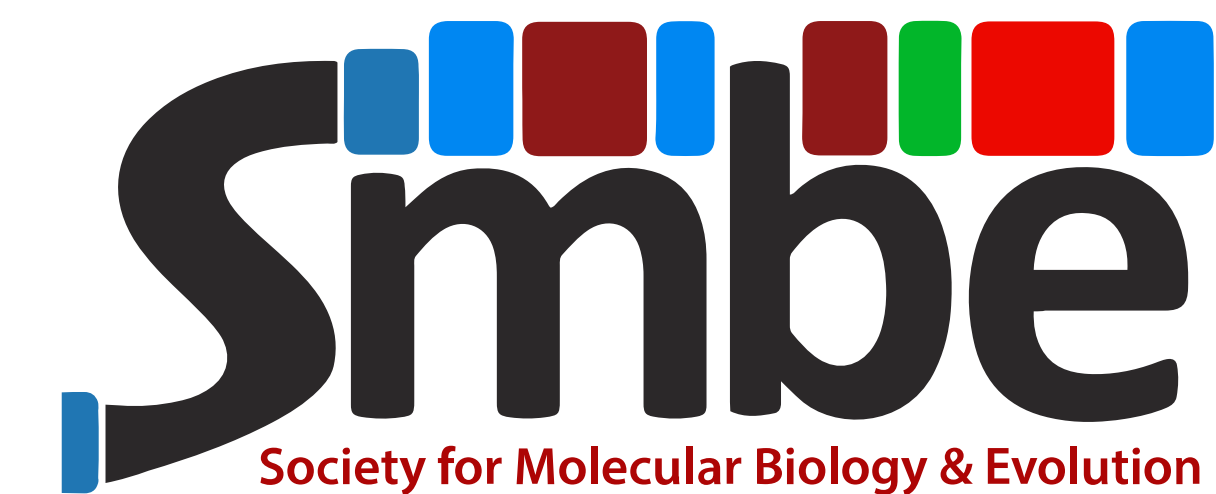
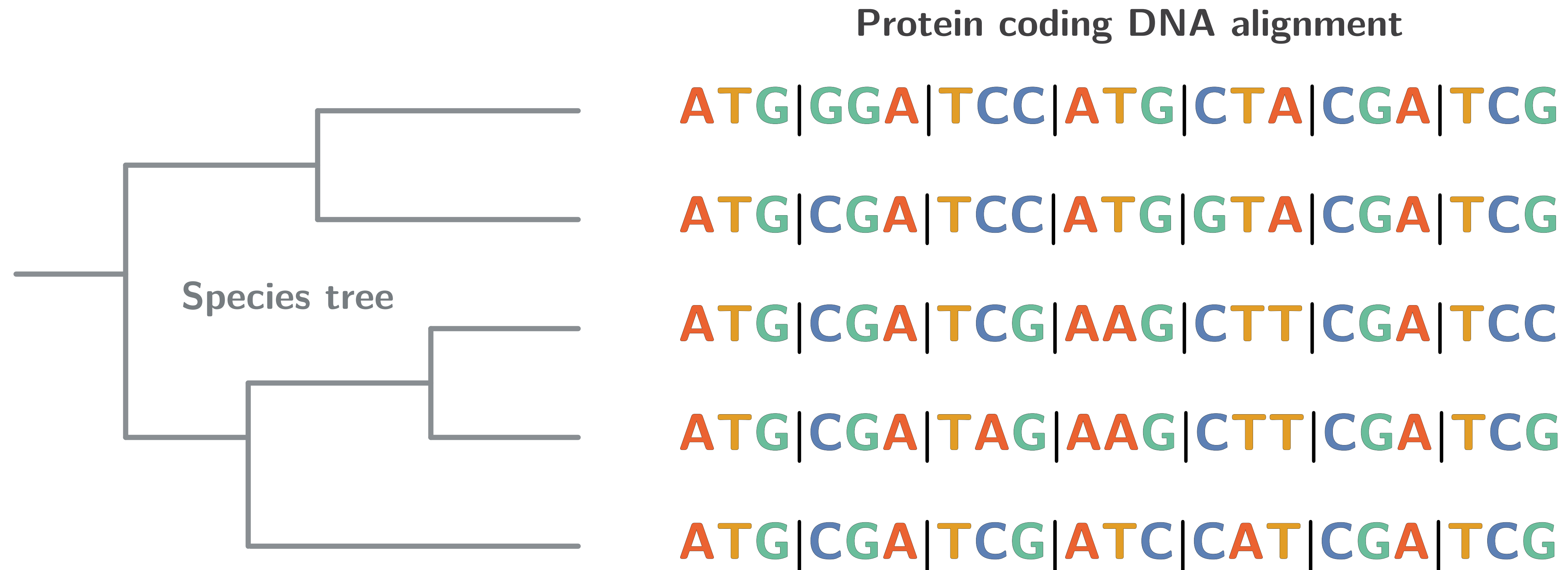


Inferring fluctuating population size and selection with phylogenetic codon models

Thibault Latrille - PhD student
Nicolas Lartillot - Research director

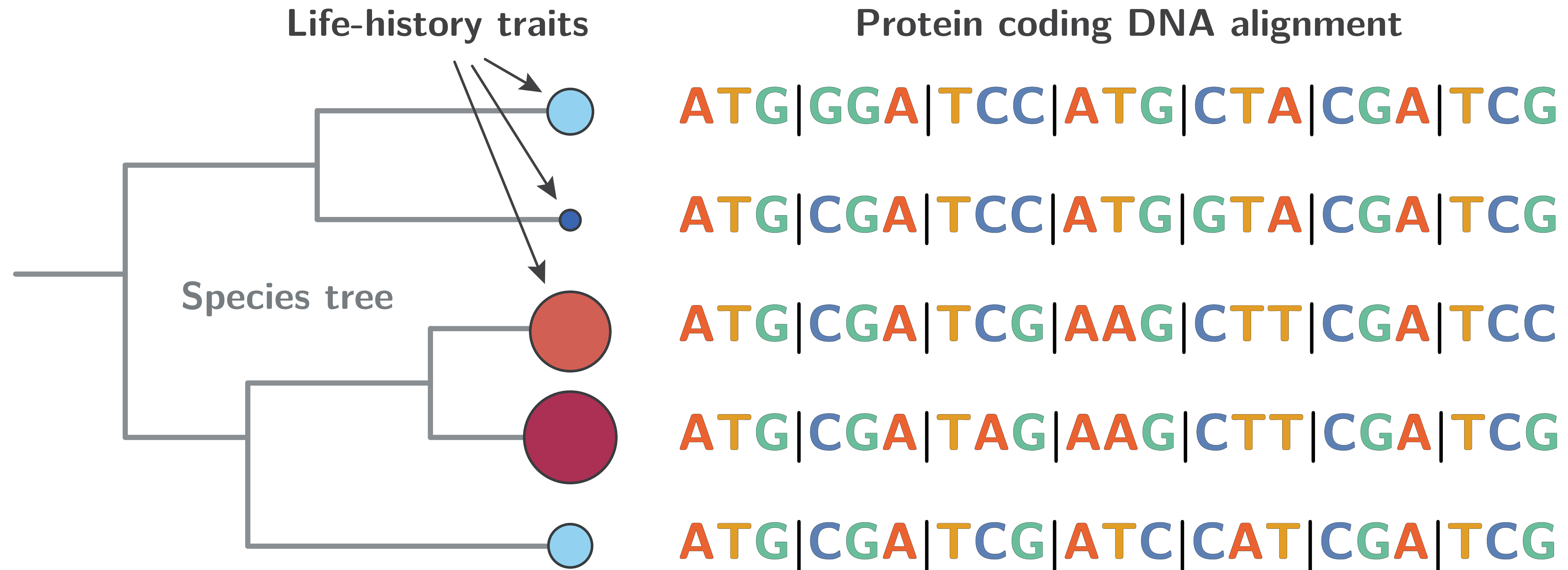


Theoretical and empirical motivations



- Can we quantify selection acting on protein coding and long term changes in effective population size, from a DNA alignment?
- Do we have enough signal for empirical estimation of selection and drift?

Theoretical and empirical motivations



- How are species life-history traits (longevity, maturity, weight, size, ...) related to population-genetics parameters (effective population size, mutation rate,...)?
- Can we reconstruct changes in effective population size instead of using dN/dS as a proxy?

Stearns (1972); Lartillot & Poujol (2011); Weber et al (2014); Figuet (2016); James (2018).

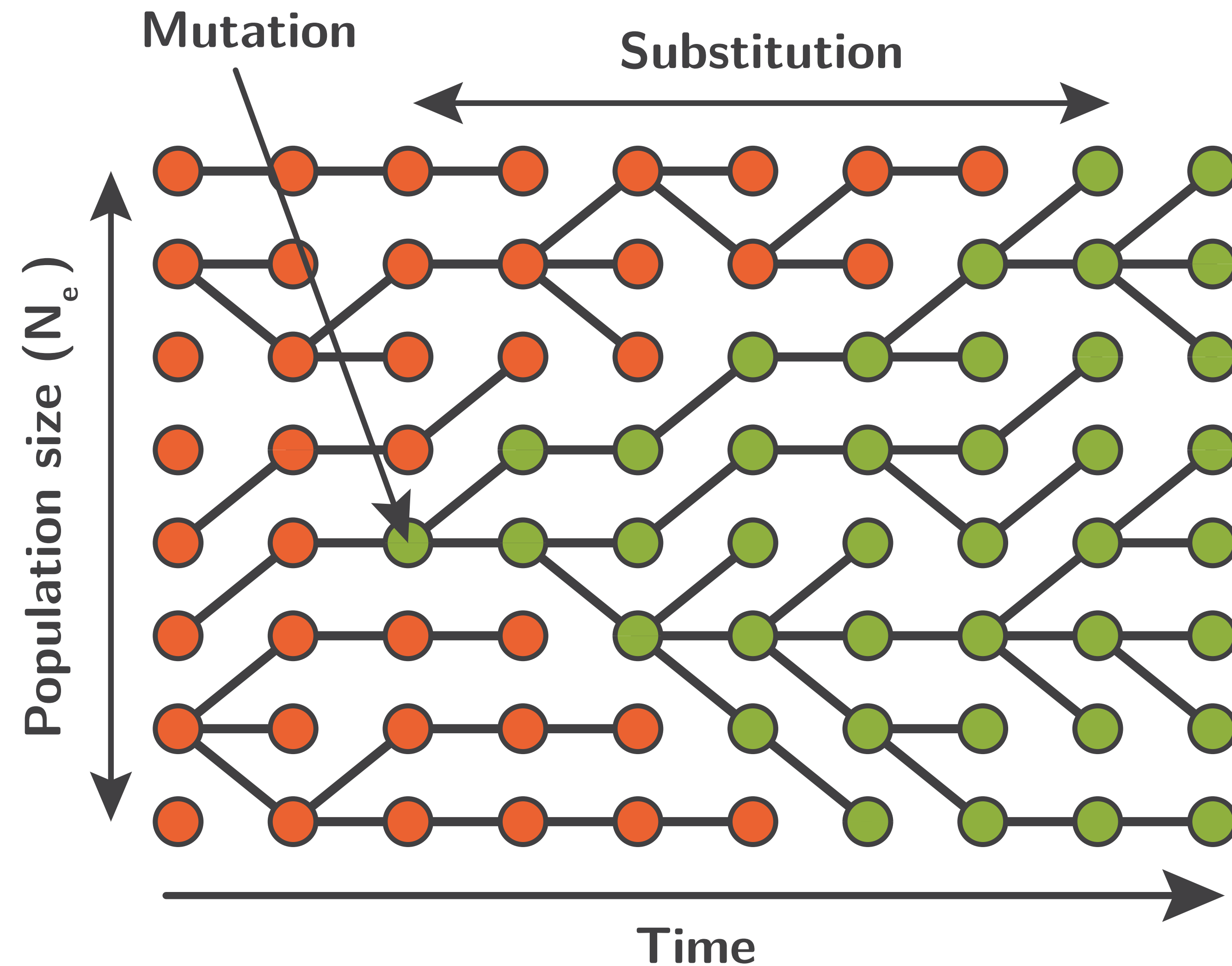
Mapping substitutions along the species tree



- **Synonymous** mutations (not changing the protein) are assumed to be neutral.
- **Non-synonymous** mutations (changing the protein) are assumed to be under selective pressure at the amino-acid level.

Stearns (1972); Lartillot & Poujol (2011); Weber et al (2014); Figuet (2016); James (2018).

Substitution is a mutation that reached fixation



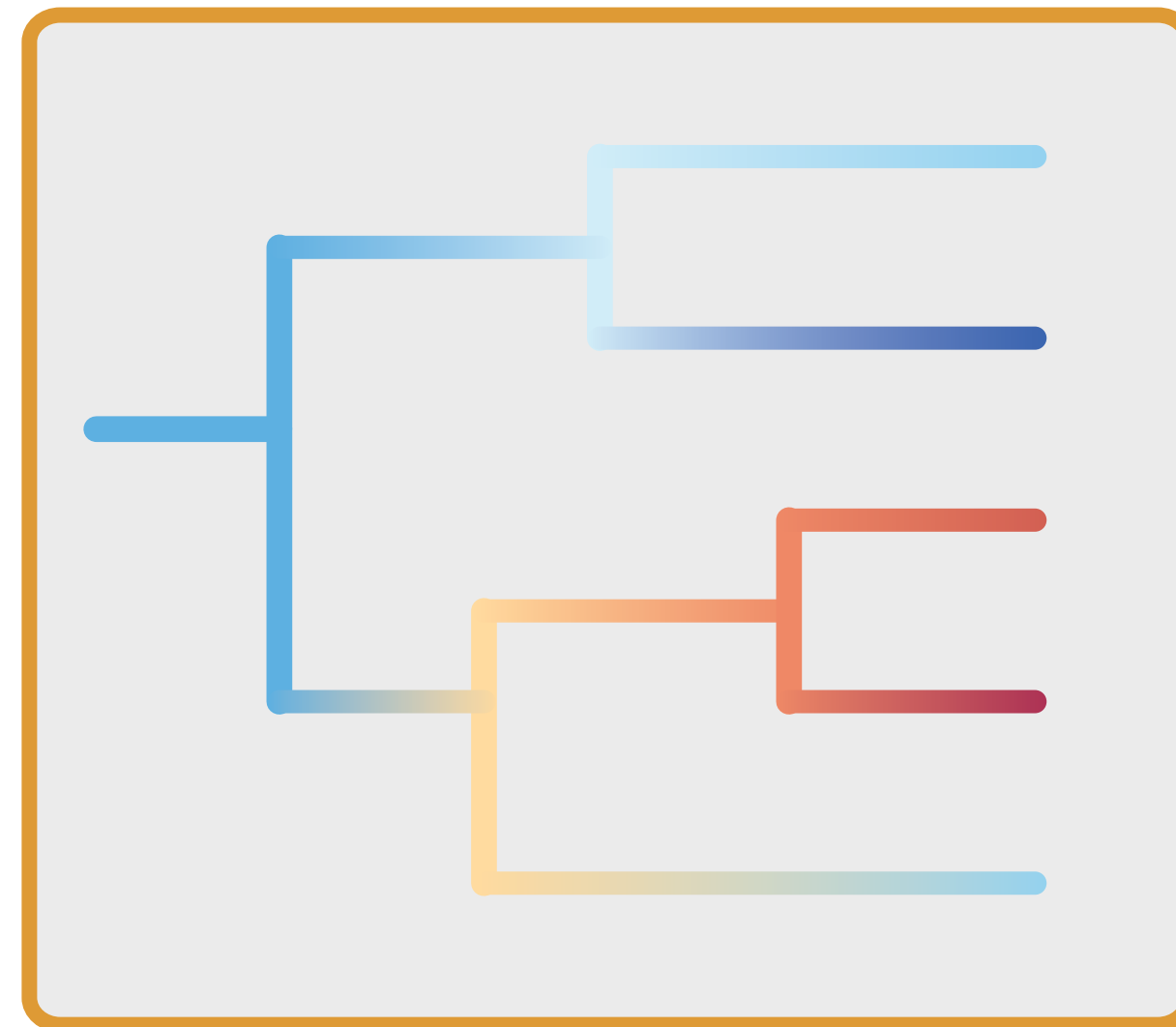
$$Q_{\text{ATT} \rightarrow \text{ATG}} = \mu_{\text{T} \rightarrow \text{G}} \frac{4N_e [F_{\text{Met}} - F_{\text{Ile}}]}{1 - e^{4N_e [F_{\text{Ile}} - F_{\text{Met}}]}}$$

- $Q_{\text{ATT} \rightarrow \text{ATG}}$ is the substitution rate from codon **ATT** to **ATG**.
- N_e is the effective population size.
- $\mu_{\text{T} \rightarrow \text{G}}$ is the mutation rate from nucleotide **T** to **G**.
- F_{Ile} (F_{Met}) is the fitness of Isoleucine (Methionine).

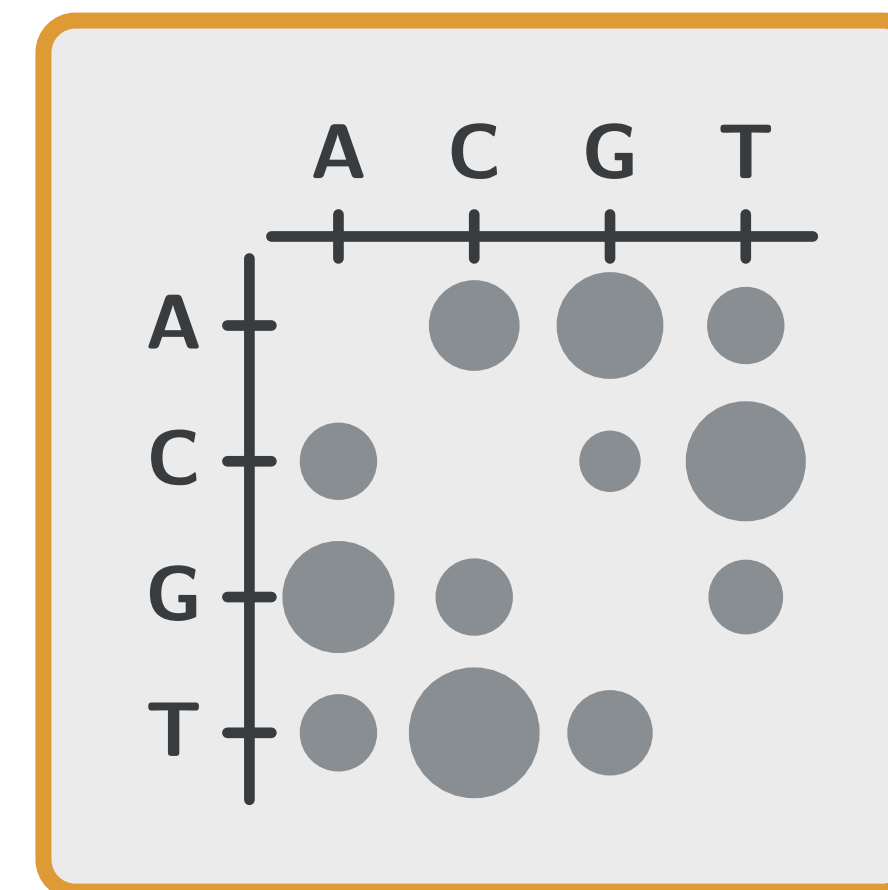
Kimura (1983); Ohta (1992); McClandish (2016).

Mechanistic mutation-selection codon models

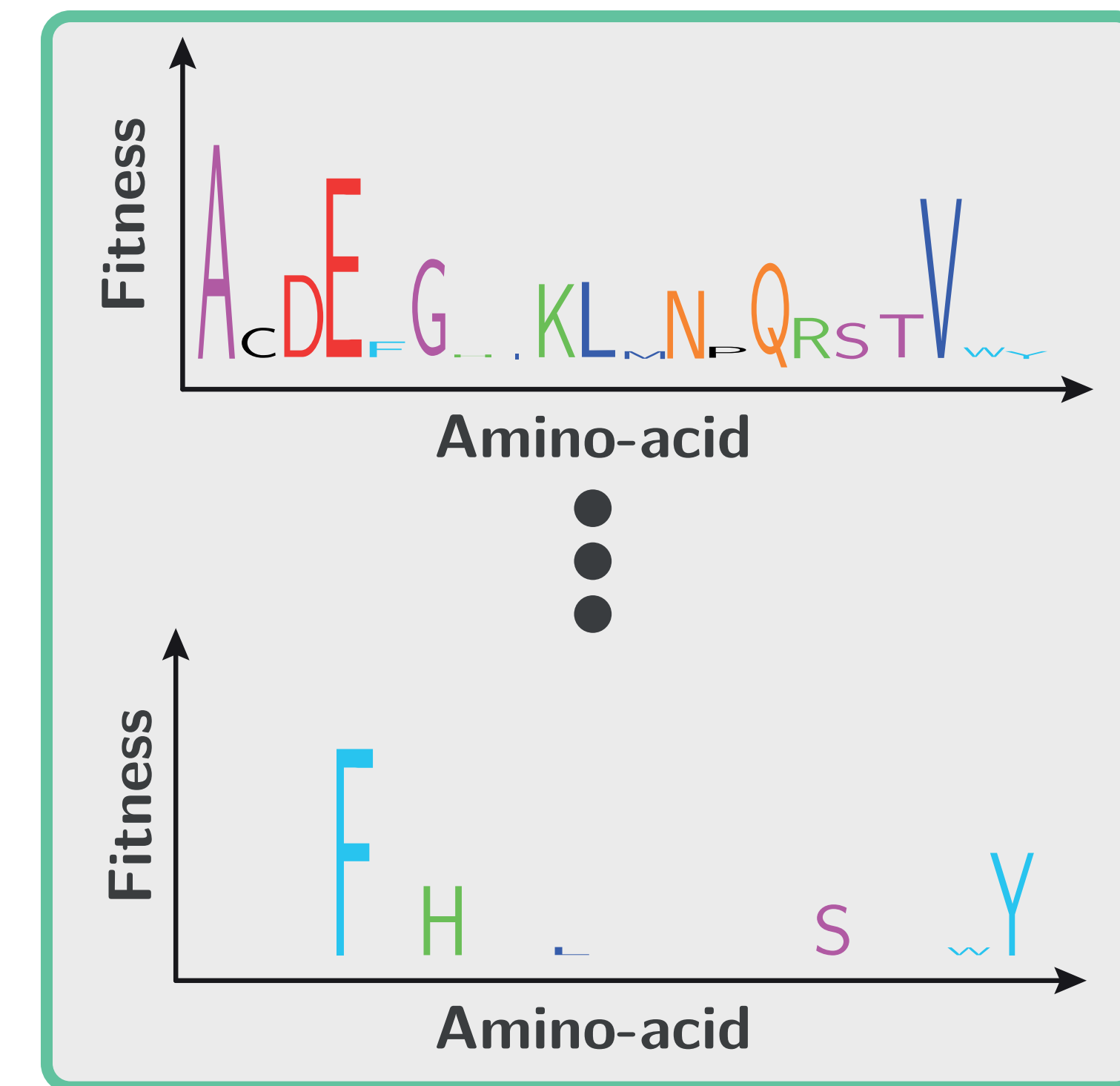
Mutation rate per
unit of time
(each branch of the tree)



Relative mutation
rate between
nucleotides



Fitness profile of 20 amino-acids
(each site of the alignment)

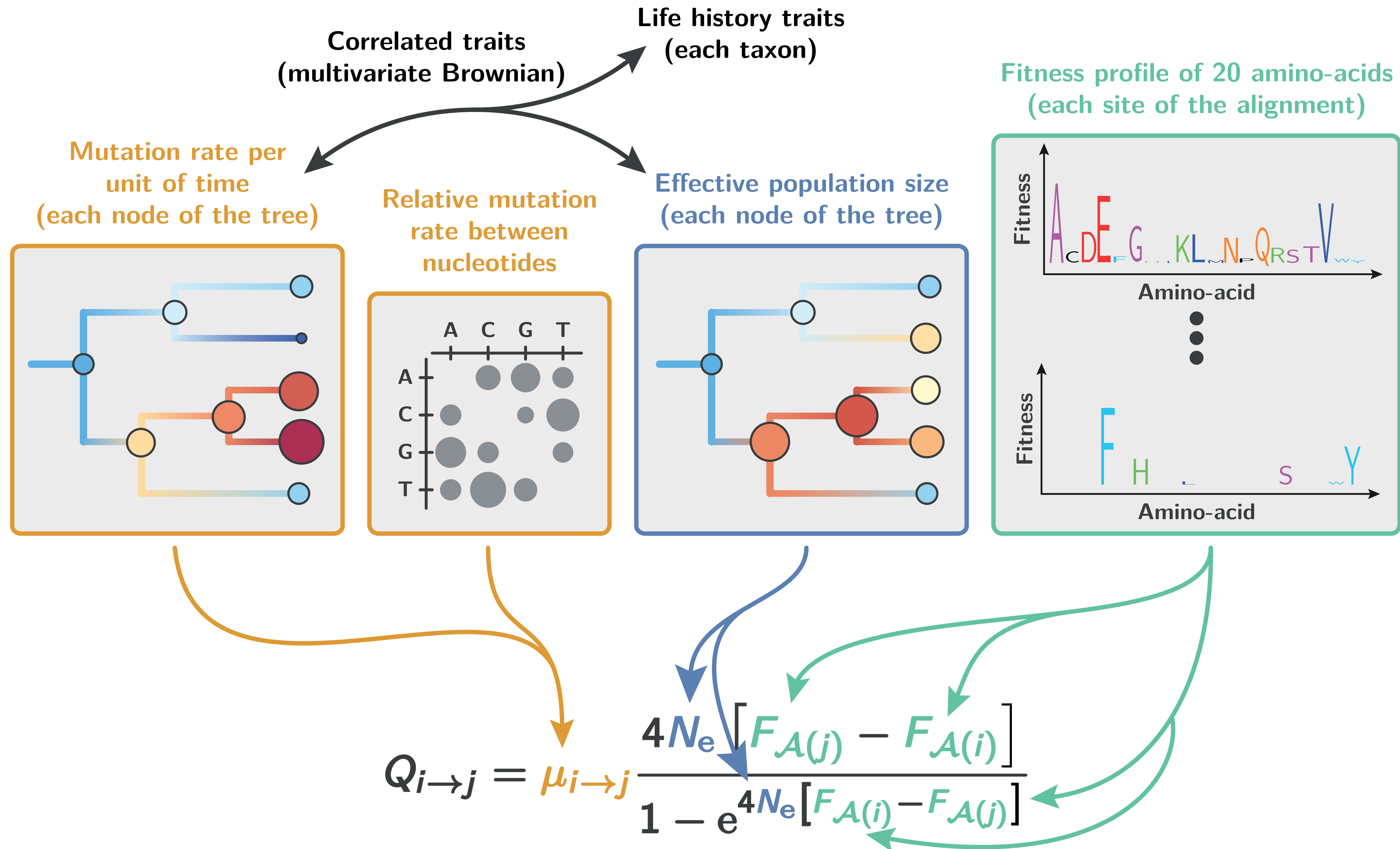


Effective population size
is considered constant

$$Q_{i \rightarrow j} = \mu_{i \rightarrow j} \frac{4N_e [F_{\mathcal{A}(j)} - F_{\mathcal{A}(i)}]}{1 - e^{4N_e [F_{\mathcal{A}(i)} - F_{\mathcal{A}(j)}]}}$$

Halpern & Bruno (1998); Tamuri *et al* (2014); Rodrigue & Lartillot (2016).

Reconstructing changes in effective population size



<https://github.com/bayesiancook/bayescodetree/chronogram>

What need to be estimated?

Tree:

- Age for each internal node of the dated tree.

Mutation:

- Mutation rate (per unit of time) for each node of the tree.
- Nucleotide relative rate matrix.

Selection:

- 20 Amino-acid fitnesses for each profile category (K categories).
- Which profile category (1..K) for each site of the alignment.

Drift:

- Population size for each node of the tree.

Traits:

- Life-history-traits for each node of the tree.
- Covariance matrix between traits (molecular and life-history).

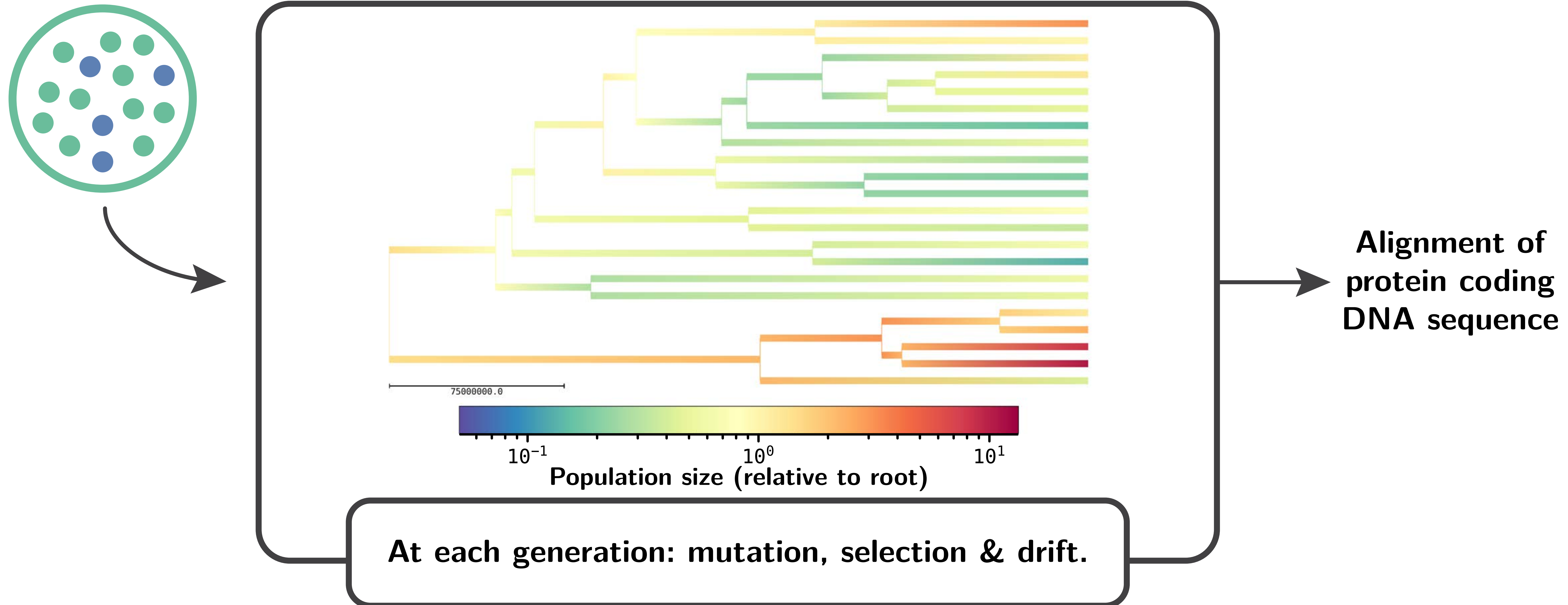
Wright-Fisher simulator along the phylogeny

Initial population size of 5000 individuals.

Mutation rate of 10^{-8} per generation per site.

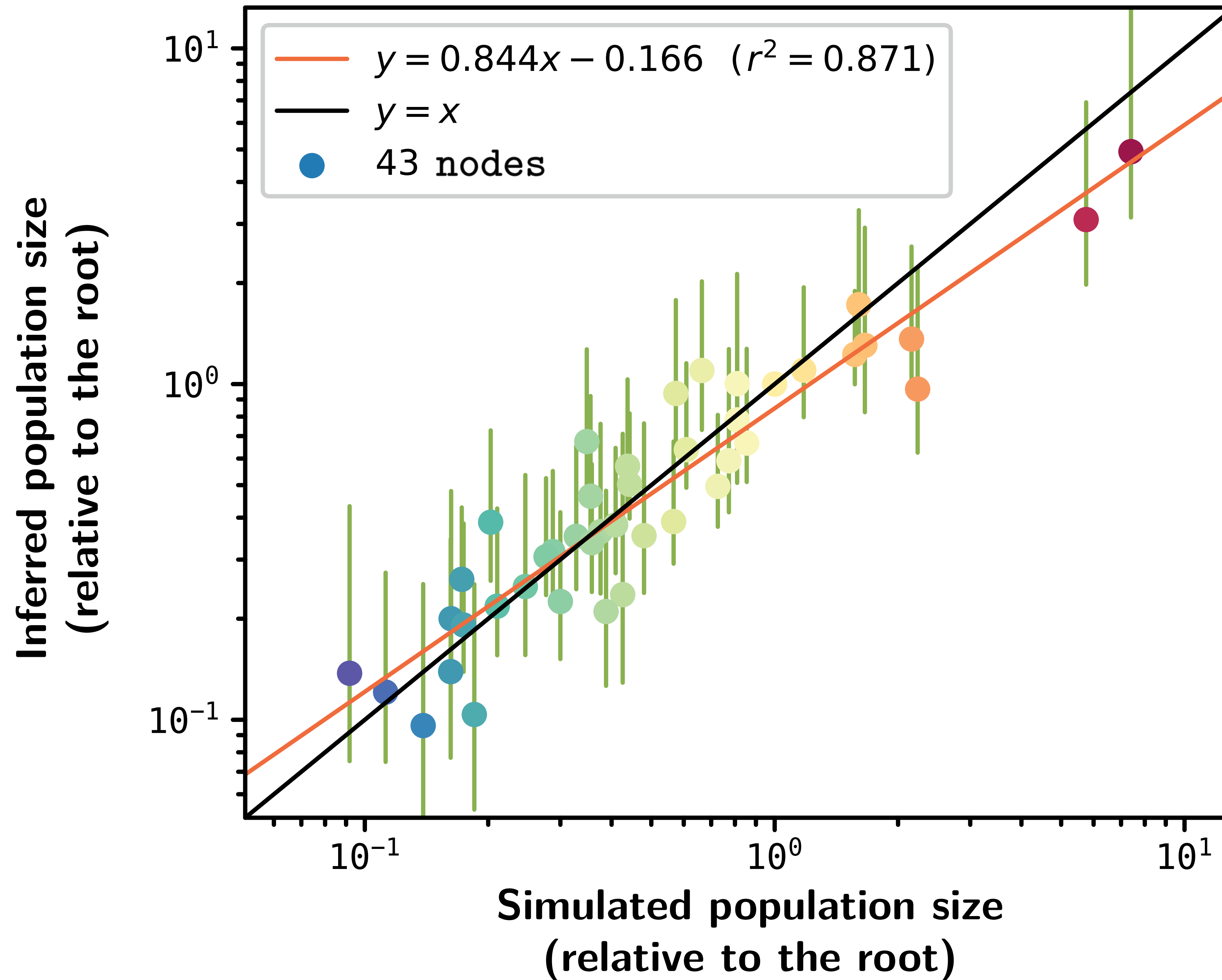
Generation time of 20 years.

Sequence of 15 000 codon sites (50 exons).



<https://github.com/ThibaultLatrille/SimuEvol>

Inference against simulated data (1/2)



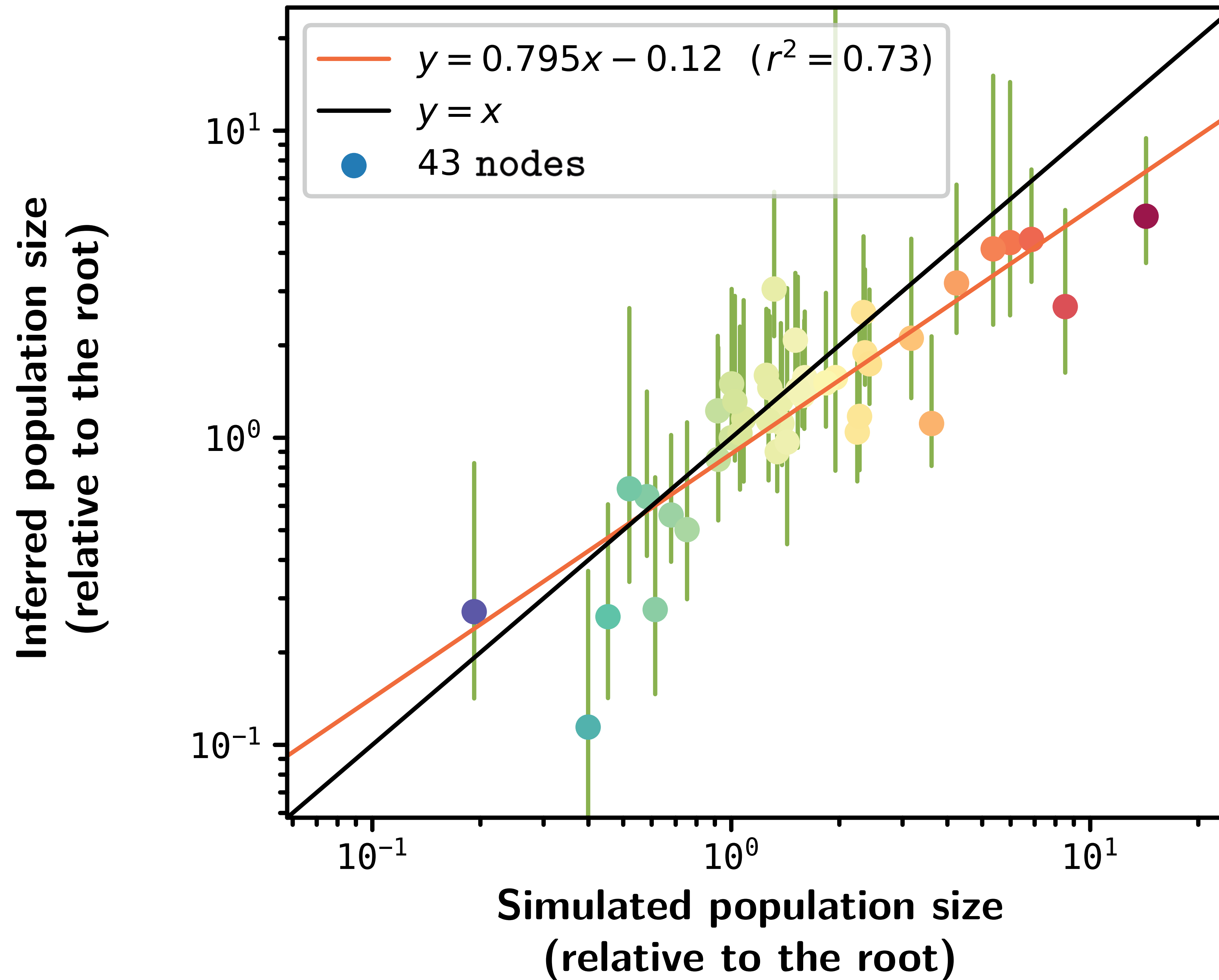
Taken into account by the simulator:

- Long term fluctuation of population size, mutation rate per time, generation time;
- Fitness landscape for each site;
- Finite population size;
- Allele hitching;

Not taken into account by the simulator:

- Species tree \neq gene tree;
- Epistasis;
- Fluctuating fitness landscape;
- Biased gene conversion;
- Selection on codon usage;

Inference against simulated data (2/2)



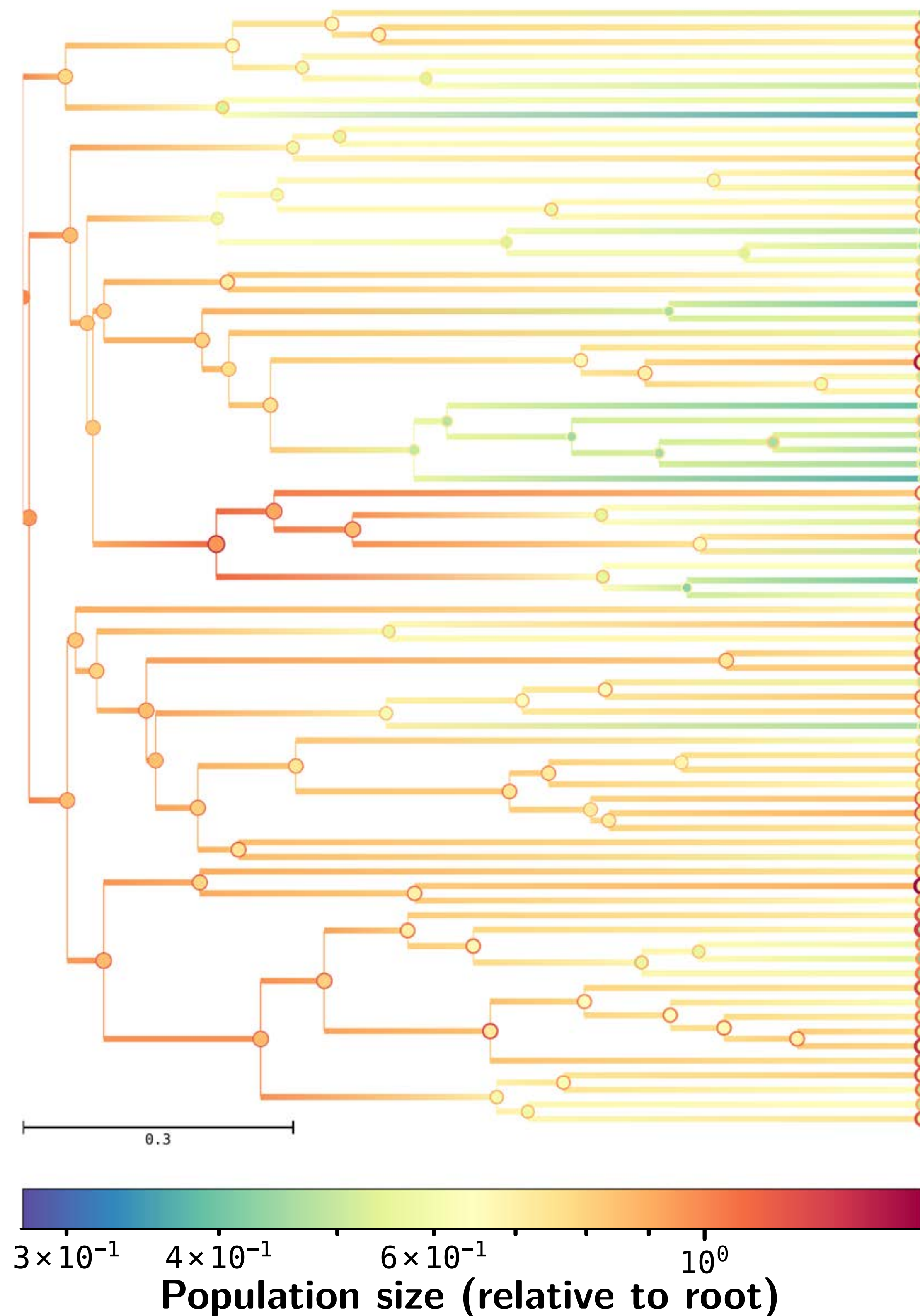
Taken into account by the simulator:

- Long term fluctuation of population size, mutation rate per time, generation time;
- Fitness landscape for each site;
- Finite population size;
- Allele hitching;
- **Short term fluctuation of population size;**

Not taken into account by the simulator:

- Species tree \neq gene tree;
- Epistasis;
- Fluctuating fitness landscape;
- Biased gene conversion;
- Selection on codon usage;

Inference with mammalian empirical data



Traits correlation

	Population size (relative)	Mutation rate (per time)	Maximum longevity (yrs)	Adult weight (g)	Female maturity (days)
Population size (relative)		0.6	-0.63	-0.55	-0.53
Mutation rate (per time)	0.6		-0.87	-0.85	-0.8
Maximum longevity (yrs)	-0.63	-0.87		0.84	0.85
Adult weight (g)	-0.55	-0.85	0.84		0.81
Female maturity (days)	-0.53	-0.8	0.85	0.81	

Correlation coefficient

- Concatenated random sample of 24 highly conserved coding sequences (>99% coverage) from OrthoMam database.
- Life-history traits extracted from AnAge database.

Tacutu (2013), Scornavacca (2019)

Take home message

- Fluctuating population size and selection can be inferred from protein coding DNA sequences using a phylogenetic approach.
- In mammals, population size correlates negatively with longevity, weight and maturity, and positively with mutation rate.
- The mechanistic mutation-selection model can be extended by taking into account polymorphism within species.
- Which mechanism could explain such a low variance of population size observed in empirical data? Epistasis, fluctuating selection, short-term fluctuation of population size...

Thanks

Nicolas Lartillot, Vincent Lanore, Philippe Veber, Nicolas Rodrigue,
Bastien Boussau, Florian Benitiere.

All lab members.

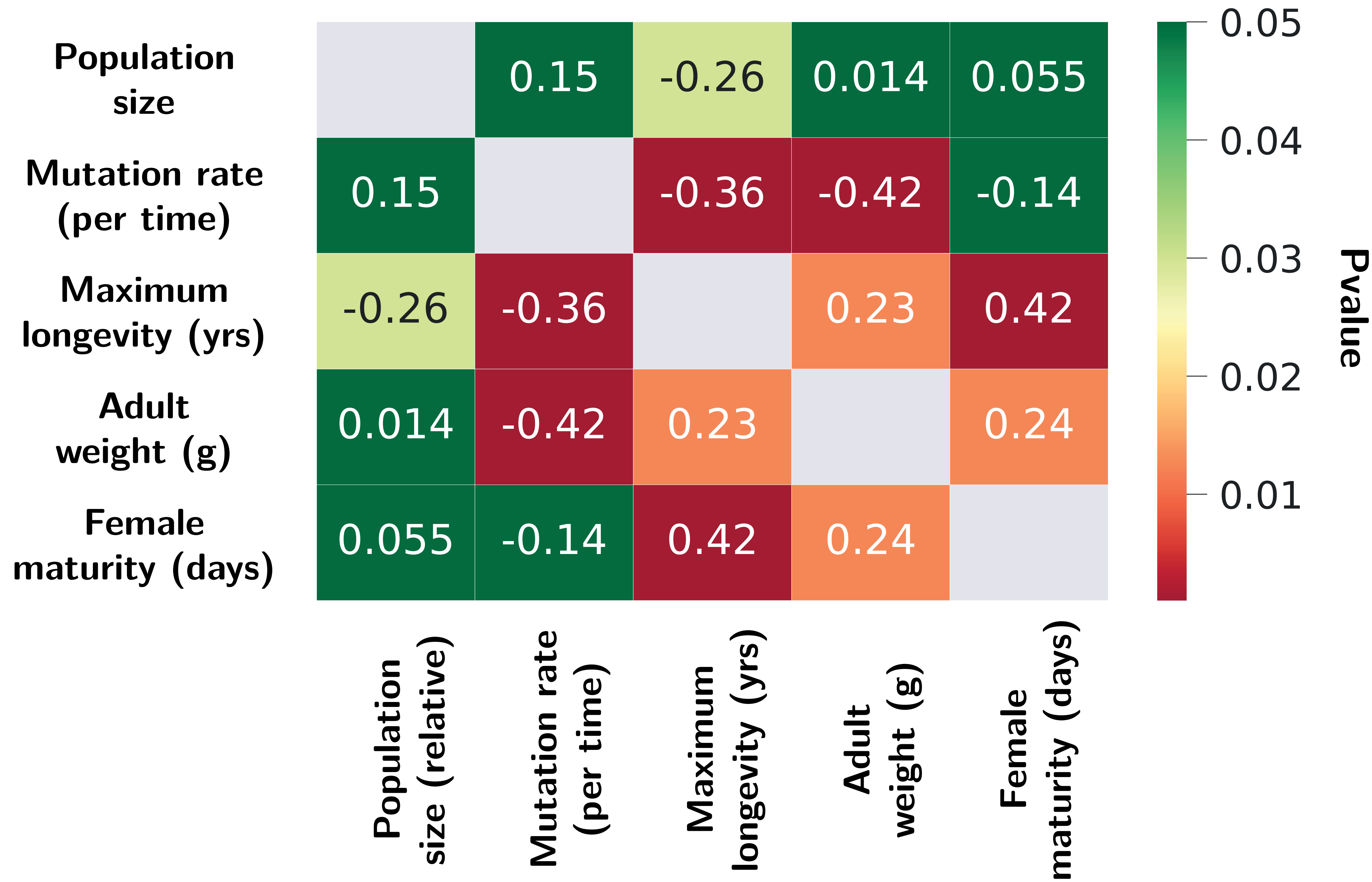


You for your attention!

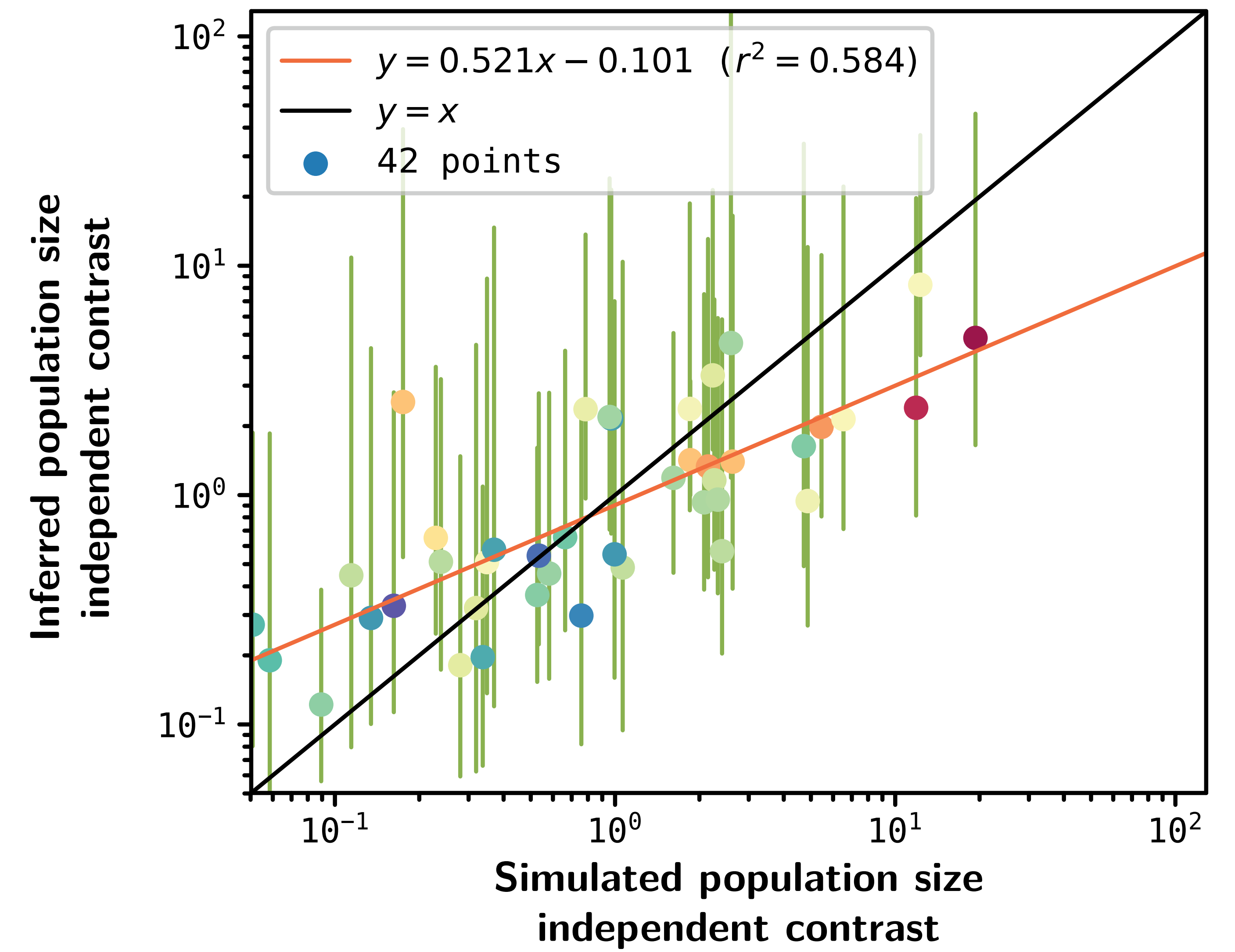
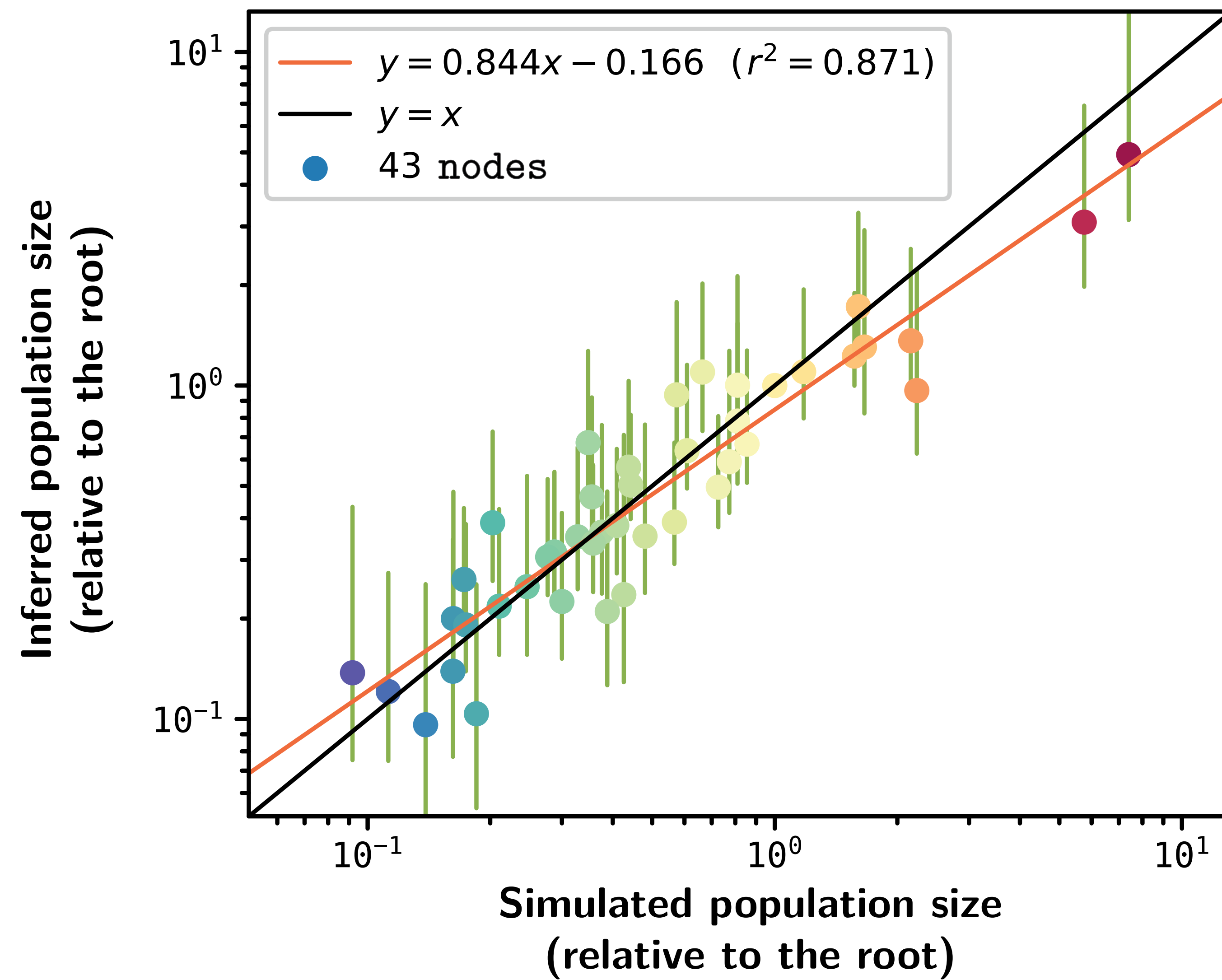
thibault.latrille@univ-lyon1.fr
@phylogenetrips

<https://github.com/ThibaultLatrille/SimuEvol>
<https://github.com/bayesiancook/bayescode>

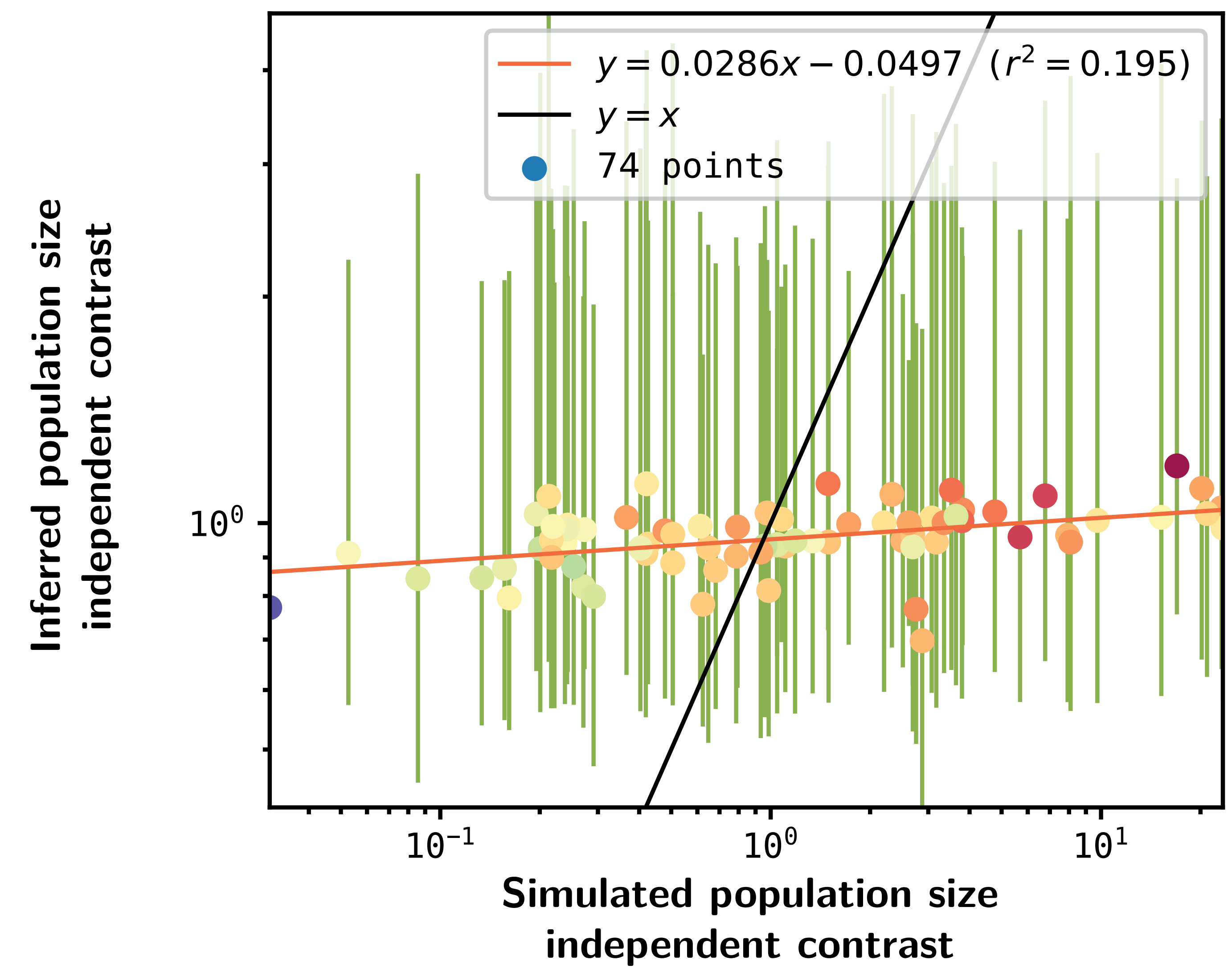
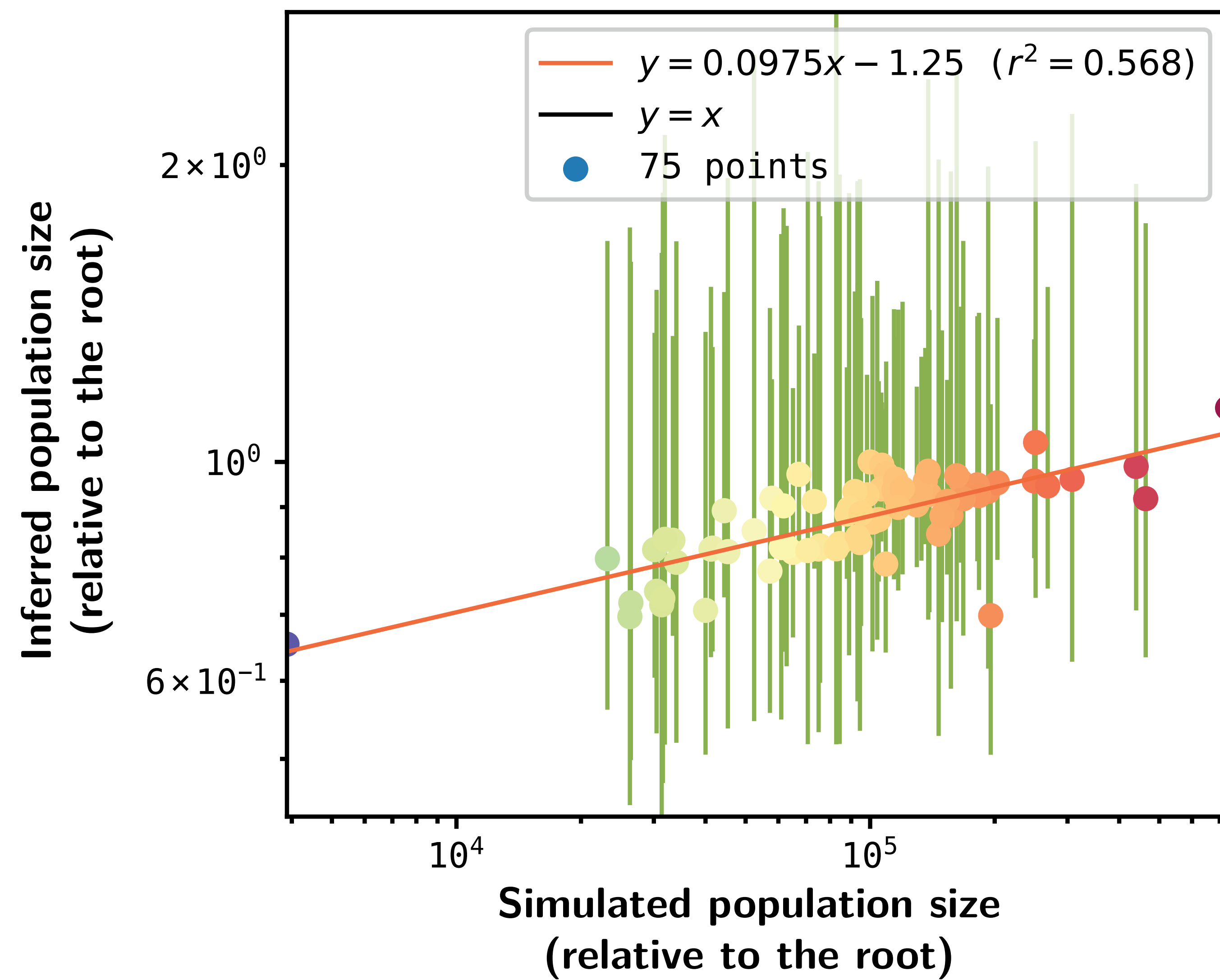
Partial correlation matrix



Comparing simulation and inference using independent contrast

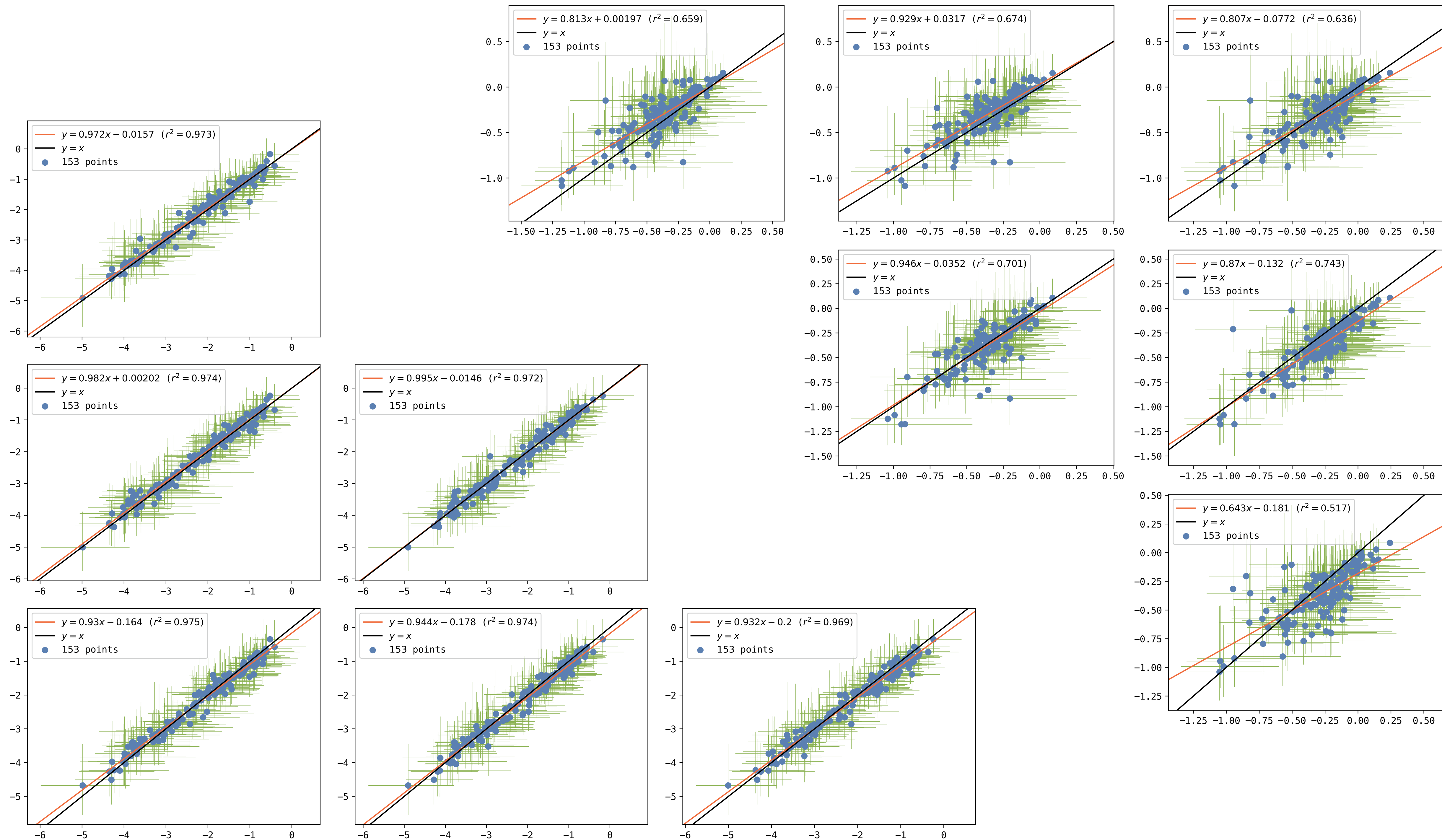


Estimating population size in the presence of epistasis



Mammalian experiment repeatability

Effective population size



Mutation rate per unit of time