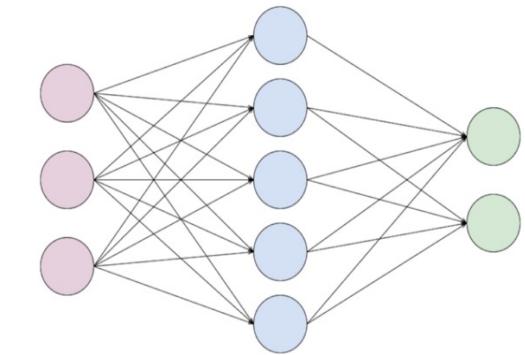
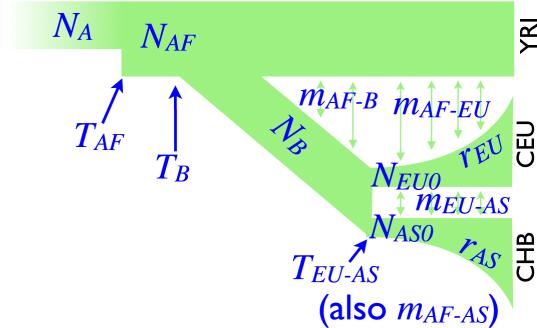
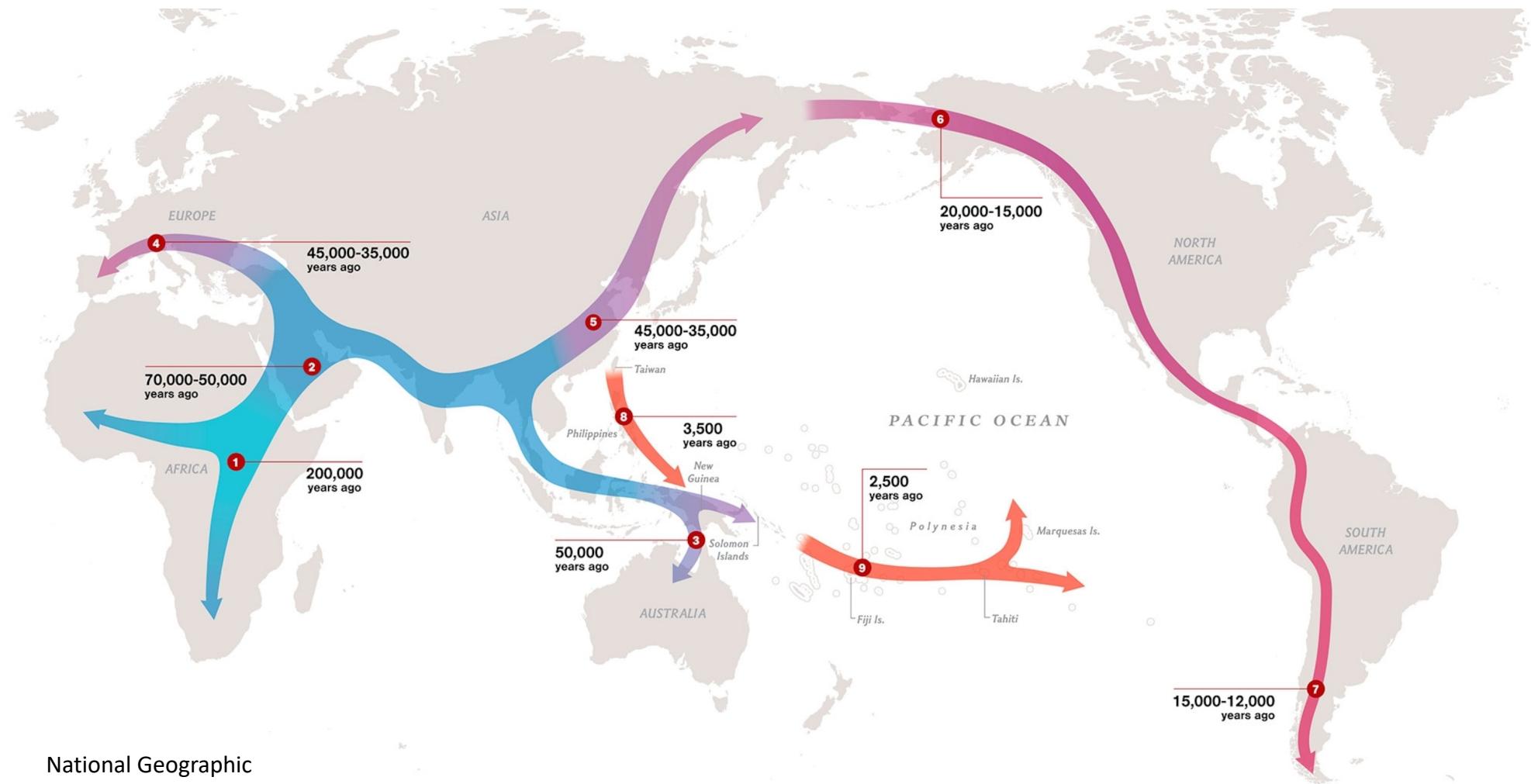


Computationally Efficient Demographic Inference with Supervised Machine Learning



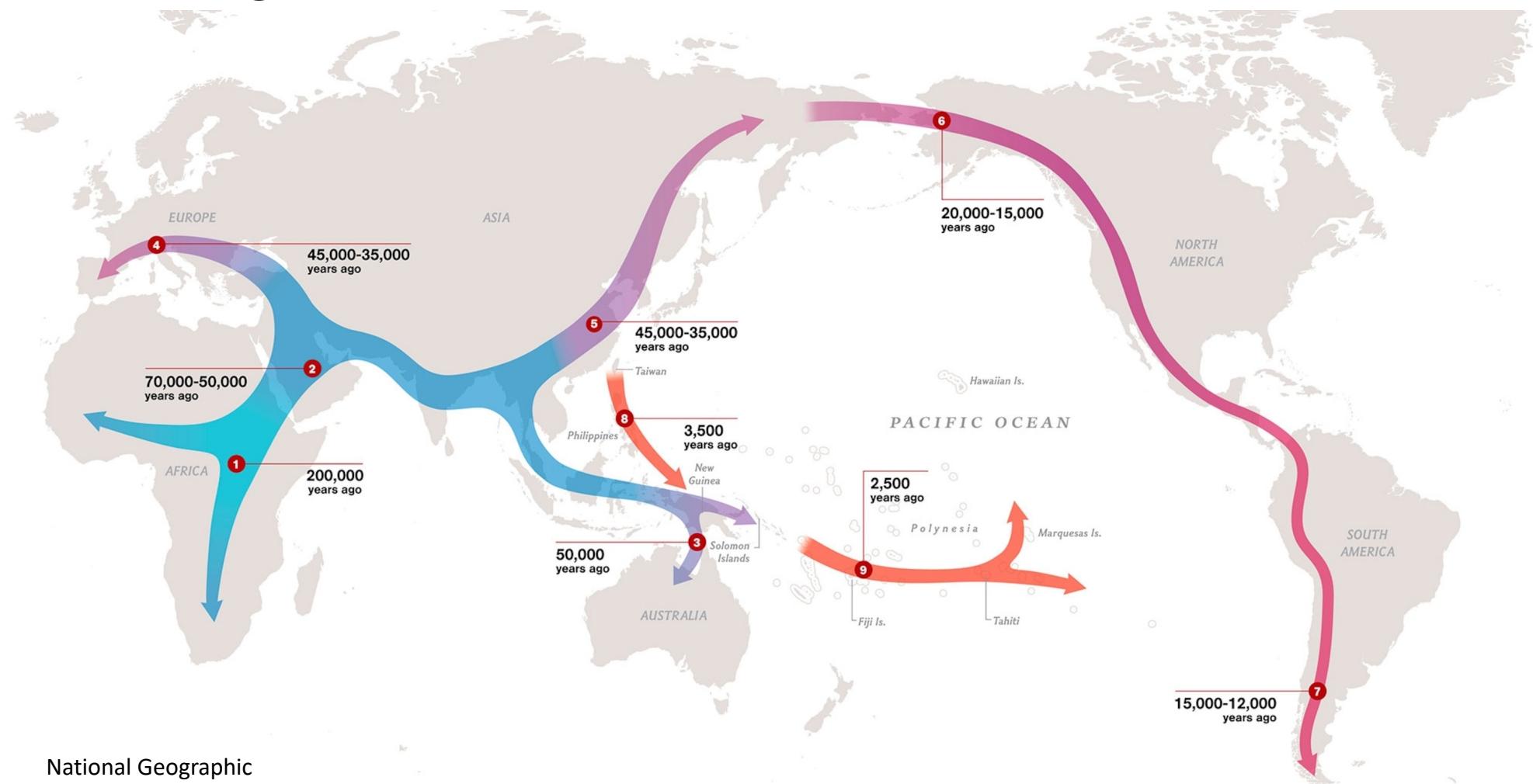
Linh N. Tran, Connie Sun, Mathews Sajan, Ryan Gutenkunst
University of Arizona

Demographic history inference



Demographic history inference – Motivation

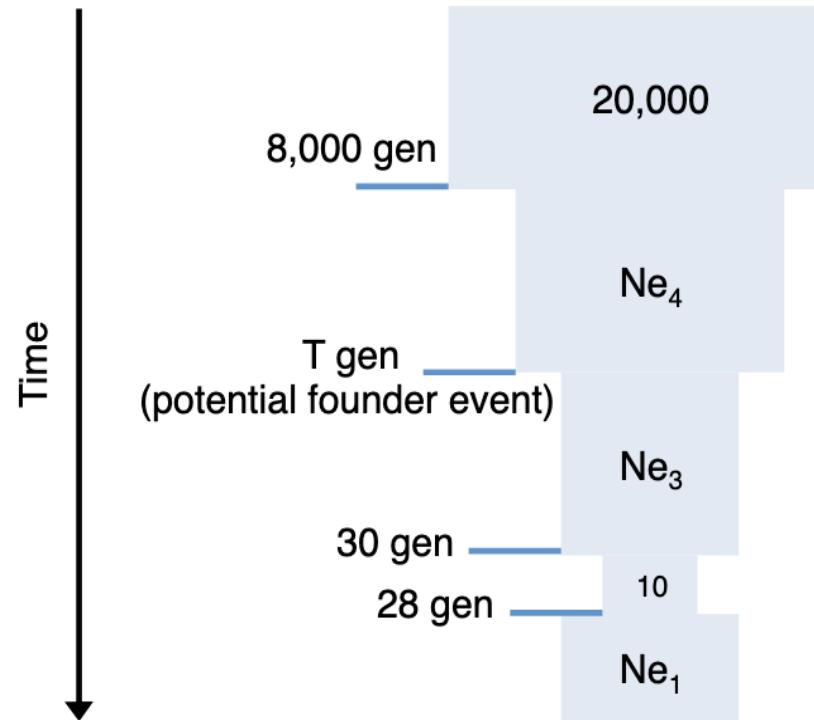
- Understand population history:
bottlenecks, gene flow, etc.



Demographic history inference – Motivation

National Geographic

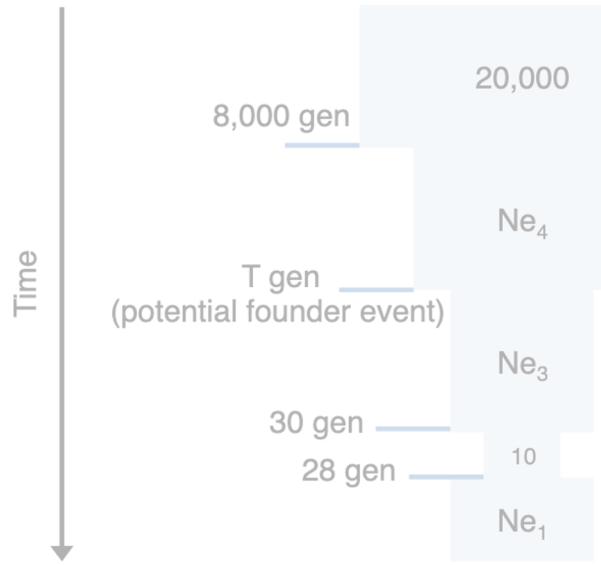
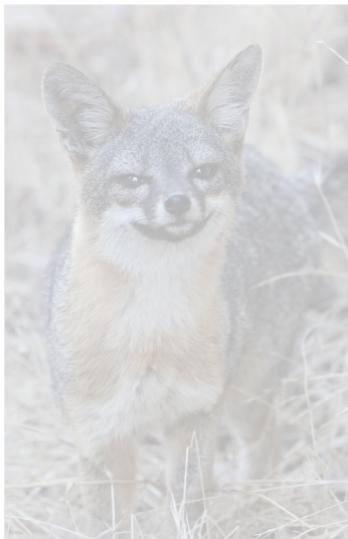
- Understand population history: bottlenecks, gene flow, etc.
- Conservation: historical genetic diversity of endangered species



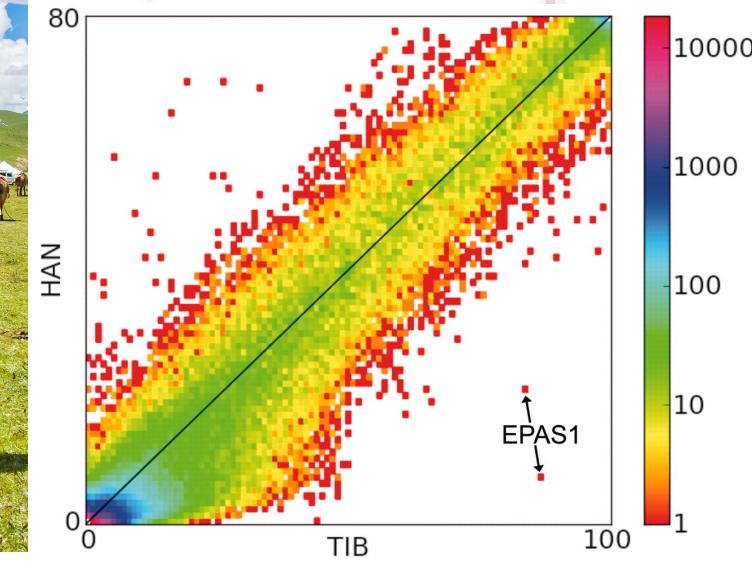
Demographic history inference – Motivation

National Geographic

- Understand population history: bottlenecks, gene flow, etc.
- Conservation: historical genetic diversity of endangered species
- Selection: sets neutral background for detecting adaptation



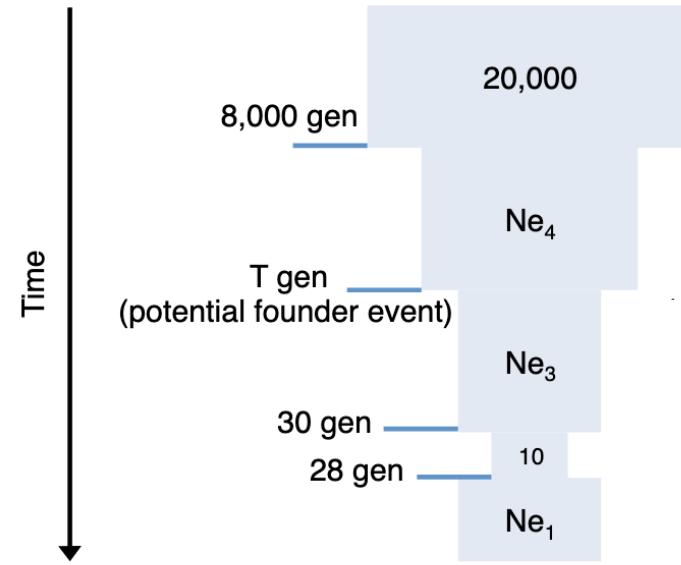
Robinson et al. (2016) *Current Biology*



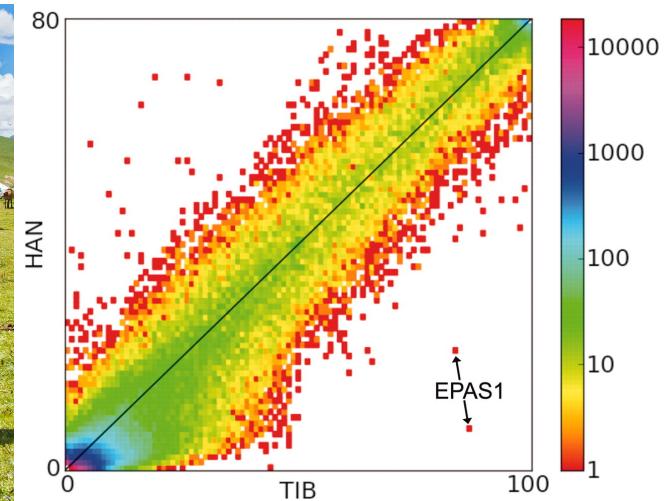
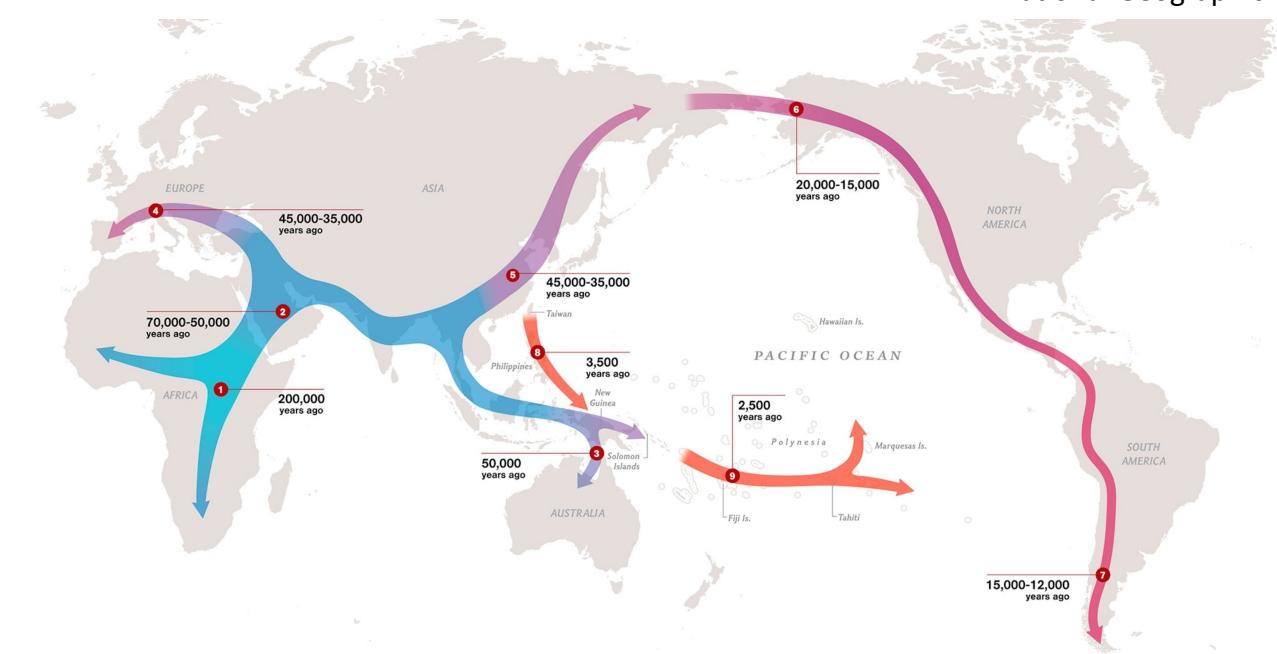
Yi et al. (2010) *Science*

Demographic history inference – Motivation

- Understand population history: bottlenecks, gene flow, etc.
- Conservation: historical genetic diversity of endangered species
- Selection: sets neutral background for detecting adaptation

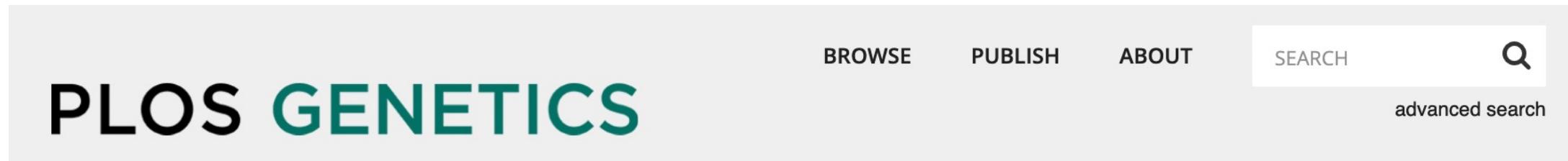


Robinson et al. (2016) *Current Biology*



Yi et al. (2010) *Science*

Diffusion approximation for demographic inference (dadi)



The image shows the header of a PLOS GENETICS research article. At the top right are navigation links: BROWSE, PUBLISH, ABOUT, SEARCH (with a magnifying glass icon), and advanced search. The title "PLOS GENETICS" is prominently displayed in large, bold, black and teal letters. Below the title are status icons: OPEN ACCESS (padlock) and PEER-REVIEWED (document). The article type is RESEARCH ARTICLE. To the right, there are four metrics in a grid: 1,409 Save (dark teal), 1,286 Citation (light teal), 55,478 View (dark teal), and 8 Share (light teal).

Inferring the Joint Demographic History of Multiple Populations from Multidimensional SNP Frequency Data

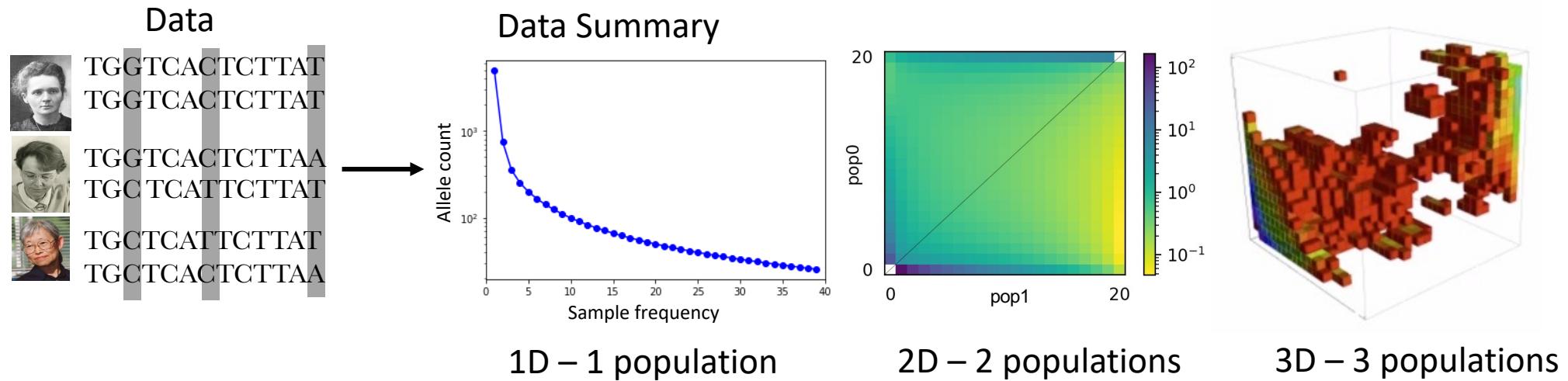
Ryan N. Gutenkunst , Ryan D. Hernandez, Scott H. Williamson, Carlos D. Bustamante

Published: October 23, 2009 • <https://doi.org/10.1371/journal.pgen.1000695>

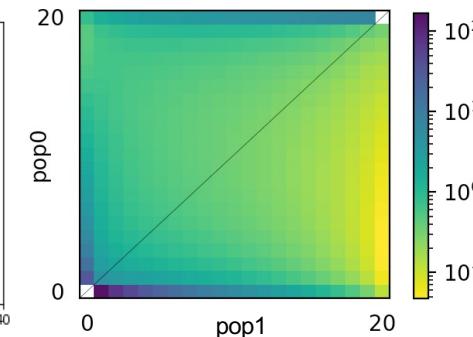
≡  Groups

☆ dadi-user 391 members

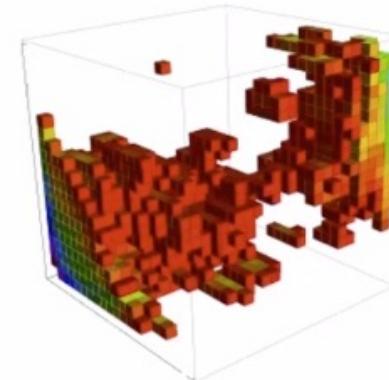
Diffusion approximation for demographic inference (dadi)



1D – 1 population

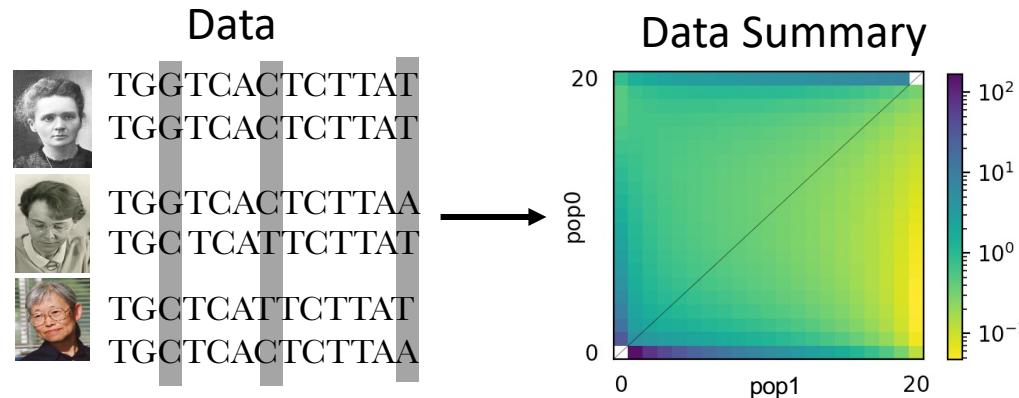


2D – 2 populations

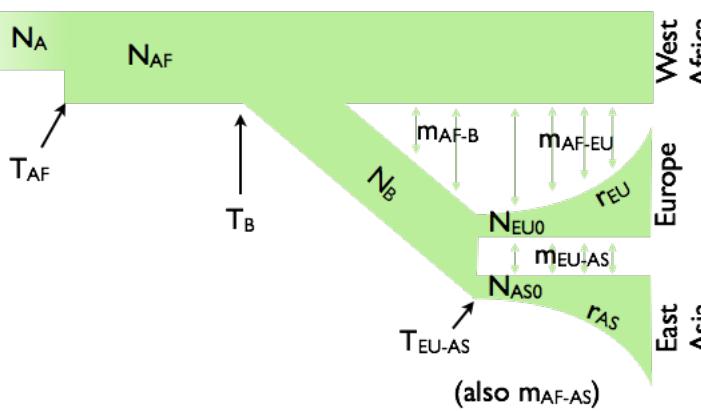


3D – 3 populations

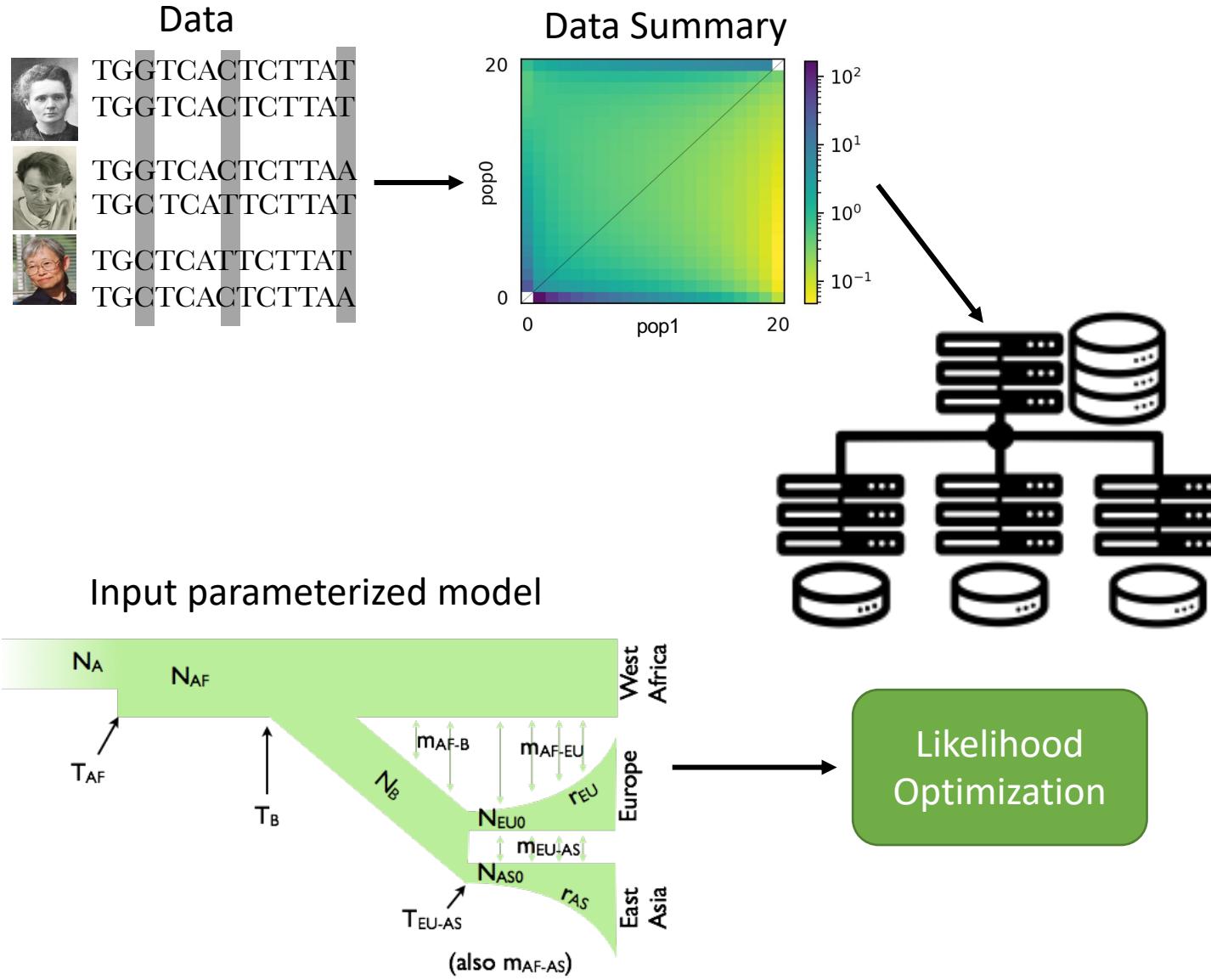
Diffusion approximation for demographic inference (dadi)



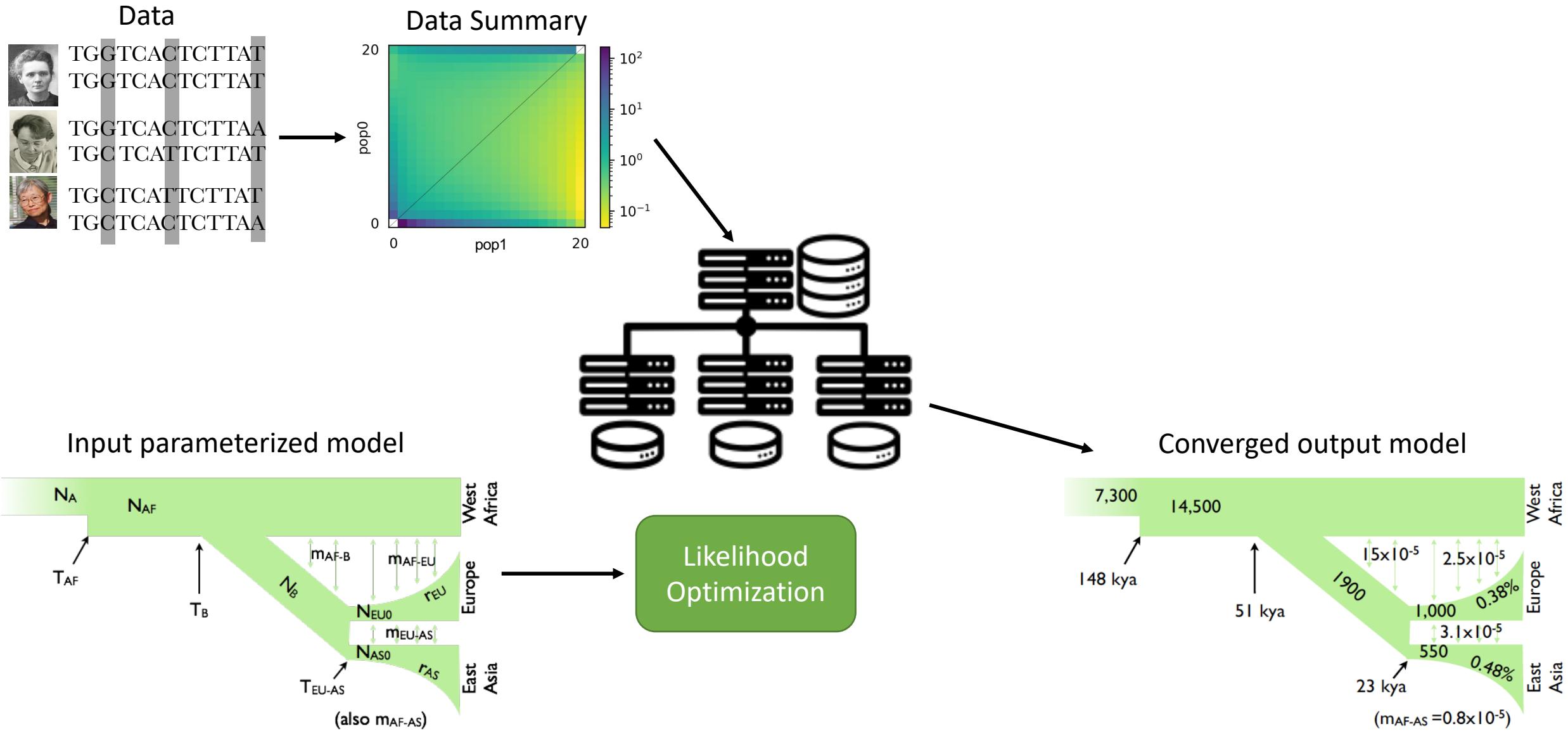
Input parameterized model



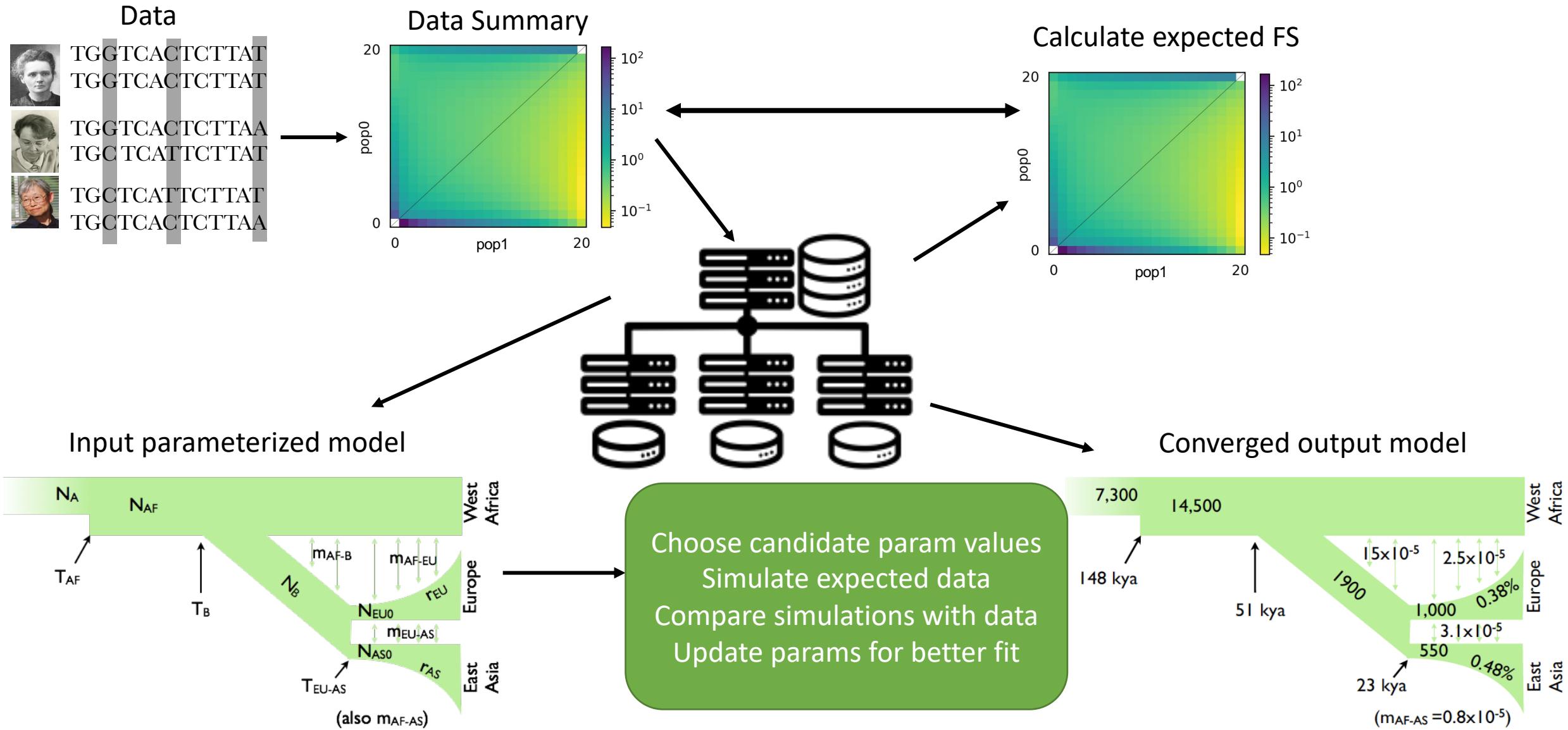
Diffusion approximation for demographic inference (dadi)



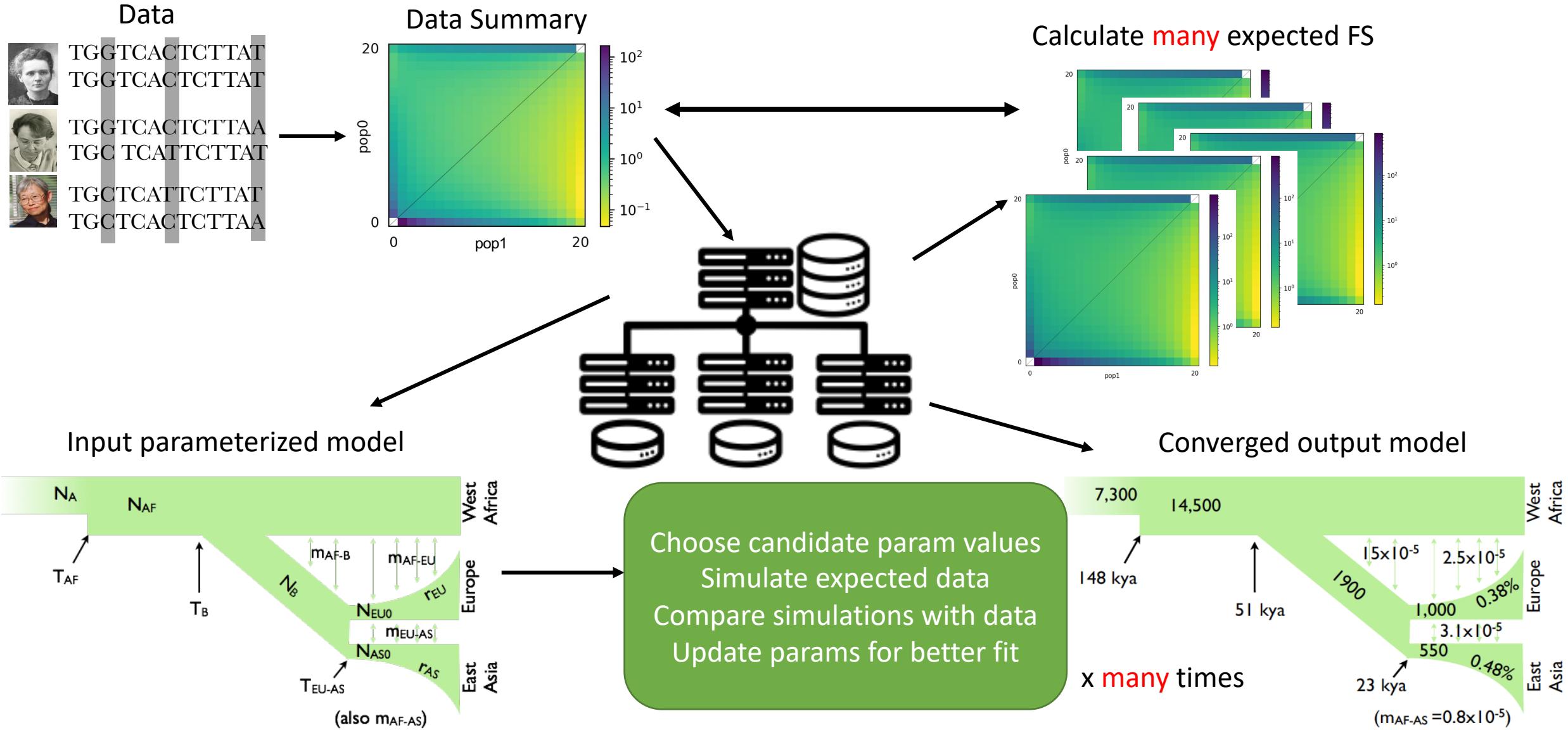
Diffusion approximation for demographic inference (dadi)



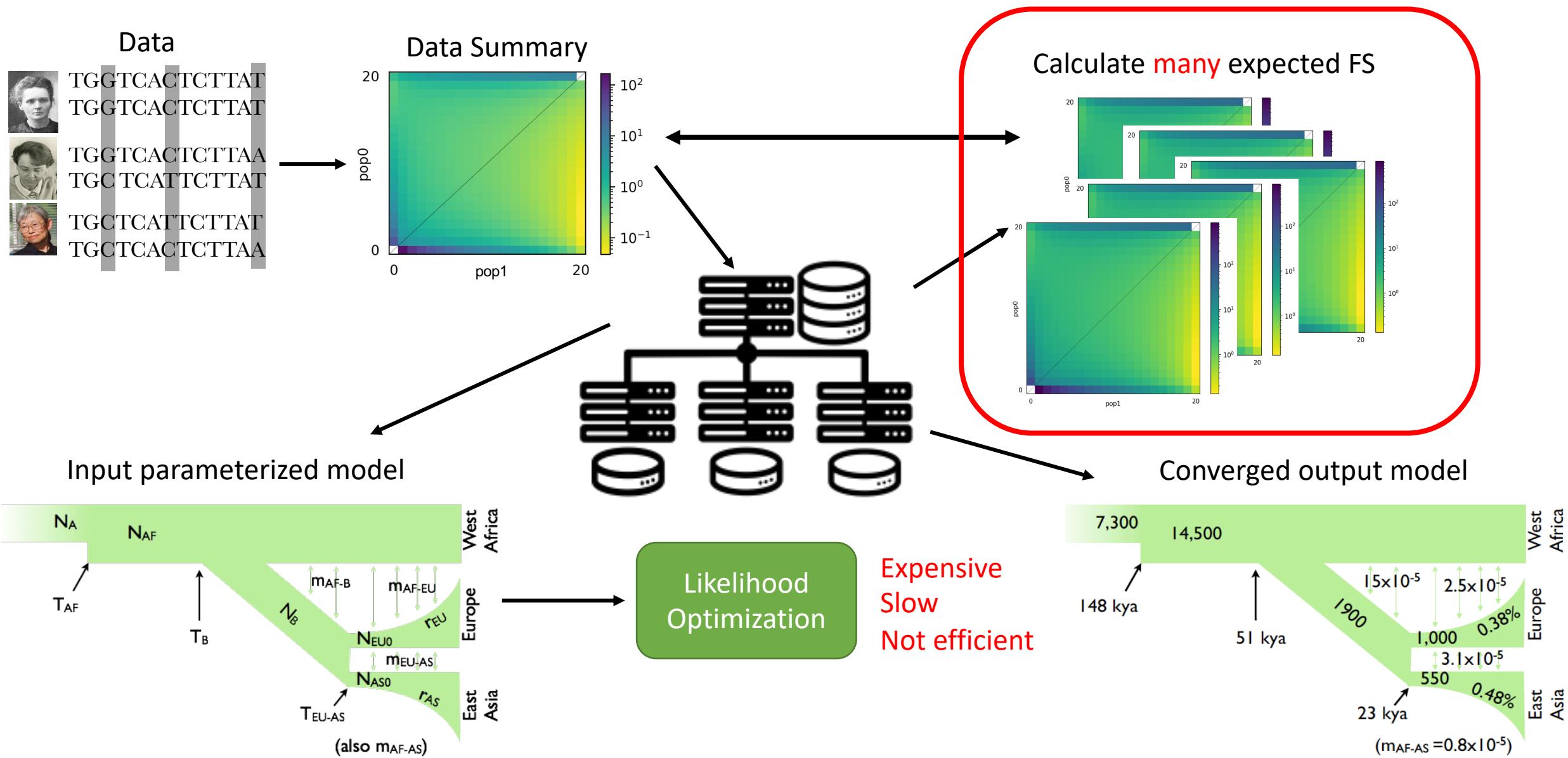
Diffusion approximation for demographic inference (dadi)



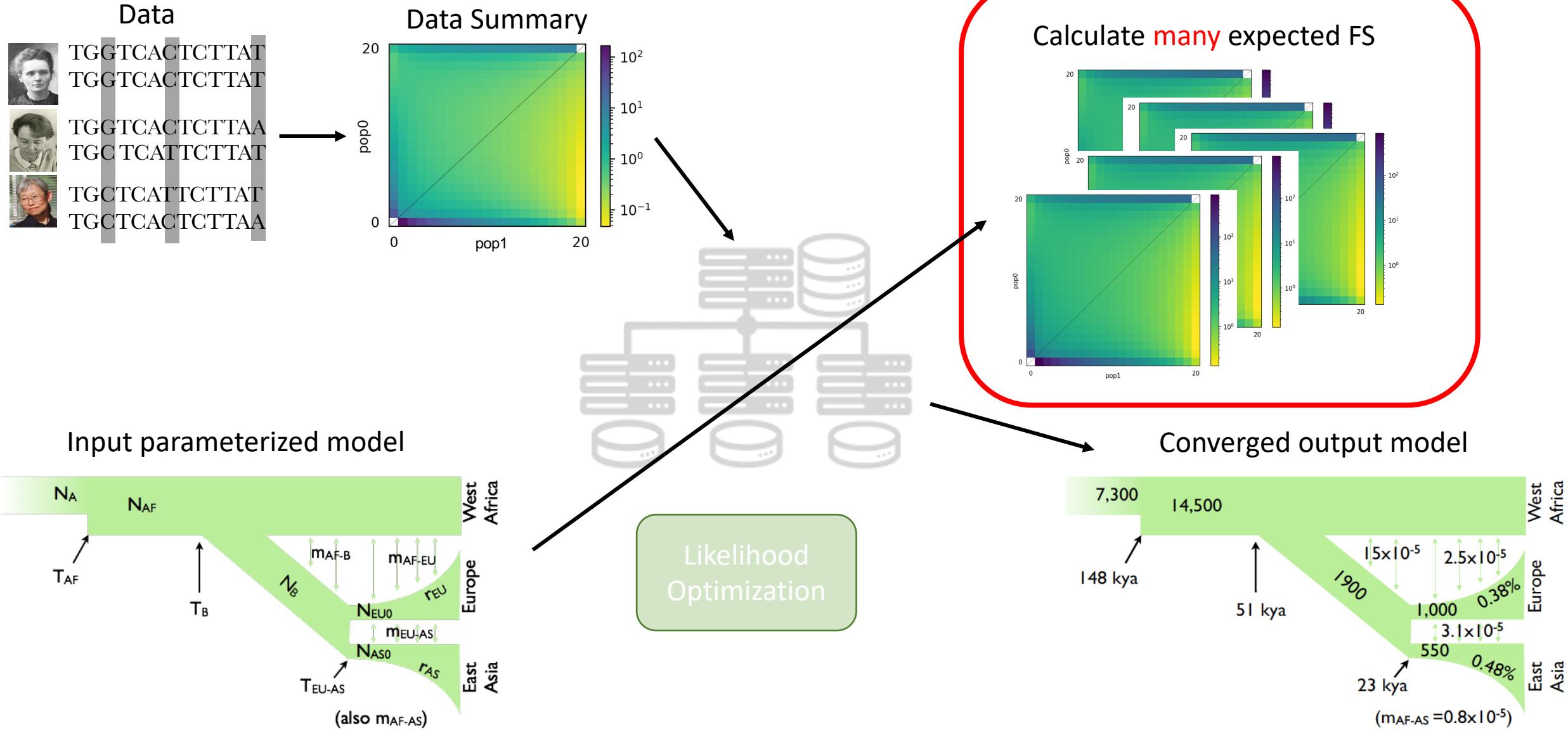
Diffusion approximation for demographic inference (dadi)



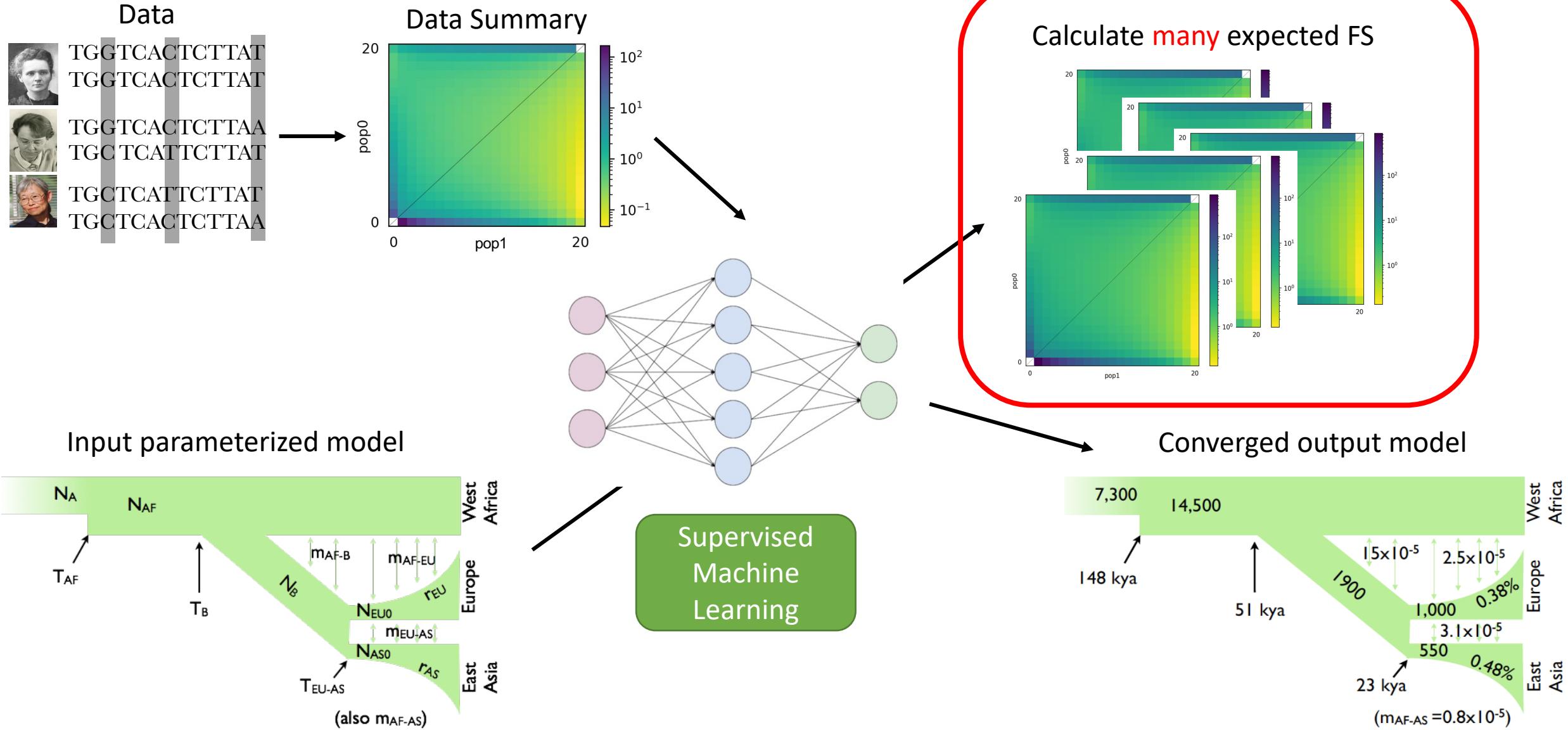
Diffusion approximation for demographic inference (dadi)



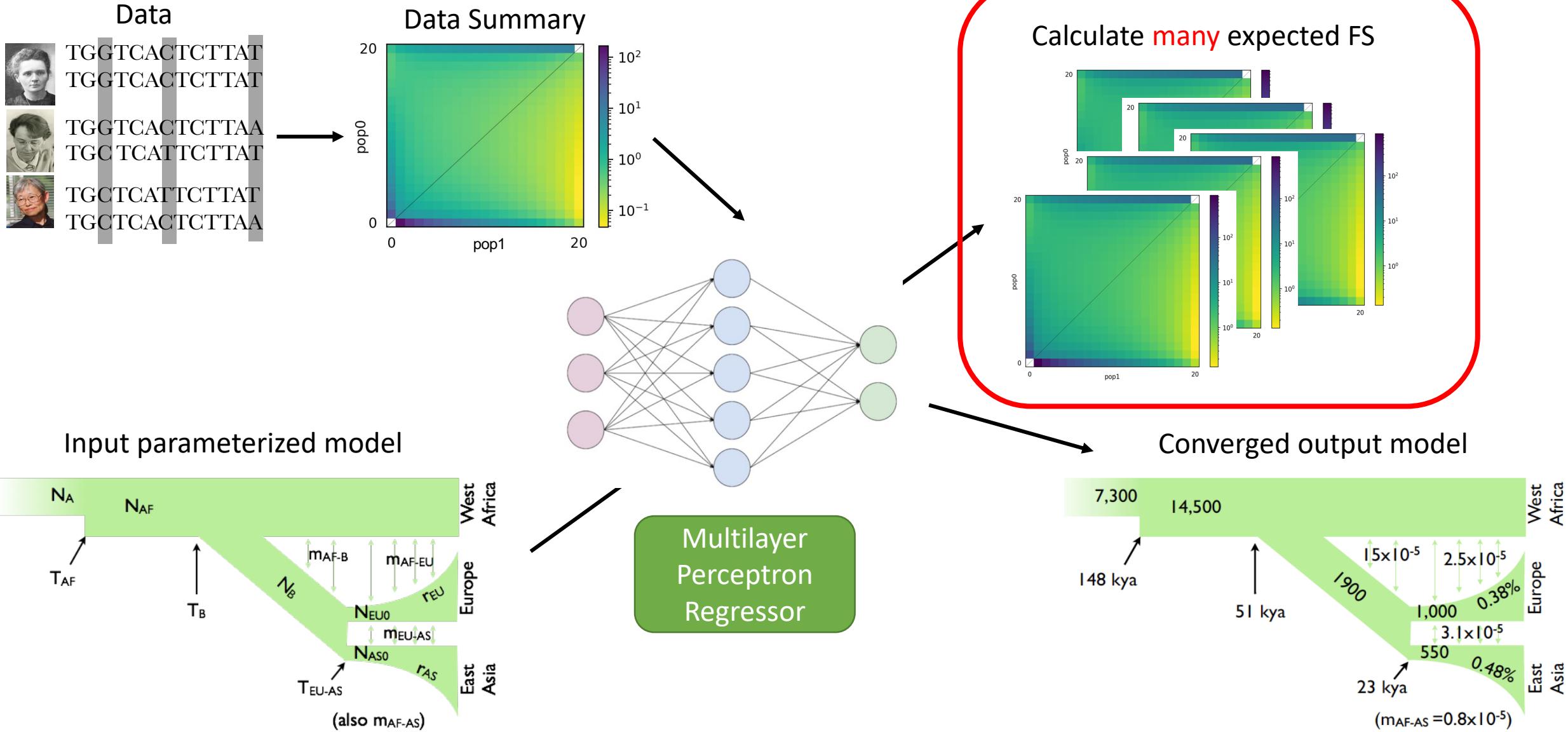
Diffusion approximation for demographic inference (dadi)



Diffusion approximation for demographic inference (dadi)



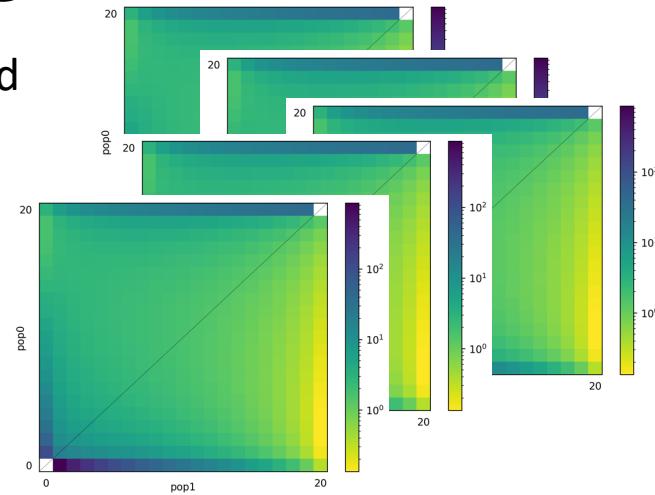
Diffusion approximation for demographic inference (dadi)



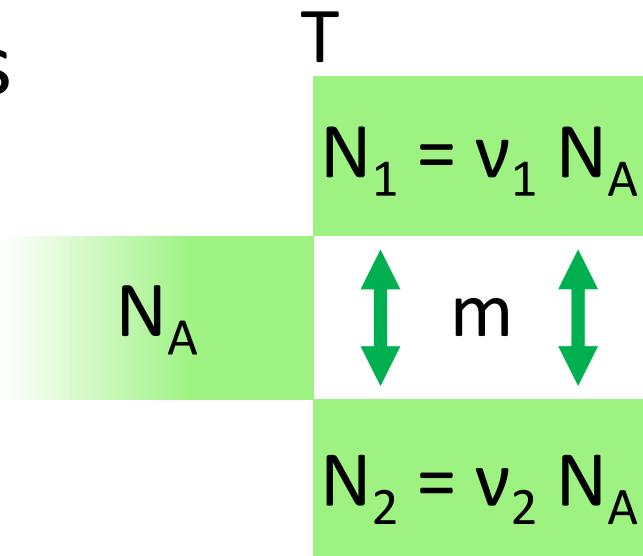
dadi-machine learning workflow

Training data

dadi-simulated
spectra
(1000)



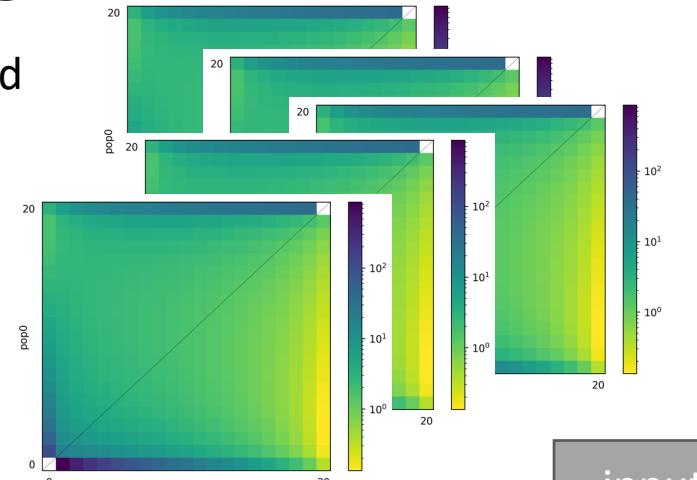
Labels
parameter
values



dadi-machine learning workflow

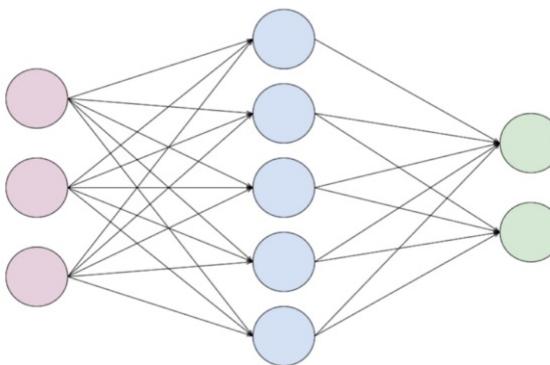
Training data

dadi-simulated
spectra
(1000)



MLP Regressor
algorithm

input



Labels

parameter
values

$$\begin{matrix} & T \\ N_1 = v_1 N_A & \end{matrix}$$

N_A

$\updownarrow m \updownarrow$

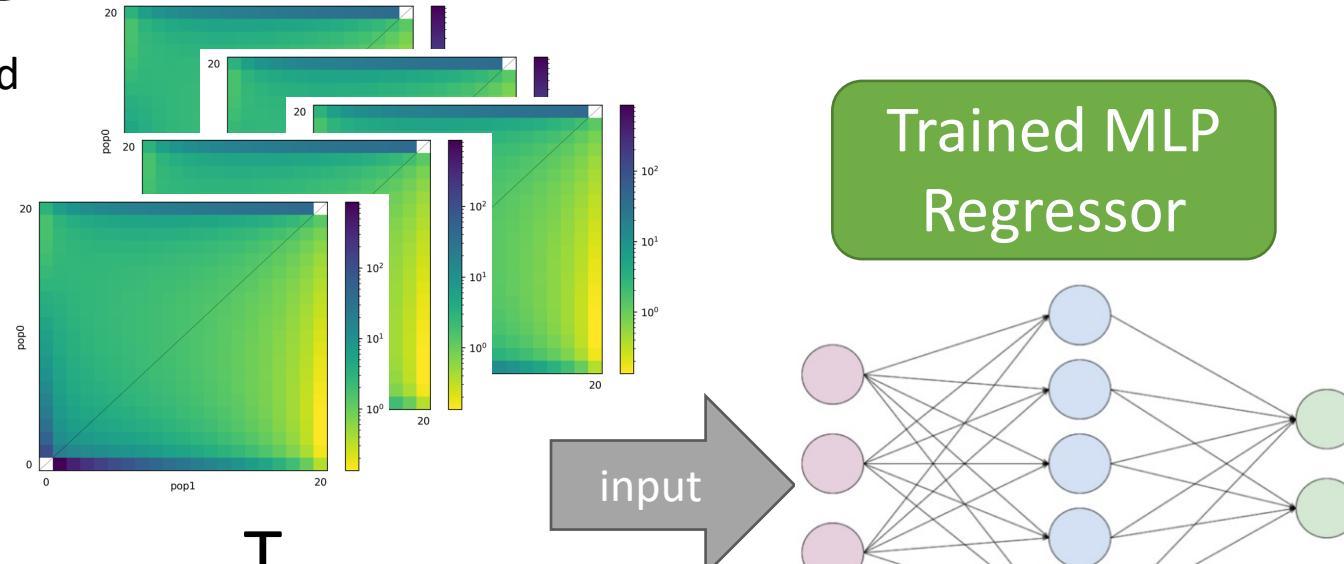
$$\begin{matrix} & \\ N_2 = v_2 N_A & \end{matrix}$$

Learns:
Finds pattern
of association

dadi-machine learning workflow

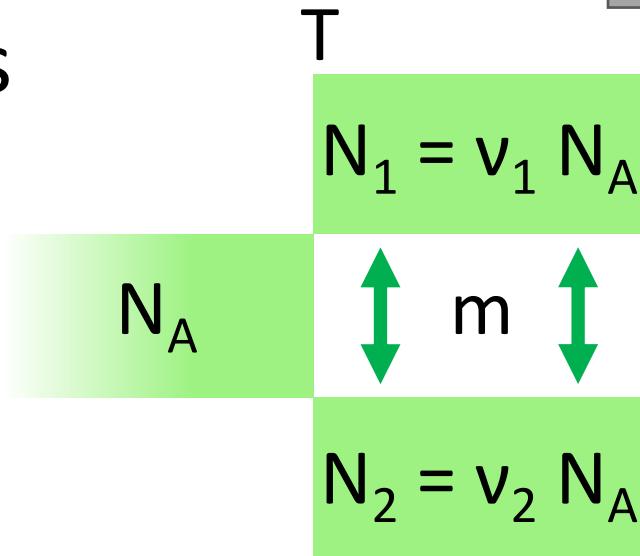
Training data

dadi-simulated
spectra
(1000)



Labels

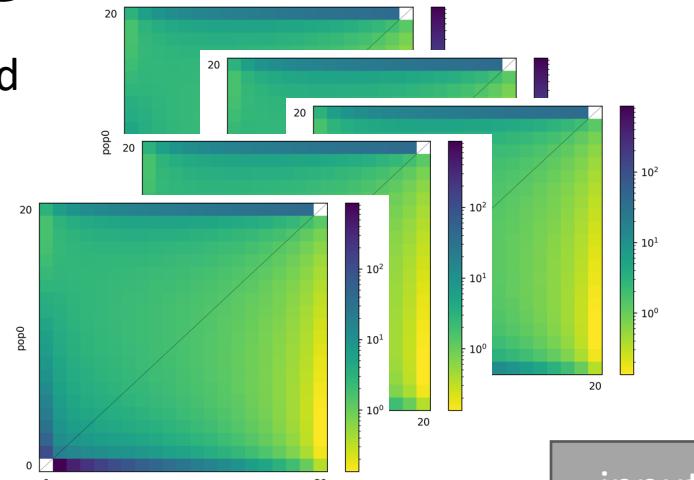
parameter
values



dadi-machine learning workflow

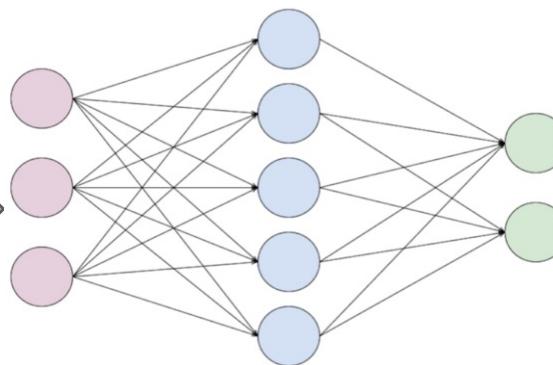
Training data

dadi-simulated
spectra
(1000)



Trained MLP
Regressor

input



Labels

parameter
values

$$\begin{matrix} & T \\ N_1 = v_1 N_A & \end{matrix}$$

N_A

$\updownarrow m \updownarrow$

$$\begin{matrix} & \\ N_2 = v_2 N_A & \end{matrix}$$

Strength of frequency spectra simulation with dadi:

- Faster and less expensive than coalescent simulations
- Organism-agnostic and transferability

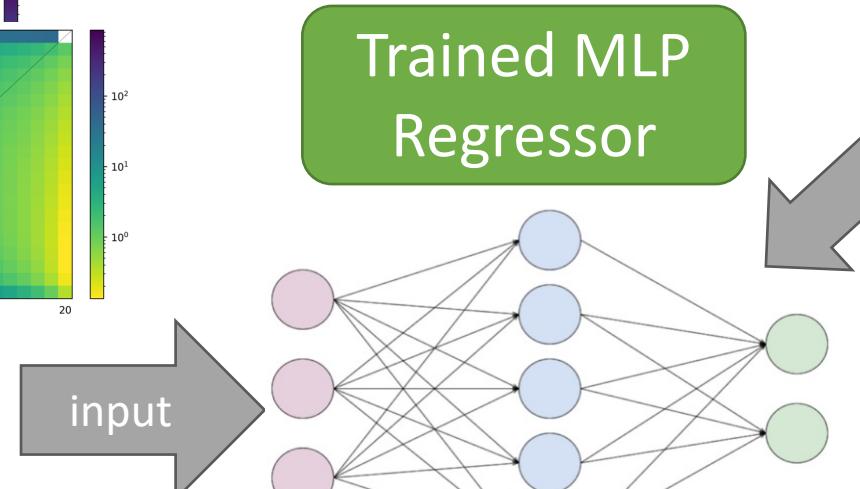
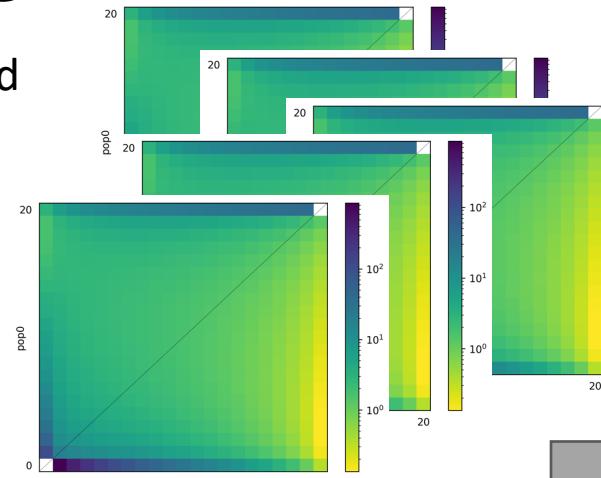
Weakness:

- Ignores linkage

dadi-machine learning workflow

Training data

dadi-simulated
spectra
(1000)

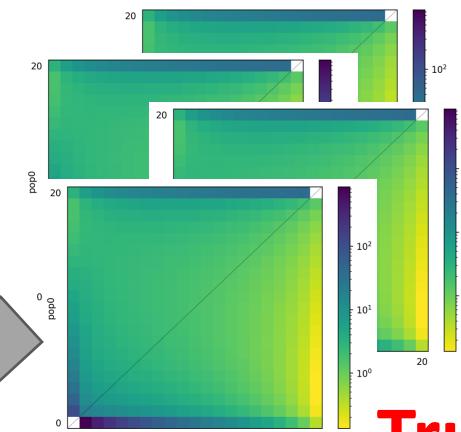


Labels

parameter
values

$$\begin{array}{c} T \\ \hline N_1 = v_1 N_A \\ \hline N_A \quad \updownarrow m \quad \updownarrow \\ \hline N_2 = v_2 N_A \end{array}$$

Test data



dadi- and
msprime-
simulated
spectra
(100-200)

True params:
set aside for
comparison

Strength of frequency spectra simulation with dadi:

- Faster and less expensive than coalescent simulations
- Organism-agnostic and transferability

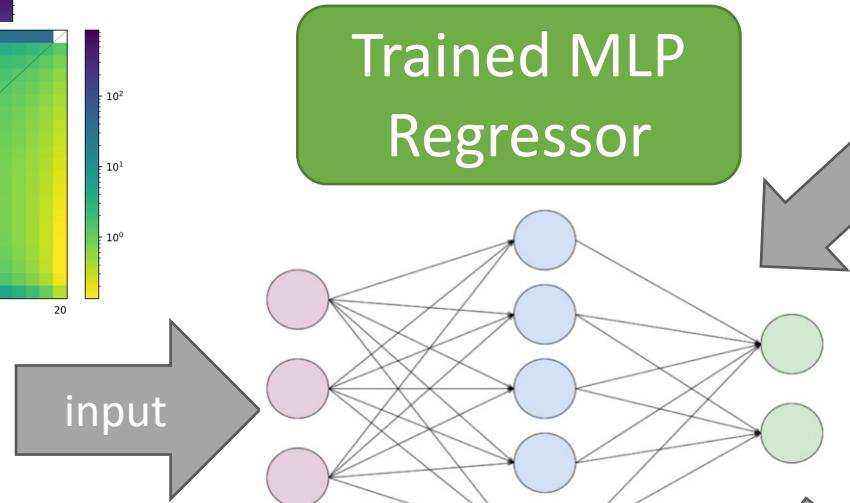
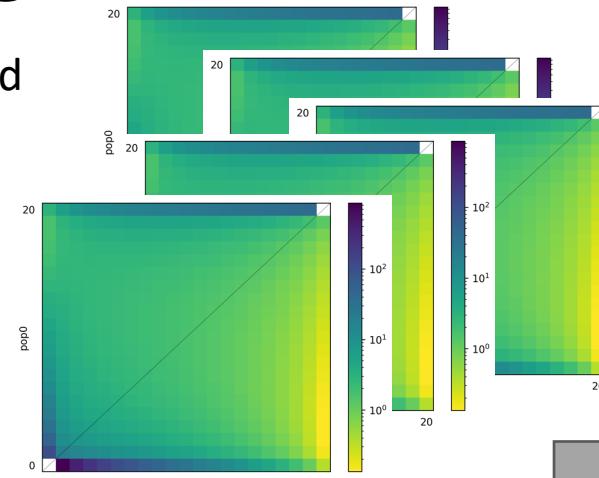
Weakness:

- Ignores linkage

dadi-machine learning workflow

Training data

dadi-simulated
spectra
(1000)



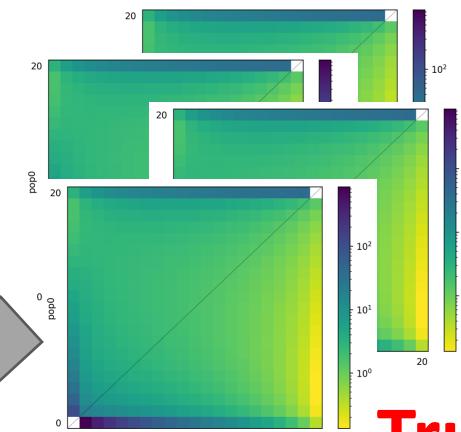
Trained MLP
Regressor

Labels

parameter
values

$$\begin{matrix} & T \\ N_1 = v_1 N_A & \\ \uparrow & \downarrow \\ N_A & m \\ \uparrow & \downarrow \\ N_2 = v_2 N_A & \end{matrix}$$

Test data



dadi- and
msprime-
simulated
spectra
(100-200)

True params:
set aside for
comparison

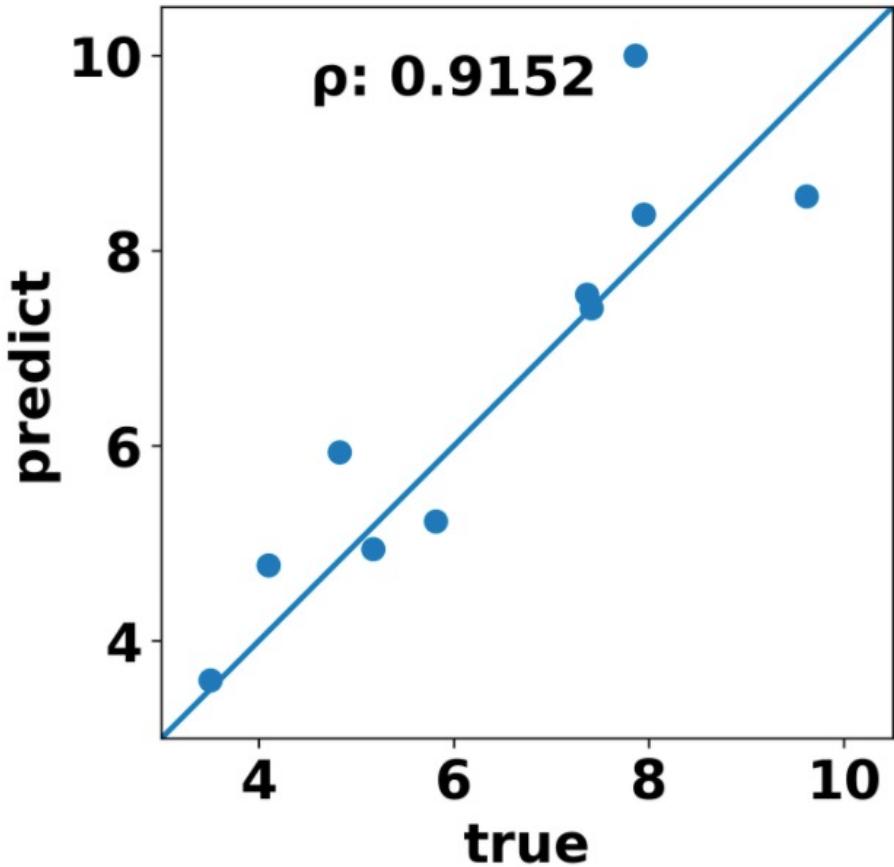
Prediction

param values for
each test spectrum

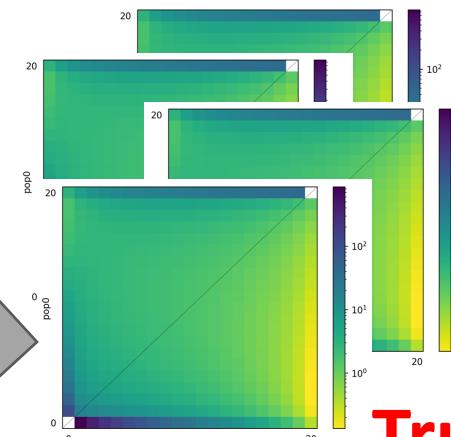
Predict :
 $v_1 = 0.1$
 $v_2 = 0.5$
 $T = 1.3$
 $m = 5$

Testing and validation

Prediction accuracy

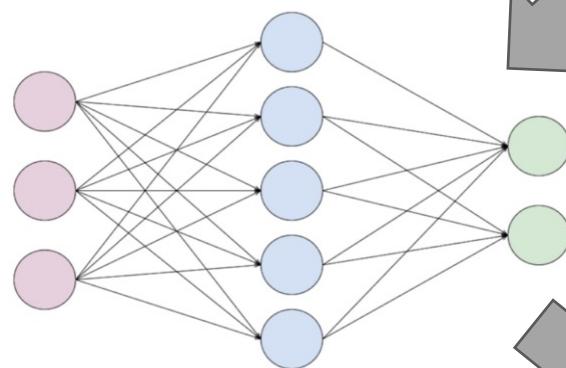


Test data



dadi- and
msprime-
simulated
spectra
(100-200)

Trained MLP
Regressor



True params:
set aside for
comparison

Prediction

param values for
each test spectrum

Predict :
 $v_1 = 0.1$
 $v_2 = 0.5$
 $T = 1.3$
 $m = 5$

Results – prediction accuracy

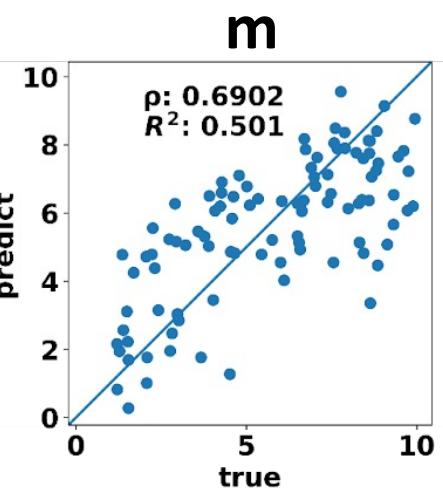
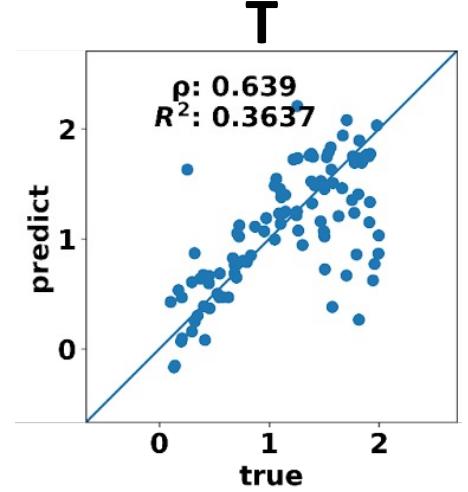
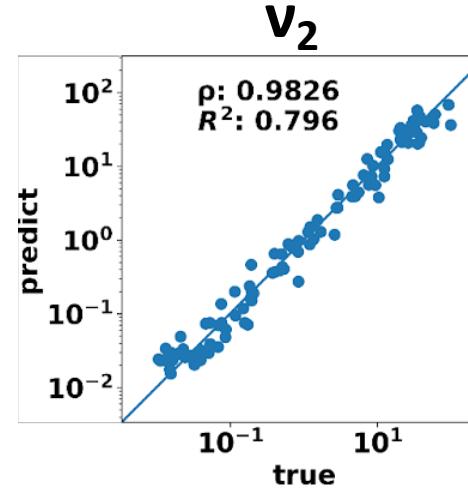
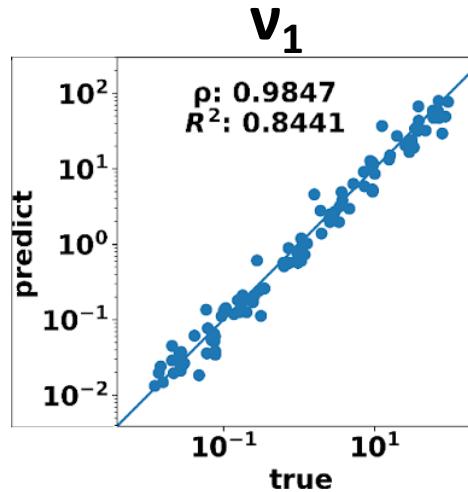
two-population
split-migration model

$$N_1 = v_1 N_A$$

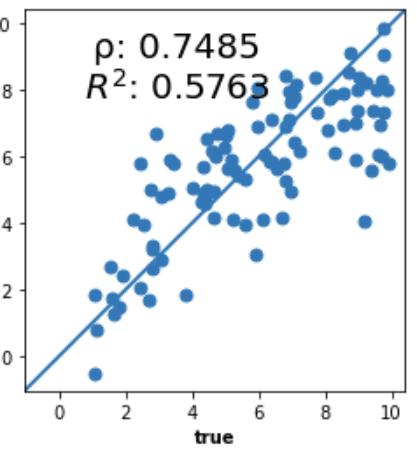
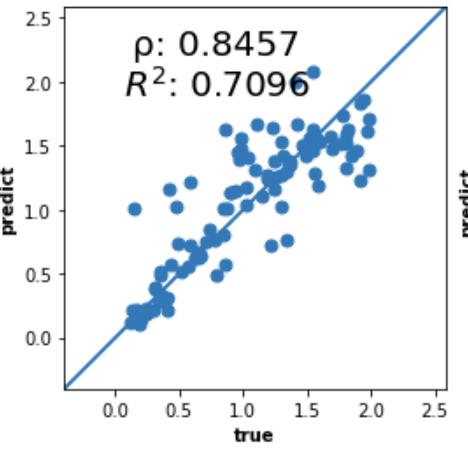
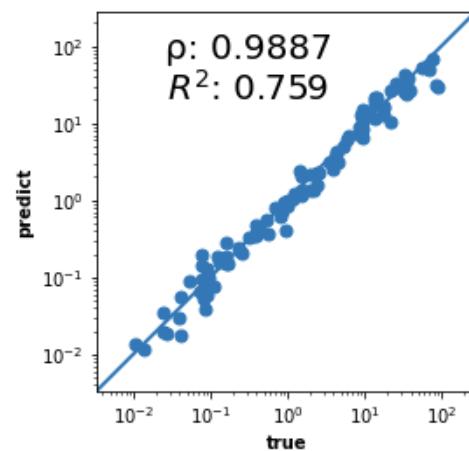
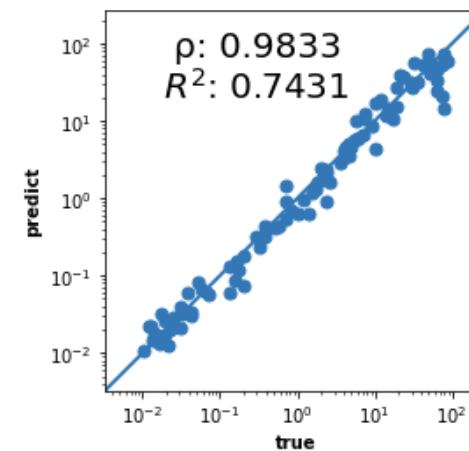
$$N_A \uparrow \downarrow m \uparrow \downarrow$$

$$N_2 = v_2 N_A$$

MLPR



dadi-simulated
data sets



msprime-simulated
data sets

Results – prediction accuracy

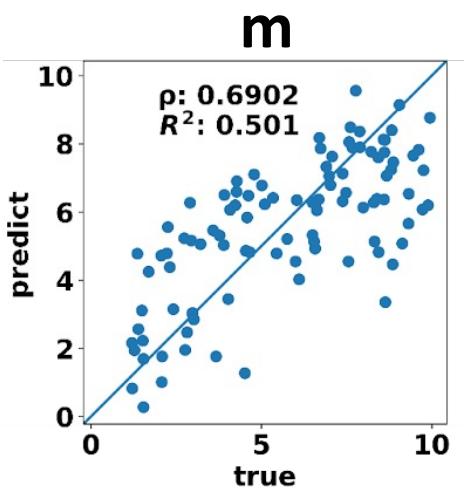
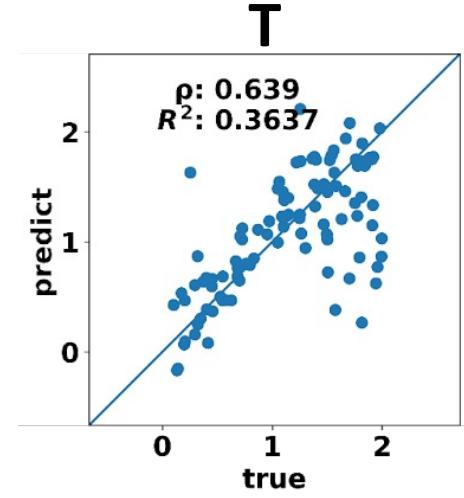
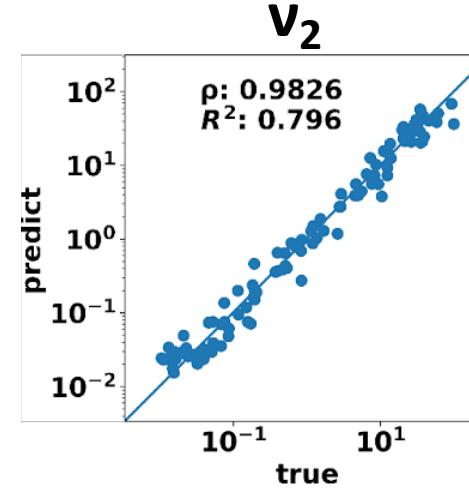
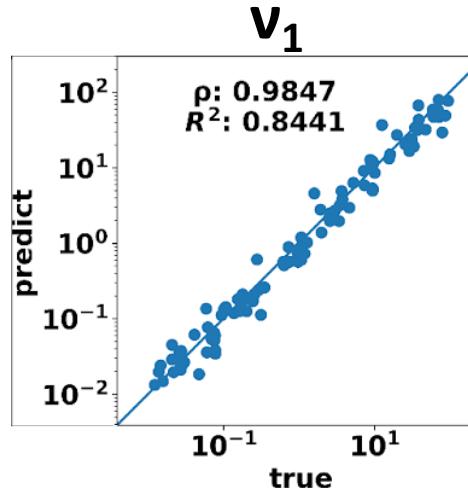
two-population
split-migration model

$$N_1 = v_1 N_A$$

$$N_A \uparrow \downarrow m \uparrow \downarrow$$

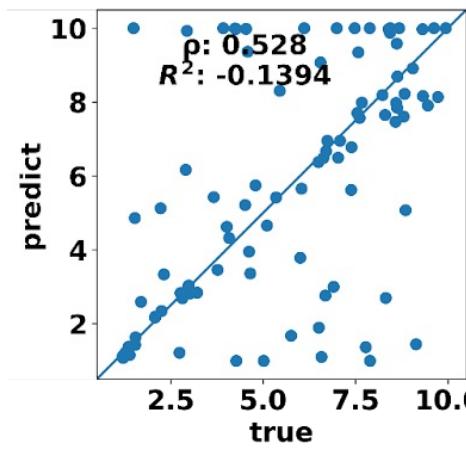
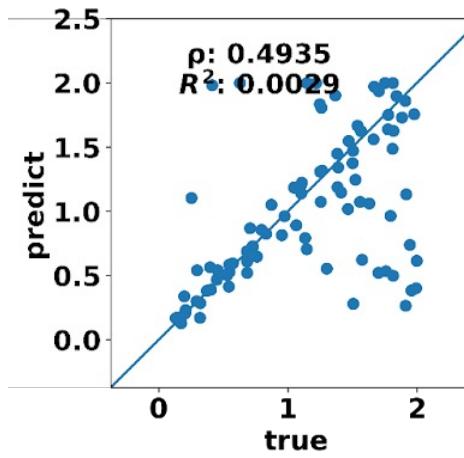
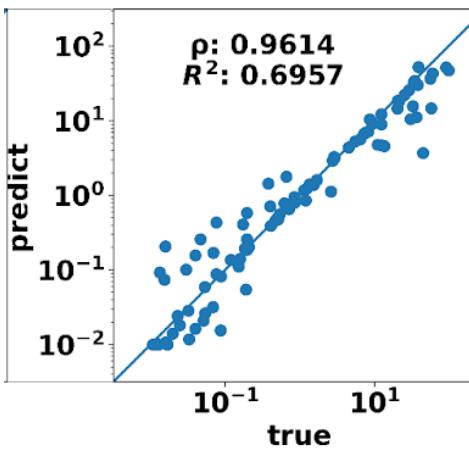
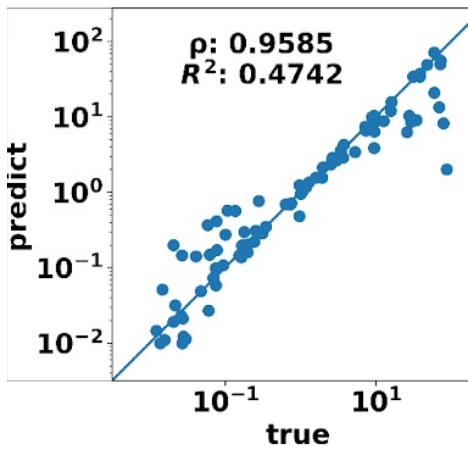
$$N_2 = v_2 N_A$$

MLPR



dadi-simulated
data sets

dadi



dadi-simulated
data sets

Results – prediction accuracy and computational efficiency

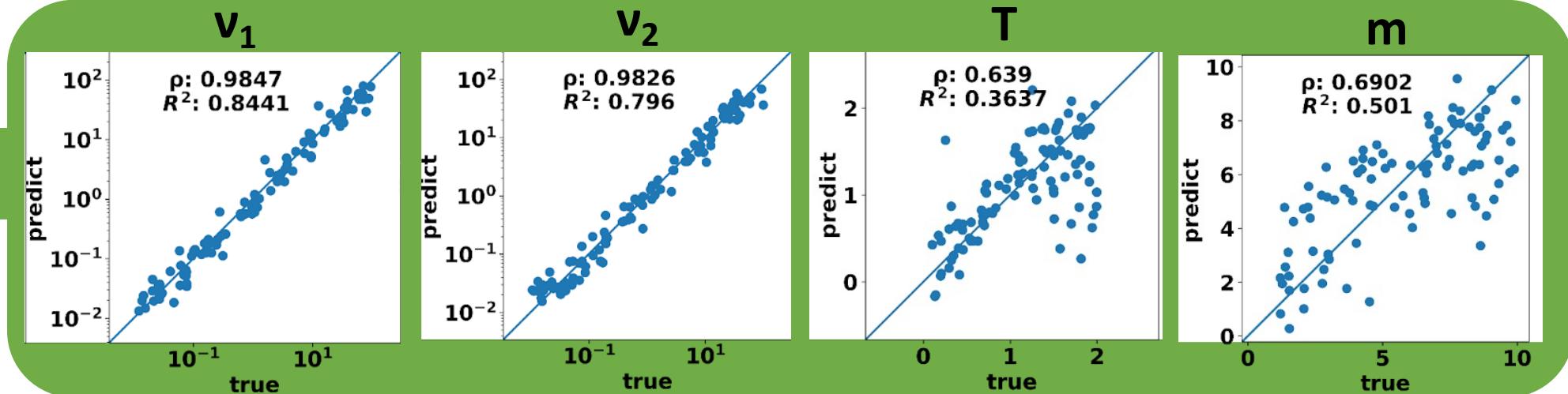
two-population
split-migration model

$$N_1 = v_1 N_A$$

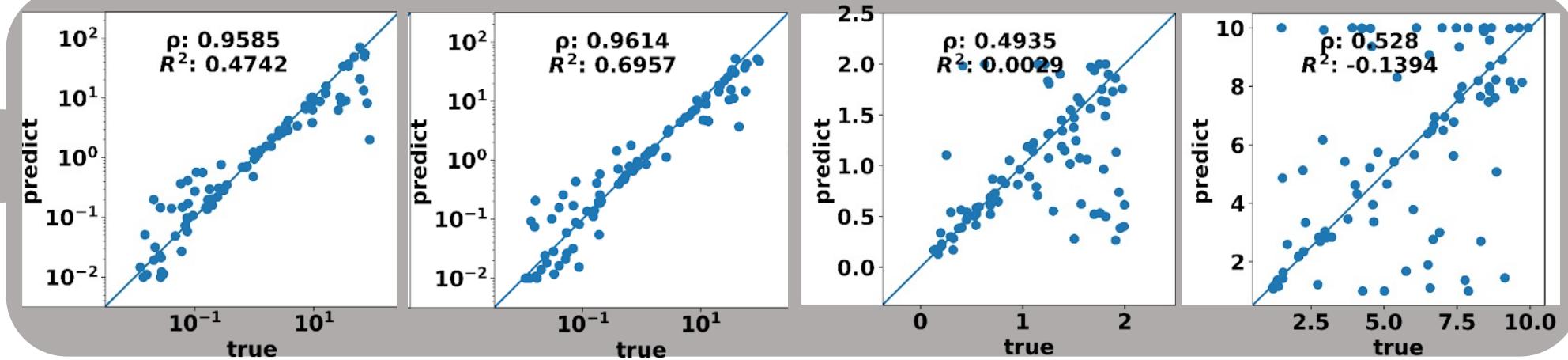
$$\uparrow \downarrow m \uparrow \downarrow$$

$$N_2 = v_2 N_A$$

MLPR



dadi



Summary

- Machine learning approach to improve the computational efficiency of a widely used likelihood-based demographic inference method
- Comparable accuracy with significantly reduced computational cost & complexity
- Robust to test simulations generated with linkage
- Accompanying uncertainty quantification method that provides prediction intervals
- Trained MLPs for frequently used demographic models and workflow will be distributed



@LnTran26

<http://gutengroup.arizona.edu>

