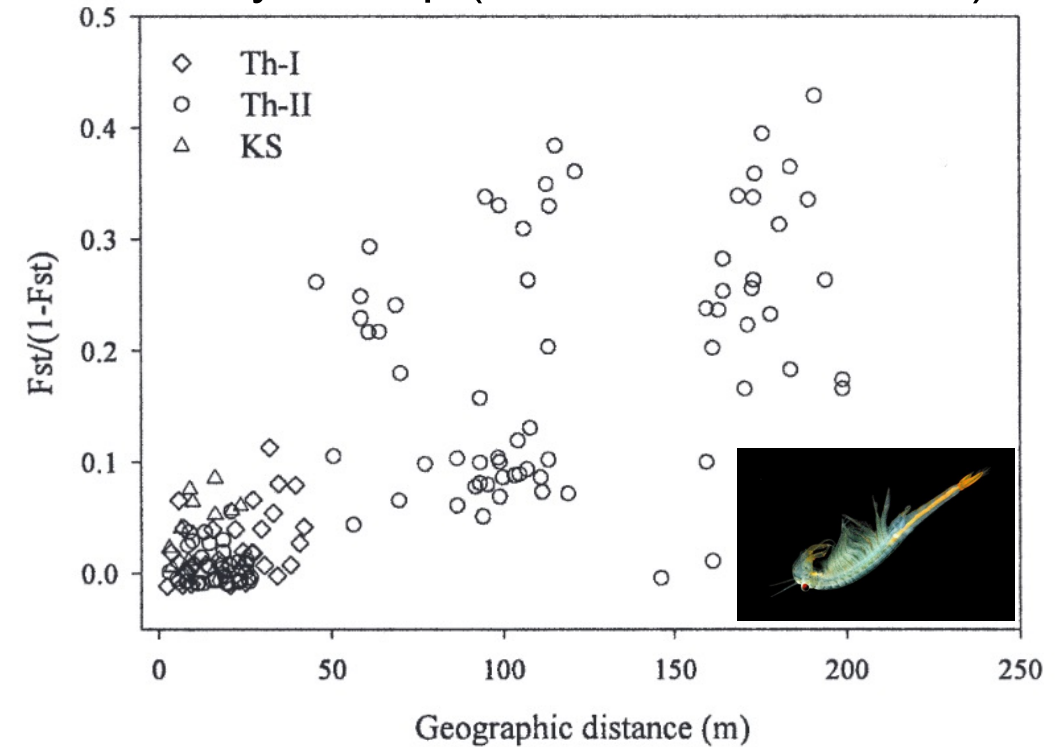


## Pattern of IBD



(initially formulated by Wright 1943)

# A few examples of IBD in different species

Humans (Ramachandran *et al* 2005)

Fairy shrimp (Hulsmans *et al* 2007)

Grey wolves (Schweizer *et al* 2016)

Long-range gene flow events can cause deviations from this pattern of IBD…

# Long-range gene flow events can cause deviations from this pattern of IBD…

…but so can a host of other <u>genetic</u> processes like fluctuating population size & assortative mating and other <u>ecological</u> processes like spatially heterogeneous landscapes & habitat fragmentation

# Long-range gene flow can be caused by:

• Human-mediated translocations



Wolves in Yellowstone

# Long-range gene flow can be caused by:

• Human-mediated translocations

• Extreme weather events



Tornado carrying a cow

# Long-range gene flow can be caused by:

- Human-mediated translocations

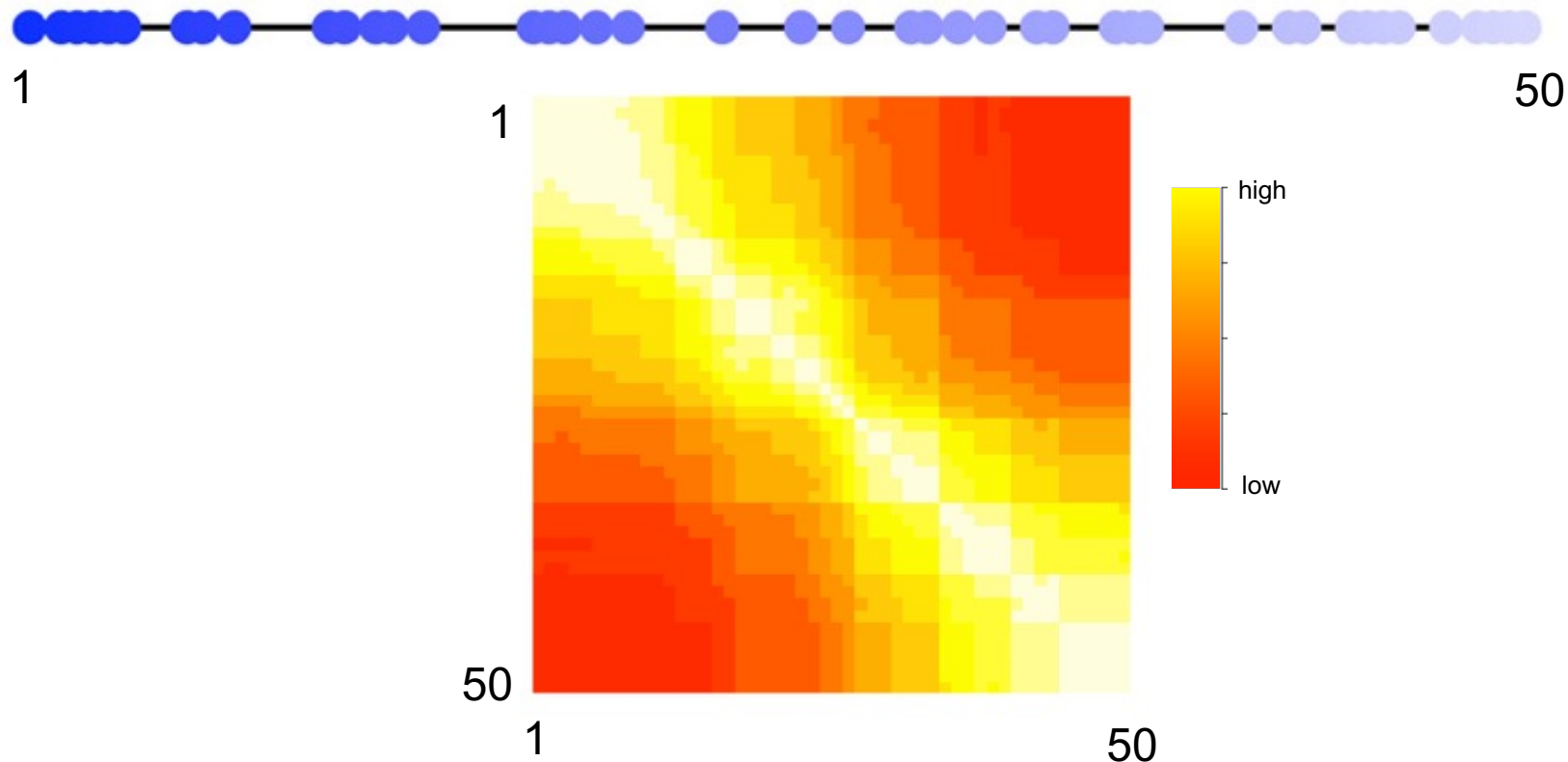- Extreme weather events

- Natural migration for greener pastures

# Outline

- Introduction + Motivation

- **Spatially heterogeneous models of isolation-by-distance (IBD)**

- Our model for long-range gene flow events (`FEEMSmix`)

- Results:
    1. Simulations
    2. North American grey wolves
    3. Afro-Eurasian panel of humans

# Spatial models of genetic variation

- Typically, revolves around modeling the **observed covariance matrix** of genotypes with an **expected covariance matrix**
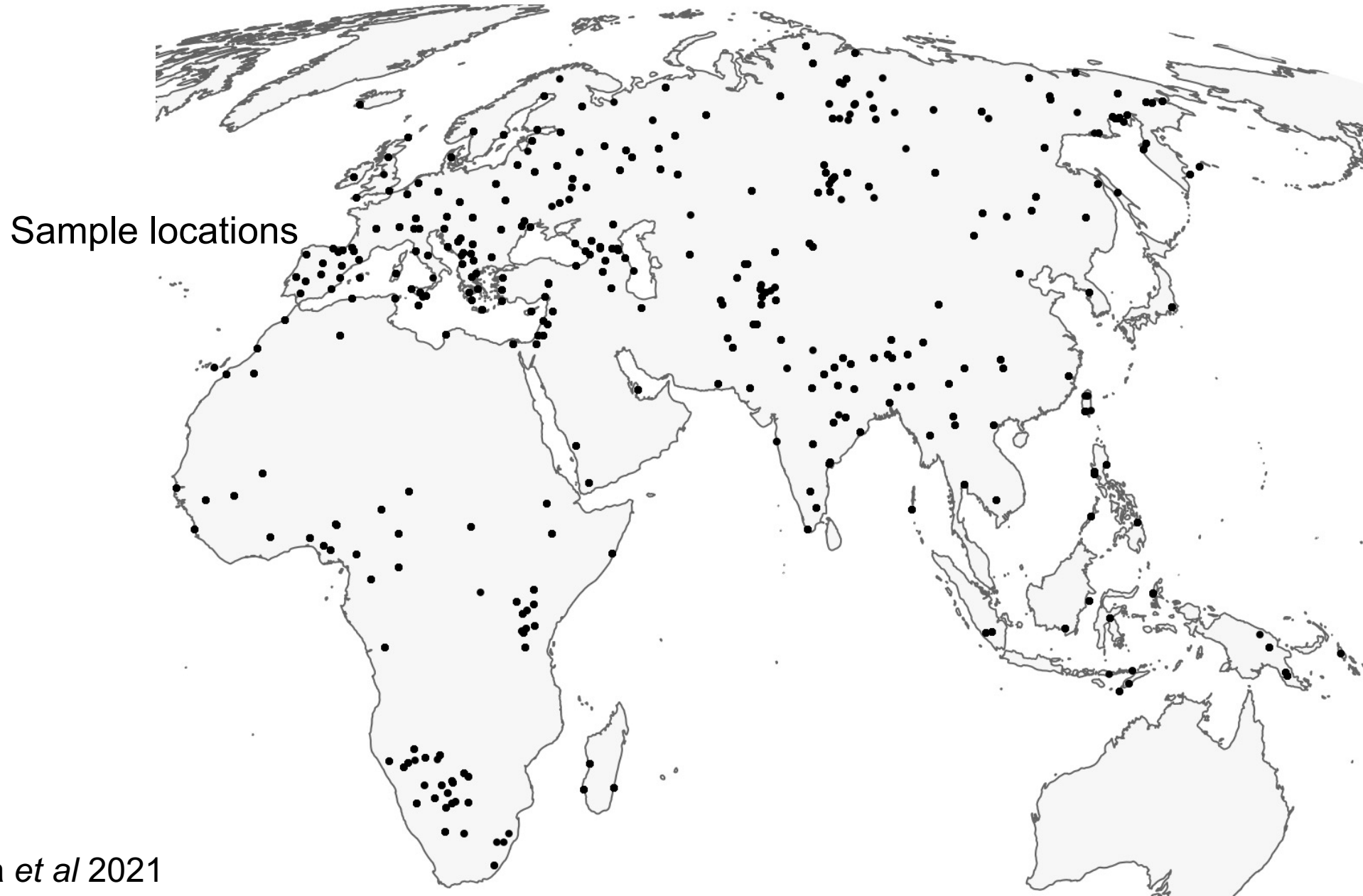
# Constructing the expected covariance matrix

1. Parametric variogram (`SpaceMix`)

2. Deep neural networks (`disperseNN` & `Locator`)

3. **Pairwise coalescent times (approximated in `EEMS`/`FEEMS` using circuit theory)**
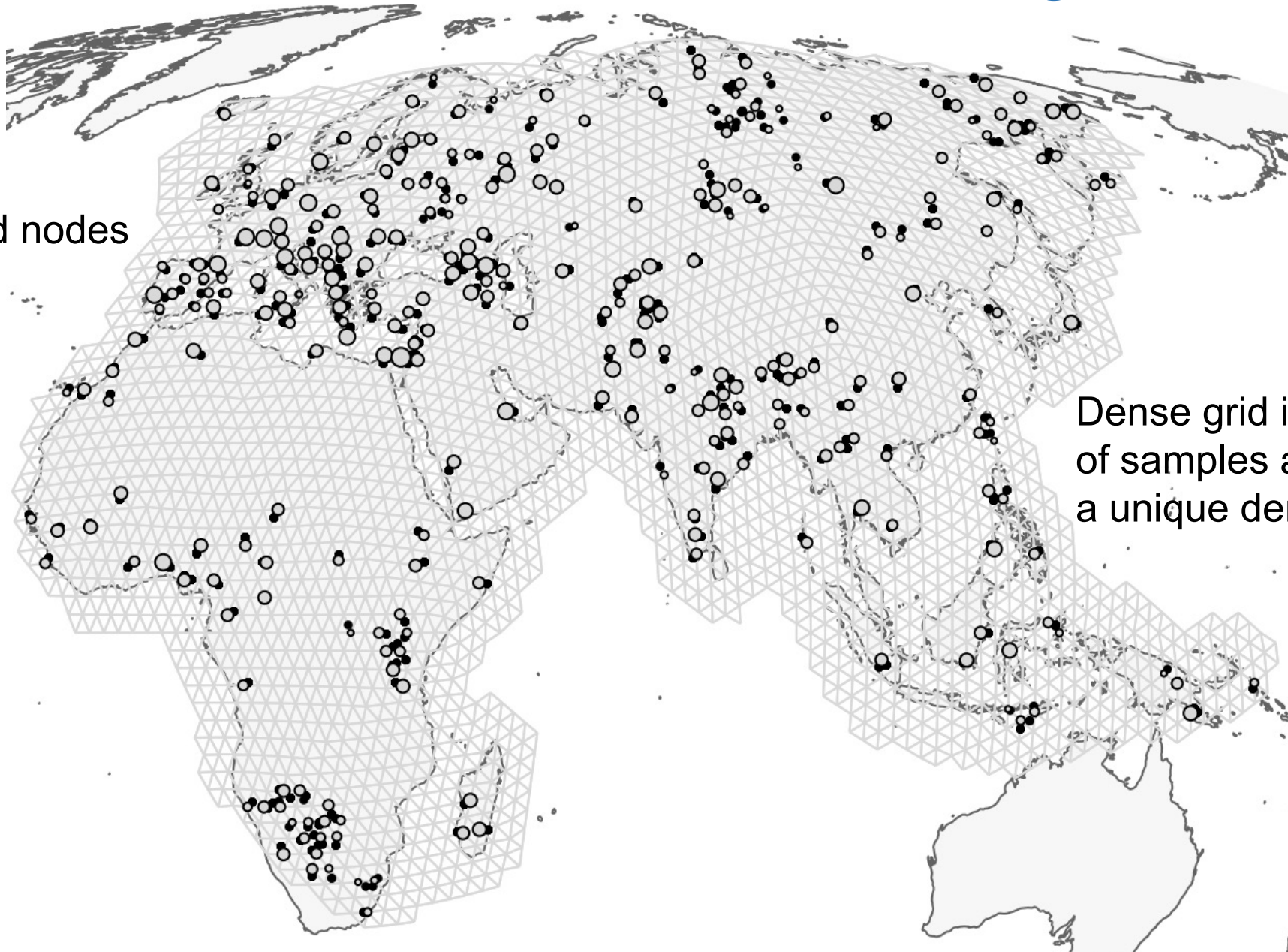
**F**ast
**E**stimation of
**E**ffective
**M**igration
**S**urfaces

1. Bradburd *et al* 2016
2. Smith *et al* 2020
   Battey *et al* 2020
3. Hanks & Hooten 2013
   Petkova *et al* 2016
   Lundgren & Ralph 2017
   Marcus, Ha *et al* 2021



16

# Brief overview of FEEMS



Sample locations

Marcus, Ha *et al* 2021

# Graph construction & spatial assignment



Observed nodes
(demes)

Dense grid in which ~80%
of samples are assigned to
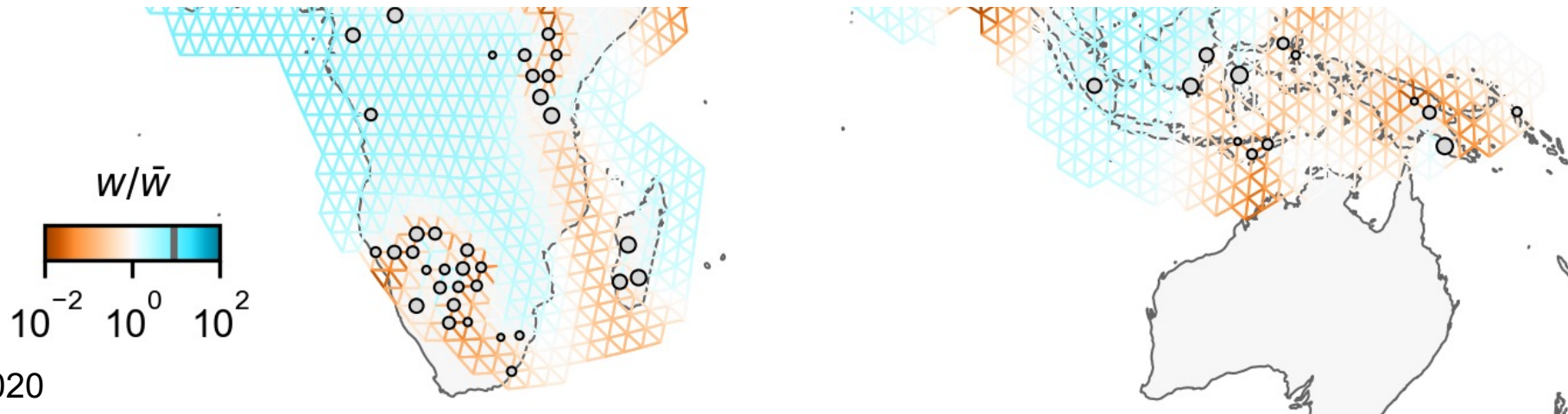a unique deme

# Penalized likelihood estimation



Effective migration rates (edge weights)
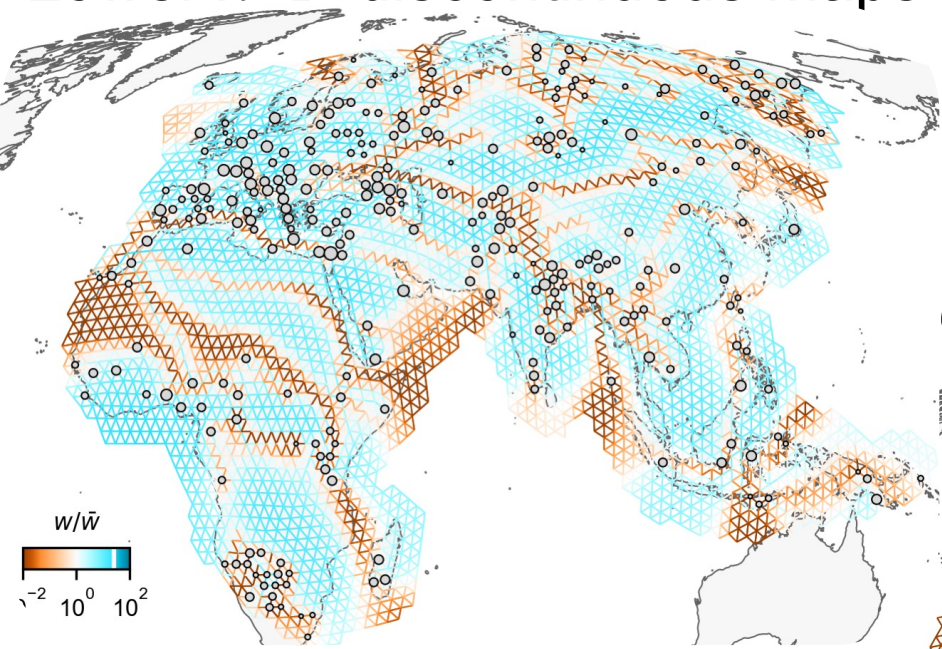
Deme-specific heterozygosity (nodes)

Lower than average "effective migration"

$w/\bar{w}$

Higher than average "effective migration"

$10^{-2}$ $10^{0}$ $10^{2}$

Data from Peter *et al* 2020

19

# Penalized likelihood estimation



Assumptions:

- Allele frequency in a deme is normally distributed (unlinked SNPs)
- Individuals are exchangeable within a deme
- **Symmetric, time-stationary migration rates**

$w/\bar{w}$

$10^{-2}$  $10^{0}$  $10^{2}$

Data from Peter *et al* 2020

Lower $\lambda$ ➜ discontinuous maps

Optimal $\lambda$ (chosen by CV)

Higher $\lambda$ ➜ smoother maps

LOO-CV error

$\lambda$

$w/\bar{w}$

$10^{-2}$ $10^0$ $10^2$

21

# Previous work in estimating "long-range" gene flow events



Admixture graphs (Lipson *et al* 2013)
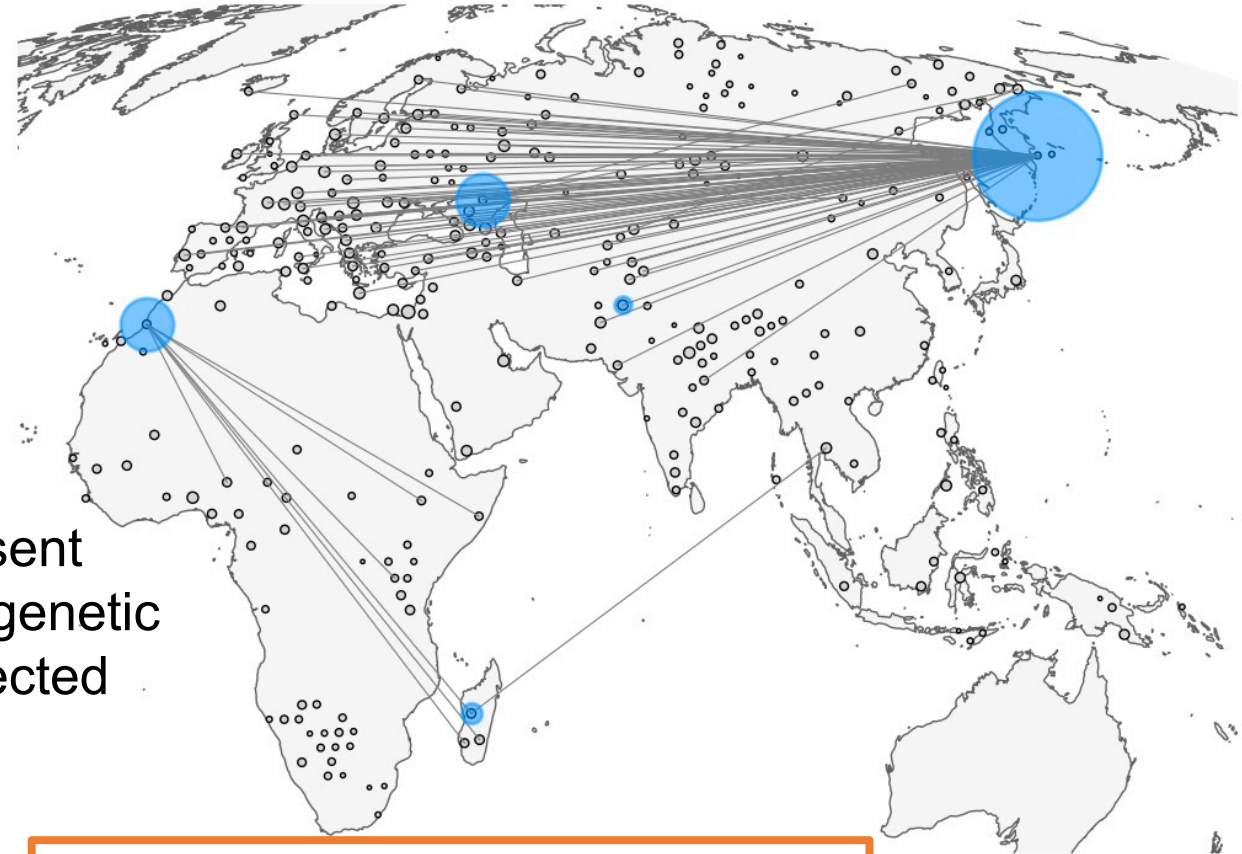
TreeMix (Pickrell & Pritchard 2012)

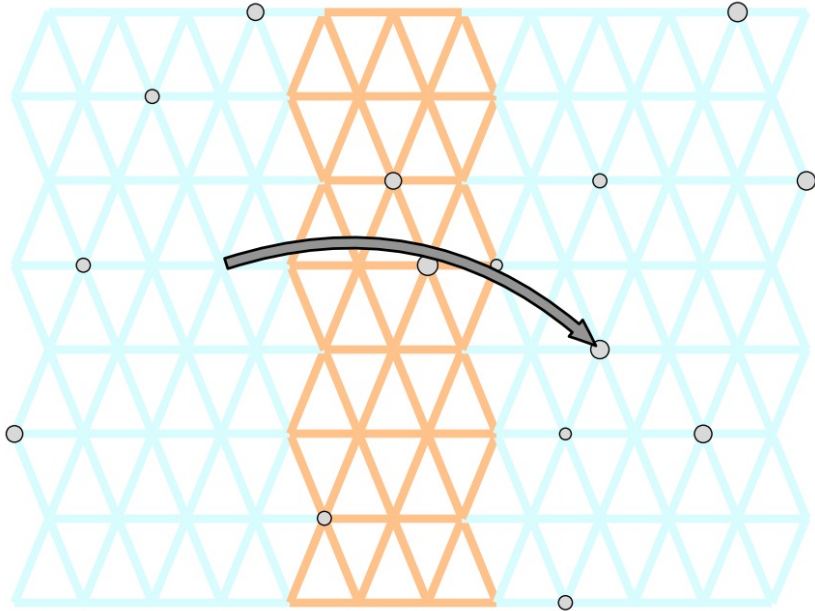SpaceMix (Bradburd *et al* 2016)

# Outline

- Introduction + Motivation

- Spatially heterogeneous models of isolation-by-distance (IBD)

- **Our model for long-range gene flow events (`FEEMSmix`)**

- Empirical results:
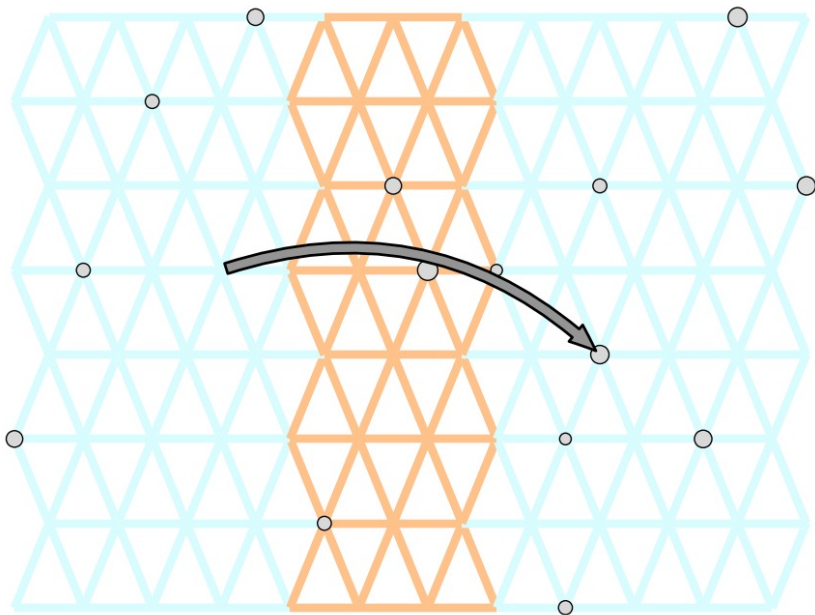  1. North American grey wolves
  2. Afro-Eurasian panel of humans

# FEEMSmix workflow



Observed (genetic) distance vs Expected (fit) distance

$R^2 = 0.884$

(red outliers represent pairs with <u>smaller</u> genetic distance than expected under the model)

$$T'_{ss} = T_{ss}$$
$$T'_{sd} = cT_{sd} + (1-c)T_{ss}$$
$$T'_{dd} = (1-c)^2 T_{dd} + 2c(1-c)T_{sd} + c^2 T_{ss}$$

# Outline

- Introduction + Motivation

- Spatially heterogeneous models of isolation-by-distance (IBD)

- Our model for long-range gene flow events (`FEEMSmix`)

- Results:
  1. Simulations
  2. North American grey wolves
  3. Afro-Eurasian panel of humans

# Brief simulation results

- 8x12 grid (only 15% sampled)
- 1-10 samples/deme
- 1,000 SNPs
- Corridor $m$ is 10x barrier $m$
- Varying population size $N$ across grid

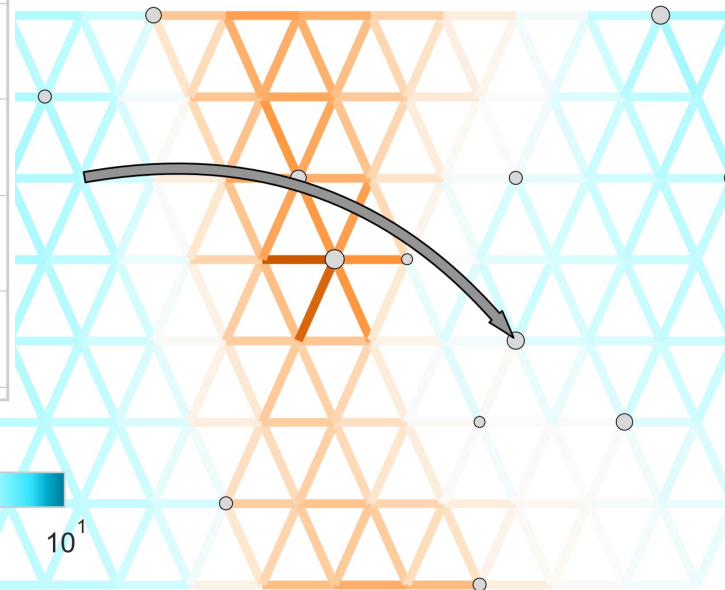$F_{ST}$ across barrier ≈ 0.1



Simulated truth
(instantaneous pulse of $c$ = 0.5)

$w/\bar{w}$

$10^{-1}$  $10^0$  $10^1$

R²=0.829

Initial FEEMS fit

Simulated truth

$F_{ST}$ across barrier ≈ 0.1

Simulated truth

R²=0.829

R²=0.918

$w/\bar{w}$

$10^{-1}$  $10^0$  $10^1$
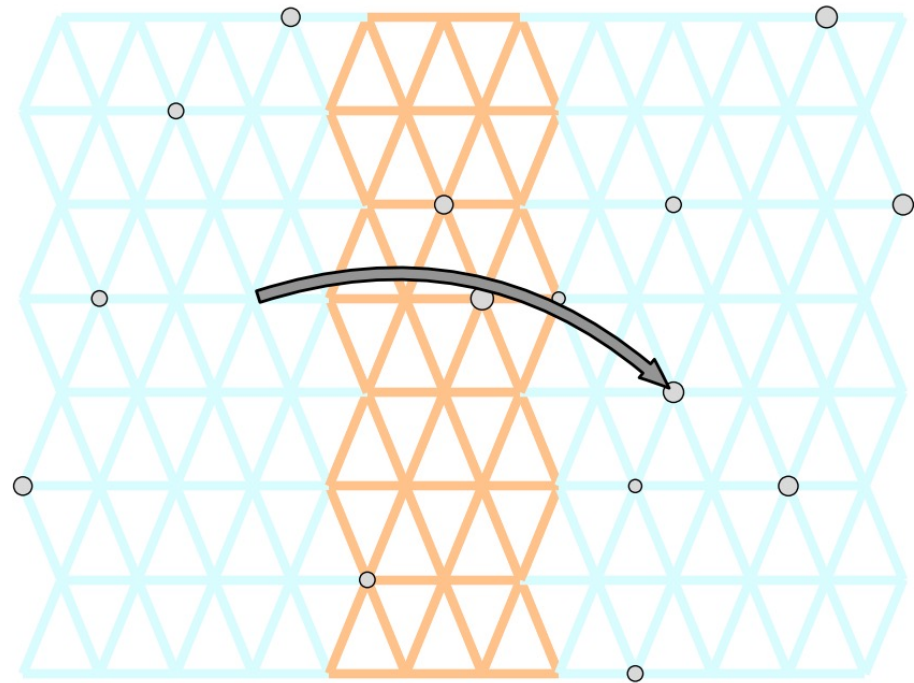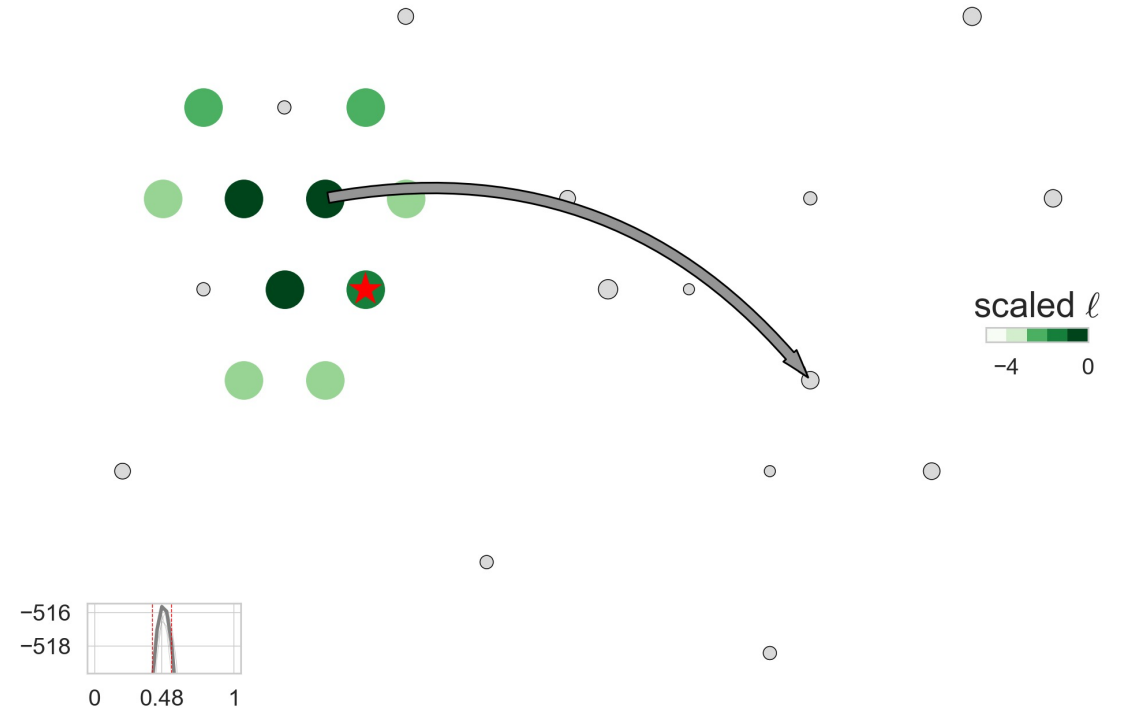
Initial FEEMS fit

FEEMSmix fit

30

# True source is within <u>two</u> log-likelihood units of MLE source



Simulated truth
($c$ = 0.5)

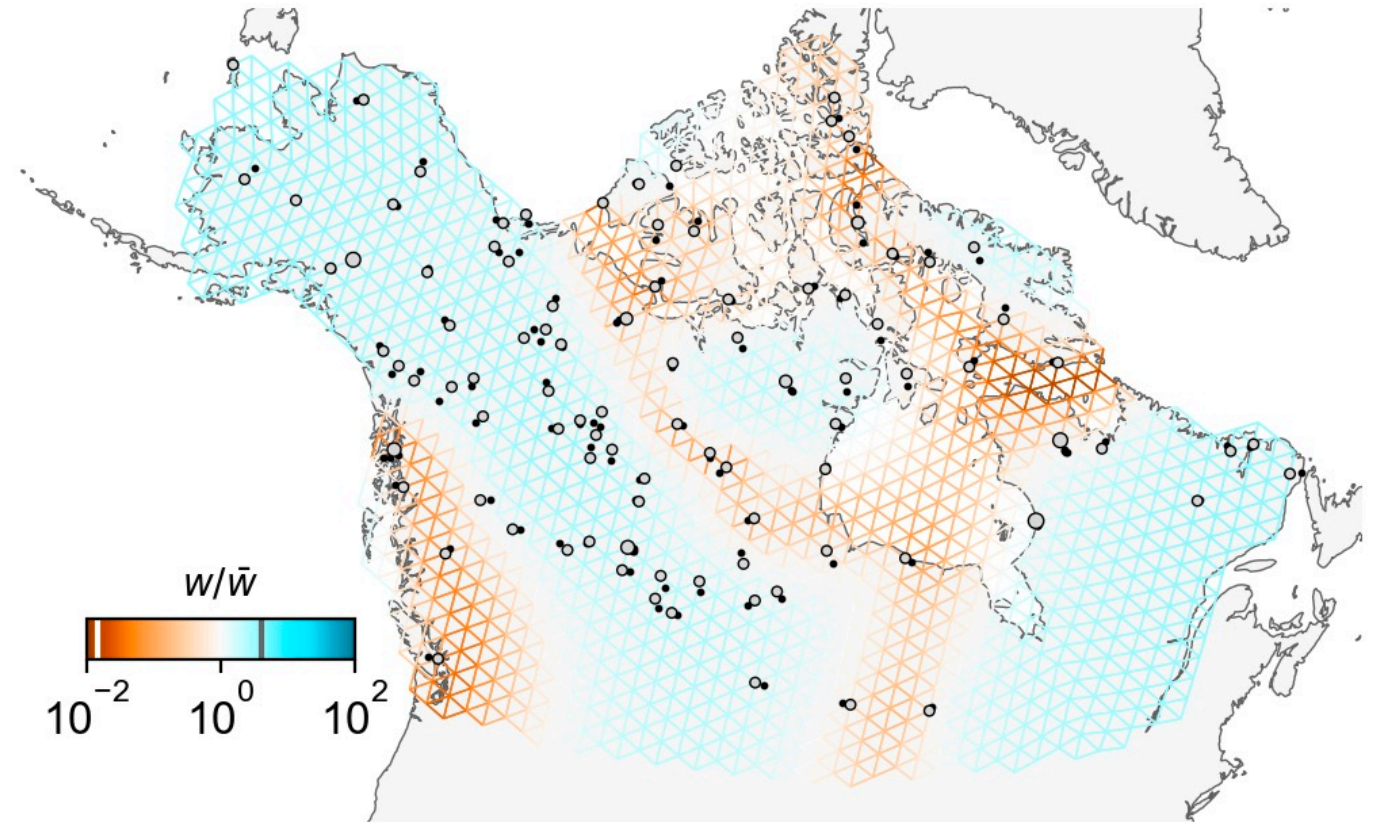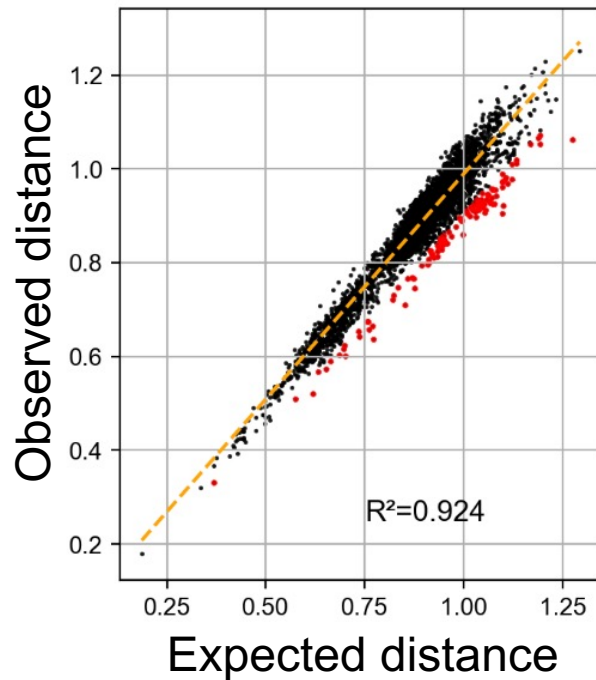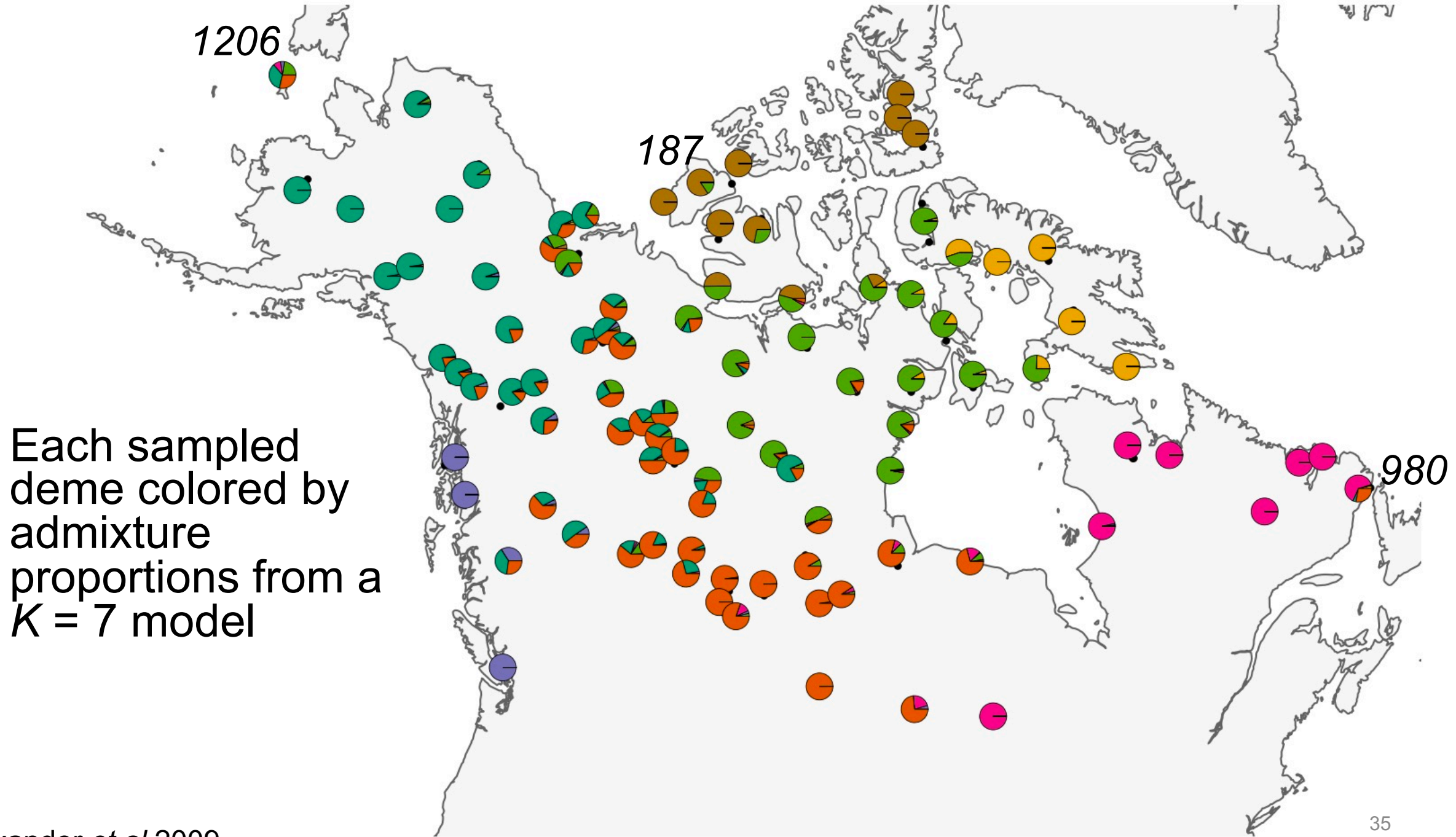Estimated `FEEMSmix` contour
(est. $c$ = 0.48)

# Outline

- Introduction + Motivation

- Spatially heterogeneous models of isolation-by-distance (IBD)

- Our model for long-range gene flow events (`FEEMSmix`)

- Results:
    1. Simulations
    2. North American grey wolves
    3. Afro-Eurasian panel of humans
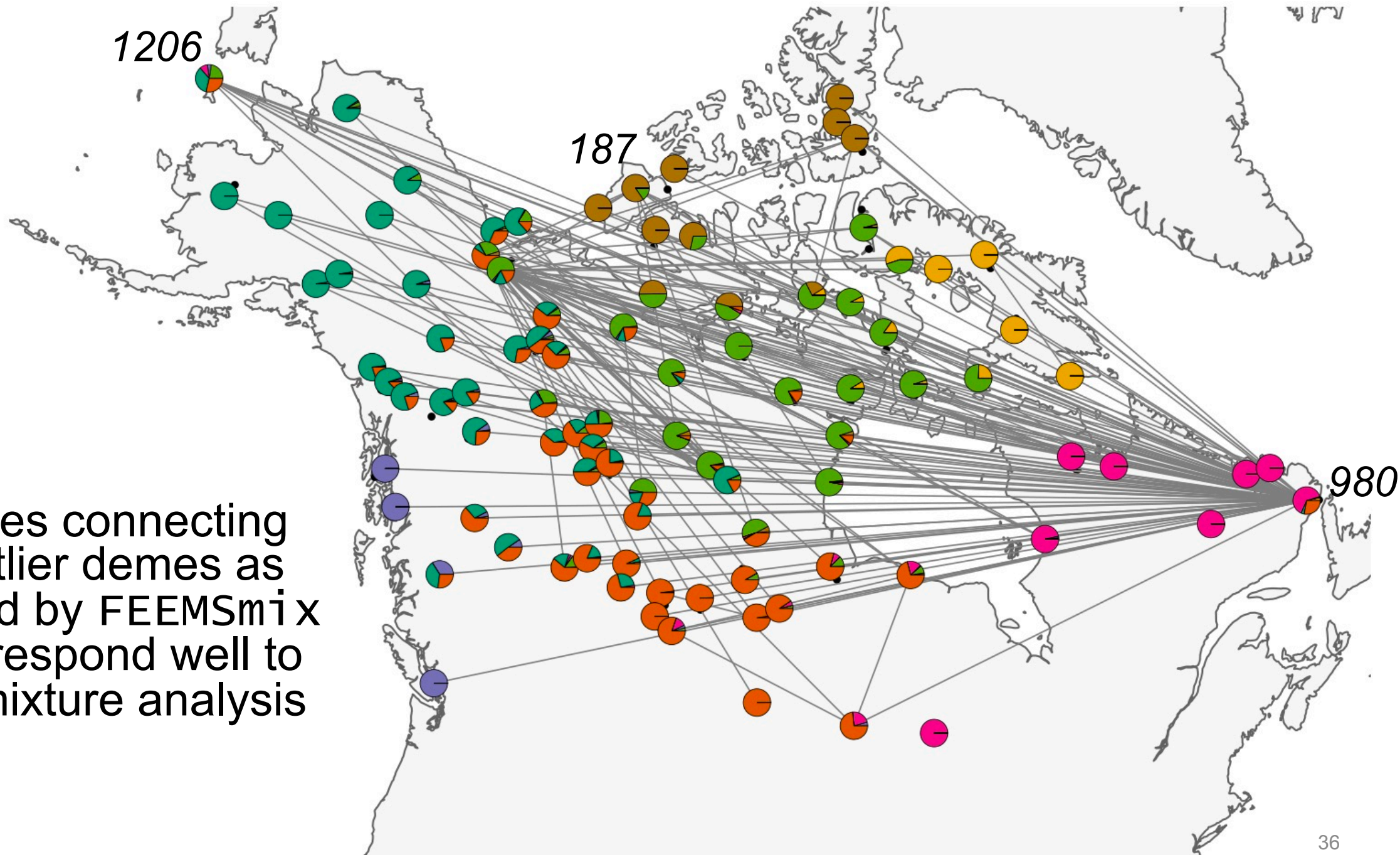
# North American grey wolves
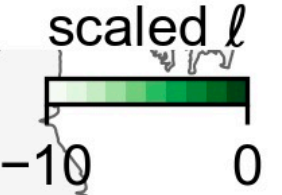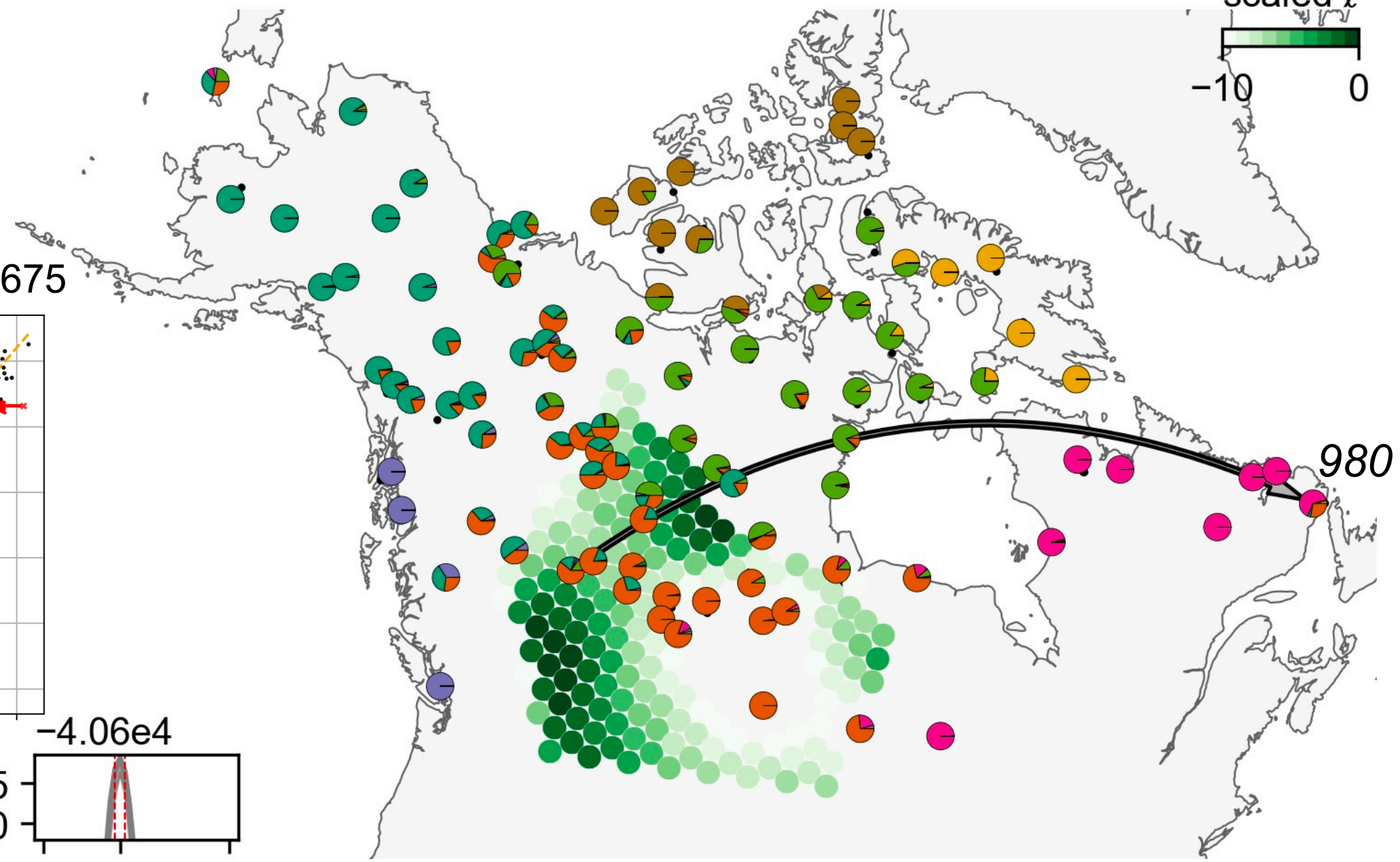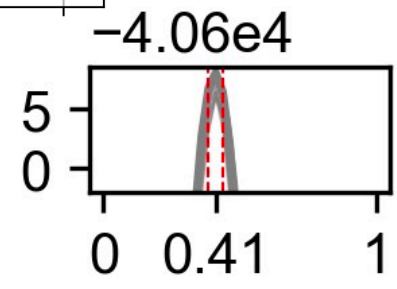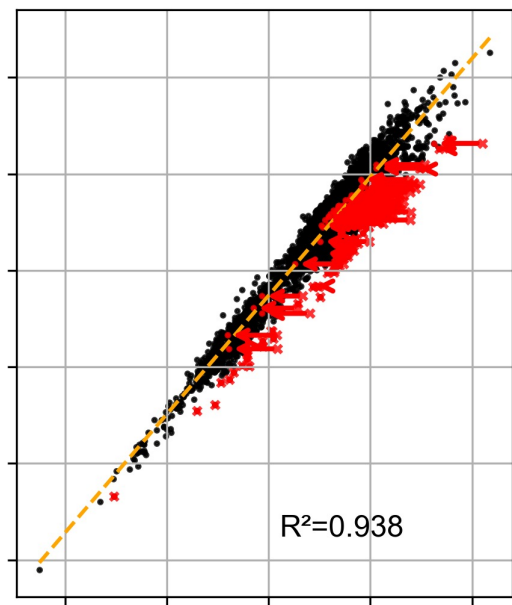
- 111 samples
- 17.8k SNPs
- 94 demes (~1 sample/deme)



Schweizer *et al* 2016

Each sampled deme colored by admixture proportions from a *K* = 7 model

*1206*

*187*

*980*

Alexander *et al* 2009

1206

187

980

Lines connecting
outlier demes as
found by `FEEMSmix`
correspond well to
admixture analysis

$\ell = $ -41070 ➜ -40675

R²=0.938

scaled $\ell$

−10          0

980

−4.06e4

37

# The full picture: wolves move around a lot

# Bonus: spatial prediction!



Median error
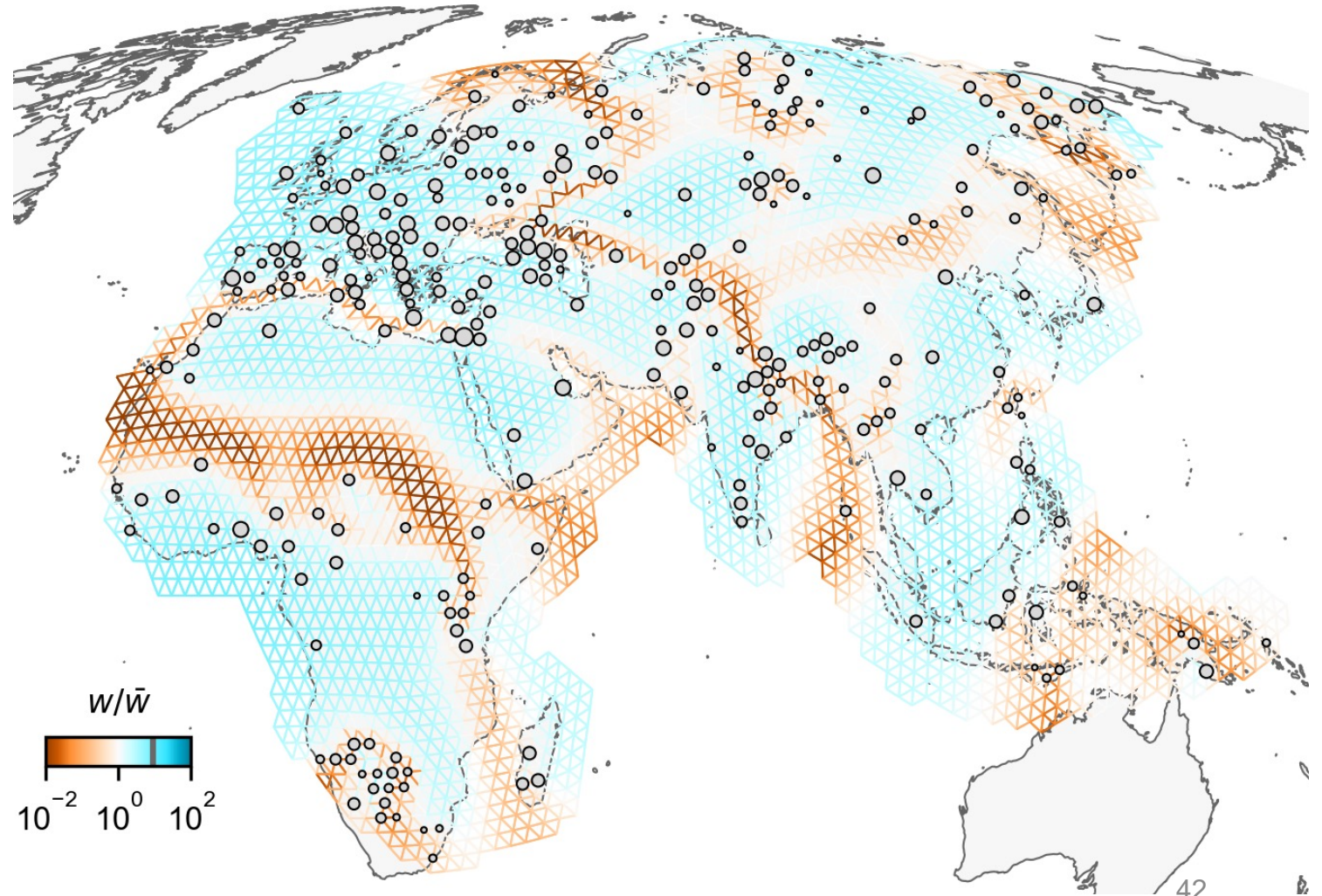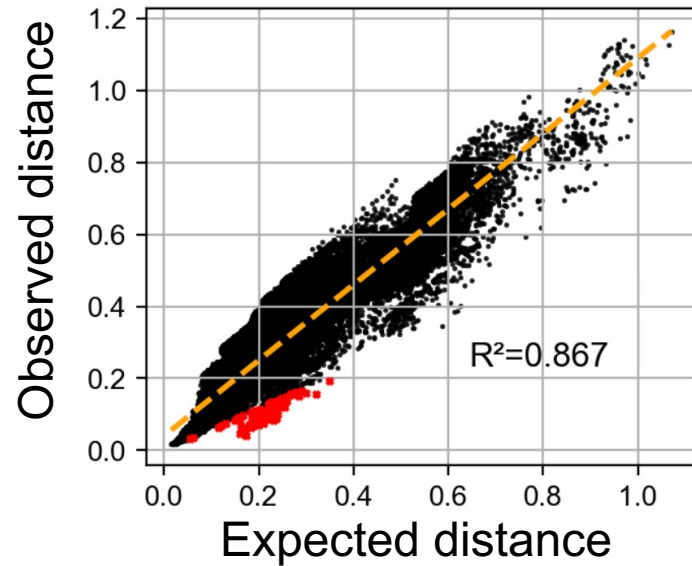~6.4 km          ~8.6 km

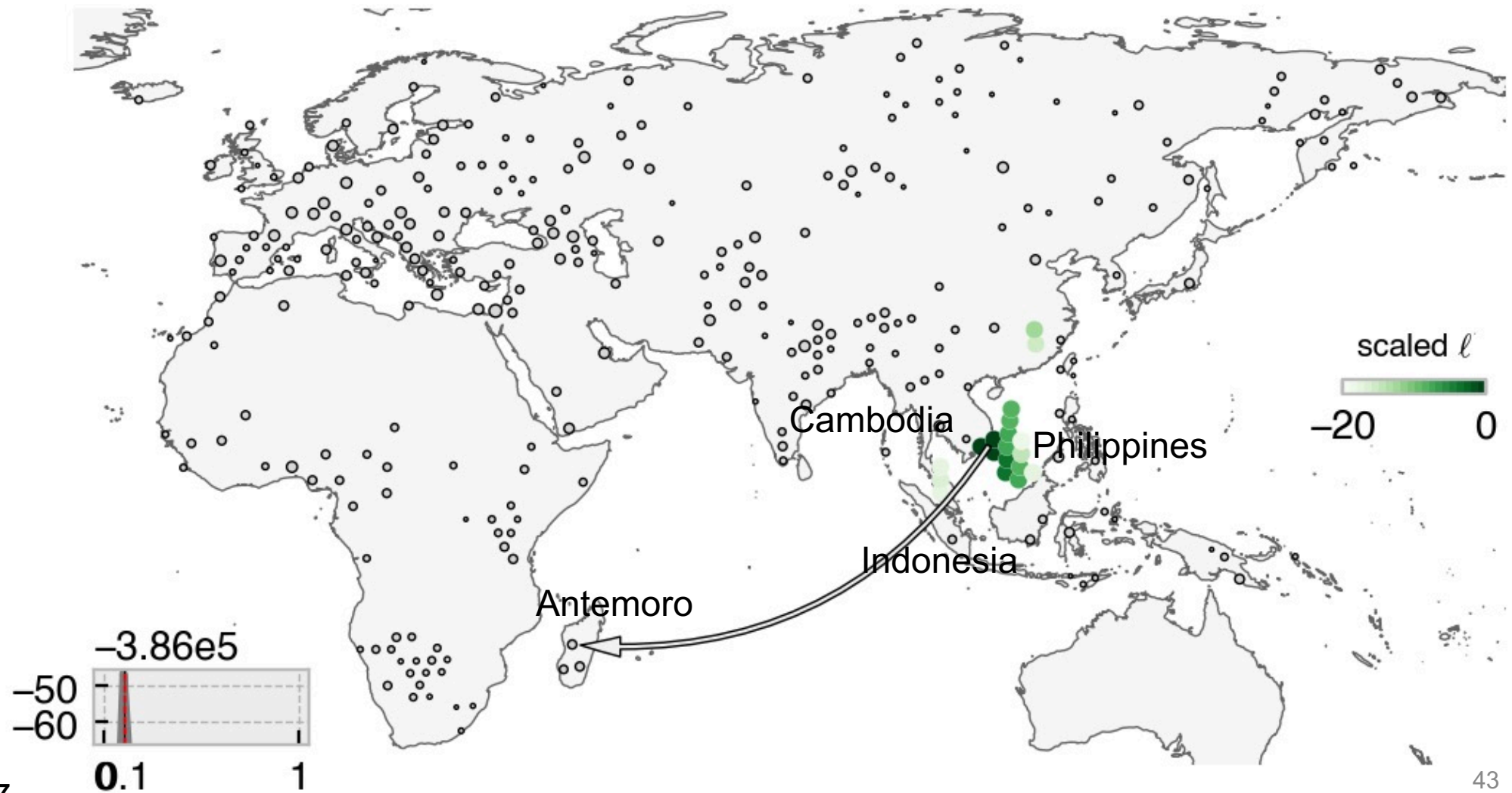Battey *et al* 2020

40

# Outline

- Introduction + Motivation

- Spatially heterogeneous models of isolation-by-distance (IBD)

- Our model for long-range gene flow events (FEEMSmix)

- Results:
    1. Simulations
    2. North American grey wolves
    3. Afro-Eurasian panel of humans
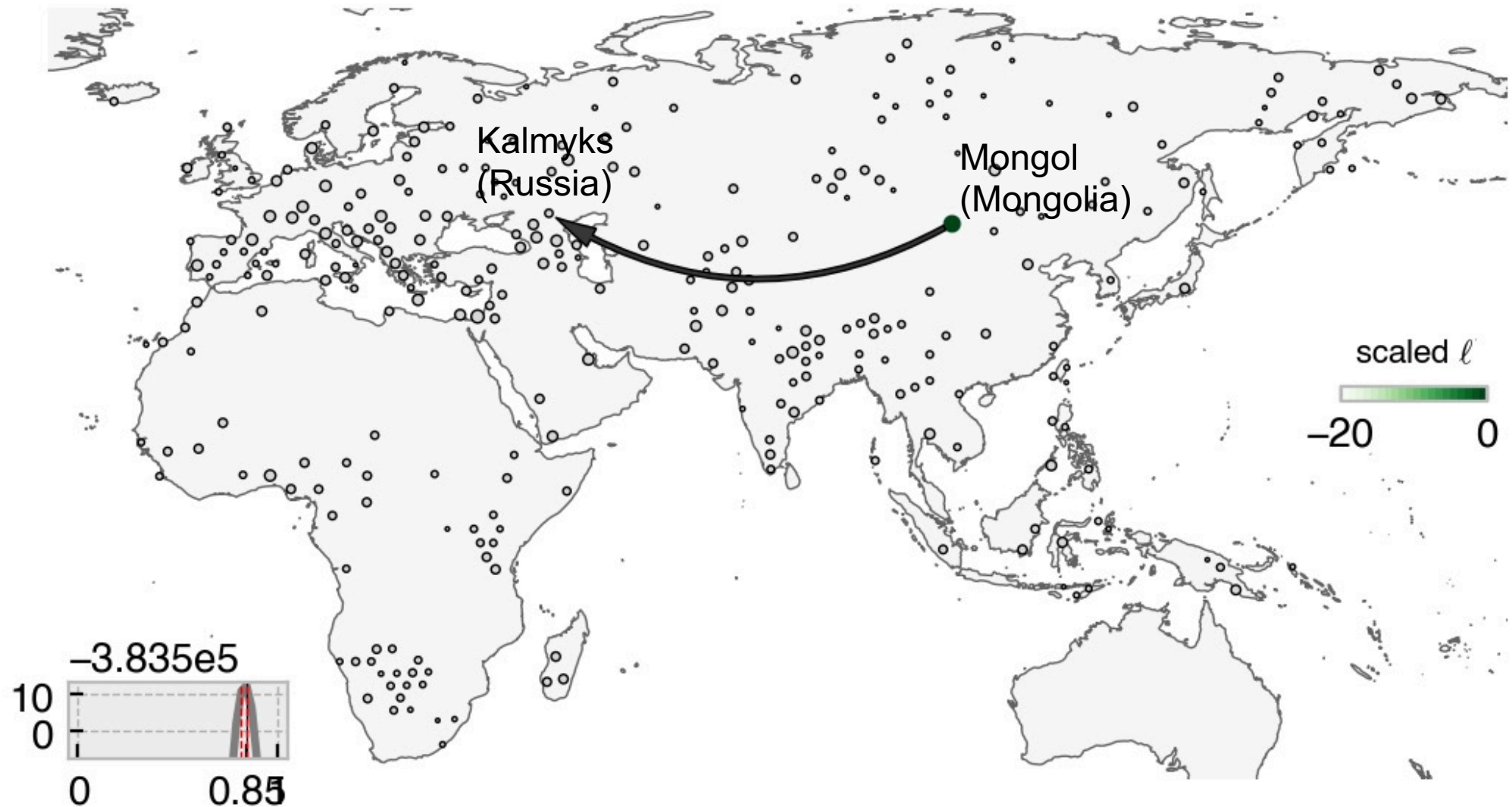
# Afro-Eurasian panel of humans

- 4,700 samples
- 20k SNPs
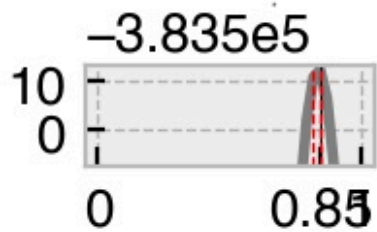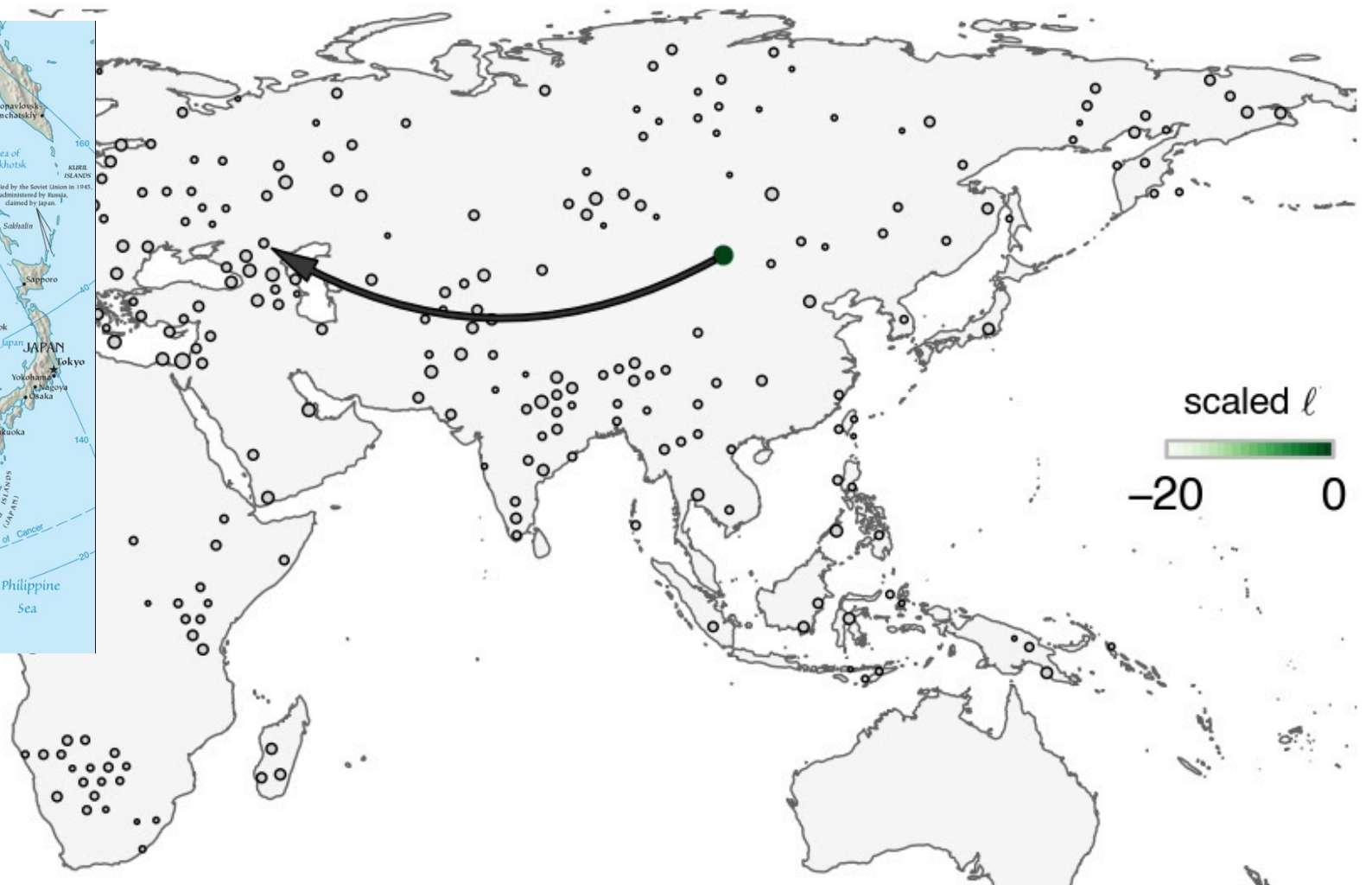- 297 demes



Peter *et al* 2020

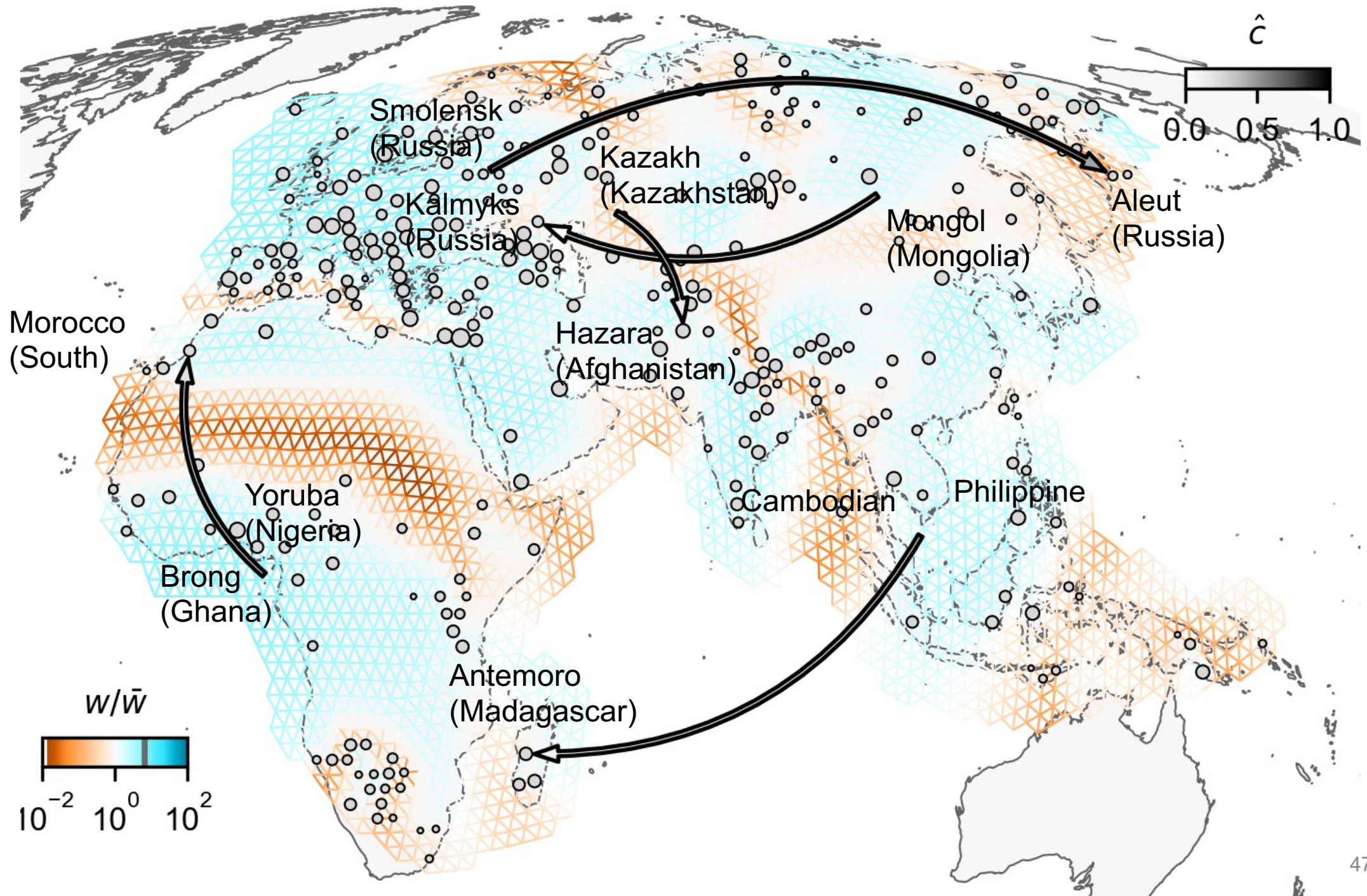# Southeast Asian ancestry source detected in Malagasy



Pierron *et al* 2017

43

# Kalmyks: only Mongolic-speaking people living in Europe

# Kalmyks: only Mongolic-speaking people living in Europe



Erdeniev 1985

Smolensk (Russia)
Kazakh (Kazakhstan)
Kalmyks (Russia)
Mongol (Mongolia)
Aleut (Russia)
Morocco (South)
Hazara (Afghanistan)
Yoruba (Nigeria)
Brong (Ghana)
Cambodian
Philippine
Antemoro (Madagascar)

$\hat{c}$
0.0  0.5  1.0

$w/\bar{w}$
$10^{-2}$  $10^{0}$  $10^{2}$

# Conclusions

- `FEEMSmix`: a method to include long-range gene flow events in `FEEMS`

- Paints a fuller picture of the spatial patterns in genetic structure

- ⚠ Caution ⚠

- Interpret value of $c$ as an informed suggestion, not as truth (e.g., if $c > 0.5$, it probably means high recent admixture)

- Reckon with uncertainty in source location (area-area vs point-point migration)