

# Population structure and demography

Prepared by Claire Mérot & Anna Tigano  
Physalia Course

# Why does population structure matter when studying adaptation?

Evolution (including adaptive evolution)

is the result of the interplay of

**Selection**

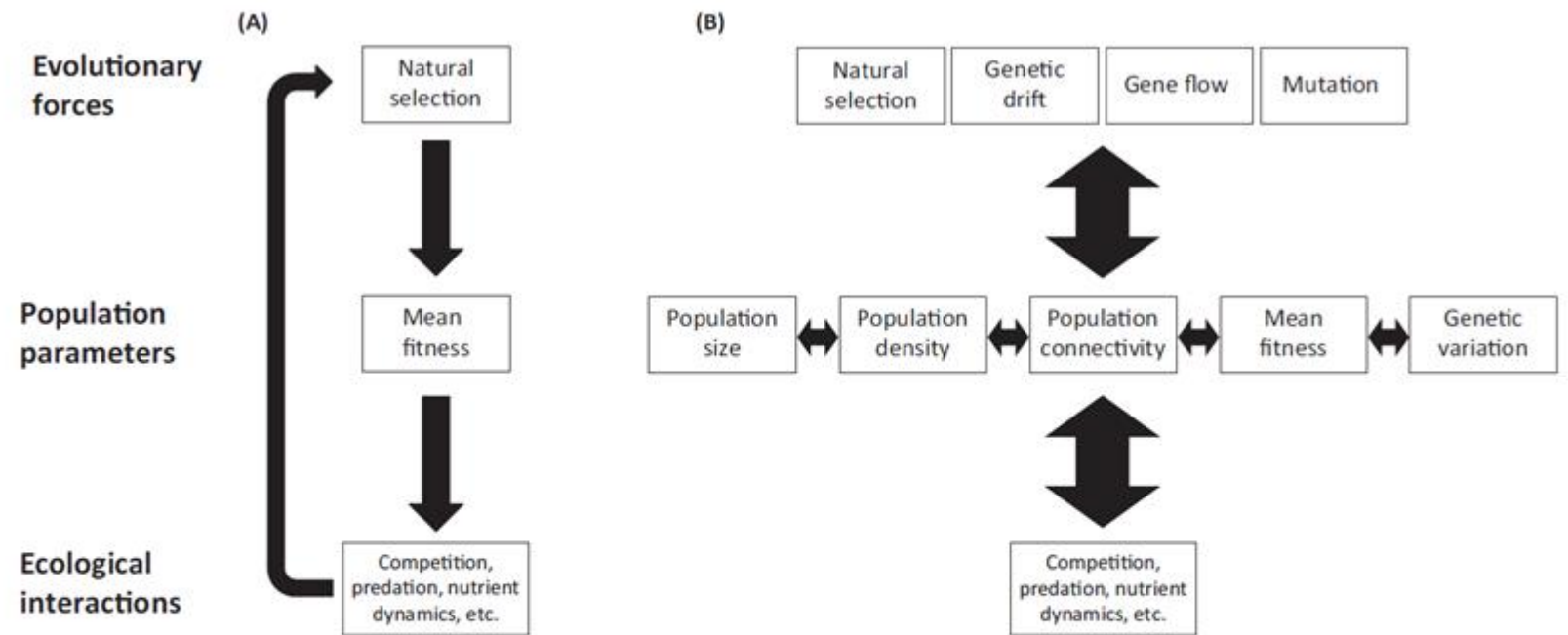
**Drift**

**Mutation**

**Gene flow** (migration + recombination)

# Evolutionary, demographic and ecological processes are inseparable

Lowe, W. H., Kovach, R. P., & Allendorf, F. W. (2017). [Population genetics and demography unite ecology and evolution](#). *Trends in Ecology & Evolution*, 32(2), 141-152.



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	<ul style="list-style-type: none"><li>. Study ecological/functional diversity</li><li>. Understand adaptative processes under divergent or balancing selection</li><li>. Identify candidate genes</li></ul>	<ul style="list-style-type: none"><li>. Understand the past history of populations</li><li>. Describe population connectivity</li><li>. Assess general genetic diversity</li></ul>

# Different loci tells a different story...

---

## AFLP

. Parapatric ecotype *Littorina saxatilis*



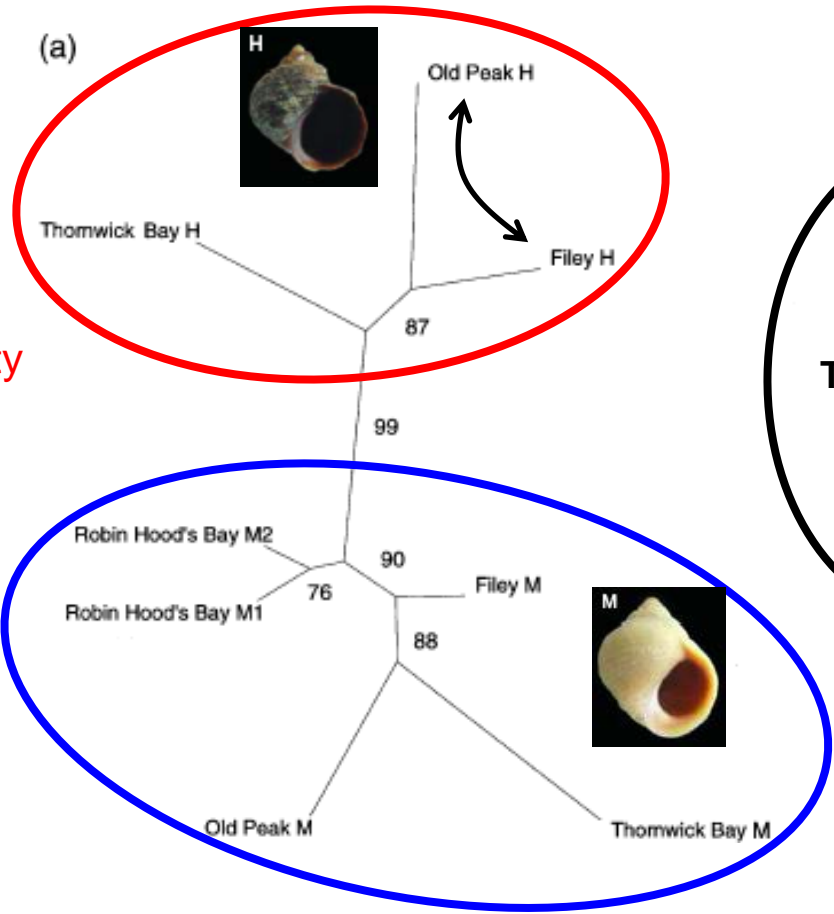
2 ecotypes  
4 regions



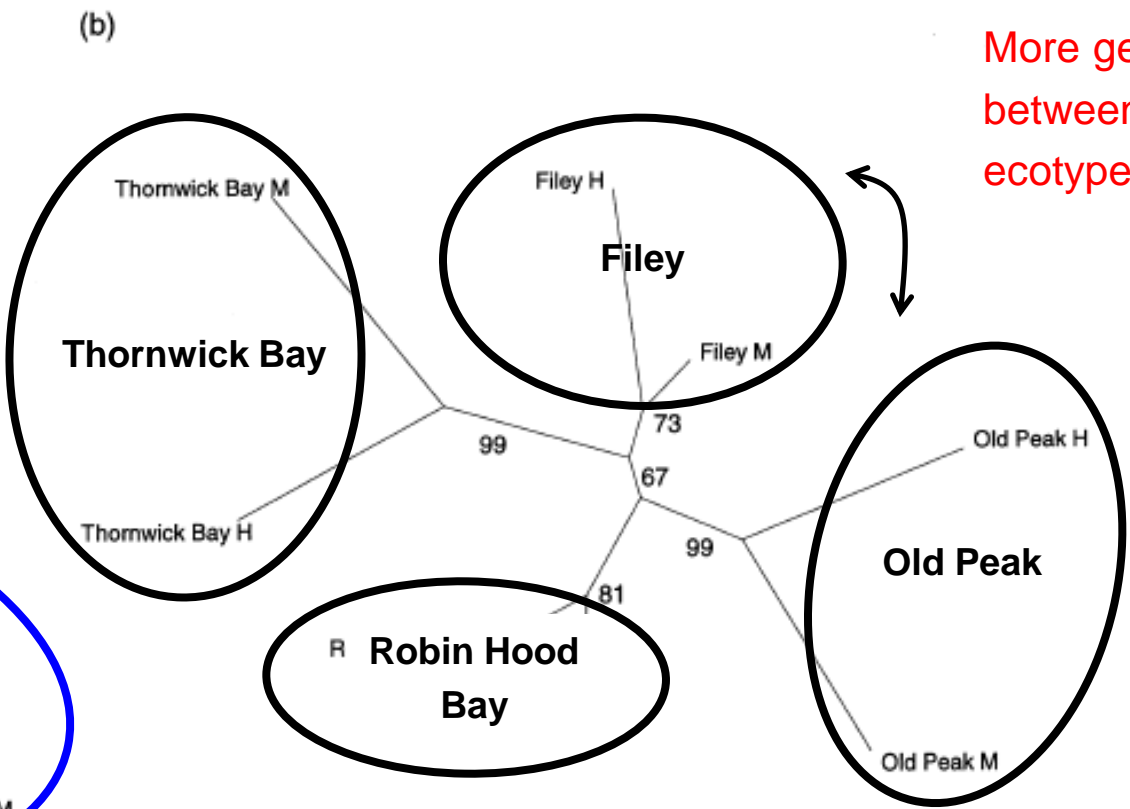
# Different loci tells a different story...

AFLP

All loci (290)



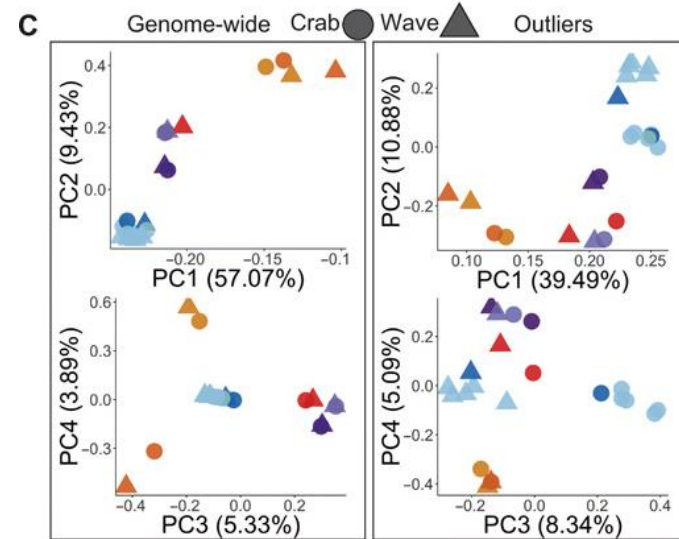
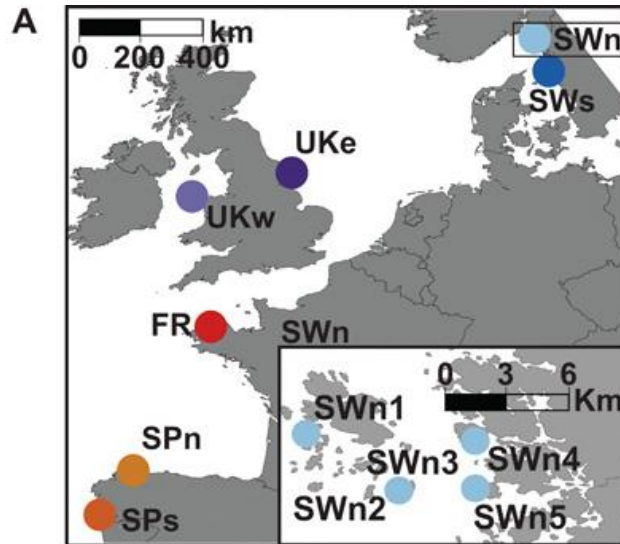
Neutral loci (275)



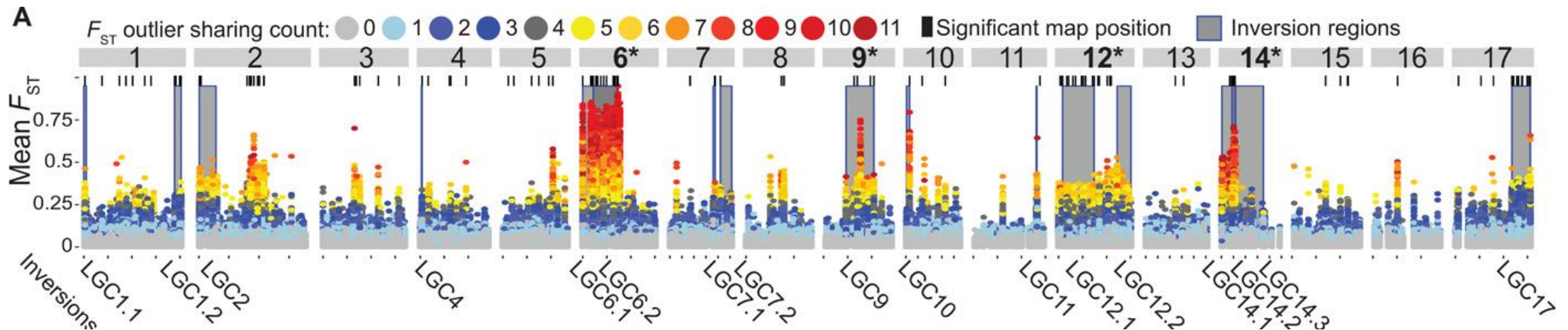
Wilding, C. S., Butlin, R. K., & Grahame, J. (2001). Differential gene exchange between parapatric morphs of *Littorina saxatilis* detected using AFLP markers. *Journal of Evolutionary Biology*, 14(4), 611-619.

# Different loci tells a different story...

## Whole-genome sequencing (pool-seq)



Morales, H. E., Faria, R., Johannesson, K., Larsson, T., Panova, M., Westram, A. M., & Butlin, R. K. (2019). Genomic architecture of parallel ecological divergence: beyond a single environmental contrast. *Science advances*, 5(12), eaav9963.



# 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.

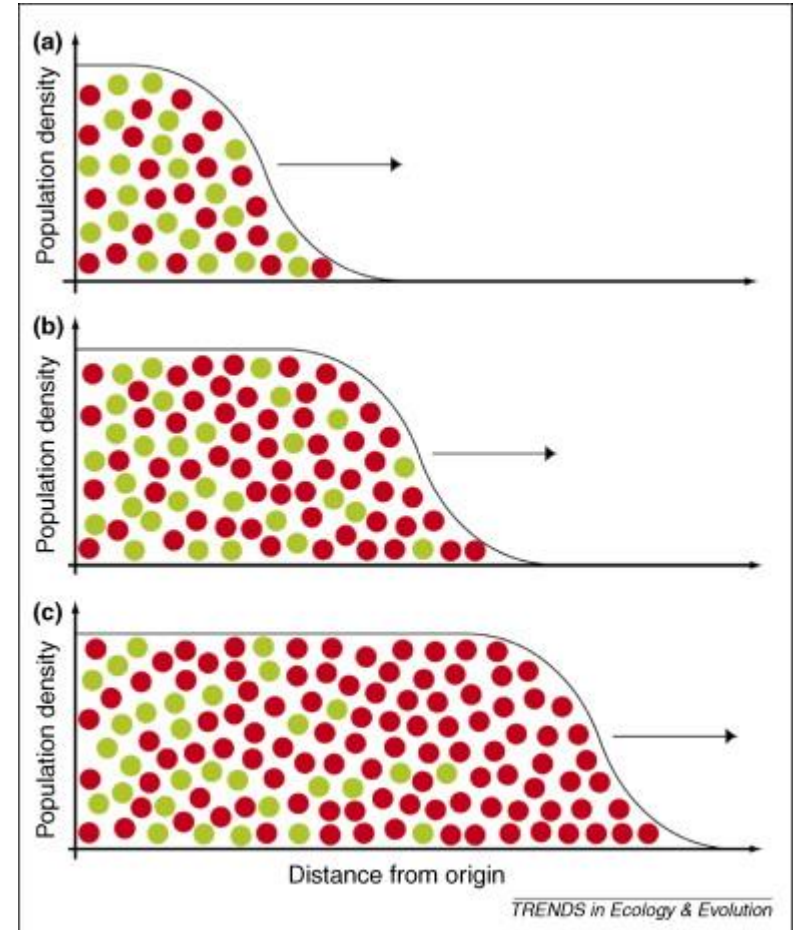
Drift is the main driver of genetic population structure, and can generate a genetic footprint similar to that of selection.



# 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.

⇒ 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



Surfing during population expansions promotes genetic revolutions and structuration

Excoffier & Ray

TREE 2008 <https://doi.org/10.1016/j.tree.2008.04.004>

# Spatial autocorrelation

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

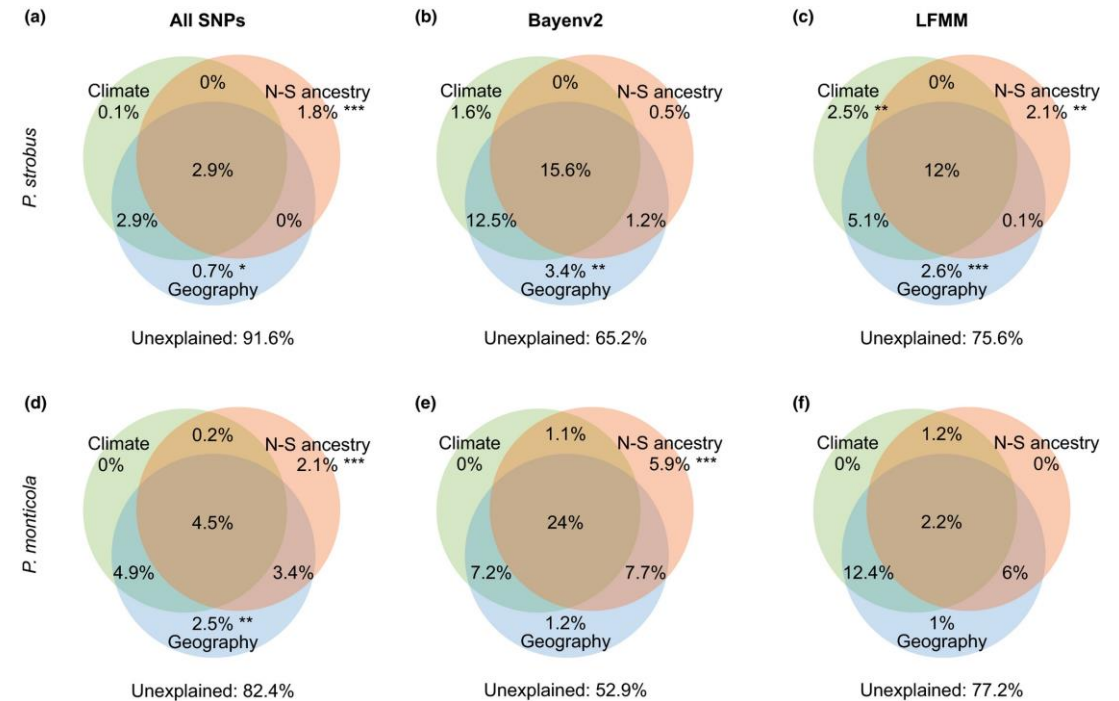
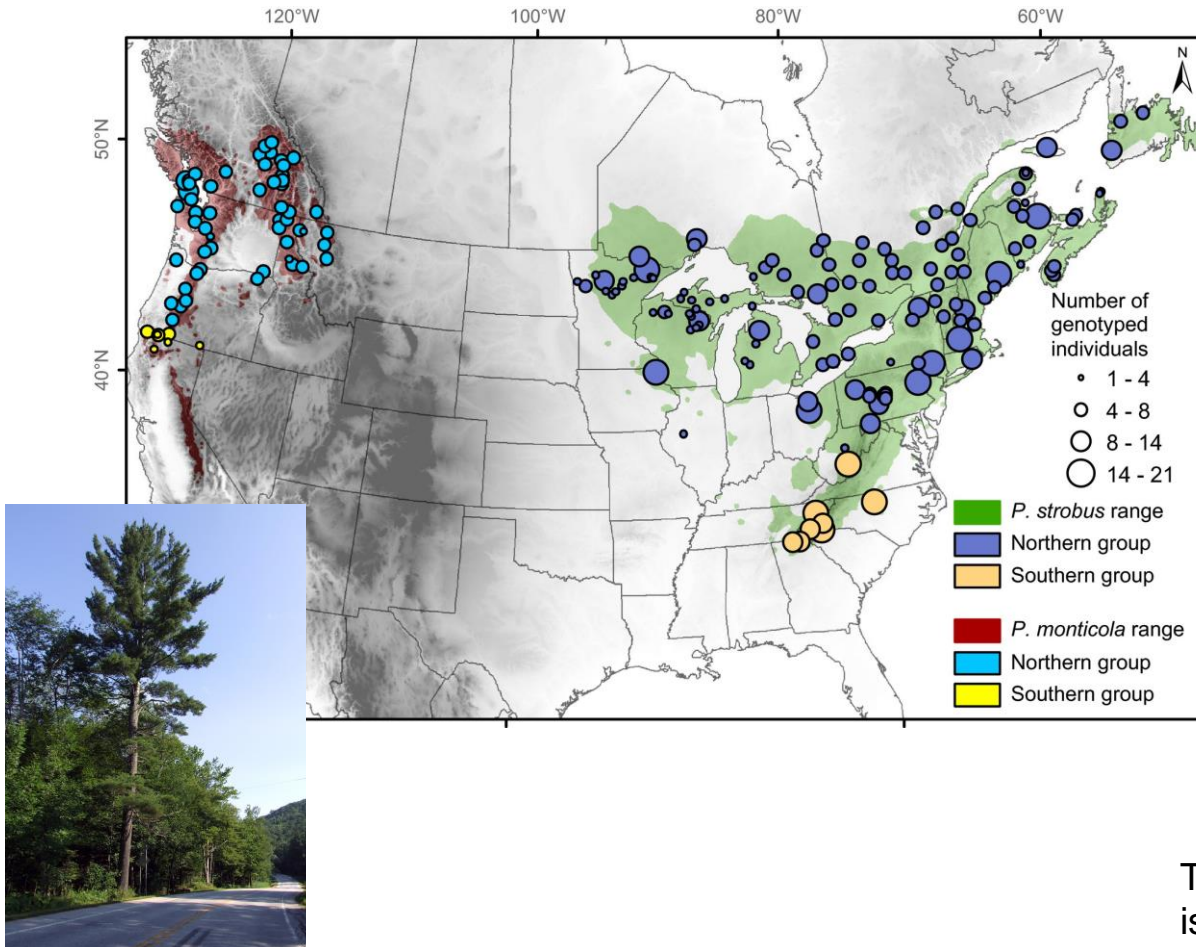
⇒ Which patterns results from the balance selection/migration? Which patterns result from the balance drift/migration?

+ 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?



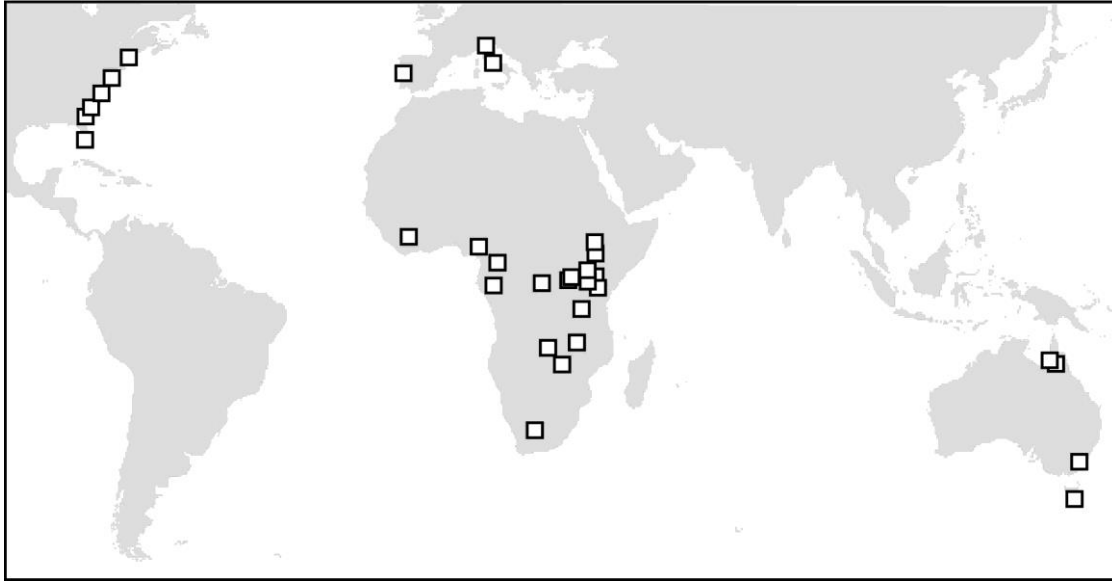
The challenge of separating signatures of local adaptation from those of isolation by distance and colonization history: The case of two white pines  
 Nadeau et al, 2016 <https://doi.org/10.1002/ece3.2550>

# Contact between different lineages / hybridization

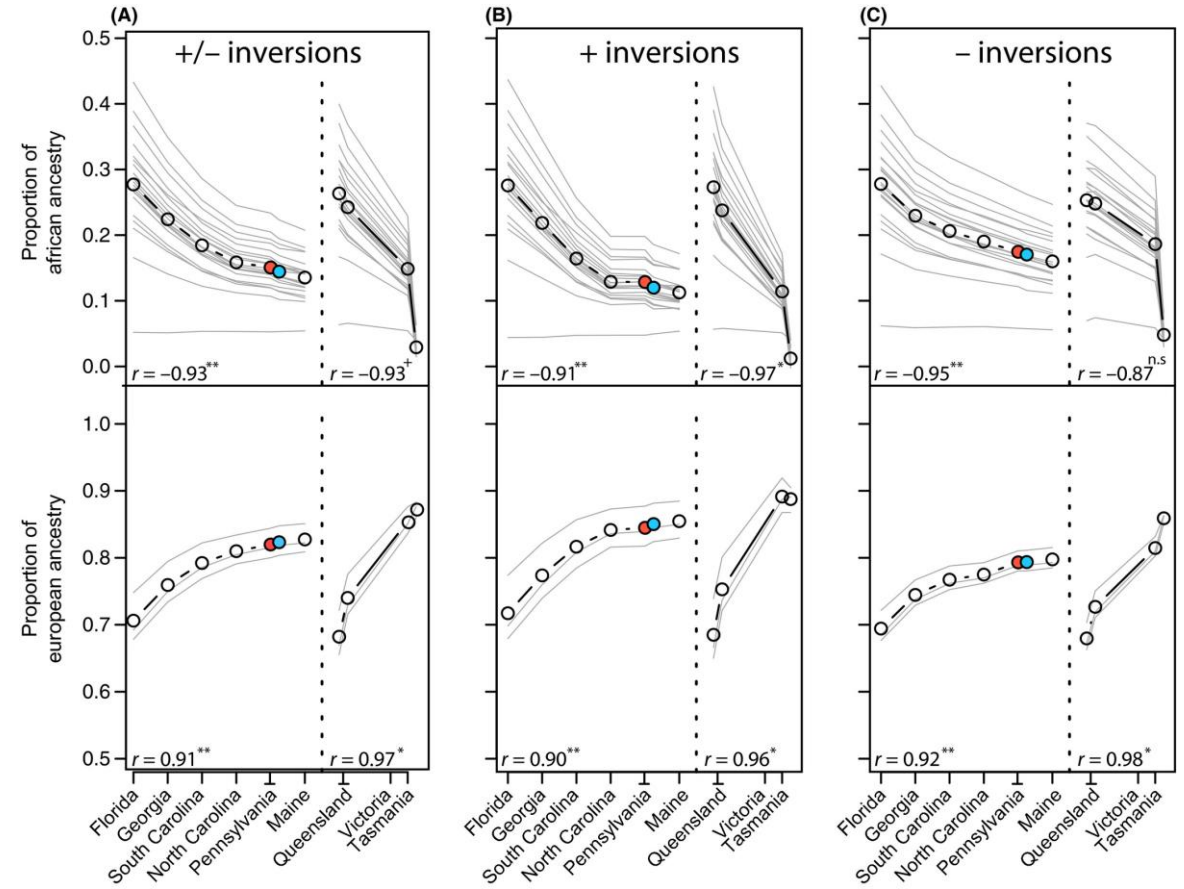
Signature of selection or of local adaptation are best detected in a context of (high) gene flow.

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
- $F_{st}$  between pairs of populations

# Principal Component Analysis (PCA)

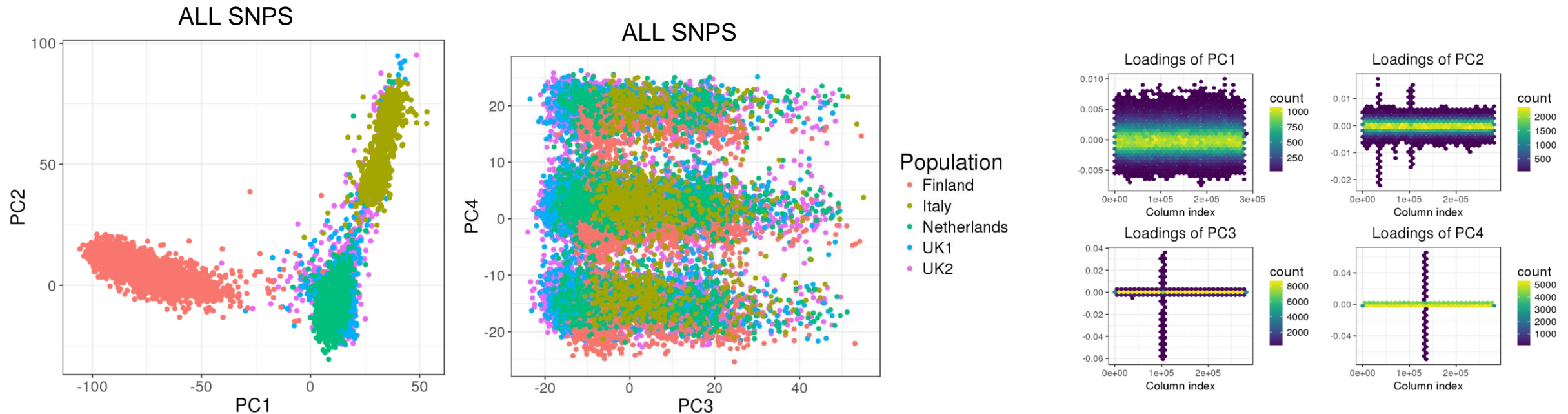
- A common statistical tool that reduces matrix complexity by identifying the eigenvectors and ordering them
- The top PCs reflect axis of genetic variation along which individuals with same ancestry, or exchanging genetic material, are more similar to each other.
- Caution: can be strongly driven by 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>

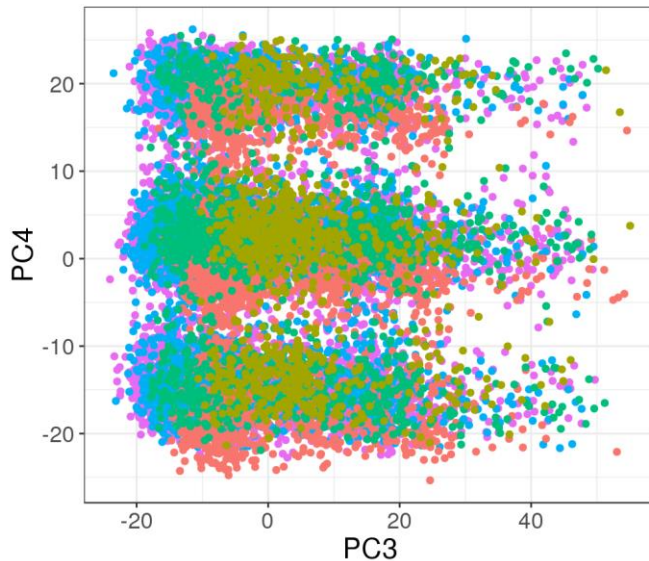


# Principal Component Analysis (PCA)

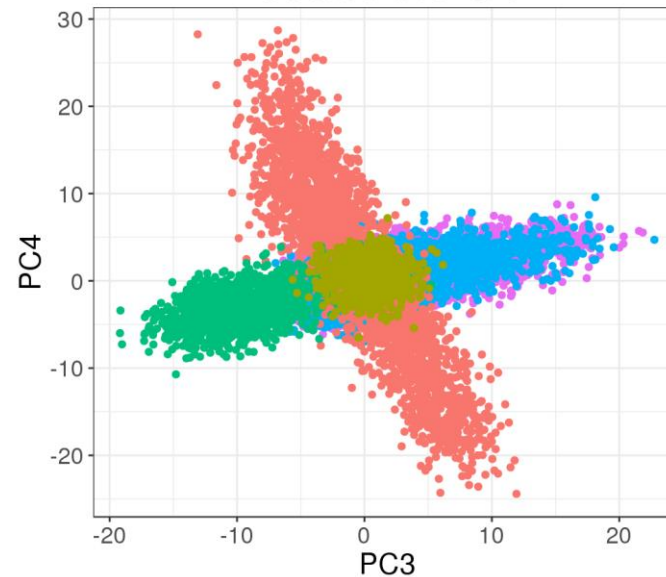
Each individual is a point with coordinates along all PCs

Each genetic marker contribute to all PCs with a different strength (loadings)

ALL SNPS

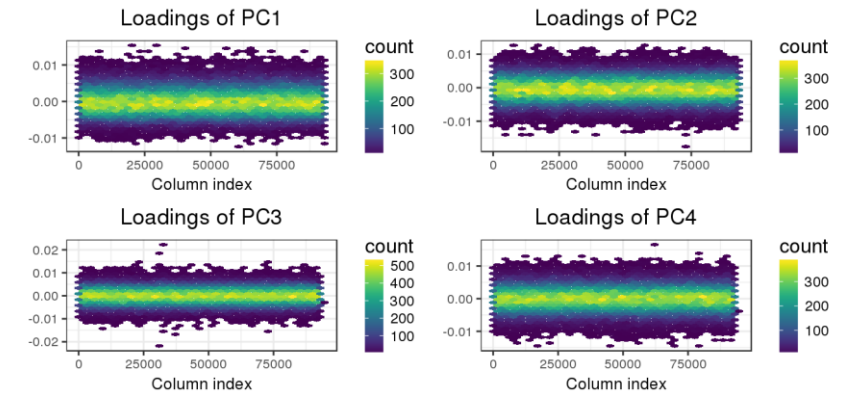


LD-pruned SNPs



Population

- Finland
- Italy
- Netherlands
- UK1
- UK2



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

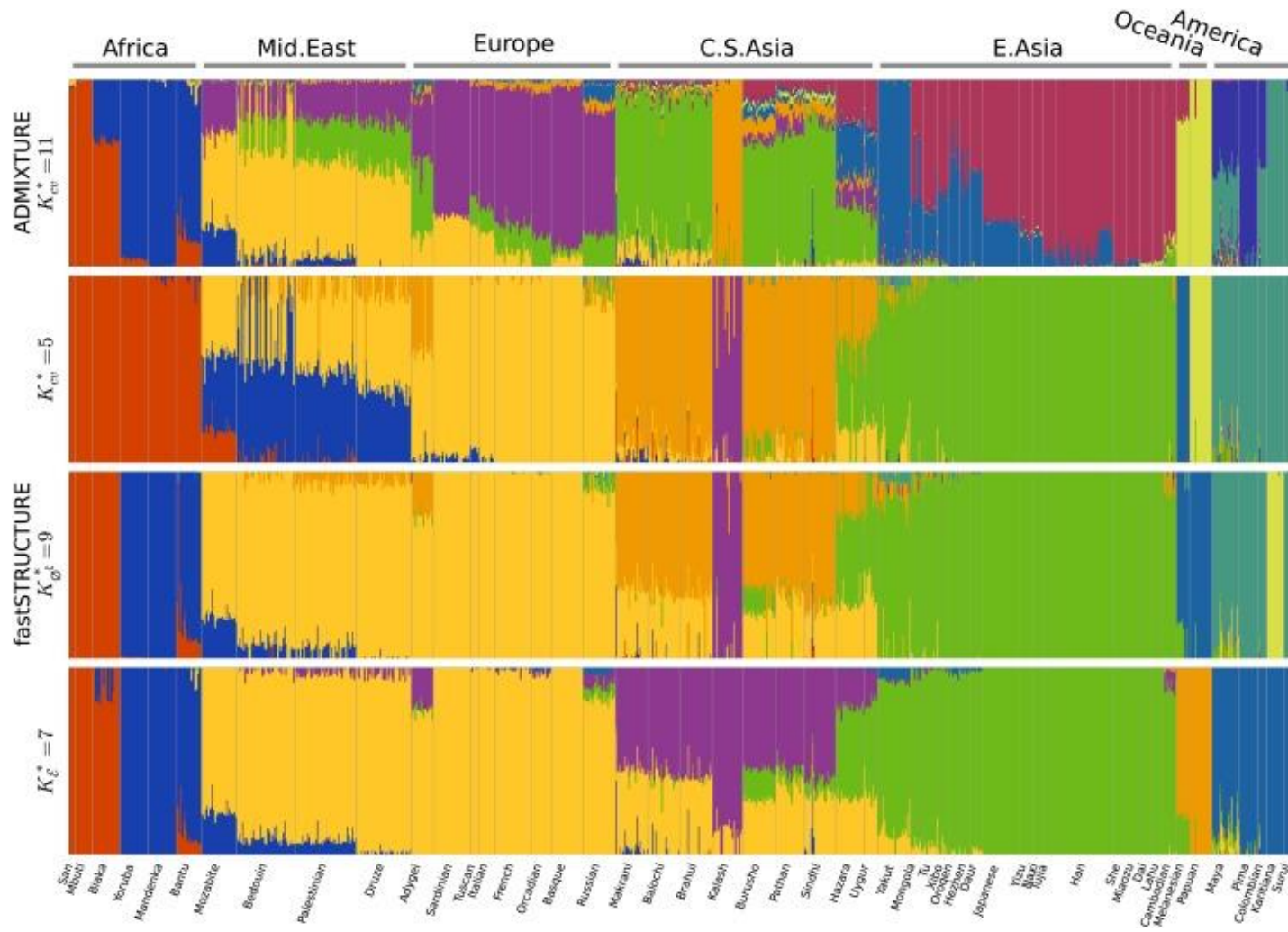
*Evaluation of model fit of inferred admixture proportions*

[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<sup>1</sup>](#), [Matthew Stephens<sup>2</sup>](#), [Jonathan K Pritchard<sup>3</sup>](#)  
[10.1534/genetics.114.164350](https://doi.org/10.1534/genetics.114.164350)

The advantage of unsupervised/semi-supervised methods:  
=> Other surprises!!

Species lineage substructure

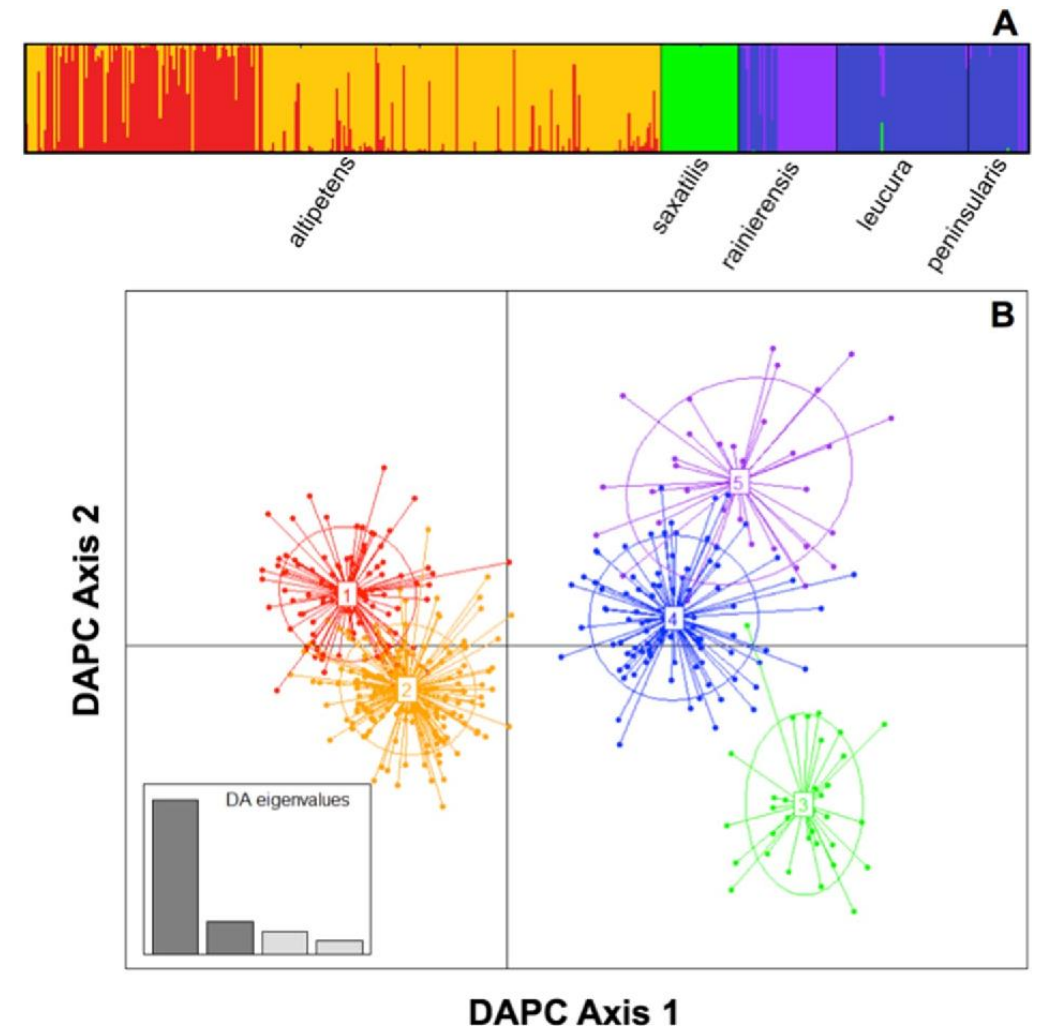
Hybridisation

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>

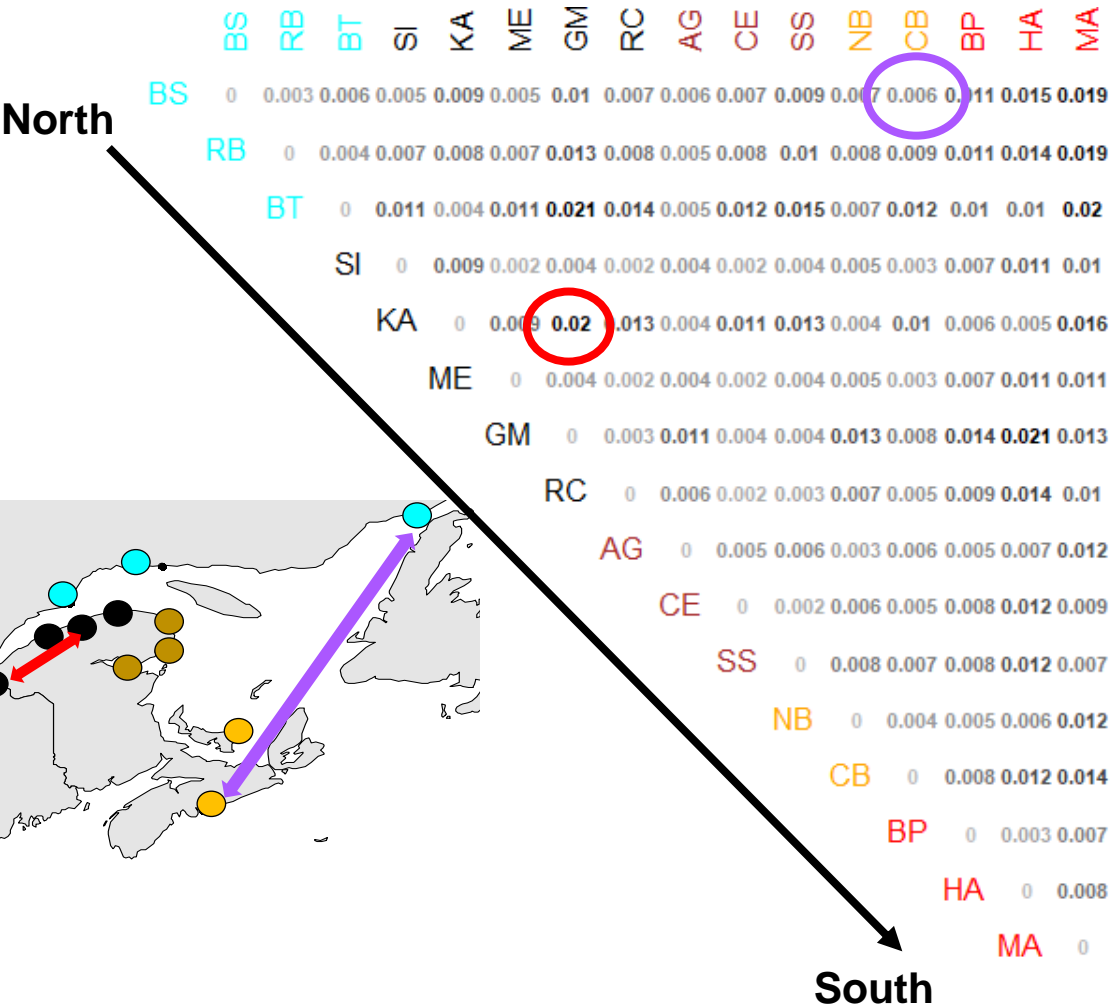


# Pairwise Fst

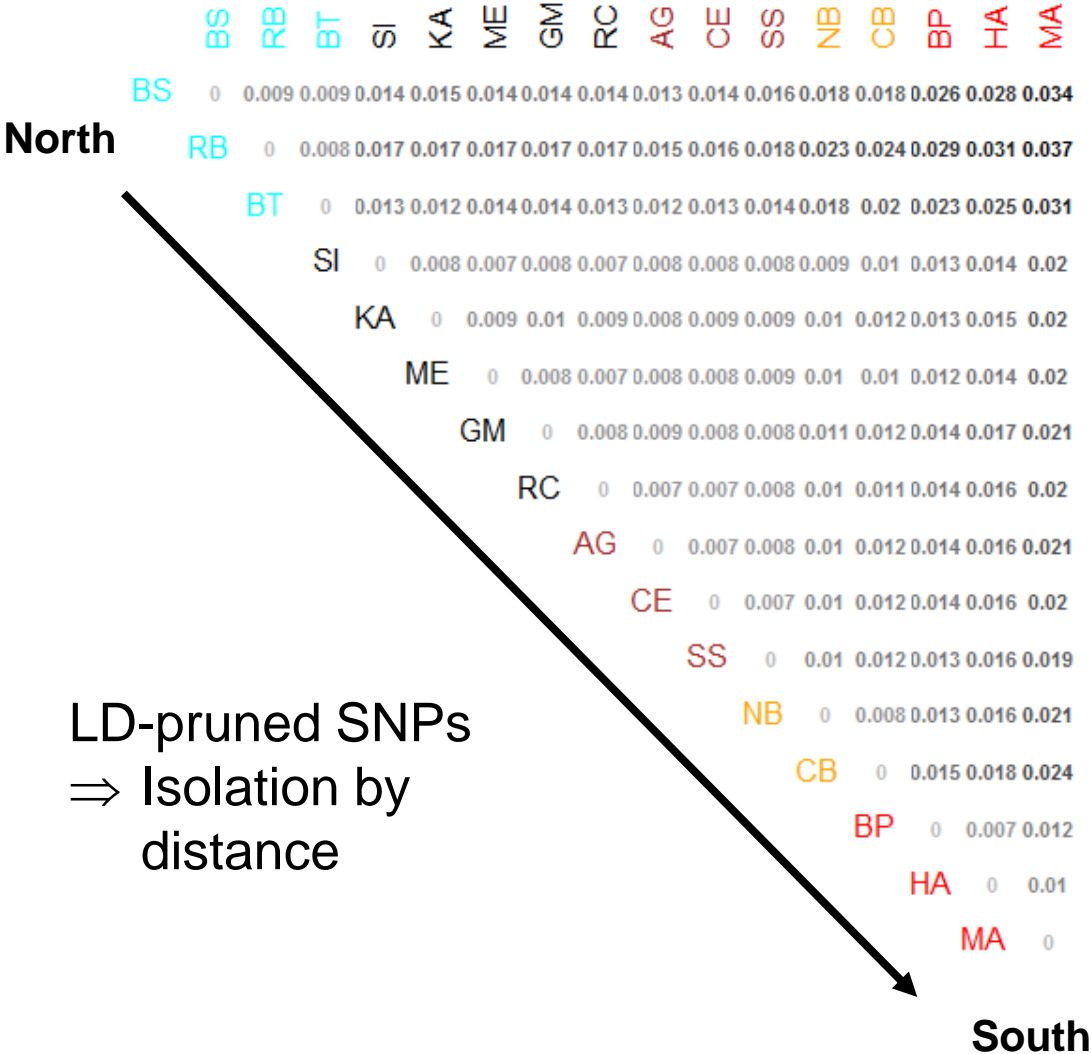
- Fst can only be computed between two groups.
- When sampling several populations, we will be interested in Fst between all pairs of individuals.
- A measure of genetic distance between all populations (does it correlate with ecological distances? Geographic distances? Etc?)
- Again, likely better on 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





# Case study

ORIGINAL ARTICLE

WILEY **MOLECULAR ECOLOGY**

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

Guillermo Friis<sup>1</sup>  | Guillermo Fandos<sup>2</sup>  | Amanda J. Zellmer<sup>3</sup>  |  
John E. McCormack<sup>3,4</sup>  | Brant C. Faircloth<sup>5</sup>  | Borja Milá<sup>1</sup> 



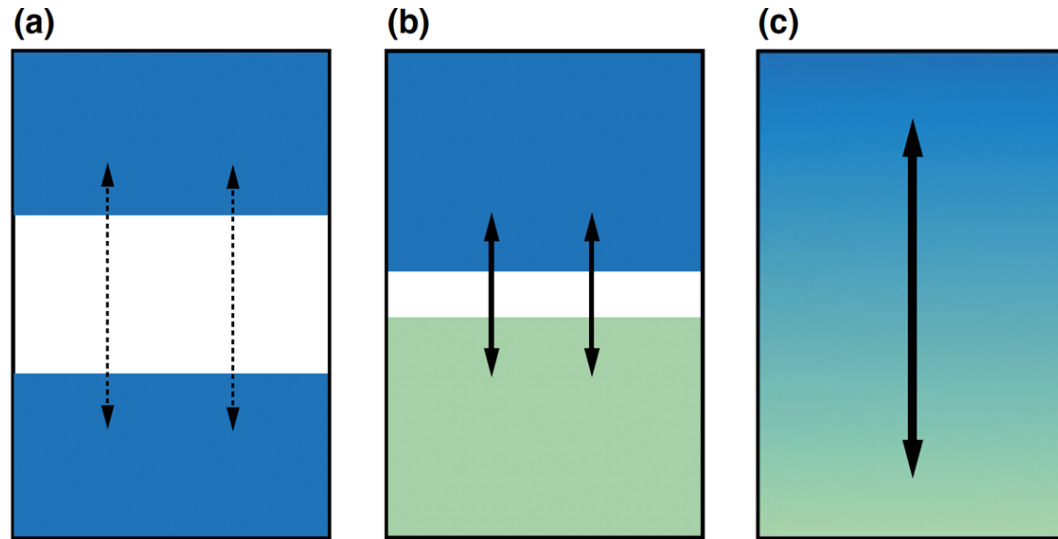


# Set expectations

(a) Geographically isolated populations in similar habitats.

(b) Parapatric populations in ecologically divergent habitats.

(c) Population continuum across a selective gradient



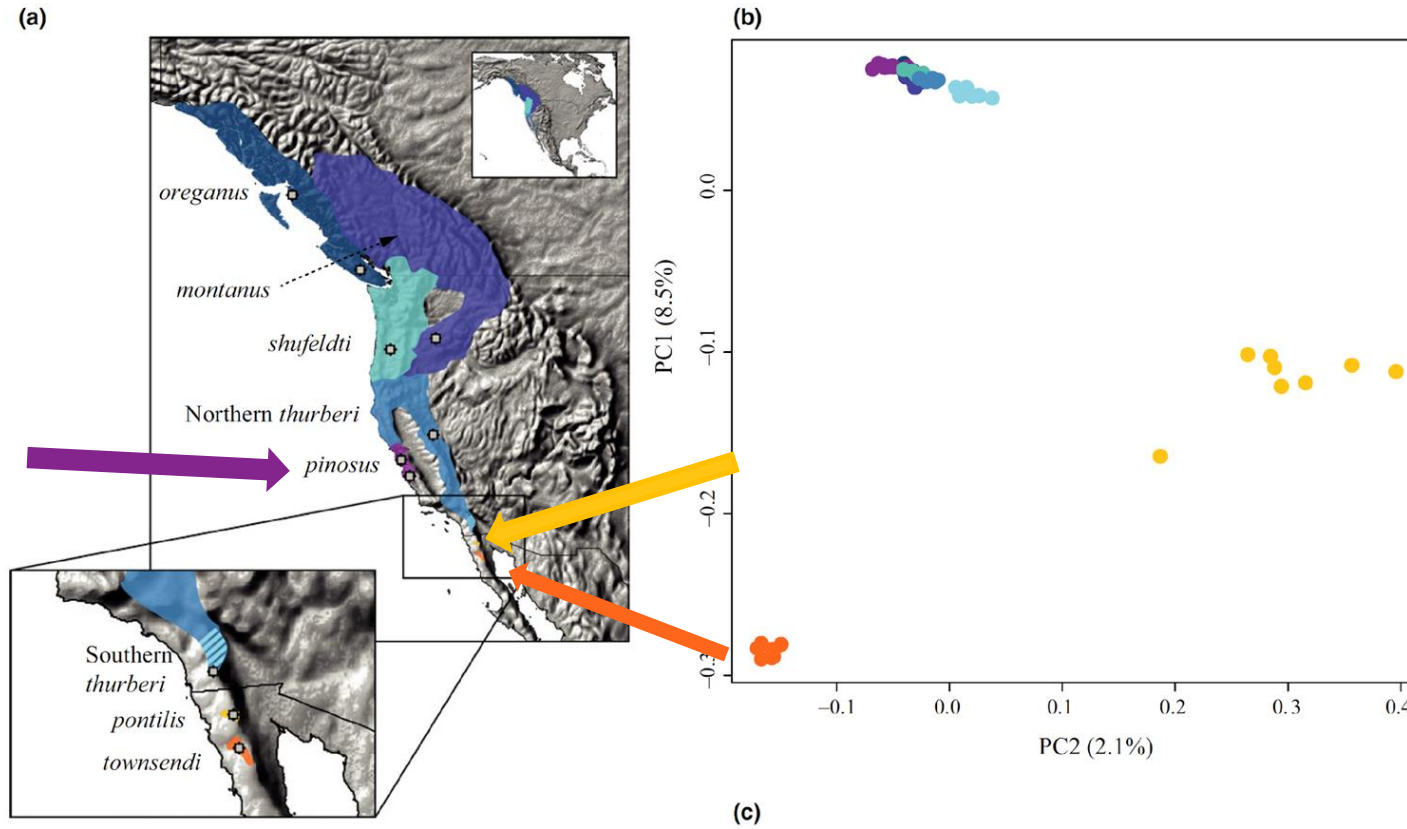
Habitat types:  
Gene flow:  
Neutral divergence:  
Adaptive divergence:

ONE  
LOW  
HIGH  
LOW

TWO  
MODERATE  
MODERATE  
HIGH

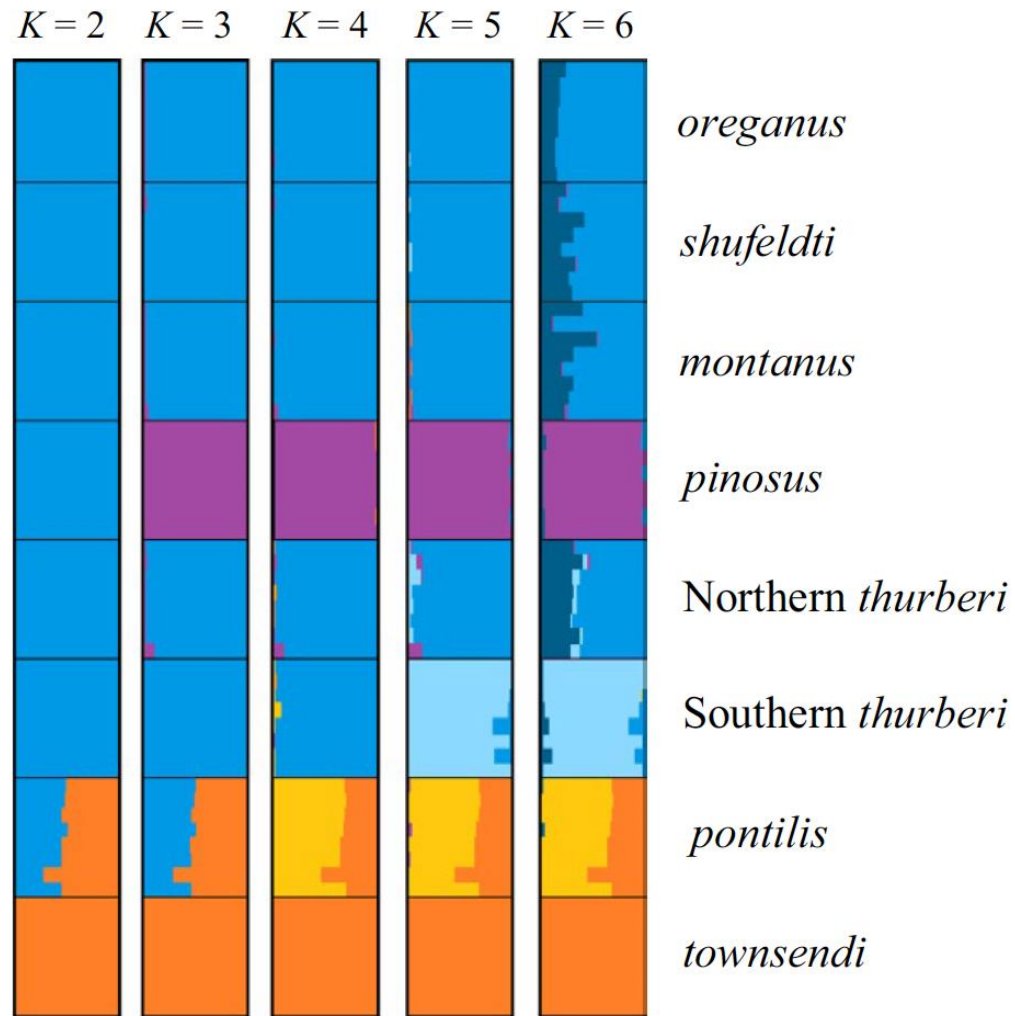
GRADIENT  
HIGH  
LOW  
HIGH

# PCA for data exploration



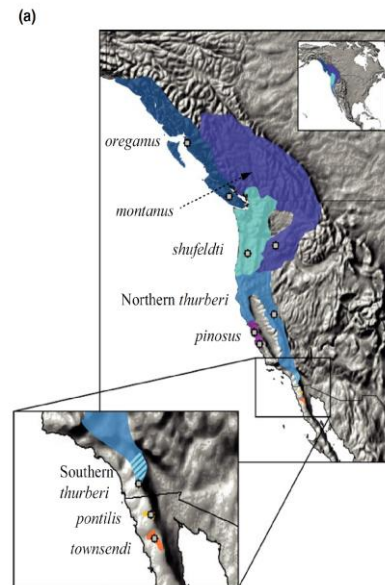
Prior-free data  
exploration

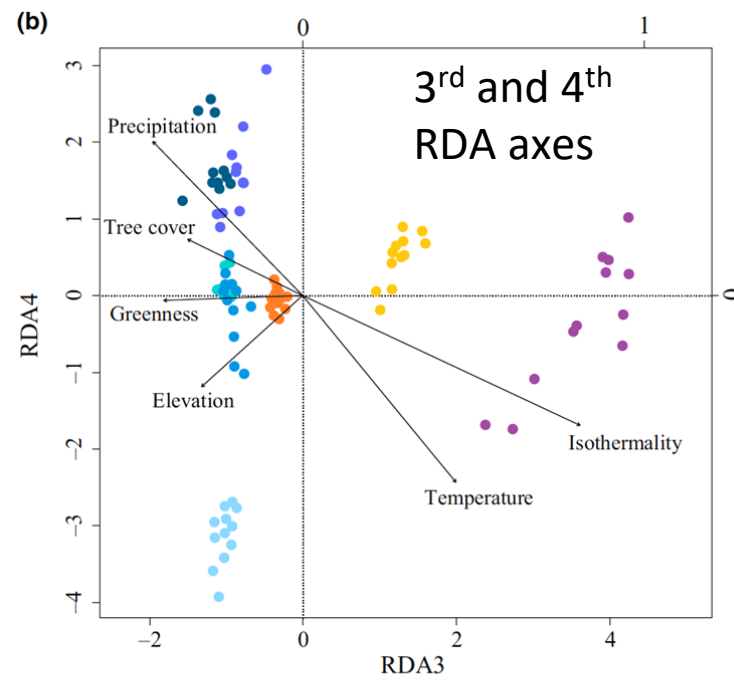
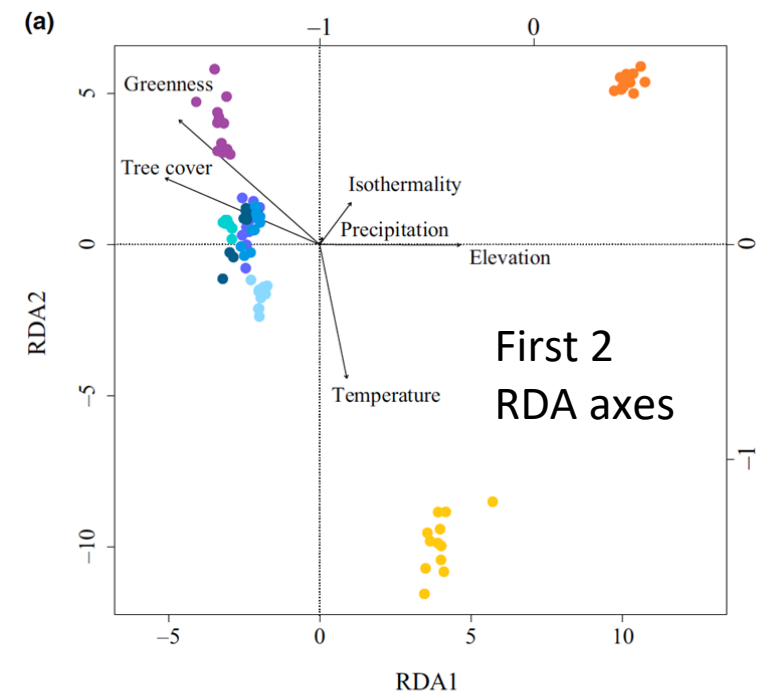
# STRUCTURE for population structure



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

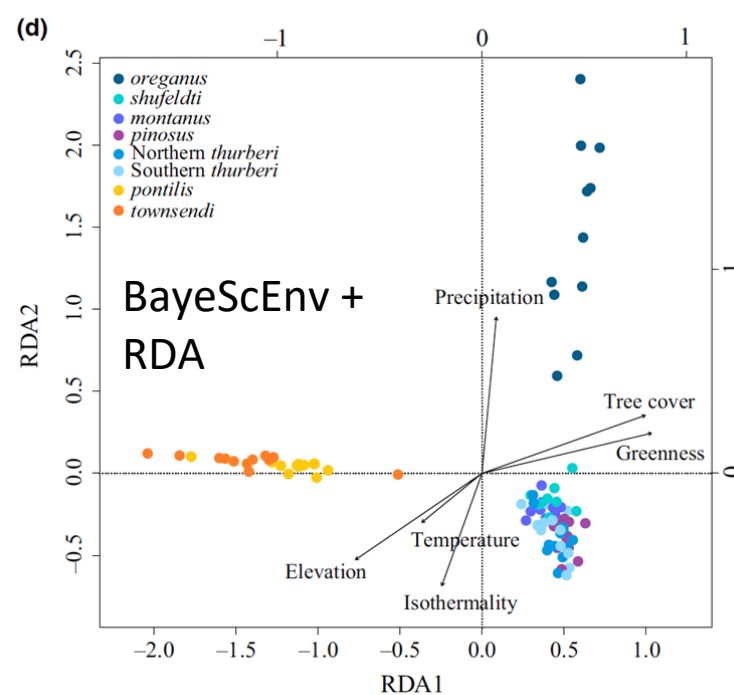
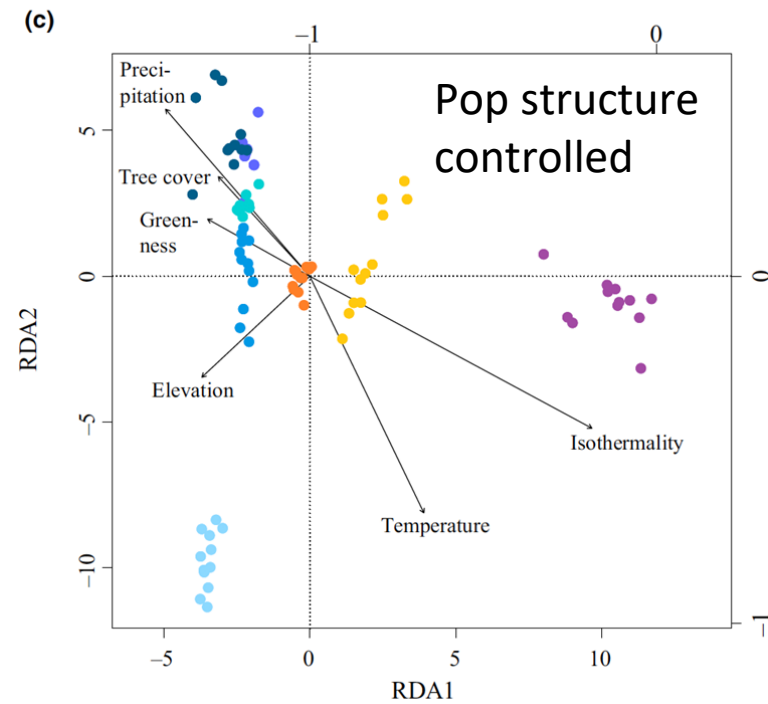
(STRUCTURE doesn't deal well with hierarchical population structure)





# GEA

Genotype-environment  
associations with multiple  
approaches

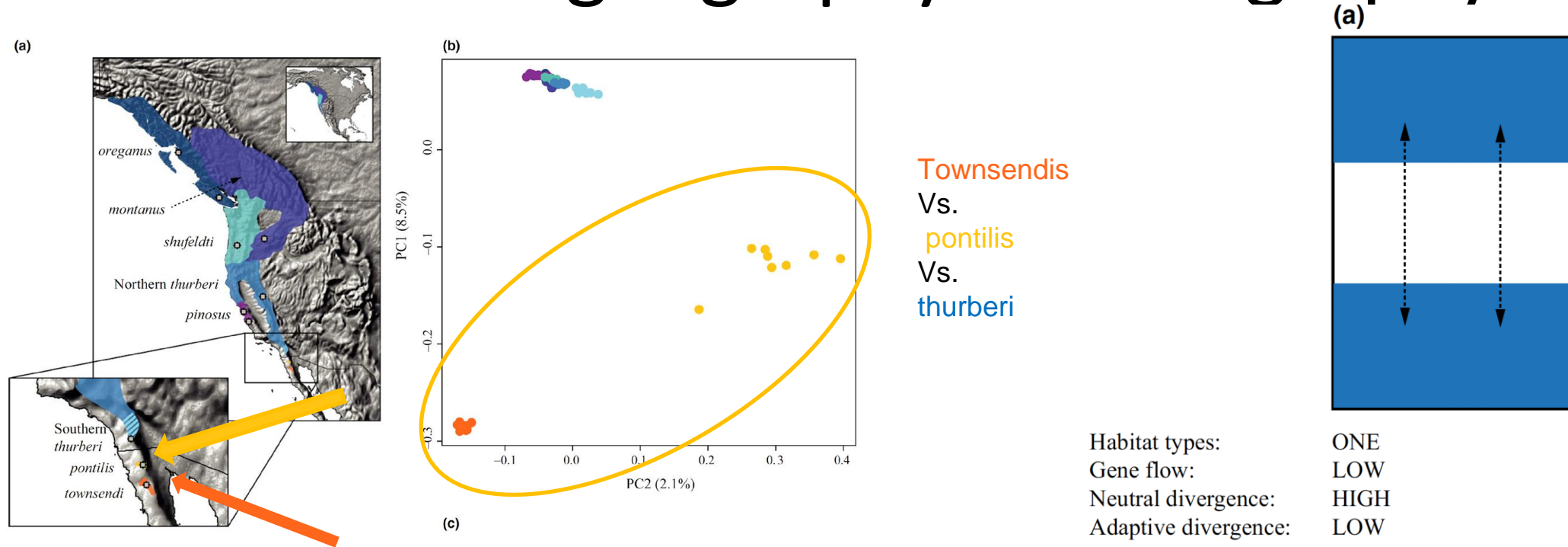


*Lecture  
tomorrow,  
we come  
back on that*

# 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

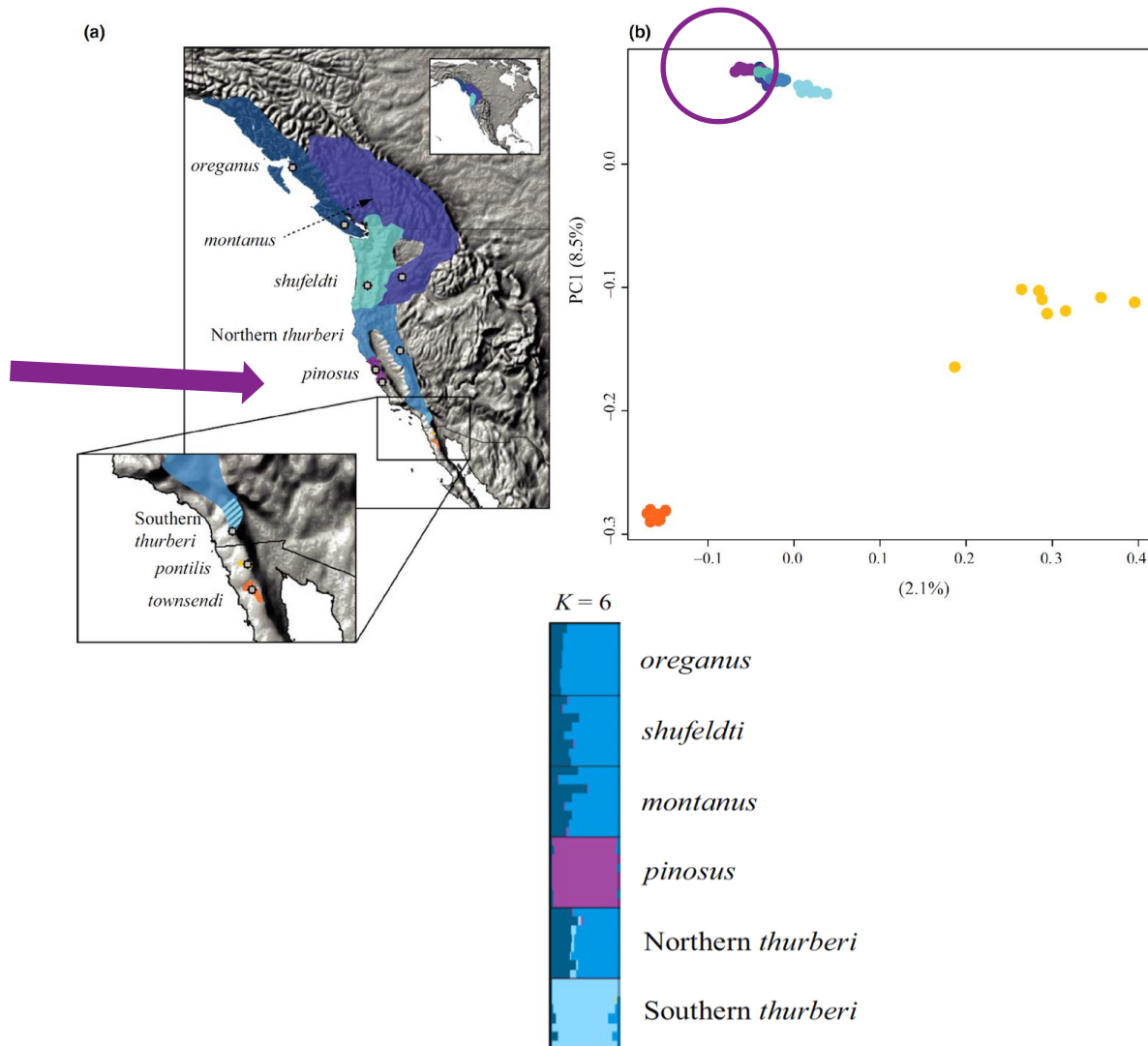


Strong population structure and weak GEA

= Isolation-by-resistance and drift

(small populations + desert between suitable habitats)

# Environment + geography + demography

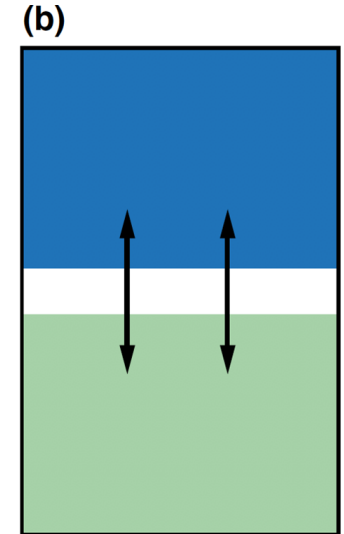


Weak pop structure and  
stronger differentiation on GEA  
= Isolation-by-adaptation?  
Or secondary contact zone?

Pinosus  
Vs.  
Thurberi

Barriers to gene flow  
due to niche  
differences?

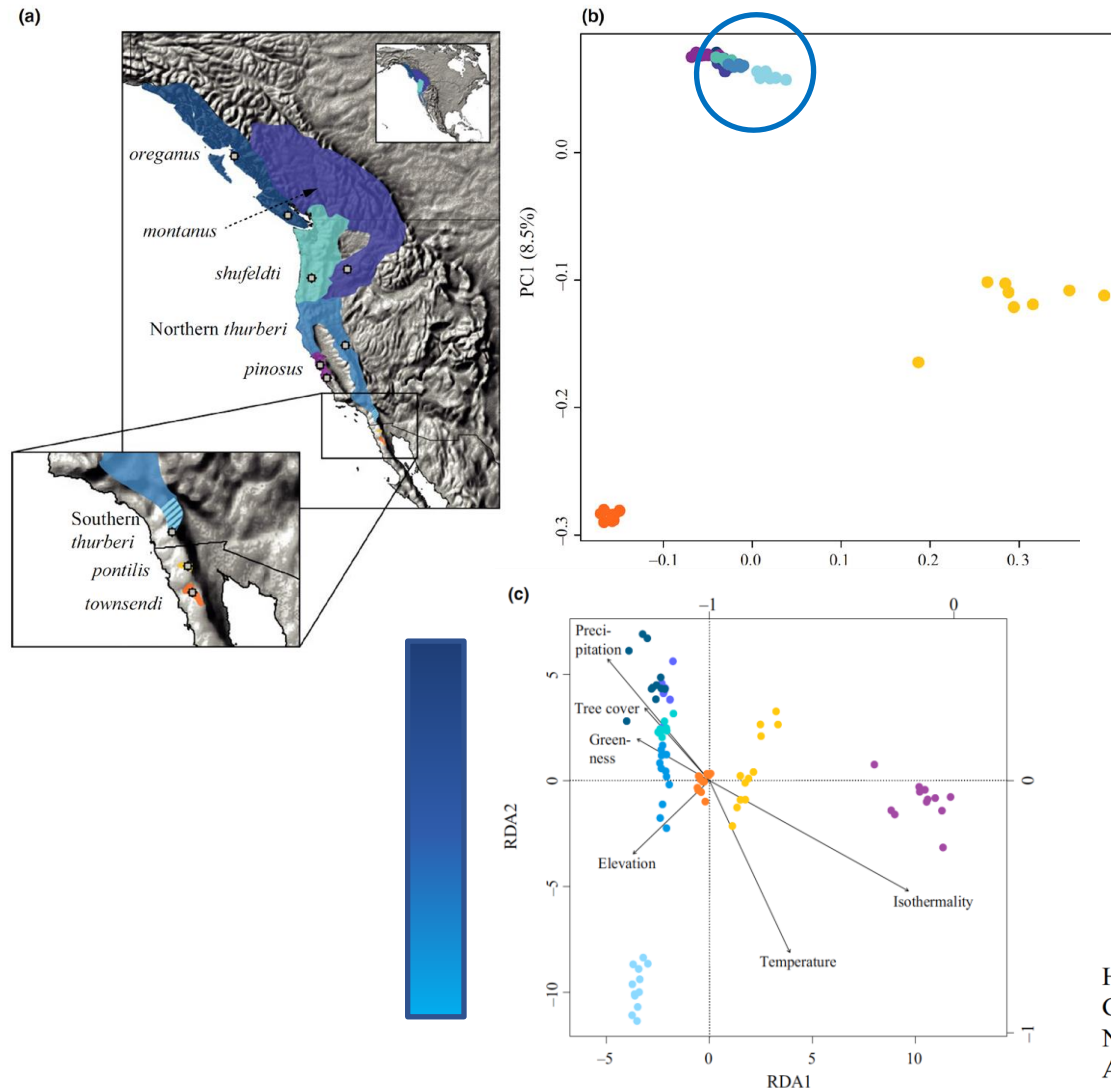
Habitat types:  
Gene flow:  
Neutral divergence:  
Adaptive divergence:



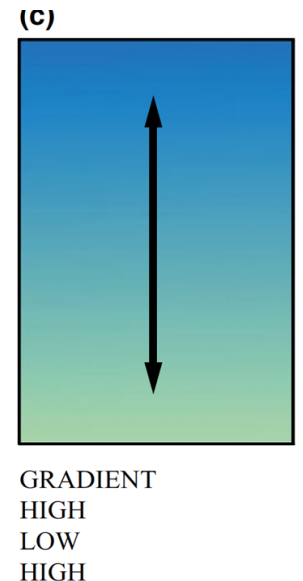
TWO  
MODERATE  
MODERATE  
HIGH



# Environment + geography + demography



No pop structure  
+ environmental associations  
⇒ Ongoing gene flow and local  
adaptation



Habitat types:  
Gene flow:  
Neutral divergence:  
Adaptive divergence:



# And help from experimental work /knowledge of natural history

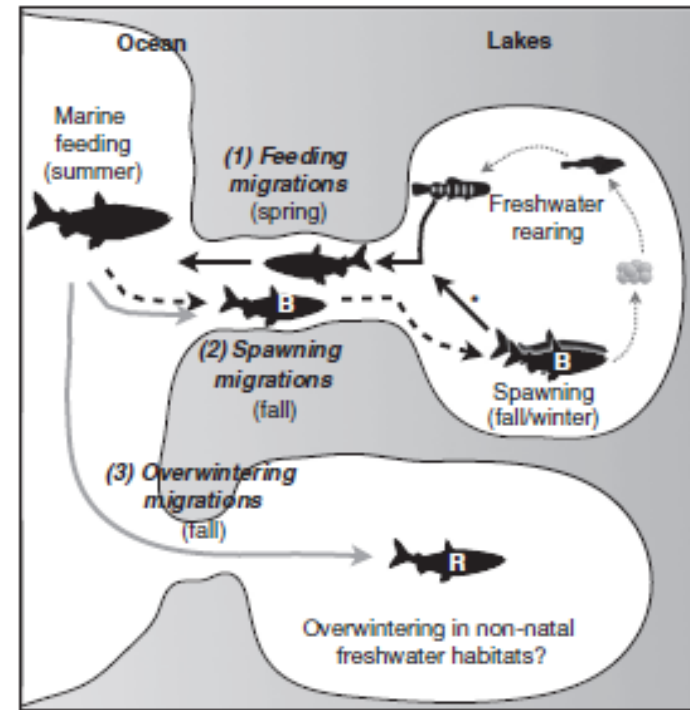
## Capture-Mark-Recapture

-> population size, dynamic and movement

## Spatial ecology

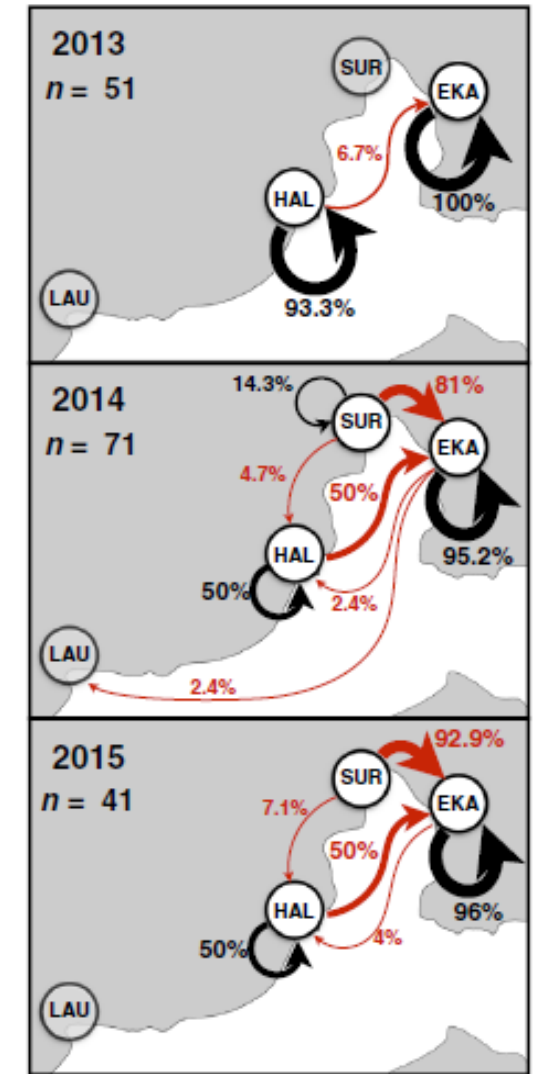
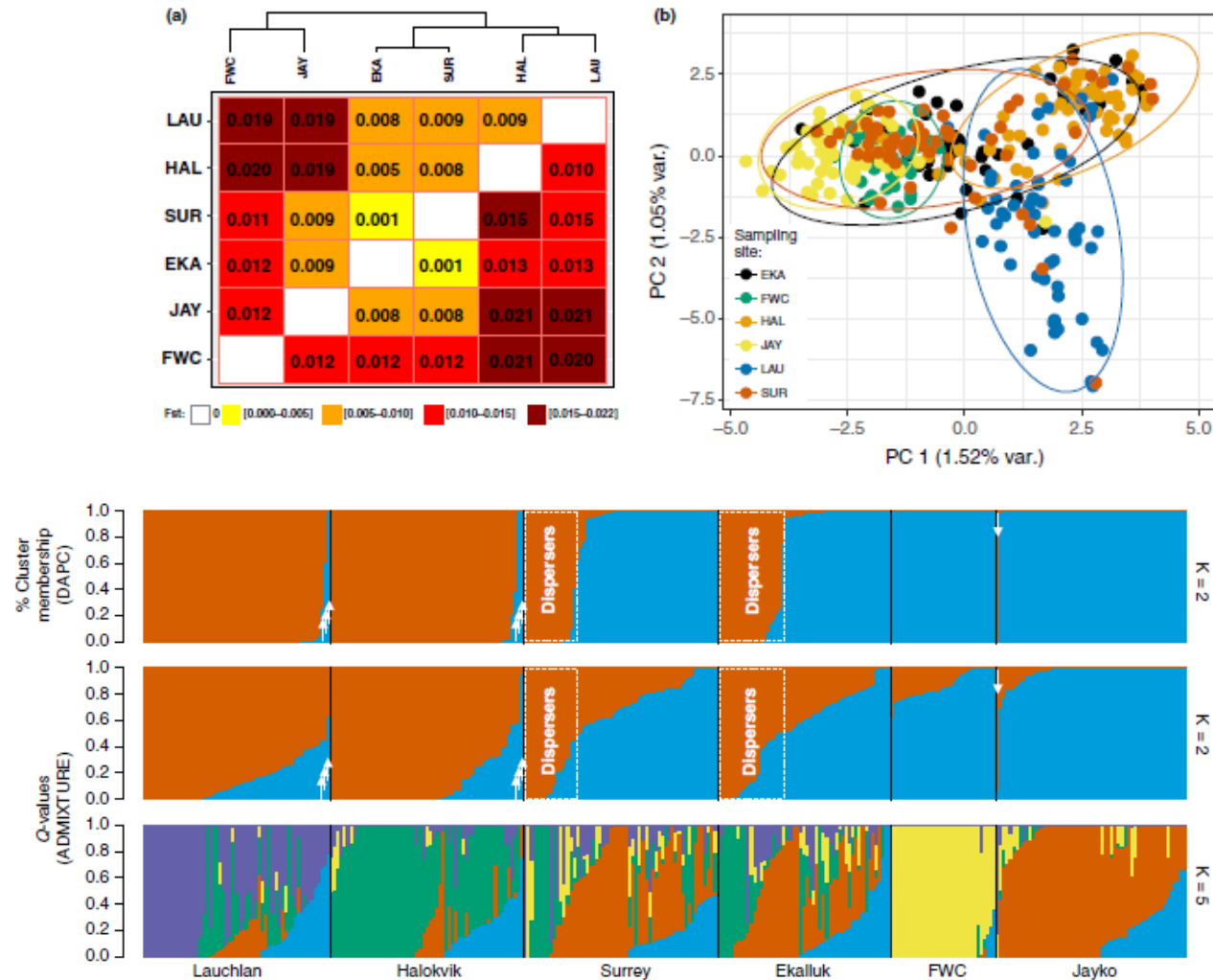
-> tracking, etc..

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



Genomics and telemetry suggest a role for migration harshness in determining overwintering habitat choice, but not gene flow, in anadromous Arctic Char

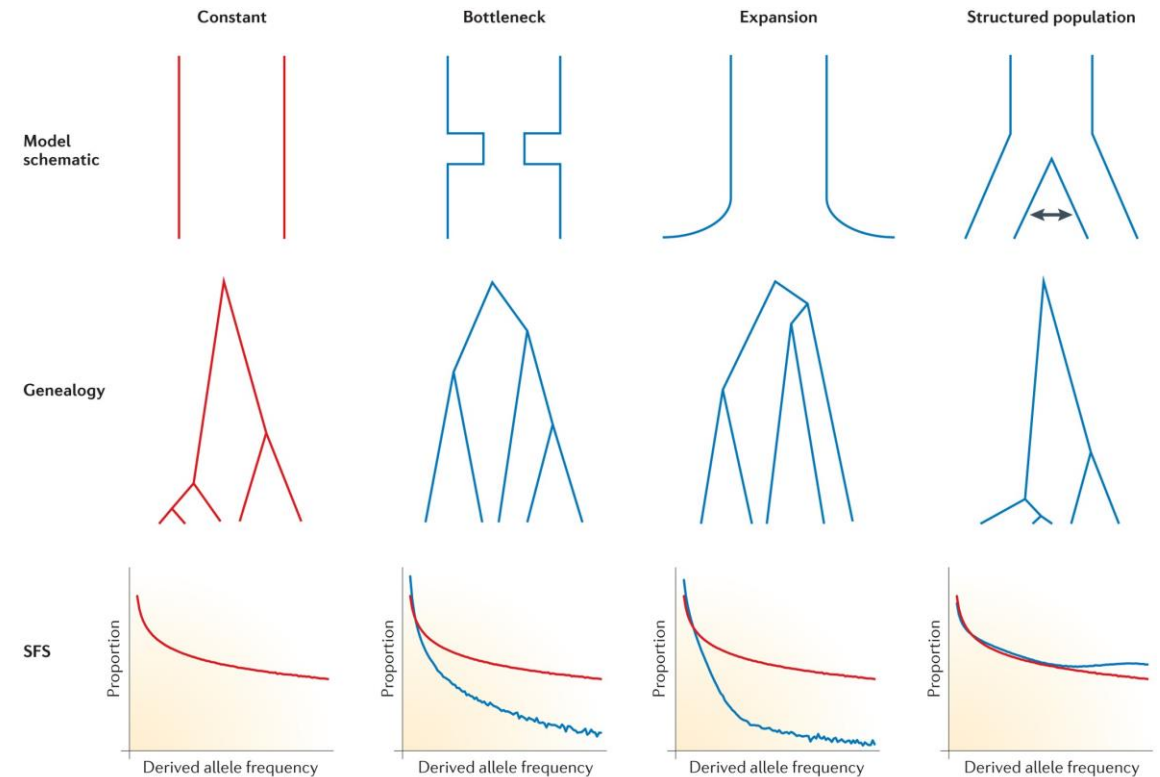
Moore, J.-S., *Molecular Ecology*, 26(24): 6784-6800

# 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



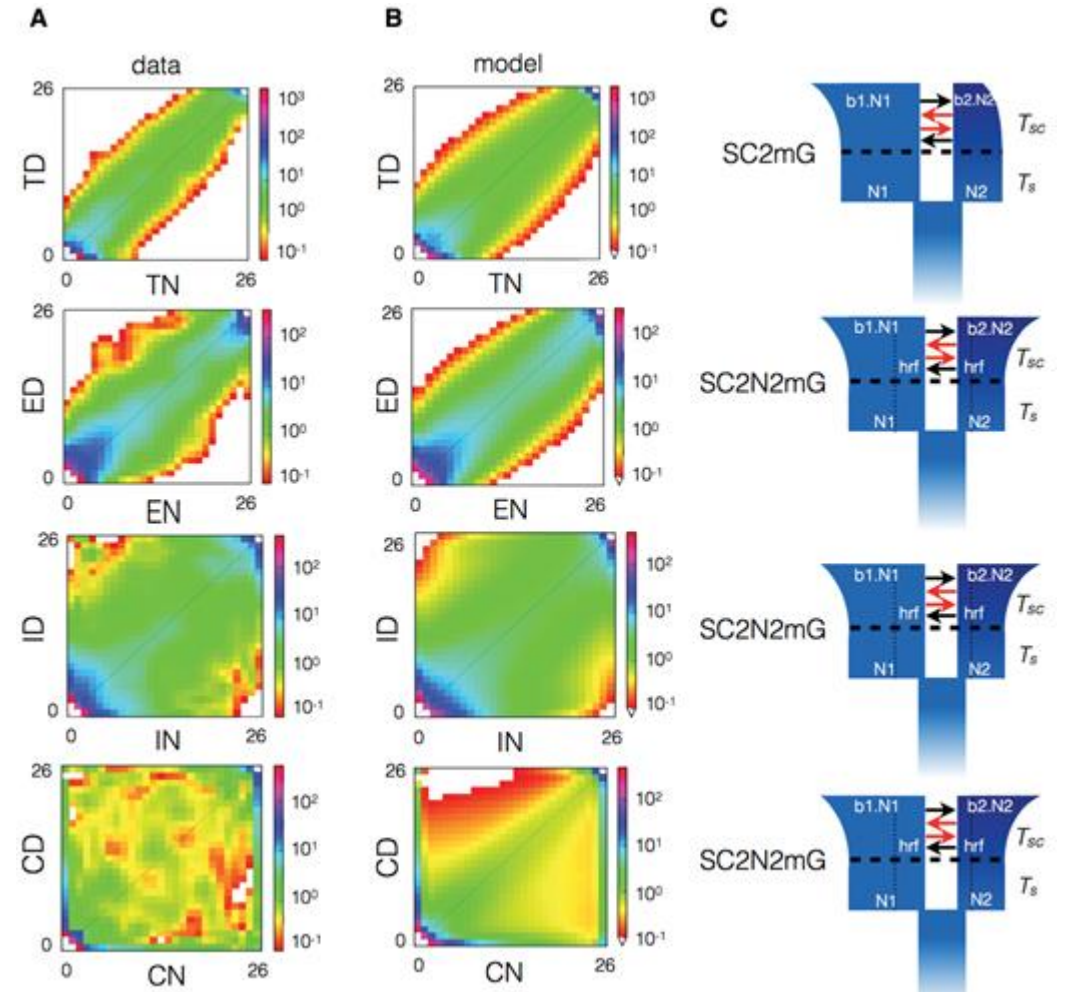
# Beyond present structure...

## 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>

