# Population structure and demography

Prepared by Claire Mérot & Anna Tigano
Physalia Course 15th September 2020

# Why does population structure matter when studying adaptation?

Evolution (including adaptive evolution)

is the result of the interplay of

**Selection**

**Drift**

**Mutation**

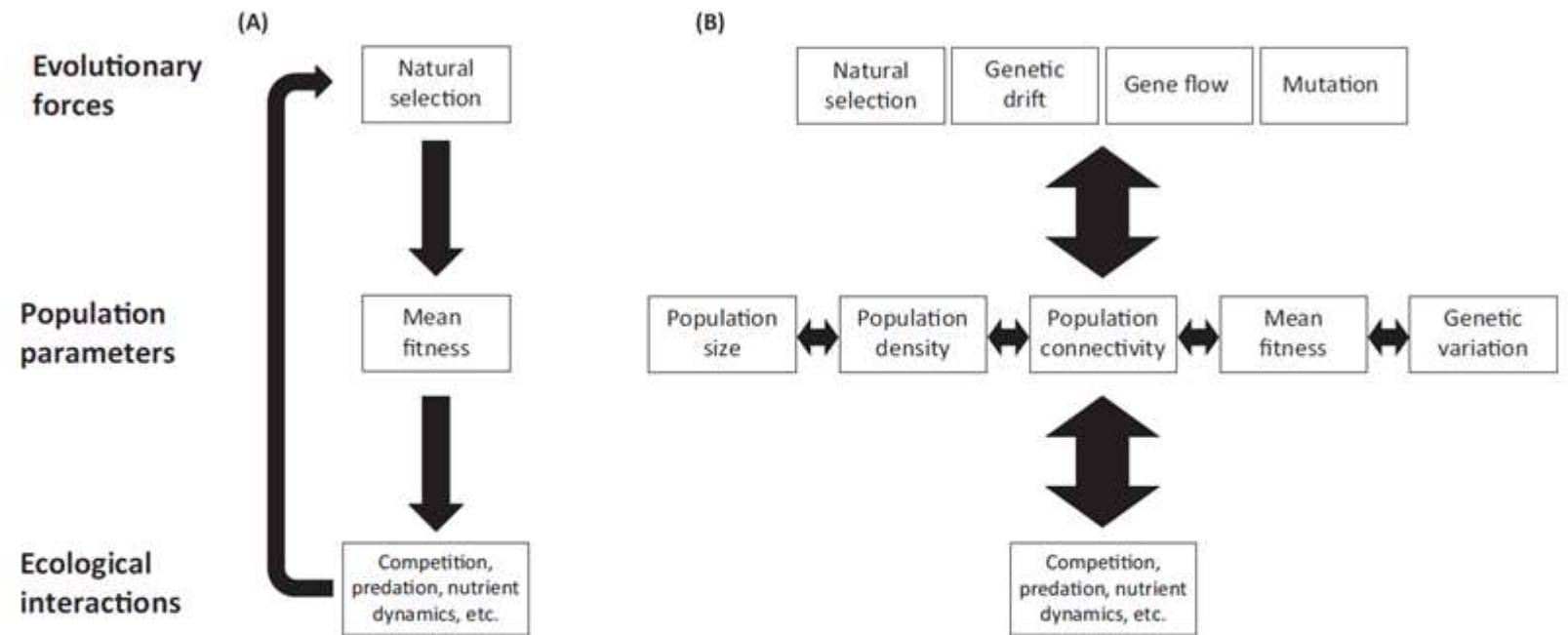**Gene flow** (migration + recombination)

# Evolutionary, demographic and ecological processes are inseparable

Population Genetics and Demography Unite Ecology and Evolution

Winsor H. Lowe,[1,*] Ryan P. Kovach,[2] and Fred W. Allendorf[1]



Trends in Ecology & Evolution

Figure 1. Evolutionary and Ecological Processes Are Inseparable. Conceptual illustration of interconnections among evolutionary forces and ecological interactions (biotic and abiotic) through population-level demographic and genetic parameters. (A) represents those interconnections emphasized in current eco-evolutionary research. (B) represents a more comprehensive model of these interconnections, including the full suite of evolutionary forces and a range of population parameters that are themselves interdependent. We build our review around population demographic parameters (size, density, connectivity), but describe key interactions with genetic parameters (mean fitness, genetic variation). We define mean fitness according to population genetics theory as the sum of the fitnesses of genotypes in a population weighted by their proportions [88], thus representing the population-level effects of local adaptation.

# Complementary objectives:

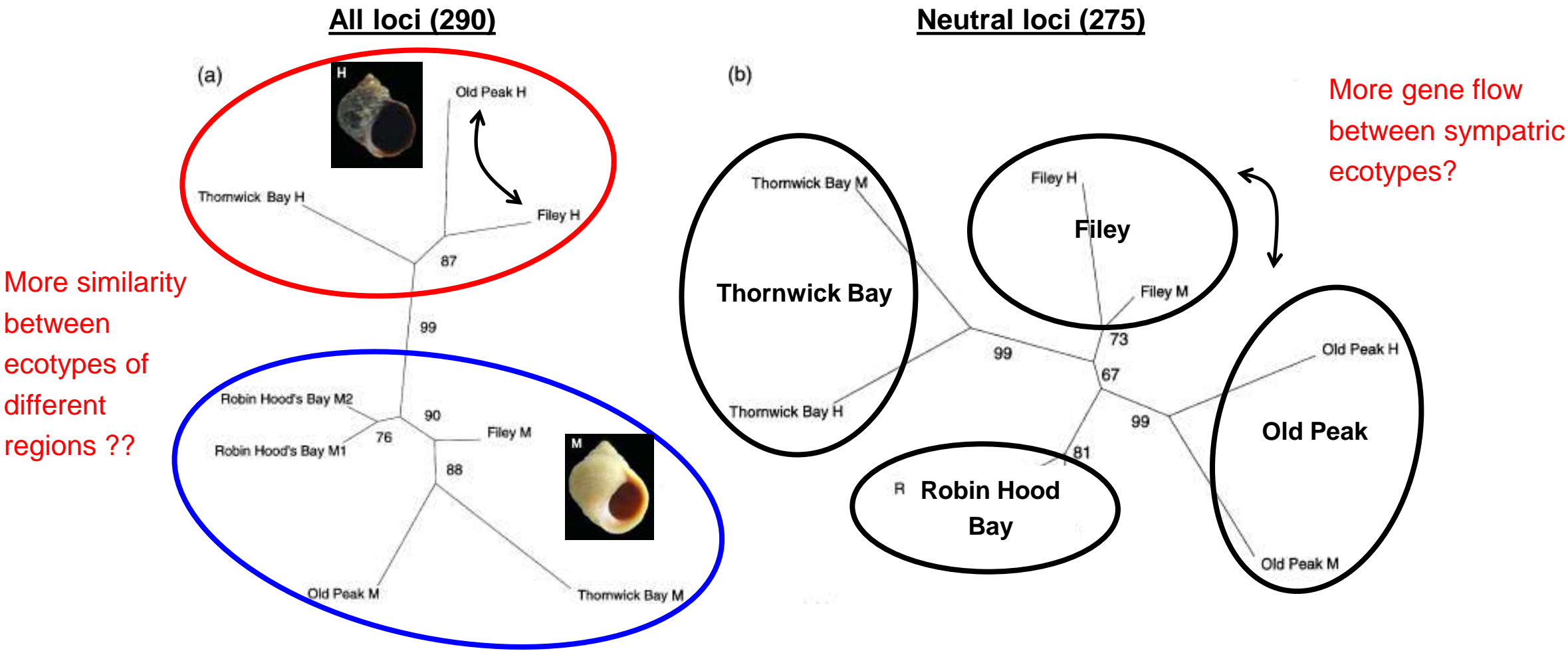|  | Study selection and adaptation | Demographic history and structure of populations |
|---|---|---|
| Actions | Focus on (putatively) adaptive loci | Focus on neutral loci |
| Use | . Study ecological/functional diversity | . Understand the past history of populations |
|  | . Understand adaptative processes under divergent or balancing selection | . Describe population connectivity |
|  |  | . Assess general genetic diversity |
|  | . Identify candidate genes |  |

# Different loci tells a different story…

. Parapatric ecotype *Littorina saxatilis*

# Different loci tells a different story…

**Exemple:**

Wilding *et al.*, 2001. *Journal of Evolutionary Biology*, 14: 611-619



**All loci (290)**

**Neutral loci (275)**

More gene flow between sympatric ecotypes?

More similarity between ecotypes of different regions ??

# Drift

= variation in allele frequency due to random processes

Drift is stronger in smaller populations and it can cause the loss or fixation of a variant due to random sampling of alleles.

Drift is the main driver of genetic population structure, and can generate a genetic footprint similar to that of selection.

# Allele surfing

Populations on the leading edge of the expansion are small, and individuals from those populations contribute disproportionately to the propagating wave of expansion.

$\Rightarrow$ Rapid drift of some alleles at the expanding edge and high differentiation in allele frequencies over the landscape for some loci, even in the absence of selection

# Spatial autocorrelation

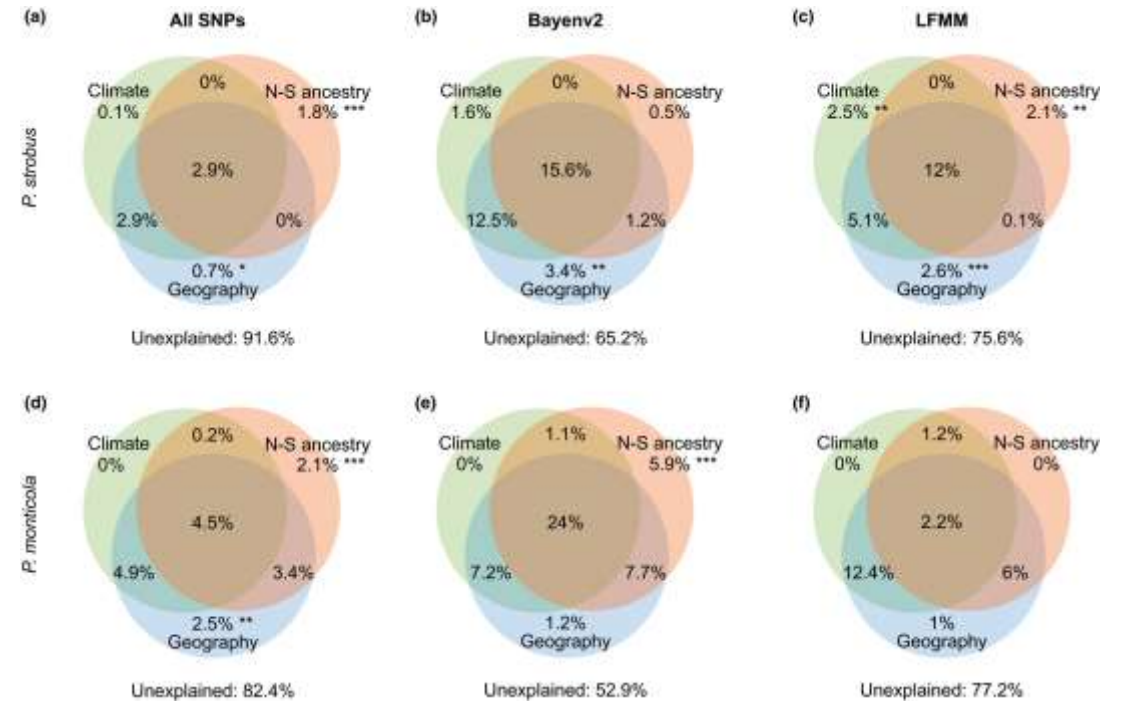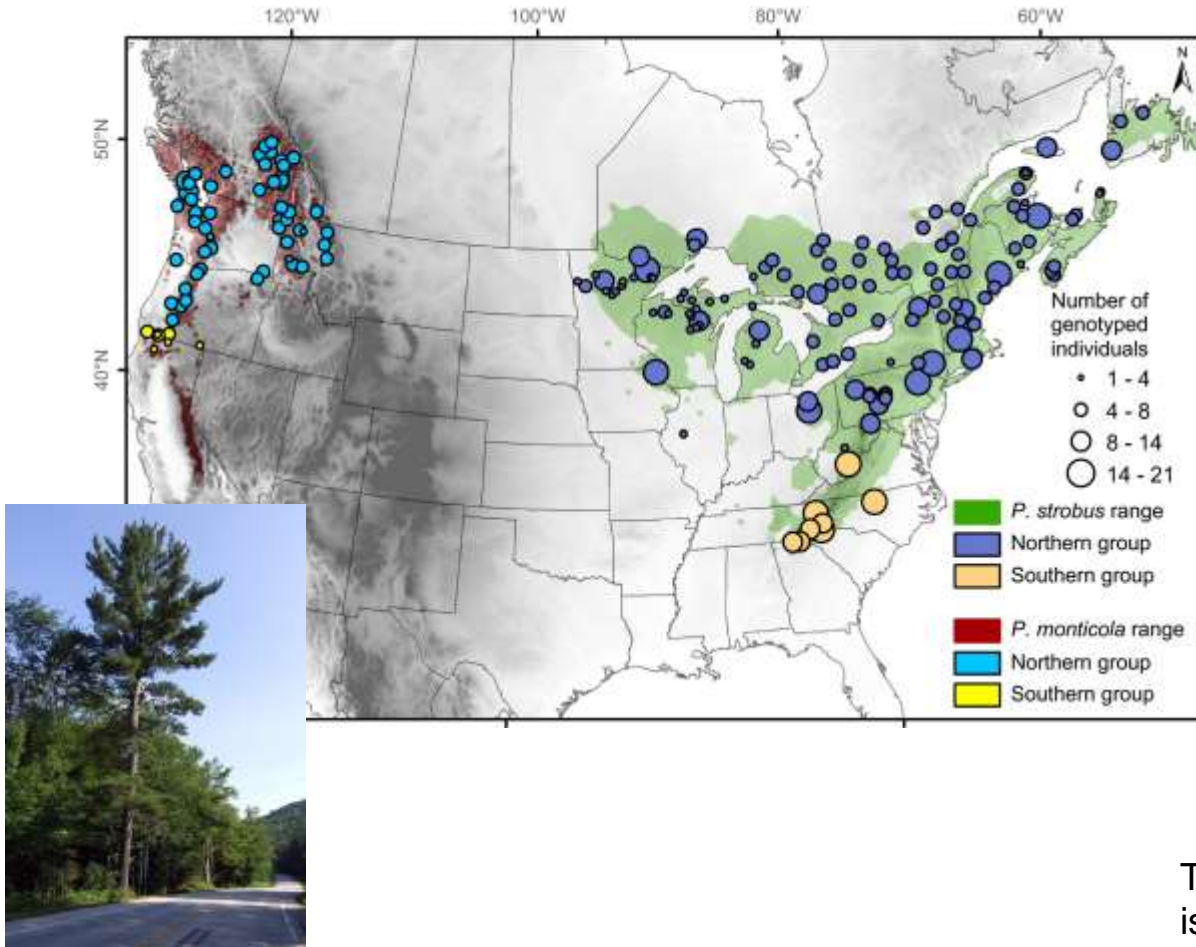Correlation between environmental variation & geographic distances
(e.g. climatic clines!)

⇒What is adaptation? What is drift?

+ Residuals of past range expansion out of glacial refugees…

⇒What is adaptation? What is the results of past history?

*Nearby locations are not statistically independent, strong correlations between neutral alleles and environmental variables are more likely to occur by chance than expected with some null models*

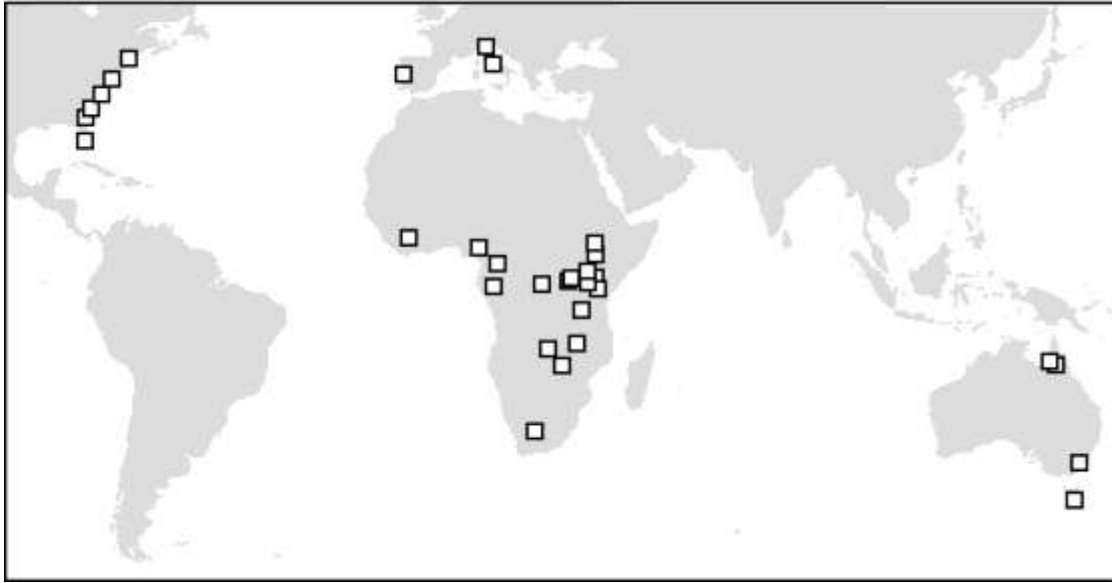# Isolation-by-distance or adaptation along a gradient … or both?



The challenge of separating signatures of local adaptation from those of isolation by distance and colonization history: The case of two white pines Nadeau et al, 2016 https://doi.org/10.1002/ece3.2550

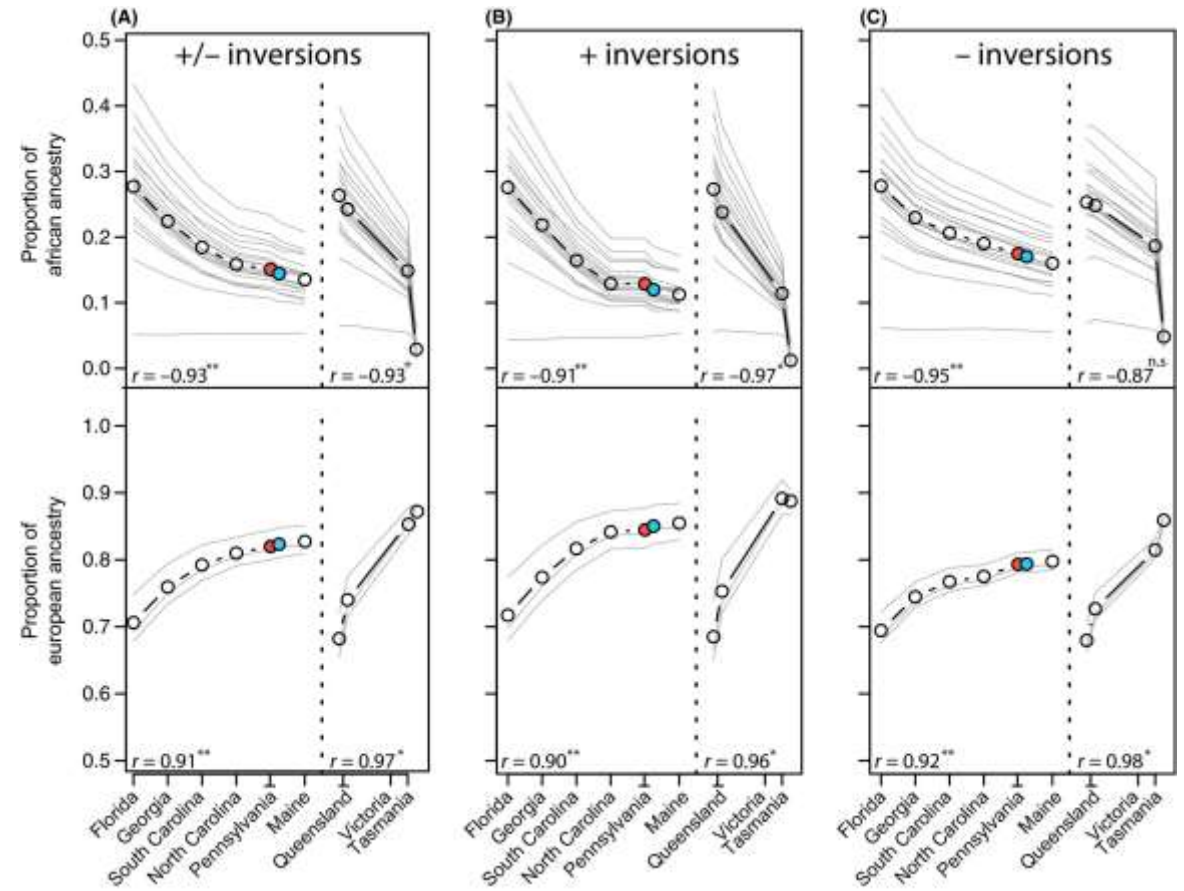# Contact between different lineages / hybridization

Signature of selection or of local adaptation are best detected in a context of (high) gene flow.

Any substructure (lineages, species, secondary contact, admixed populations) should be taken into account.

# Clinal variation or secondary contact...
# Or both?



Bergland et al, 2015 MolEcol
https://doi-org /10.1111/mec.13455

# How to characterise population structure?

Unsupervised methods:

- PCA


Semi-supervised methods (K = number of expected clusters)

- Bayesian clustering


Supervised methods (with location information for instance)

- DAPC

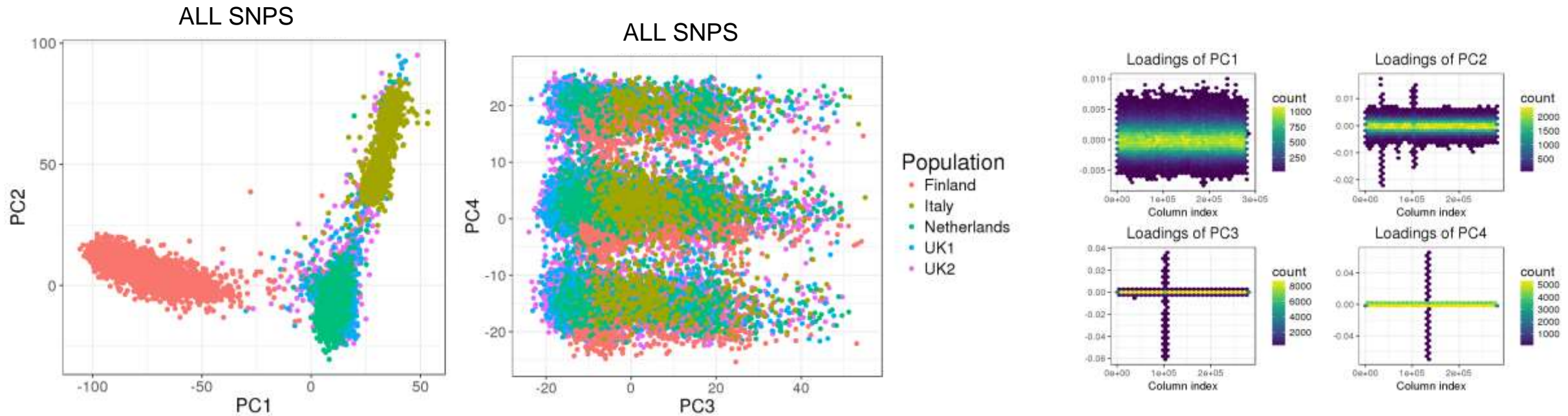- Fst between pairs of populations

# Principal Component Analysis (PCA)

- A common statistical tool that reduces matrix complexity by identifying the eigenvectors and ordering them

- The top PCs reflect axis of genetic variation along which individuals with same ancestry, or exchanging genetic material, are more similar to each other.

- Caution: can be strongly driven by few loci in linkage disequilibrium…

- For population structure purpose:
-> compare Pca on all SNPs vs. Pca on LD-pruned SNPs
-> look at loadings of the PCs: which fraction of the genome explains PC1? Explains PC2? Etc..

- There is lots of genetic variance, it can be relatively expected that even PC1 explains less than 1% of variance. (but it can also capture 20-50%... Depends on the dataset!
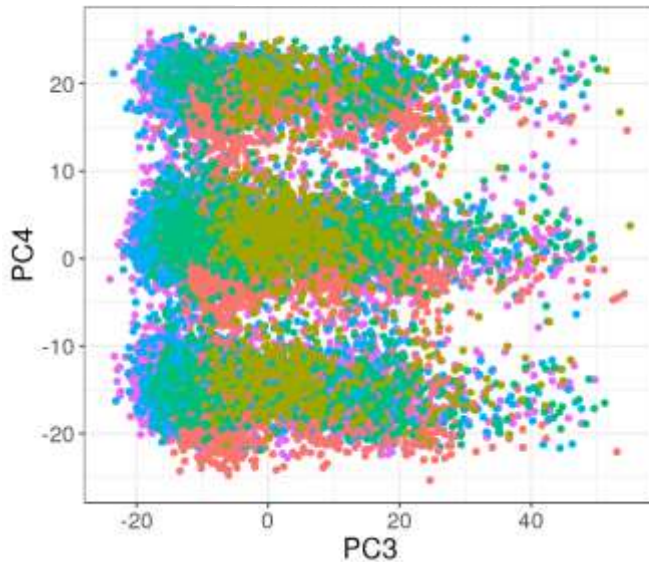
# Principal Component Analysis (PCA)

Each individual is a point with coordinates along all PCs
Each genetic marker contribute to all PCs with a different strength (loadings)



Packages *bigstatsr, bigsnpr* to remove short-range and long-range LD.
Nice tutorial about PCA for pop genomics!
Florian Privé
https://privefl.github.io/bigsnpr/articles/how-to-PCA.html
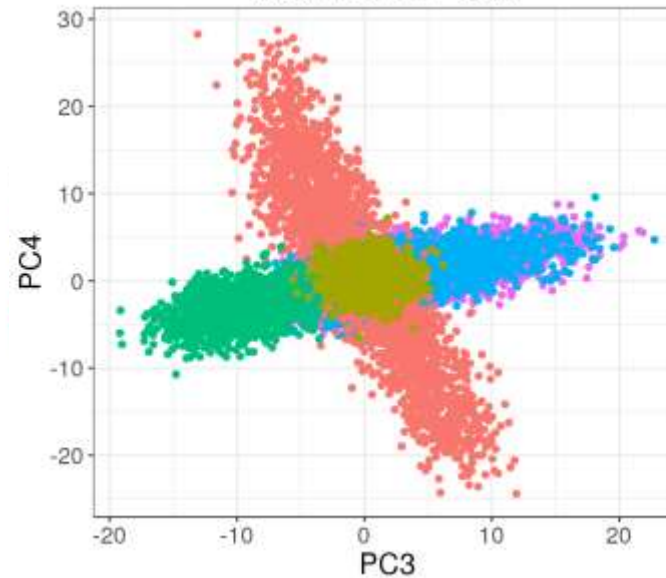
# Principal Component Analysis (PCA)

Each individual is a point with coordinates along all PCs
Each genetic marker contribute to all PCs with a different strength (loadings)
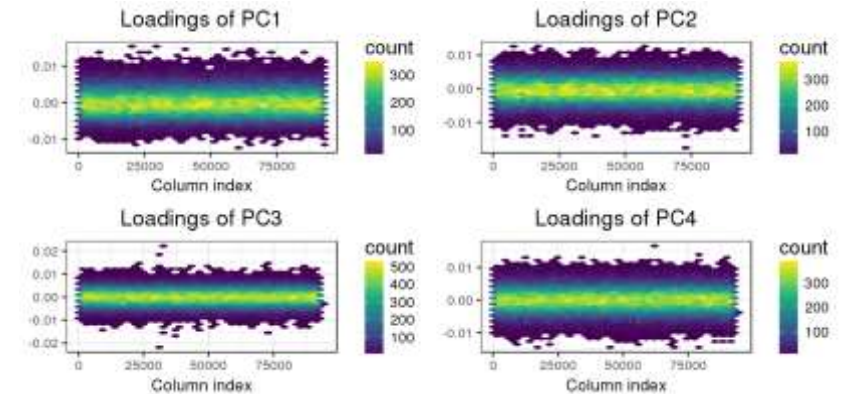


ALL SNPS



LD-pruned SNPs

Packages *bigstatsr, bigsnpr* to remove short-range and long-range LD.
Nice tutorial about PCA for pop genomics!
Florian Privé
https://privefl.github.io/bigsnpr/articles/how-to-PCA.html

# Bayesian clustering (STRUCTURE, etc..)

- Aim to sort individuals into K clusters so as to minimize departures from Hardy-Weinberg equilibrium and linkage equilibrium

- Caution: can be strongly driven by few loci in linkage disequilibrium…

- For population structure purpose:
-> compare results on all SNPs vs. results on LD-pruned SNPs
-> explore many values of K

-  Admixture or FastSTRUCTURE replace STRUCTURE for genome-wide data
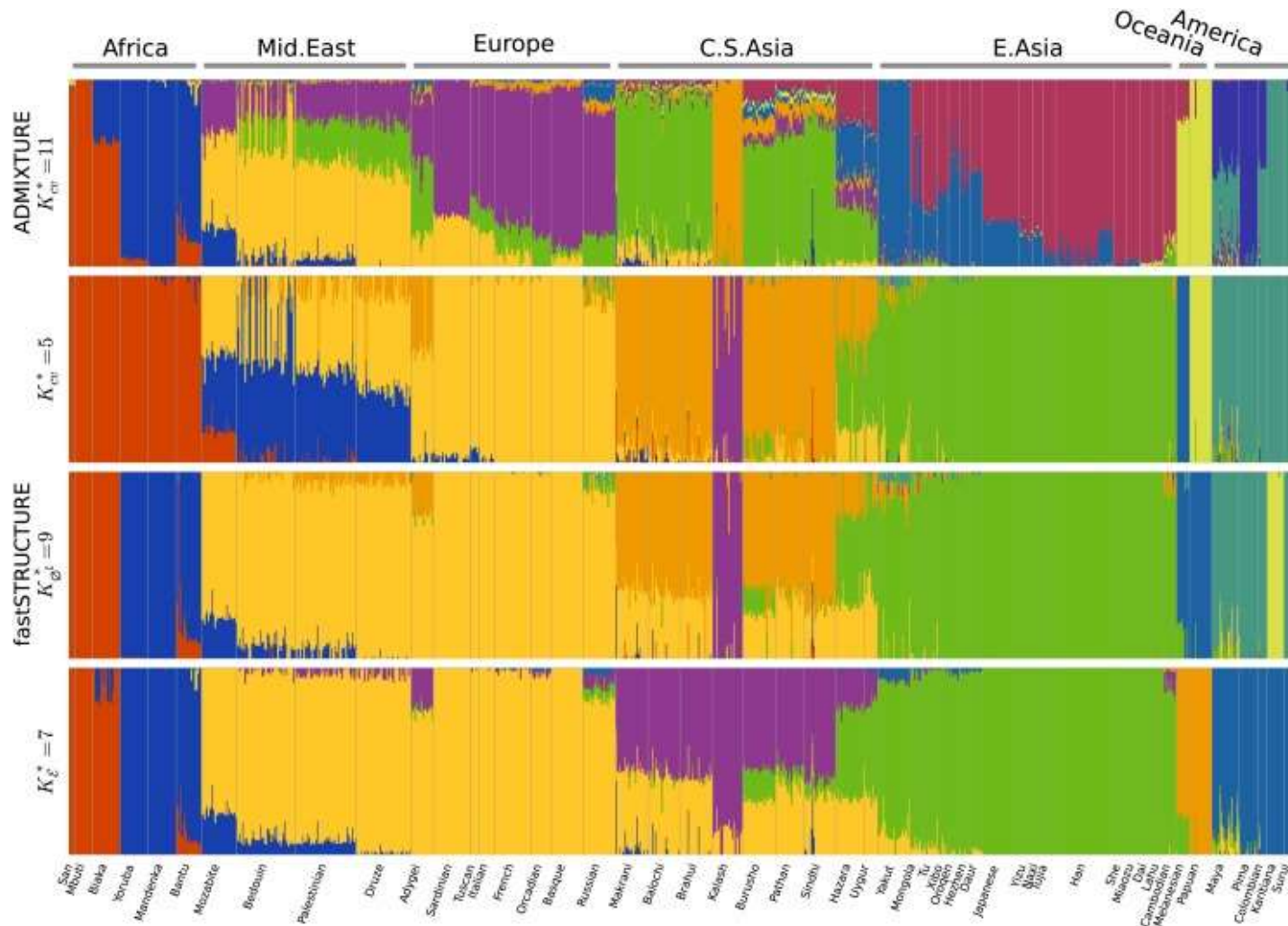
-  Evaluate the fit of the model

**Evaluation of model fit of inferred admixture proportions**
*Genís Garcia-Erill  Anders Albrechtsen*
*MER 2020*
 **https://doi.org/10.1111/1755-0998.13171**

# Bayesian clustering (STRUCTURE, etc..)



Each individual is a thin vertical line that is partitioned into *K* colored segments according to its membership coefficients in *K* clusters.

**fastSTRUCTURE: variational inference of population structure in large SNP data sets 2014 Genetics**
Anil Raj[1], Matthew Stephens[2], Jonathan K Pritchard[3]
10.1534/genetics.114.164350

# The advantage of unsupervised/semi-supervised methods:
# => Other surprises!!
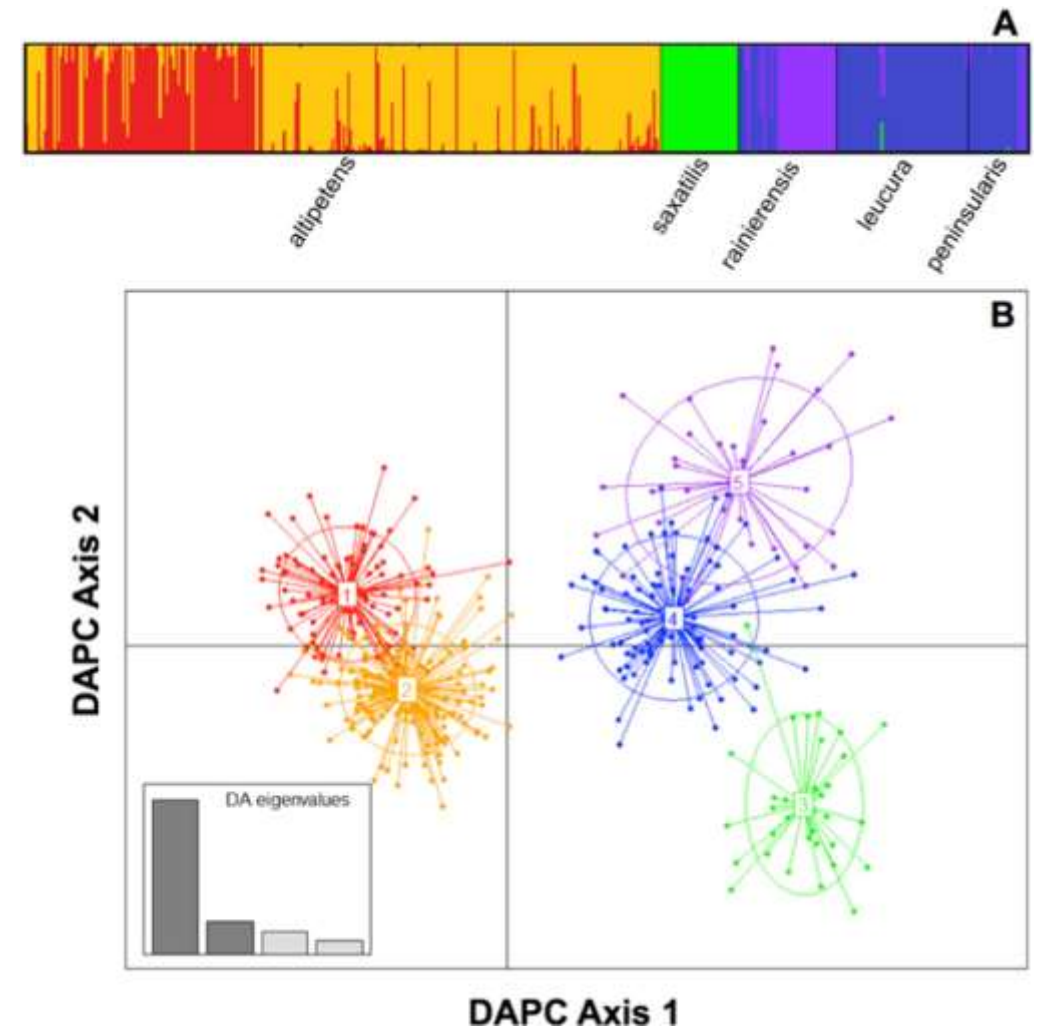
Species lineage substructure

Hybridisation

Chromosomal rearrangements….

# DAPC (discriminant PCA)

- A mix of a discriminant analysis and a PCA

- It will try very hard to find axis of variation that discriminate the groups given *a priori*

- *WARNING:* A dangerous analysis when we have much more markers (SNPs) than groups (populations)… Be well aware of not over-fitting and not-overinterpreting the output.

Miller, J.M., Cullingham, C.I. & Peery, R.M. The influence of a priori grouping on inference of genetic clusters: simulation study and literature review of the DAPC method. *Heredity* (2020). https://doi.org/10.1038/s41437-020-0348-2
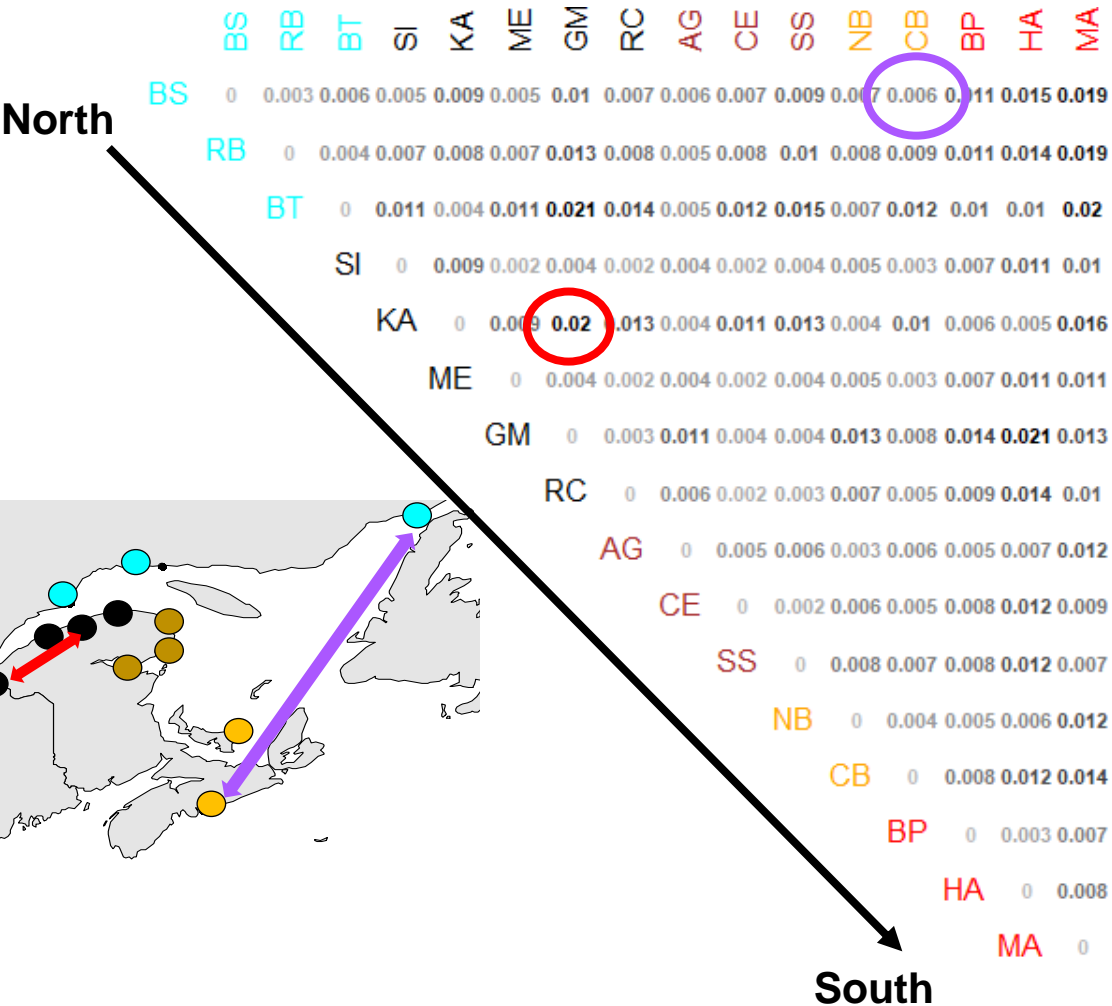


Languin et al, 2018,
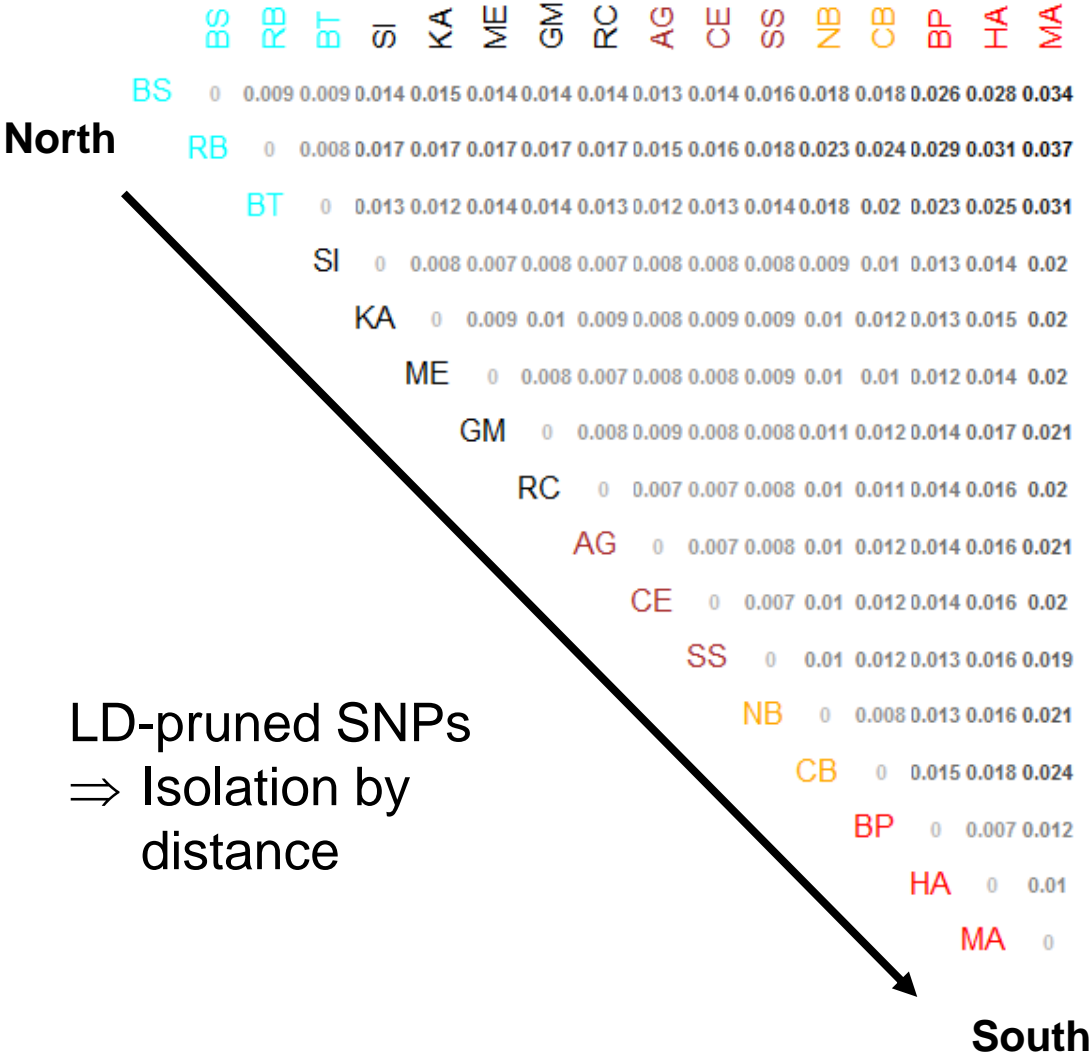Conservation genomics

# Pairwise Fst

- Fst can only been computed between two groups.

- When sampling several populations, we will be interested in Fst between all pairs of individuals.

- A measure of genetic distance between all populations (does it correlate with ecological distances? Geographic distances? Etc?)

- Again, likely better on Ld-pruned SNPs to infer neutral structure…

- Absolute values are informative… (0,000x -> high gene flow, don't bother too much about looking for structure. 0,01-0,1 -> consider carefully structure… Higher: do you really have one species?)

# Pairwise Fst



ALL SNPs

LD-pruned SNPs

LD-pruned SNPs ⇒ Isolation by distance

# Case study

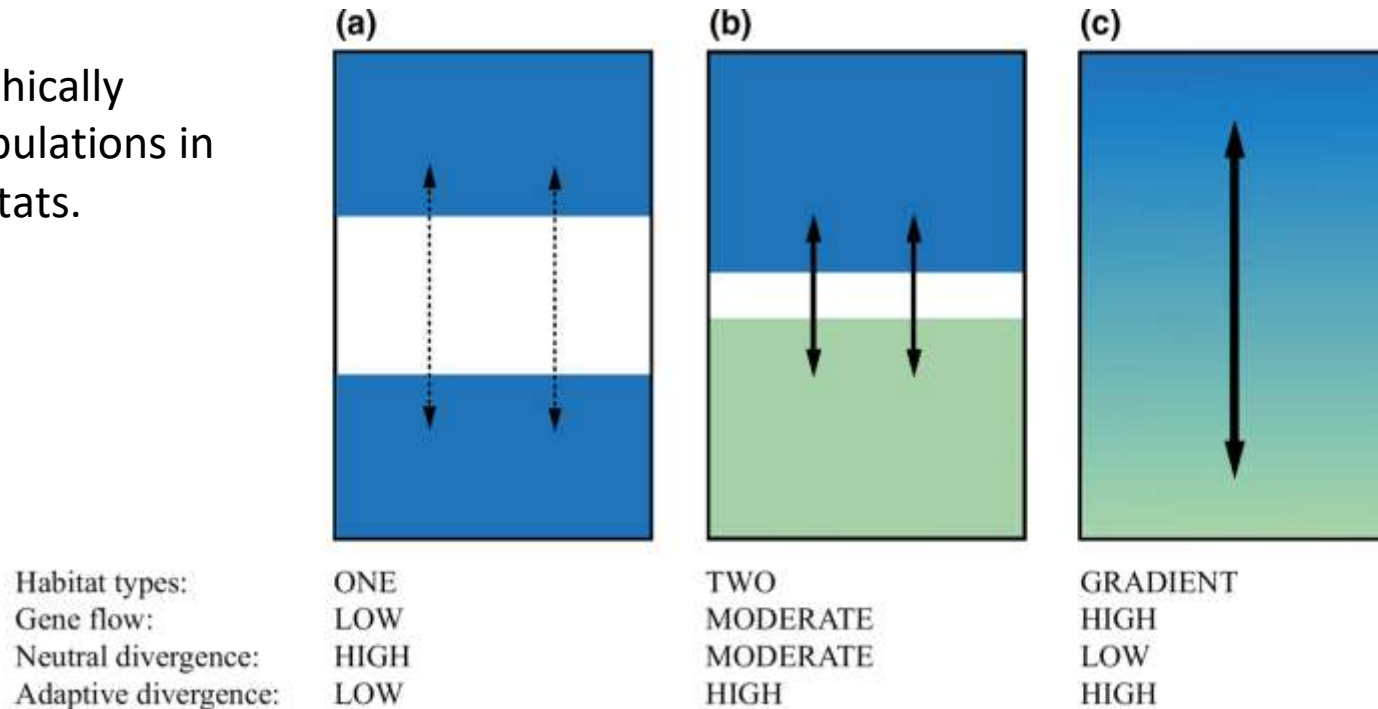## Genome-wide signals of drift and local adaptation during rapid lineage divergence in a songbird

Guillermo Friis[1] | Guillermo Fandos[2] | Amanda J. Zellmer[3] |
John E. McCormack[3,4] | Brant C. Faircloth[5] | Borja Milá[1]
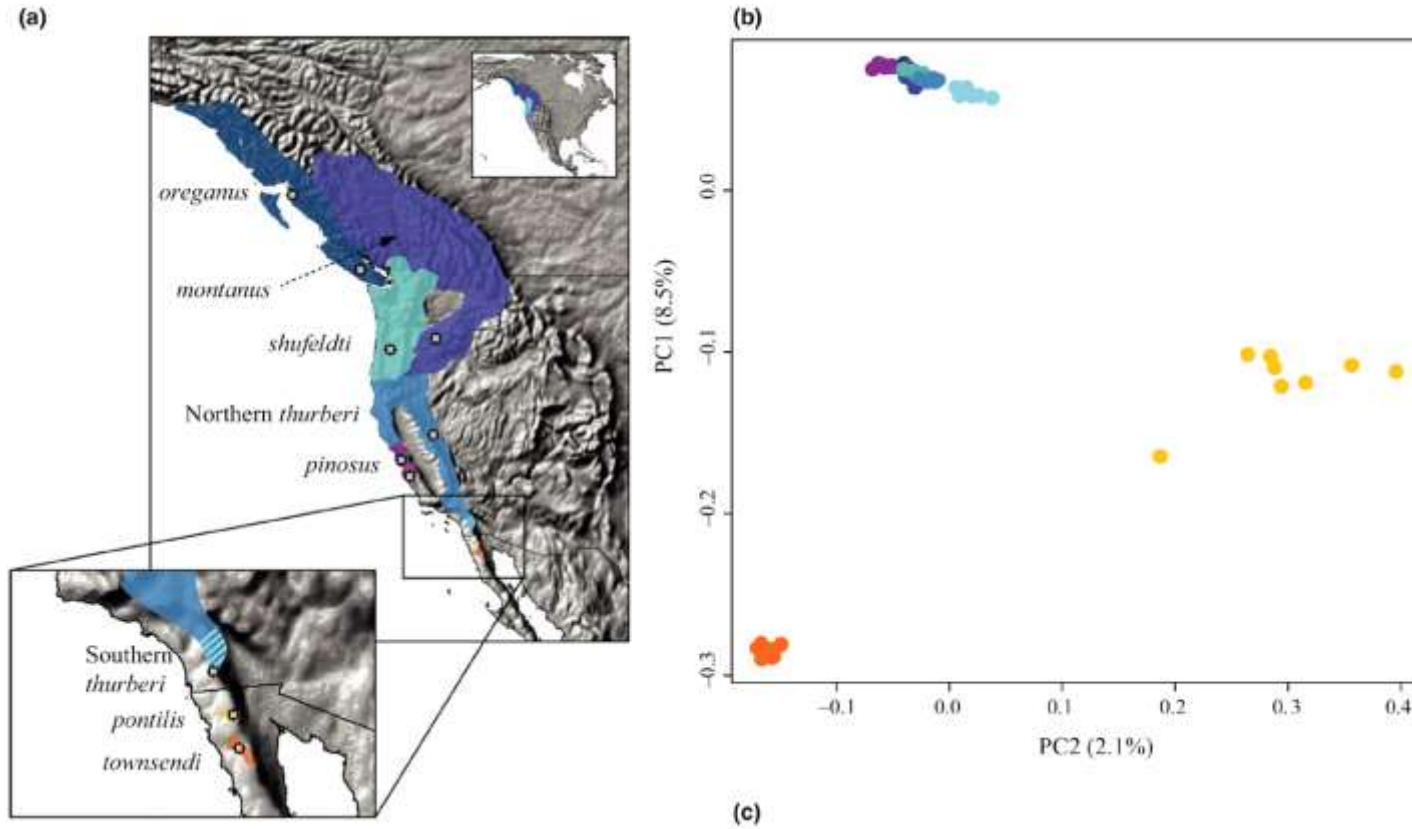
# Set expectations

(b) Parapatric populations in ecologically divergent habitats.

(a) Geographically isolated populations in similar habitats.
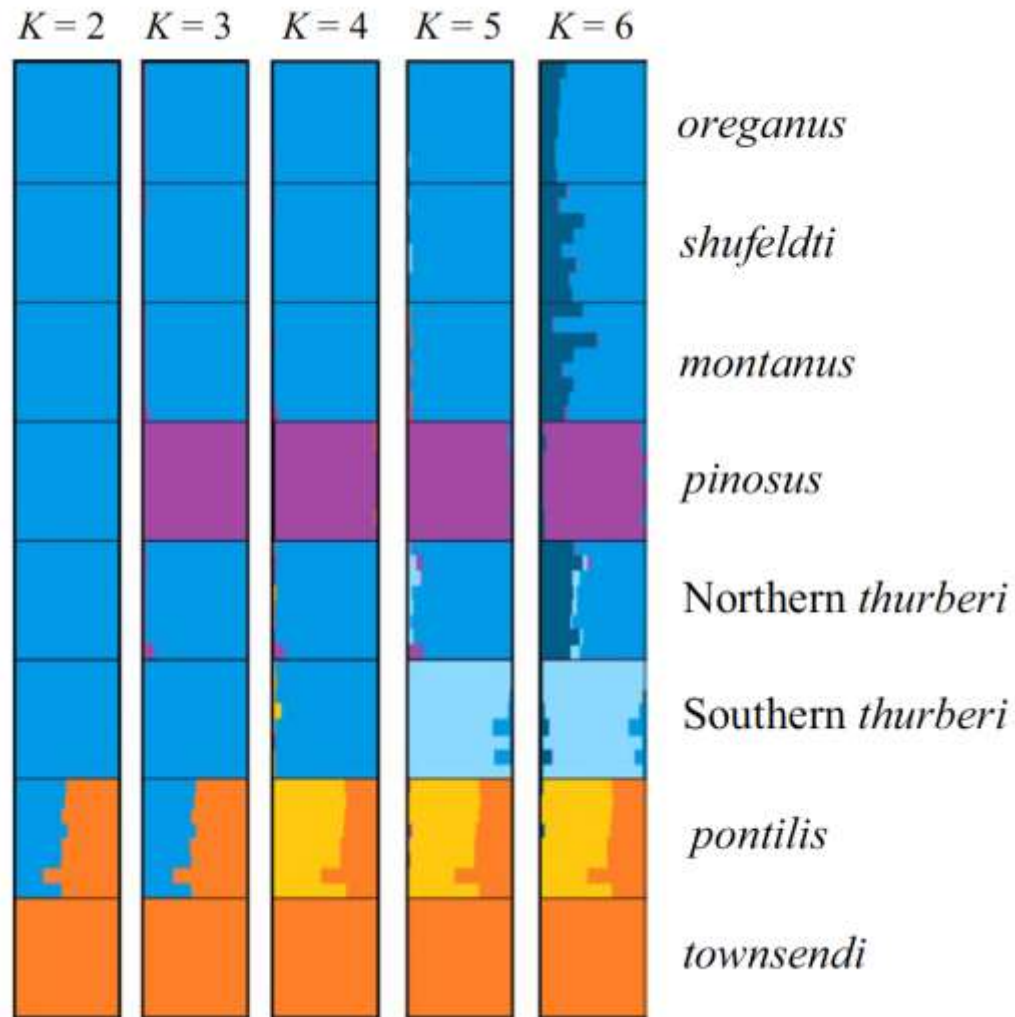
(c) Population continuum across a selective gradient



| Habitat types: | ONE | TWO | GRADIENT |
|---|---|---|---|
| Gene flow: | LOW | MODERATE | HIGH |
| Neutral divergence: | HIGH | MODERATE | LOW |
| Adaptive divergence: | LOW | HIGH | HIGH |

# PCA for data exploration



Prior-free data exploration
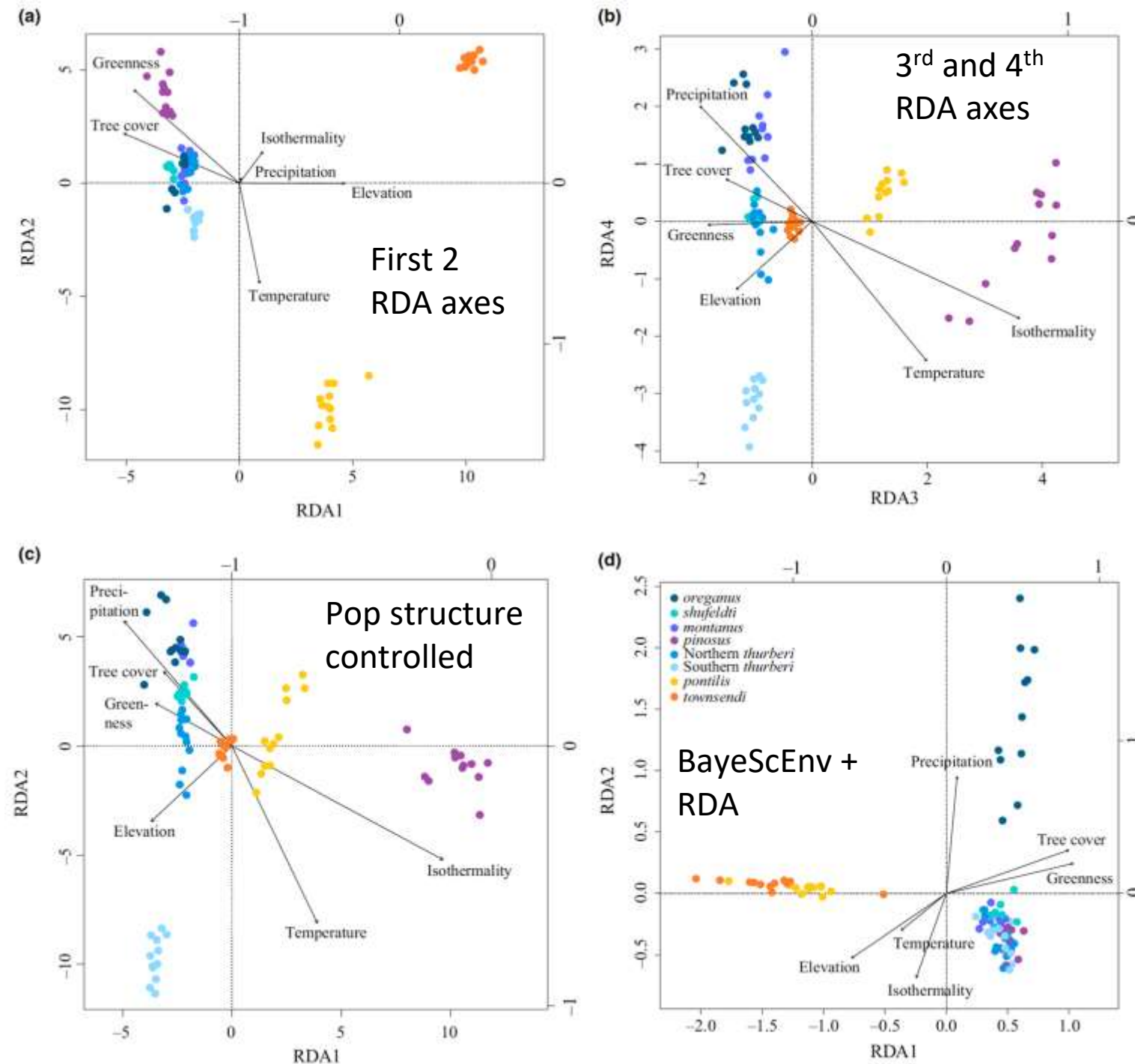
# STRUCTURE for population structure



Always examine multiple K values as more than one K could be biologically informative

(STRUCTURE doesn't deal well with hierarchical population structure)

# GEA
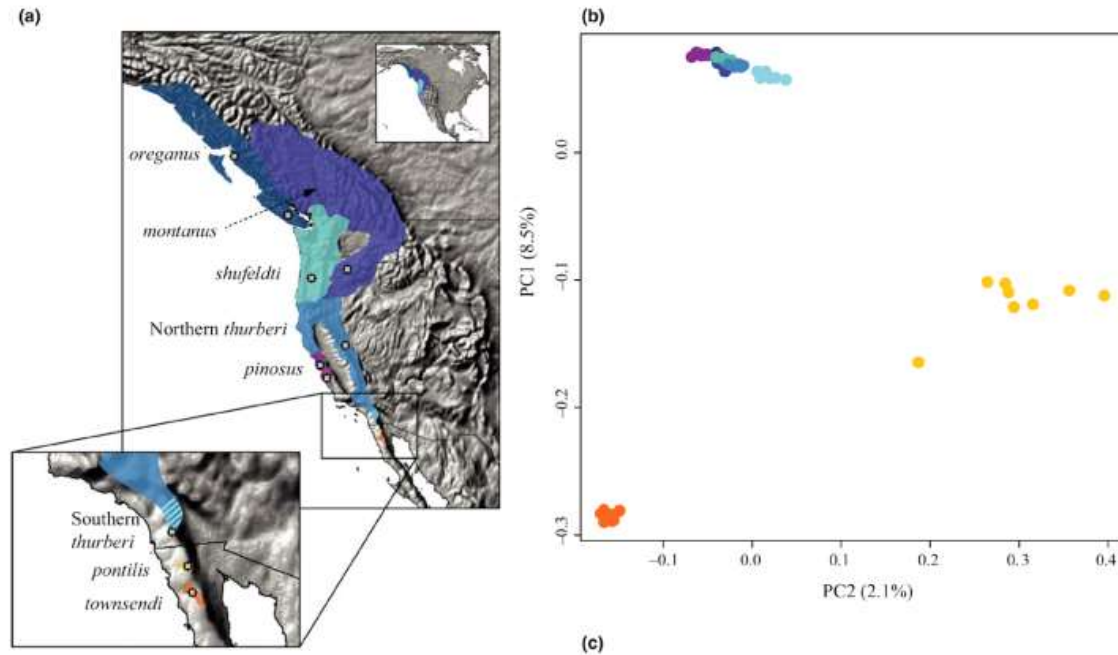
Genotype-environment associations with multiple approaches

*Lecture tomorrow, we come back on that*

(a) First 2 RDA axes

(b) 3rd and 4th RDA axes

(c) Pop structure controlled

(d) BayeScEnv + RDA

Legend (panel d):
- *oreganus*
- *shufeldti*
- *montanus*
- *pinosus*
- Northern *thurberi*
- Southern *thurberi*
- *pontilis*
- *townsendi*

# Partition of genetic variation

- Environmental variables (controlling for population structure) 1.17%

- Environmental variable + pop structure 7.41%

- 92.59%?
  - Loci under balancing selection
  - Other selective pressures
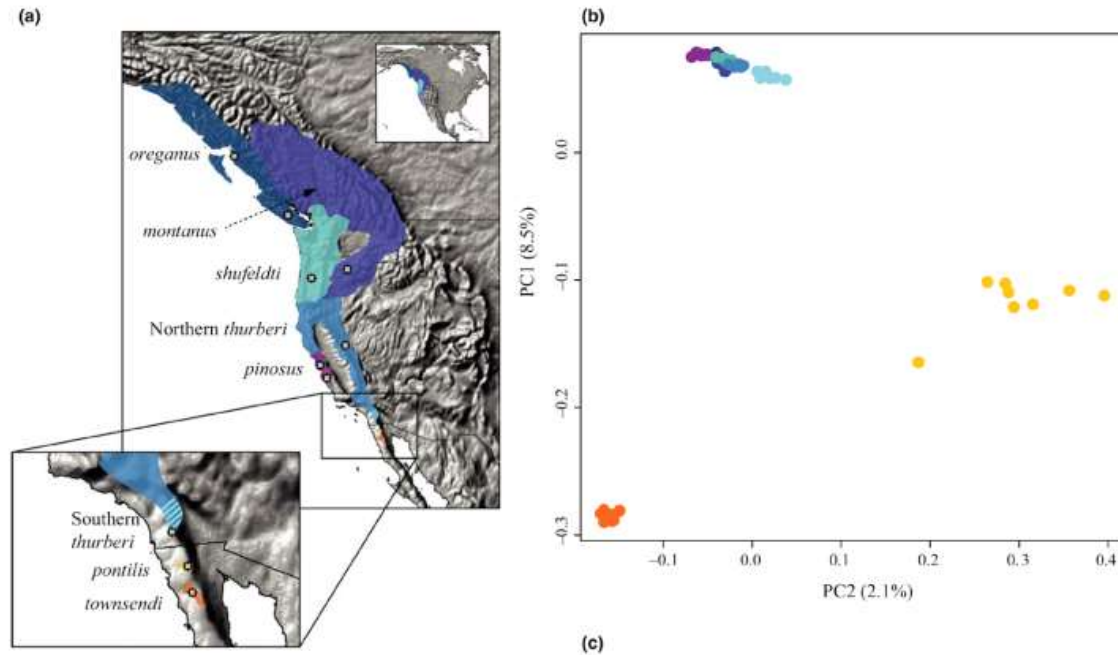  - Shared neutral variation due to relatedness and/or gene flow

# Environment + geography + demography



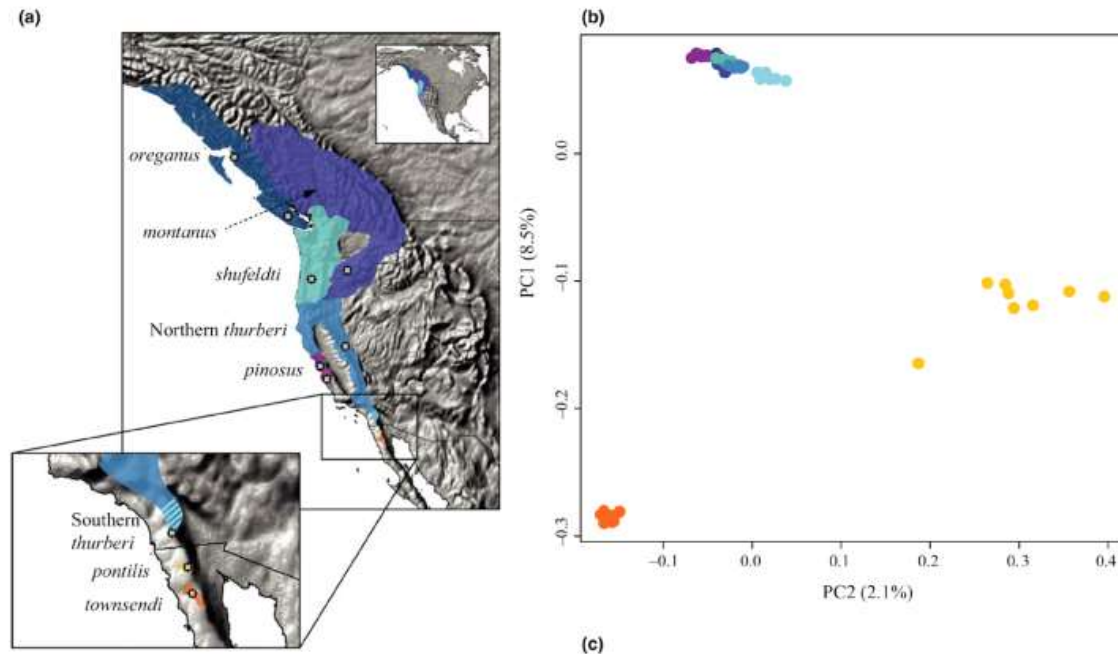Strong population structure and weak GEA

Isolation-by-resistance and drift

# Environment + geography + demography



Weak pop structure and stronger differentiation on GEA

Isolation-by-adaptation?

# Environment + geography + demography



No pop structure + environmental associations

$\Rightarrow$ Ongoing gene flow and local adaptation

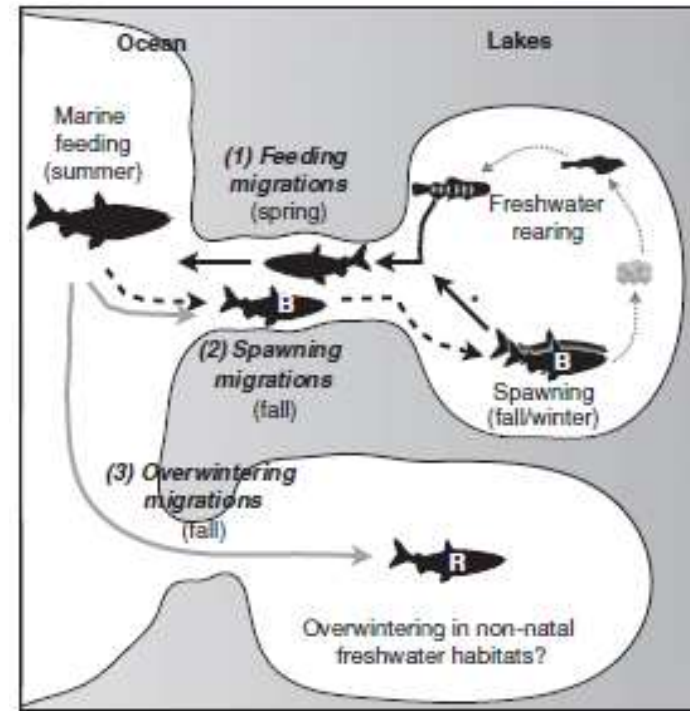# And help from experimental work /knowledge of natural history

Capture-Mark-Recapture

-> population size, dynamic and movement
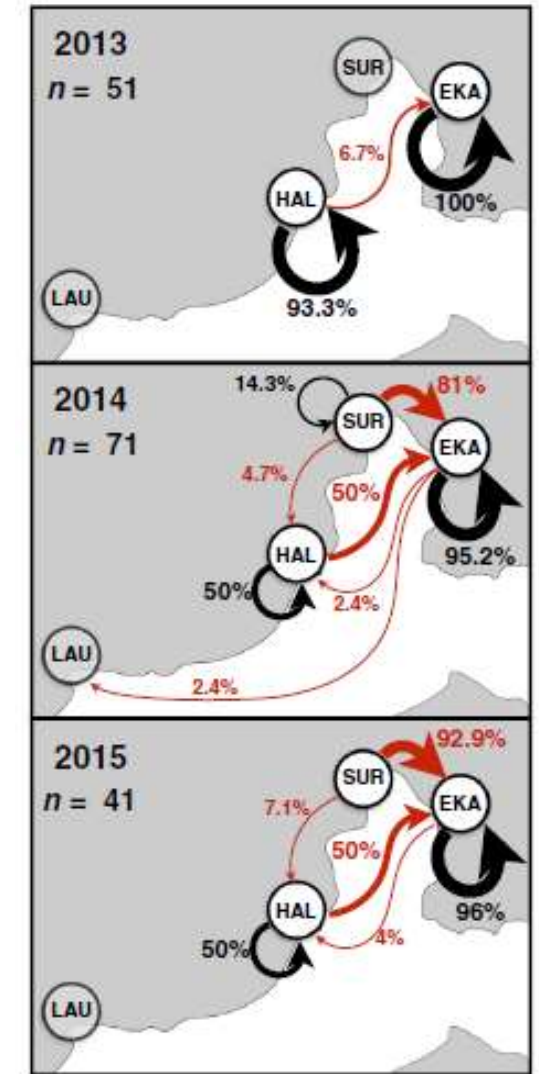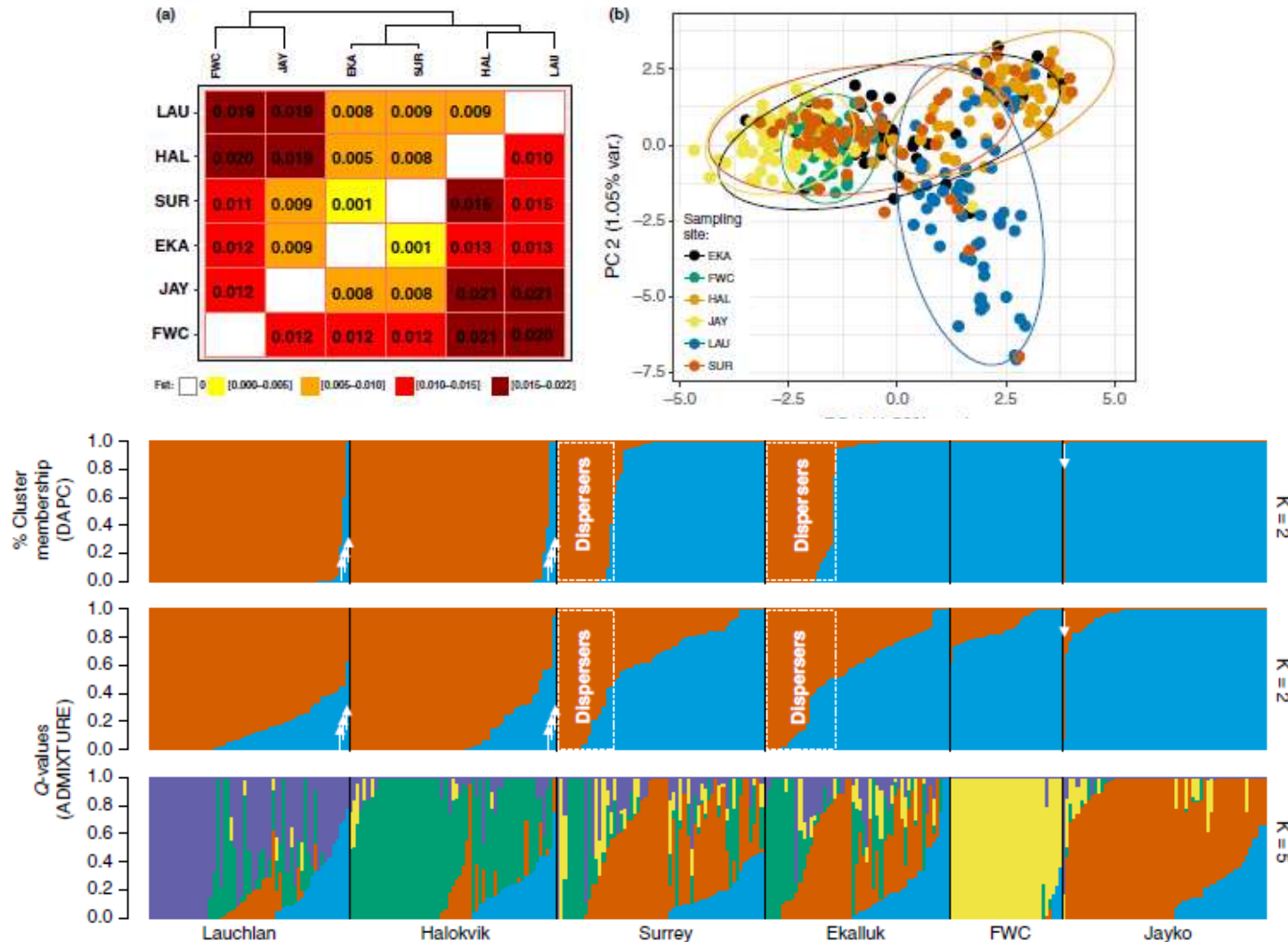
Spatial ecology

-> tracking, etc..

e.g : telemetry & genomics



**Moore, J.-S**., L.N. Harris, J. Le Luyer, B.J.G. Sutherland, Q. Rougemont, R.F. Tallman, A.T. Fisk & L. Bernatchez (2017) Genomics and telemetry suggest a role for migration harshness in determining overwintering habitat choice, but not gene flow, in anadromous Arctic Char. *Molecular Ecology,* 26(24): 6784-6800

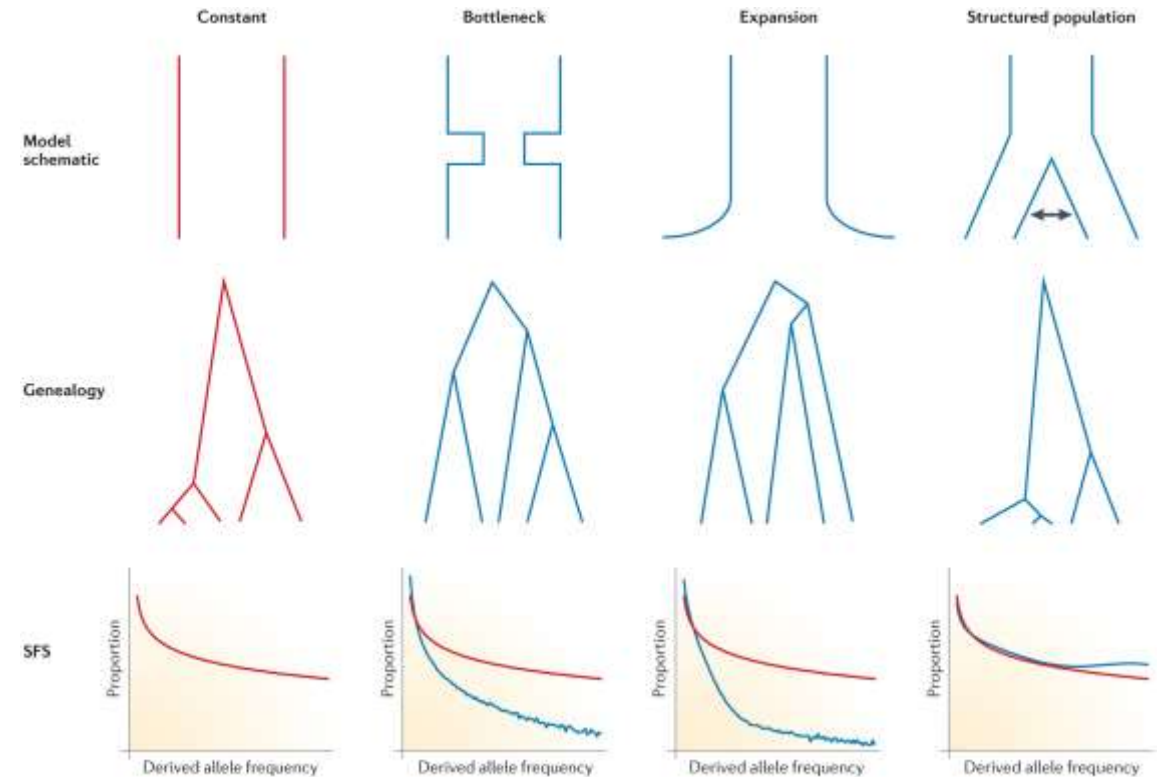# And help from experimental work /knowledge of natural history



Moore, J.-S.,. *Molecular Ecology,* 26(24): 6784-6800

# Beyond present structure…
# How to know population history and demography?

Models:

- to understand population history, bottleneck, gene flow…

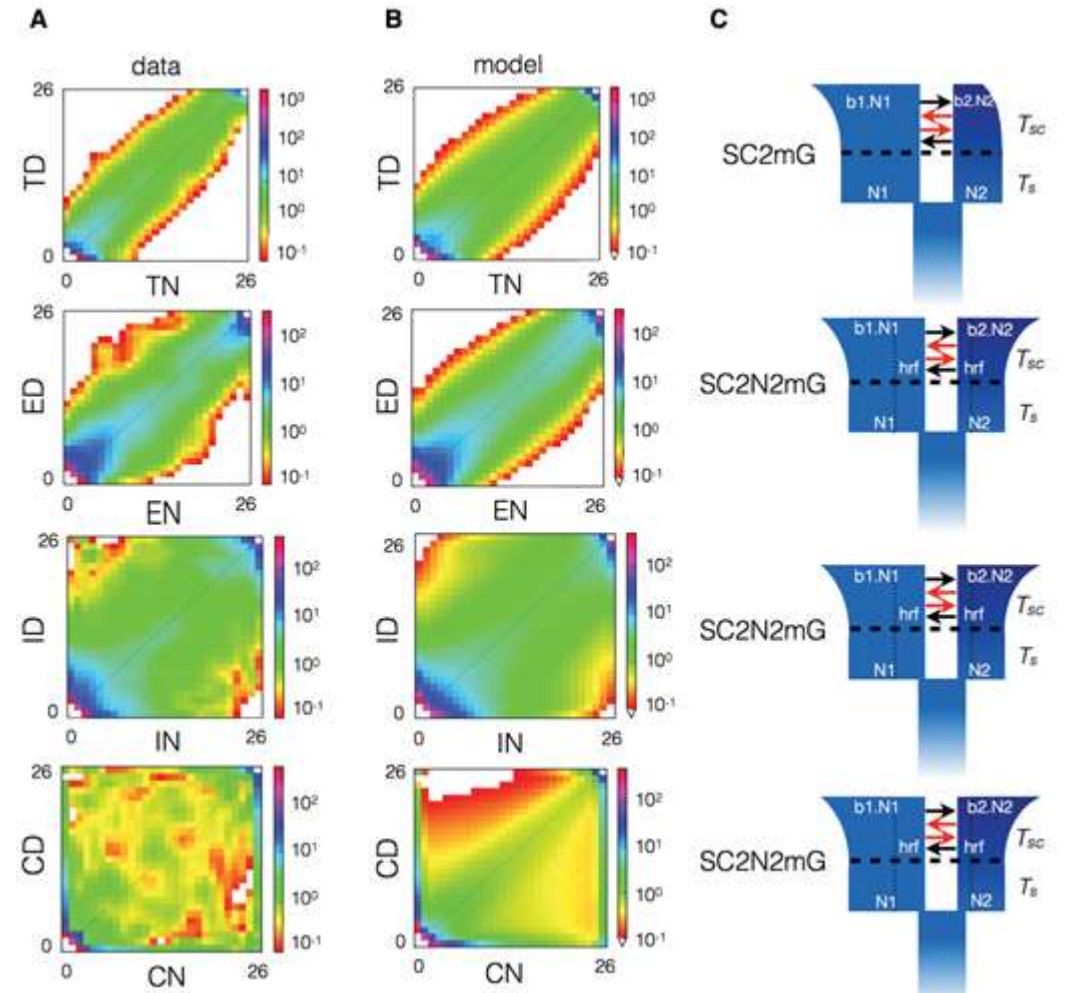- demography can set a null model against which one can look for the effect of selection



Nature Reviews | Genetics

# Beyond present structure…
# How to know population history and demography?

Based on coalescence theory

Compare SFS (site frequency spectrum) between real data and modelled data under different scenario

Common tools: dadi, FastSimCoal, ABC…

# Population structure and demography

### A good overview

Schraiber, J., Akey, J. Methods and models for unravelling human evolutionary history. Nat Rev Genet 16, 727–740 (2015). https://doi.org/10.1038/nrg4005

Nature Reviews | Genetics