

# Ecological & evolutionary genomic analyses using RADseq

2017 Workshop on Genomics  
**Český Krumlov**

Bill Cresko  
Institute of Ecology and Evolution  
Department of Biology  
University of Oregon



# Outline for today's lecture

---

RADseq for ecological and evolutionary genomics

Primer on Population Genomics

Evolutionary genomics of stickleback fish

- Population genomics of rapid adaptation
- Using long read RAD-seq for coalescent analyses
- Genome Wide Association Studies using RAD-seq

Genomically enabling the Gulf pipefish

*RAD-seq experimental and statistical considerations*

*Stacks software pipeline (this afternoon & evening)*

# January 2016 Issue

---



STUDY DESIGNS

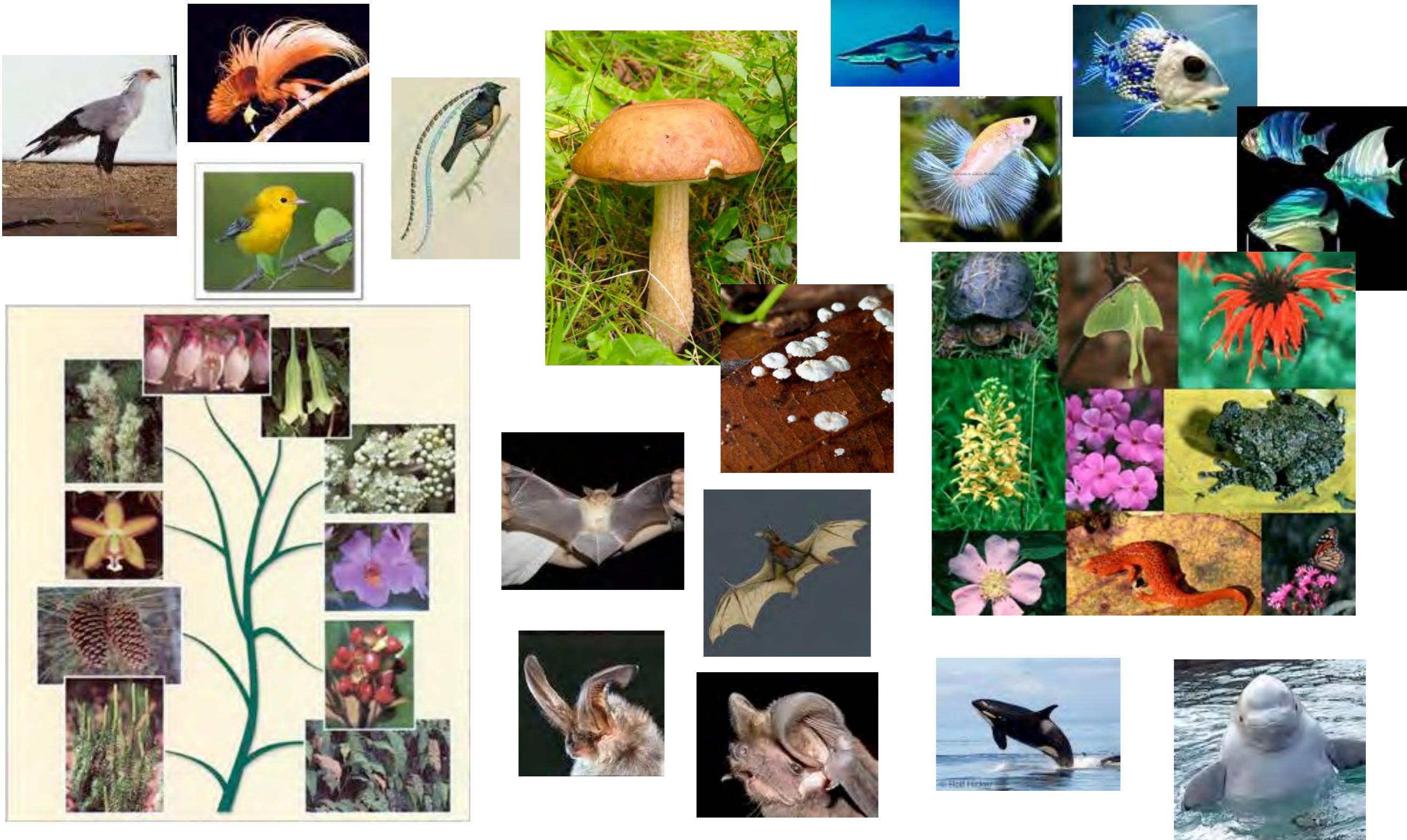
## Harnessing the power of RADseq for ecological and evolutionary genomics

---

*Kimberly R. Andrews<sup>1</sup>, Jeffrey M. Good<sup>2</sup>, Michael R. Miller<sup>3</sup>, Gordon Luikart<sup>4</sup>  
and Paul A. Hohenlohe<sup>5</sup>*

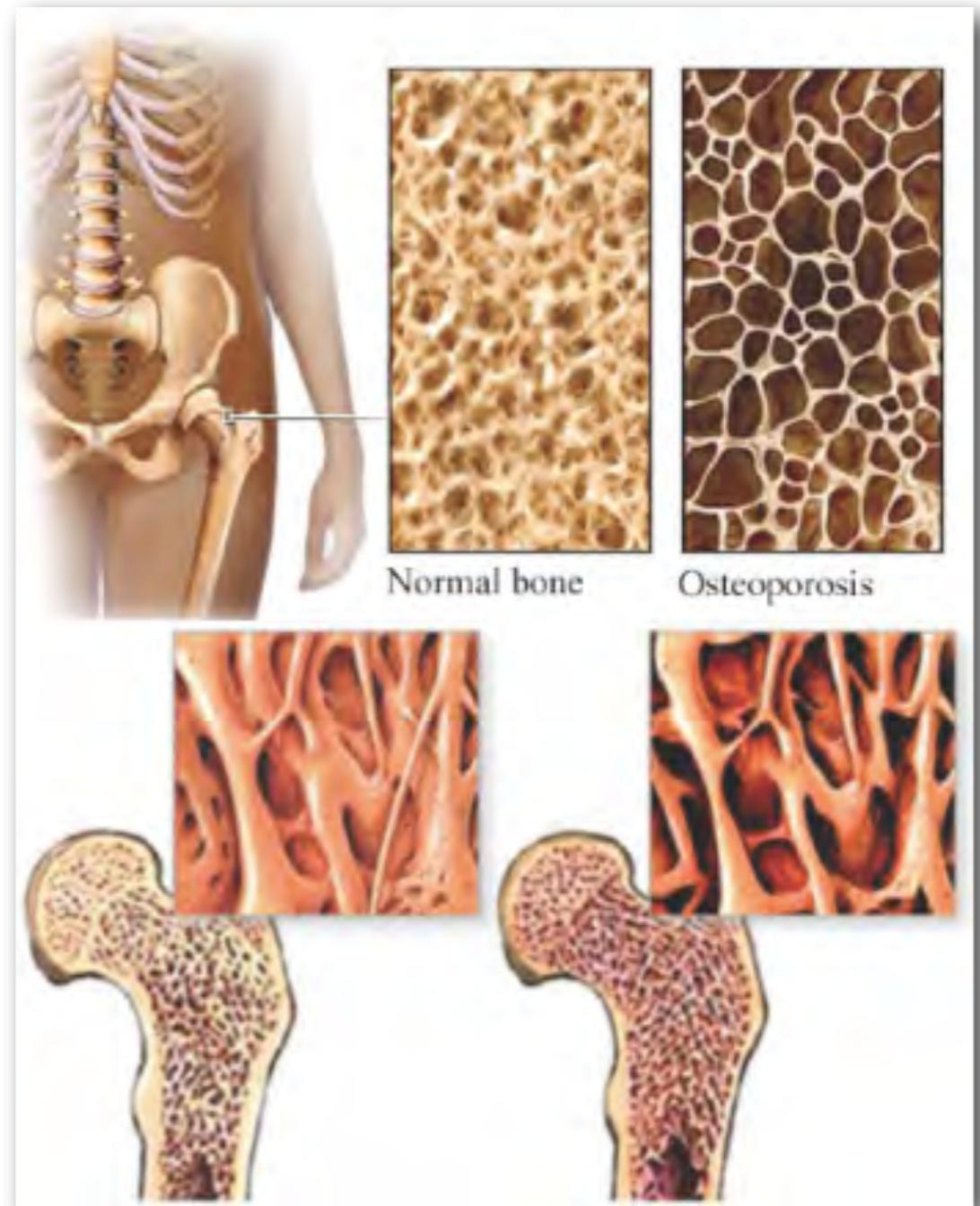
**nature**  
**REVIEWS** GENETICS

# Why do species look the way that they do?



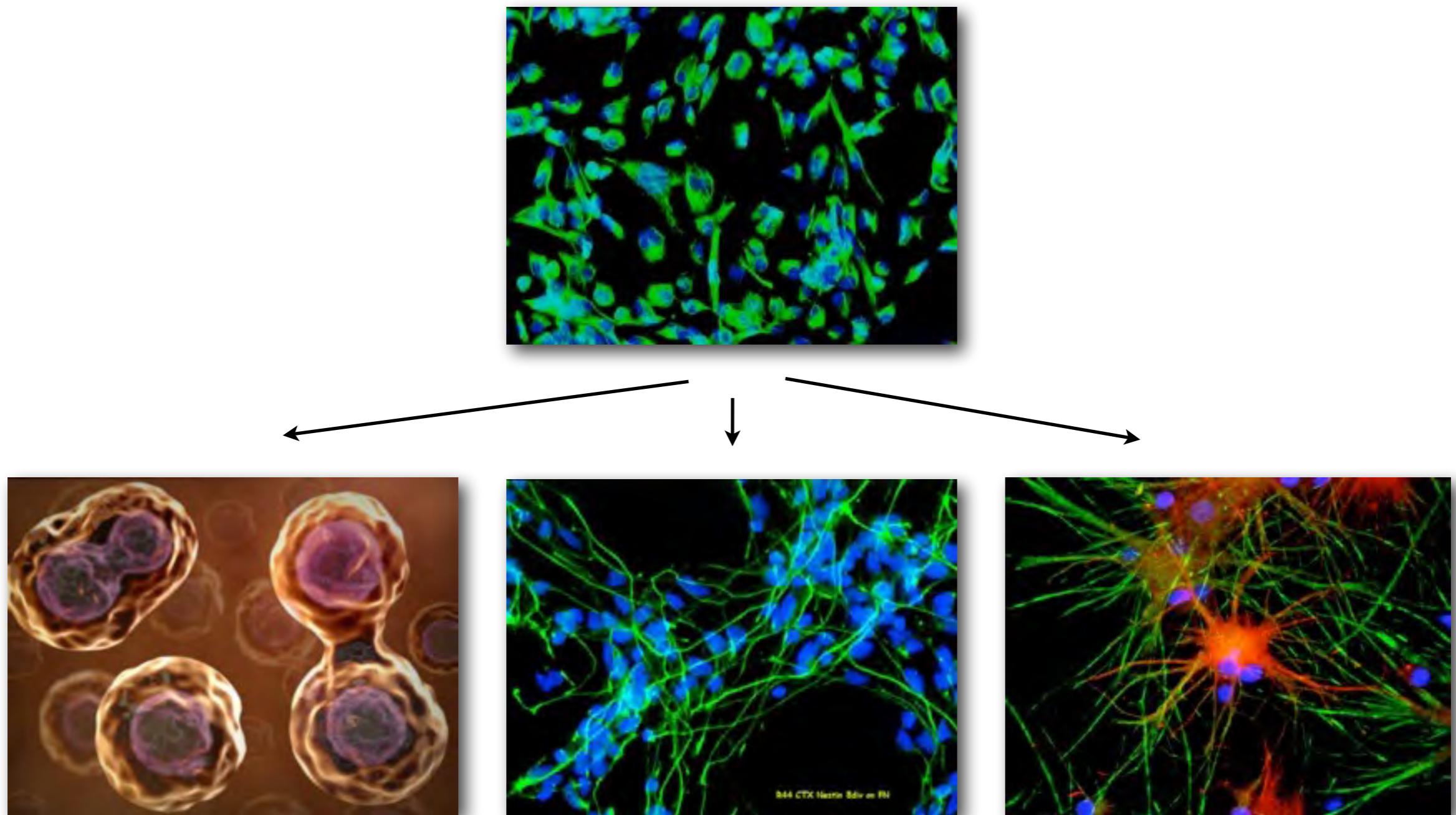
# Why do organisms vary?

---



# How is cellular functional diversity created?

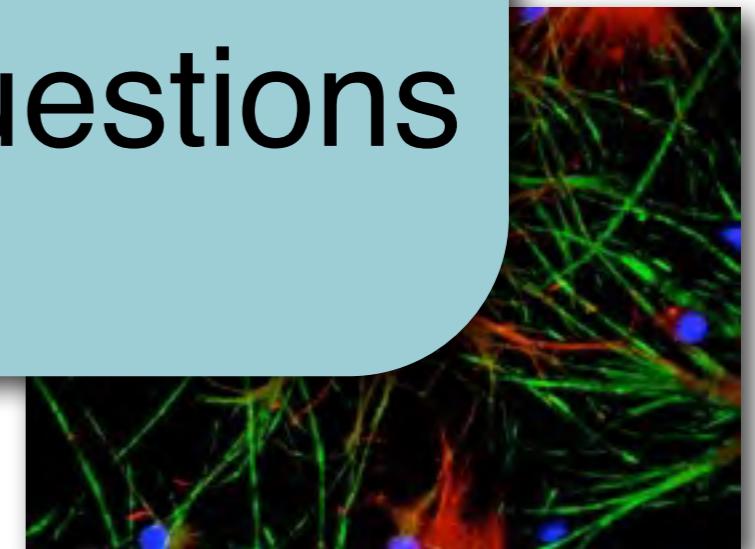
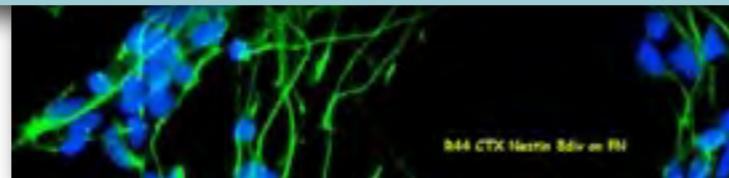
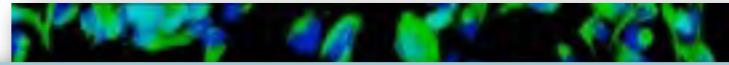
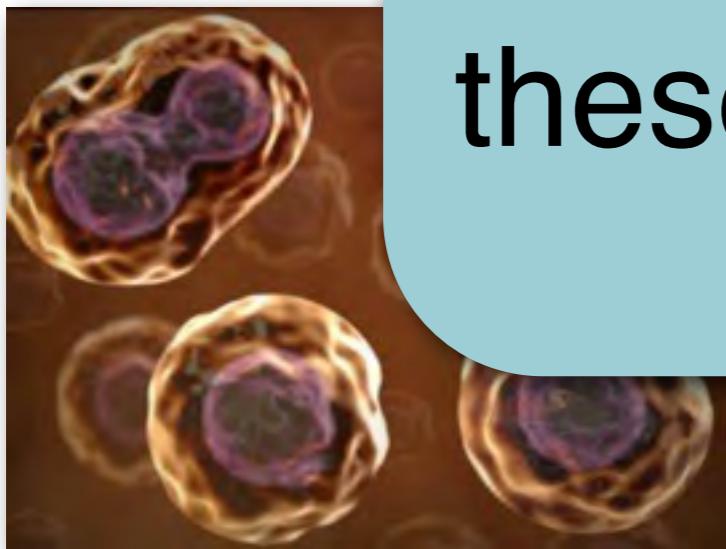
---



# How is cellular functional diversity created?

---

The \*.omics toolkit is revolutionizing our understanding of all of these biological questions

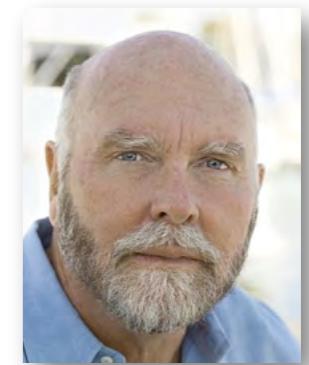




(1998)



(2000)



(2006)



(2006)



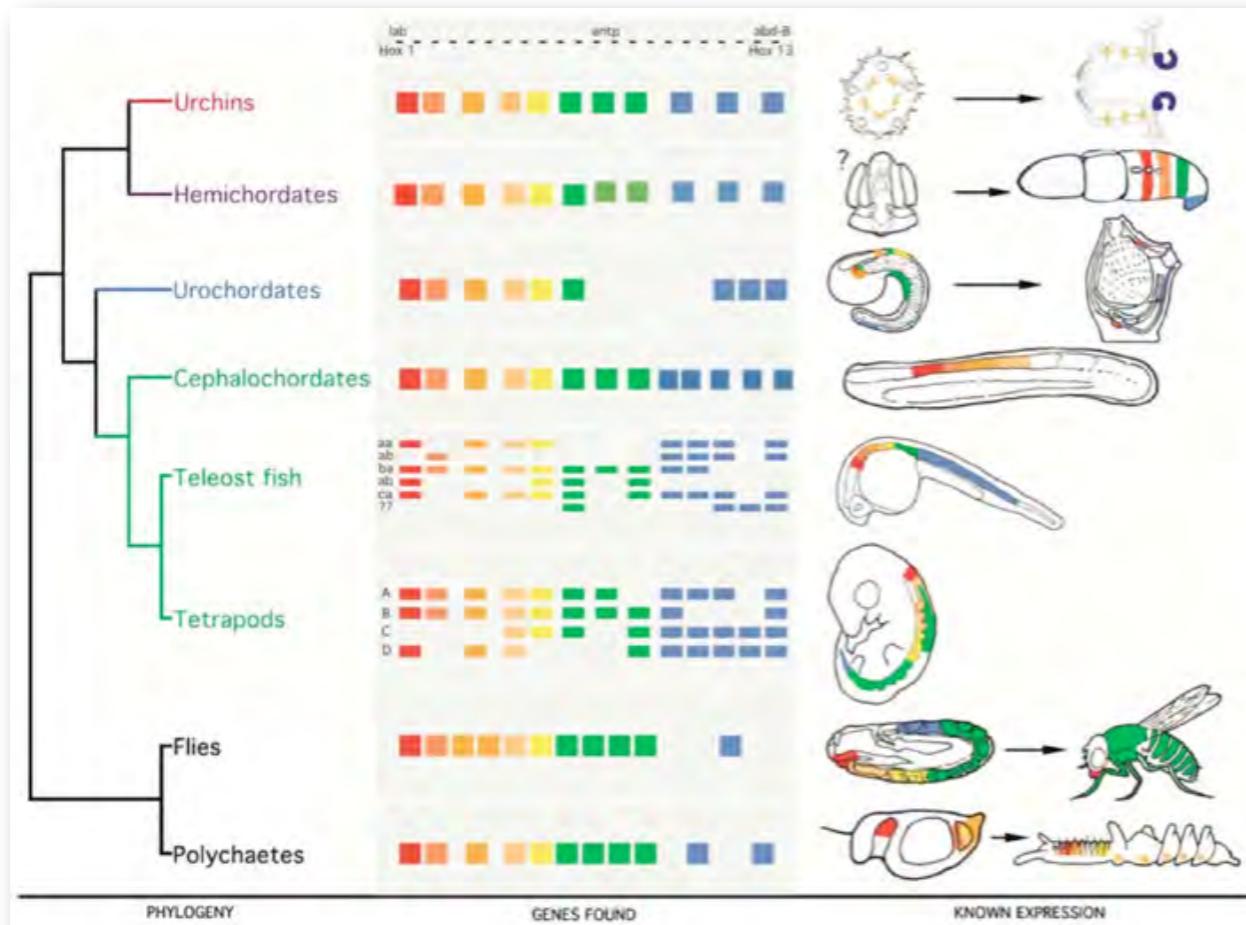
(2002)



(2007)



(2006)



# Comparative Genomics

---

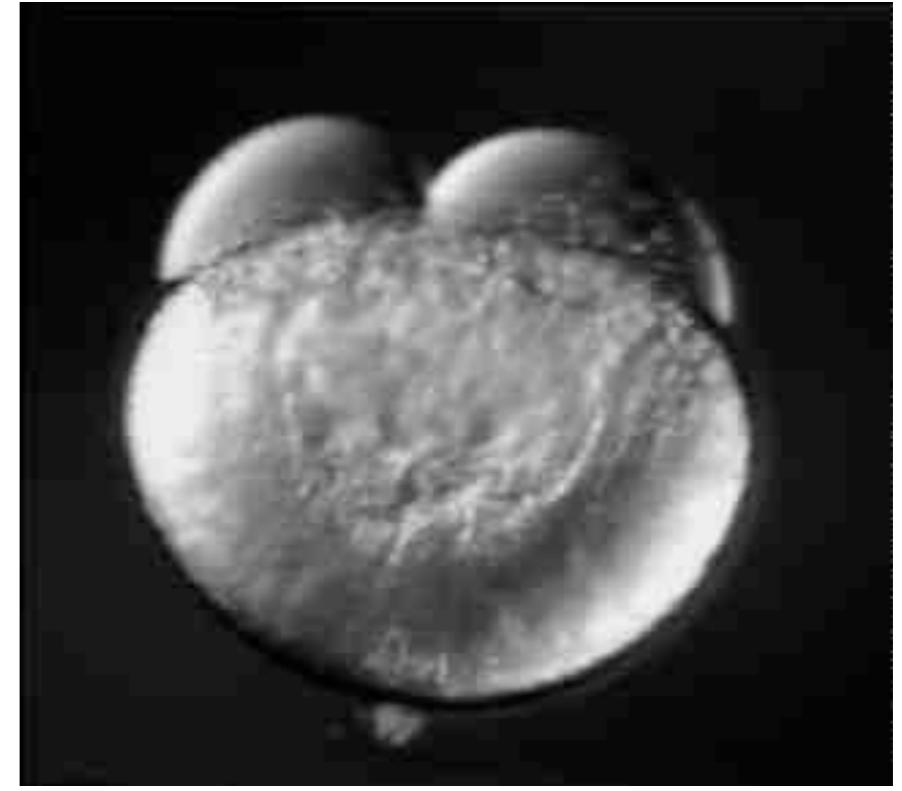
## Vertebrate zygotes or embryos



28 day human



19h zebrafish



---

## Vertebrate zygotes or embryos

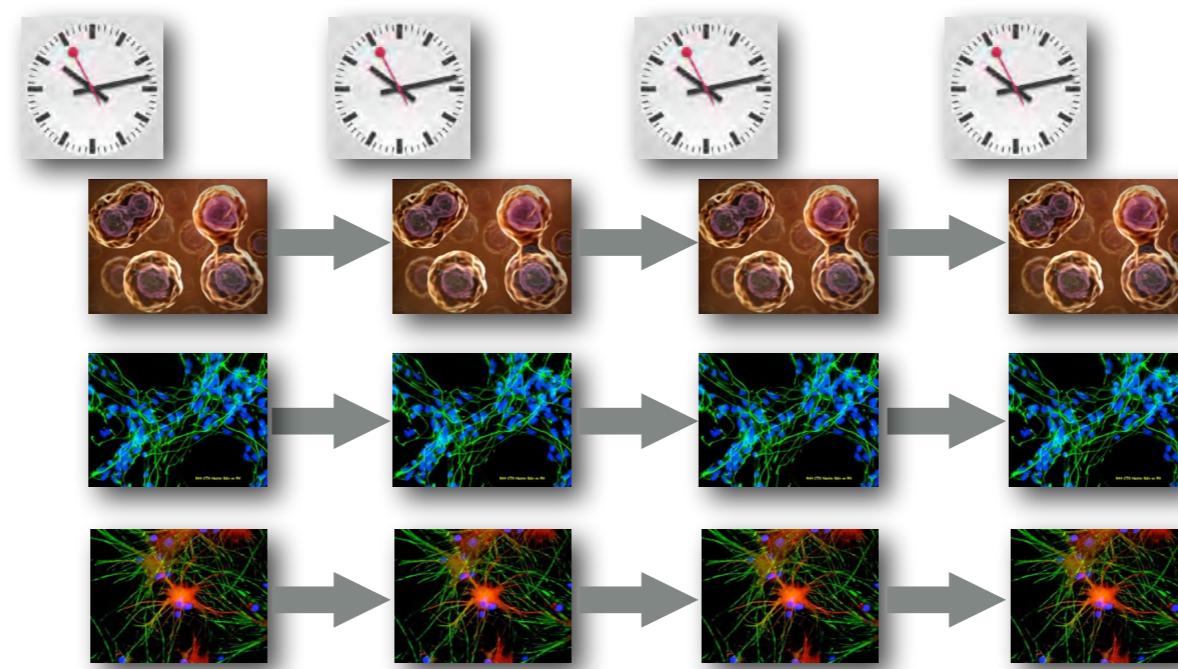
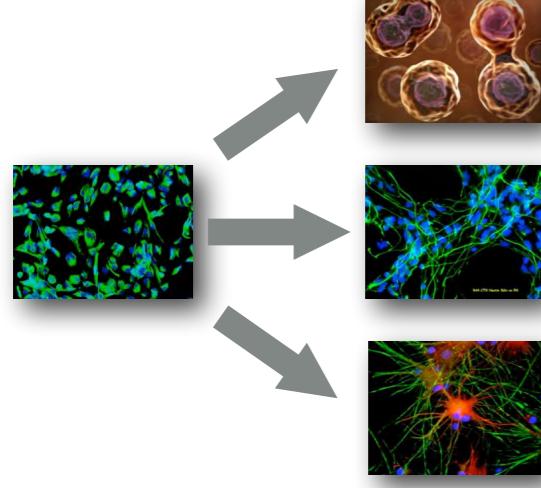


28 day human



19h zebrafish





Functional  
Genomics



# Population Genomics

# How do we ‘genomically enable’ research studies of non-model (and model) organisms?

---

1. Genetic markers & genetic maps
2. Physical maps (genomes)
3. Transcriptomes
4. Gene expression analyses
5. Epigenetic & functional analyses





Assaying genetic variation:  
Shouldn't we just sequence everything?

# Why not just sequence entire genomes??

---

- Still prohibitively expensive for many studies
- For many studies a full sequence isn't necessary
  - genomes of many organisms are organized in linkage blocks
  - well spaced markers will provide the necessary coverage
- Genetic maps are very useful in genomic studies
  - a high density genetic map can facilitate genome assembly
  - genomes may be segregating structural variation

# Alternative - Reduced representation sequencing

---

- Use restriction enzyme digestion to focus sequencing of multiple samples on homologous regions across the genome
- Simultaneous identification and typing of single nucleotide polymorphisms (SNPs) and haplotypes
- The cost is a fraction of the cost of re-sequencing the genome
- Thousands of genomes can be assayed in just a few weeks

# What is RADseq?

(Restriction-site Associated DNA)



2007

Rapid and cost-effective polymorphism identification and genotyping using restriction site associated DNA (RAD) markers

Michael R. Miller,<sup>1</sup> Joseph P. Dunham,<sup>2</sup> Angel Amores,<sup>3</sup> William A. Cresko,<sup>2</sup> and Eric A. Johnson<sup>1,\*</sup>

<sup>1</sup>Institute for Molecular Biology, University of Oregon, Eugene, Oregon 97403, USA, <sup>2</sup>Center for Ecology & Evolutionary Biology, University of Oregon, Eugene, Oregon 97403, USA, <sup>3</sup>Institute of Neuroscience, University of Oregon, Eugene, Oregon 97403, USA

2008

OPEN ACCESS Freely available online

PLOS ONE

Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers

Nathan A. Baird<sup>1,2</sup>, Paul D. Etter<sup>1,\*</sup>, Tressa S. Atwood<sup>2</sup>, Mark C. Currey<sup>3</sup>, Anthony L. Shiver<sup>1</sup>, Zachary A. Lewis<sup>1</sup>, Eric U. Selker<sup>1</sup>, William A. Cresko<sup>2</sup>, Eric A. Johnson<sup>1,2</sup>

<sup>1</sup>Institute of Molecular Biology, University of Oregon, Eugene, Oregon, United States of America, <sup>2</sup>International Maize & Wheat Improvement Center, Mexico, <sup>3</sup>The Center for Ecology and Evolutionary Biology, University of Oregon, Eugene, Oregon, United States of America

# What is RADseq?

(Restriction-site Associated DNA)



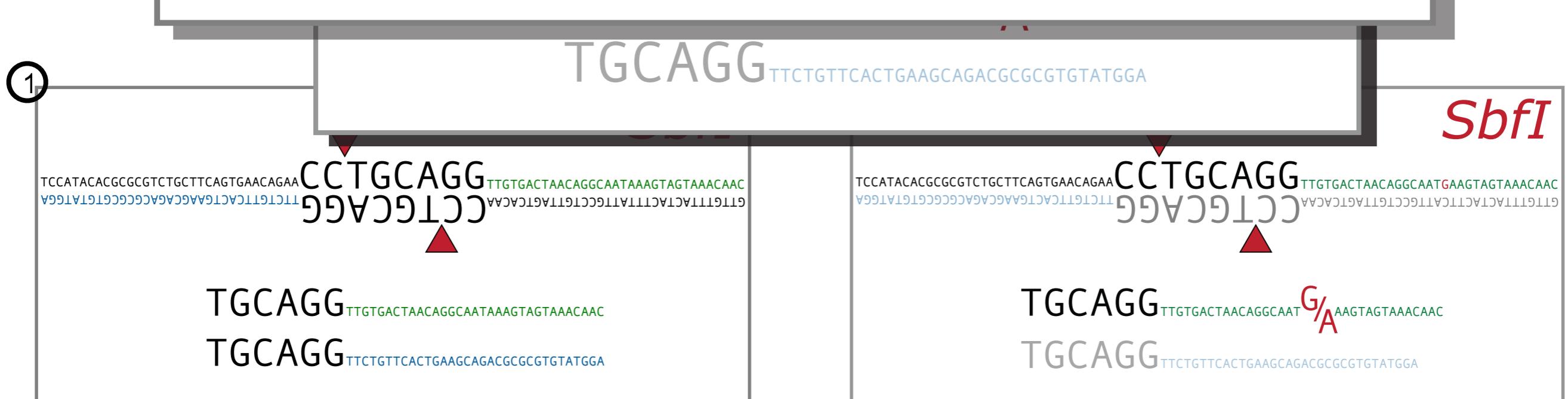
22,830 *SbfI* sites in threespine stickleback genome  
~ 45,000 RAD-Tags

HiSeq2500 Illumina Lane:

160 million reads

HiSeq4000 Illumina Lane:

350 million reads





For what types of studies can  
RADseq be useful?

# Identifying genetically distinct individuals and estimating genetic diversity

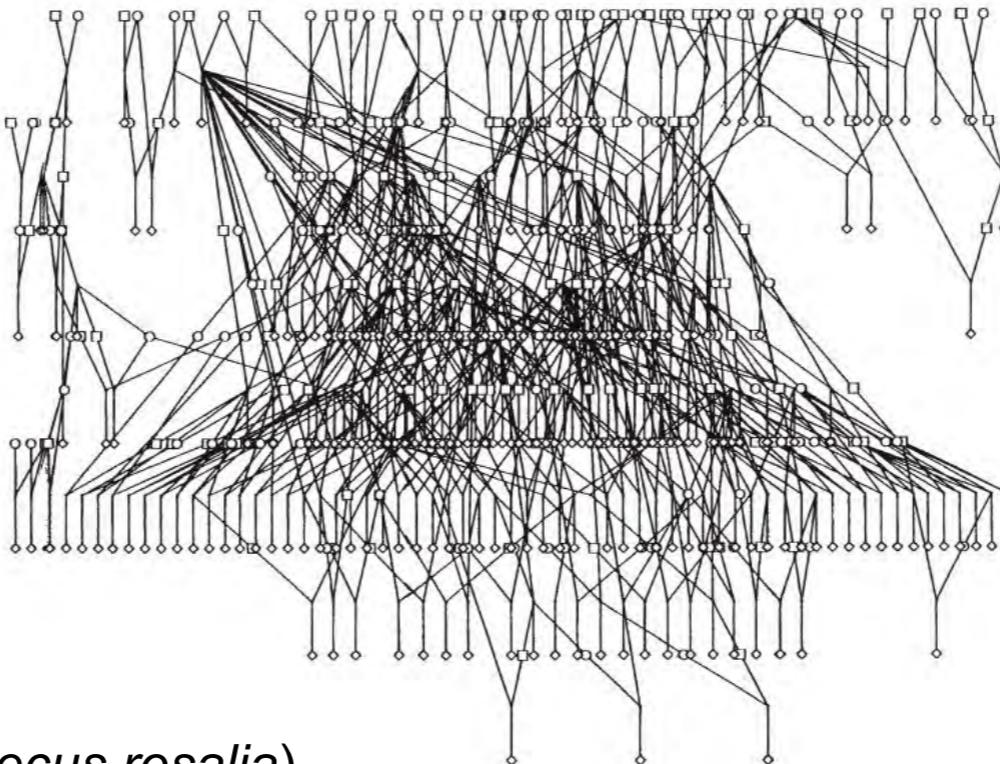
---



Quaking Aspens

# Defining the relationships among individuals and populations

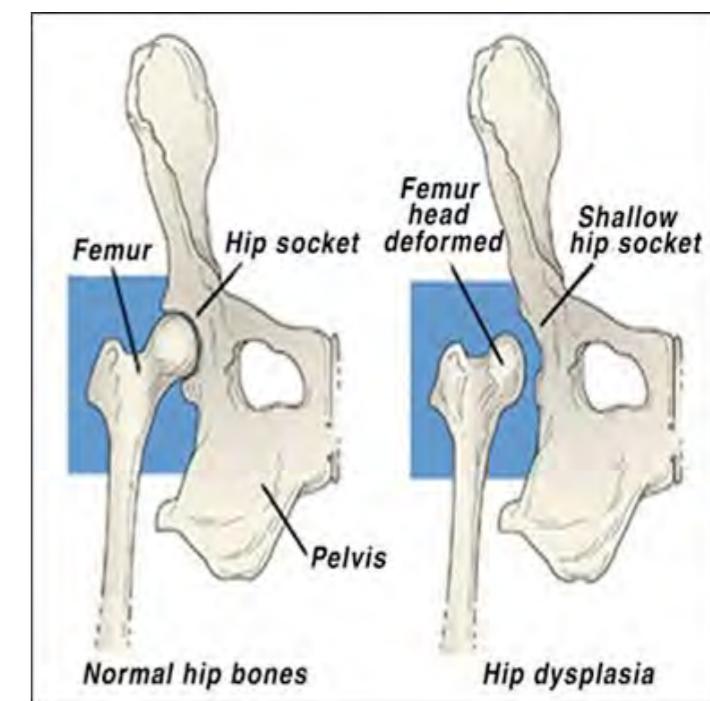
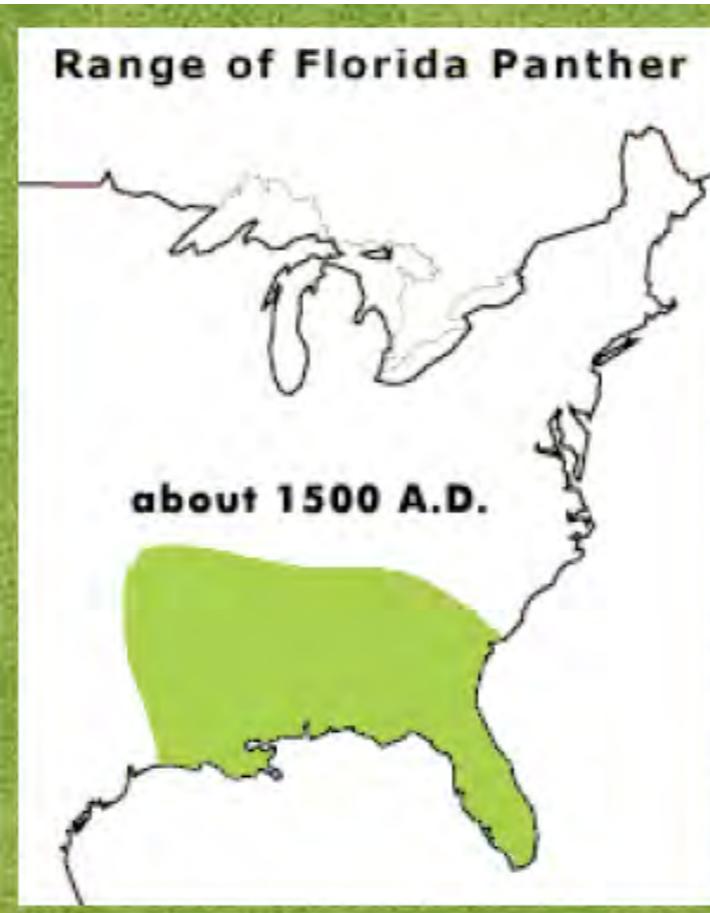
---



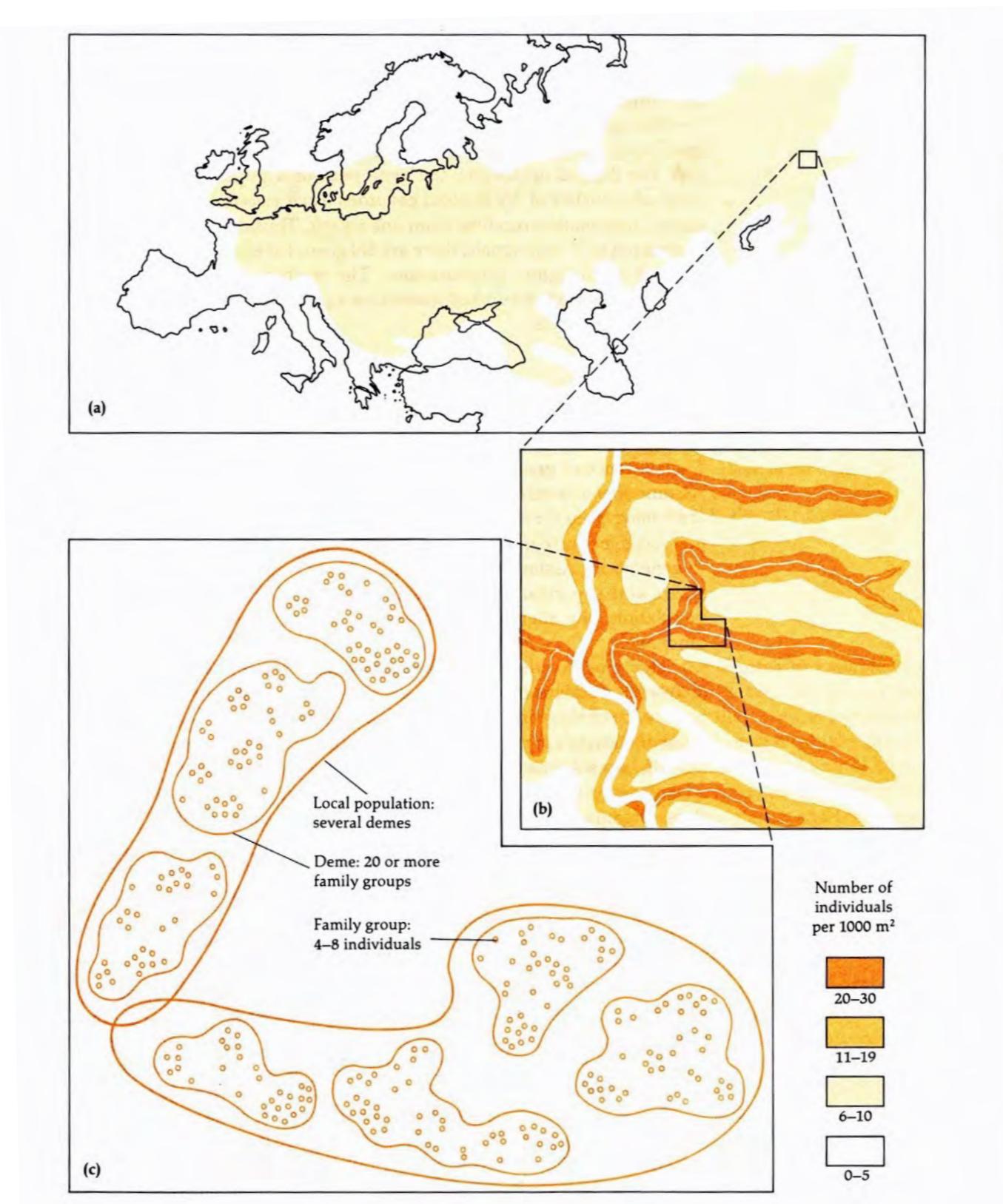
The golden lion tamarin (*Leontopithecus rosalia*).



# Precisely quantifying the amount of inbreeding in wild and captive populations



# Defining the relationships among individuals and populations



# Phylogenetic relationships

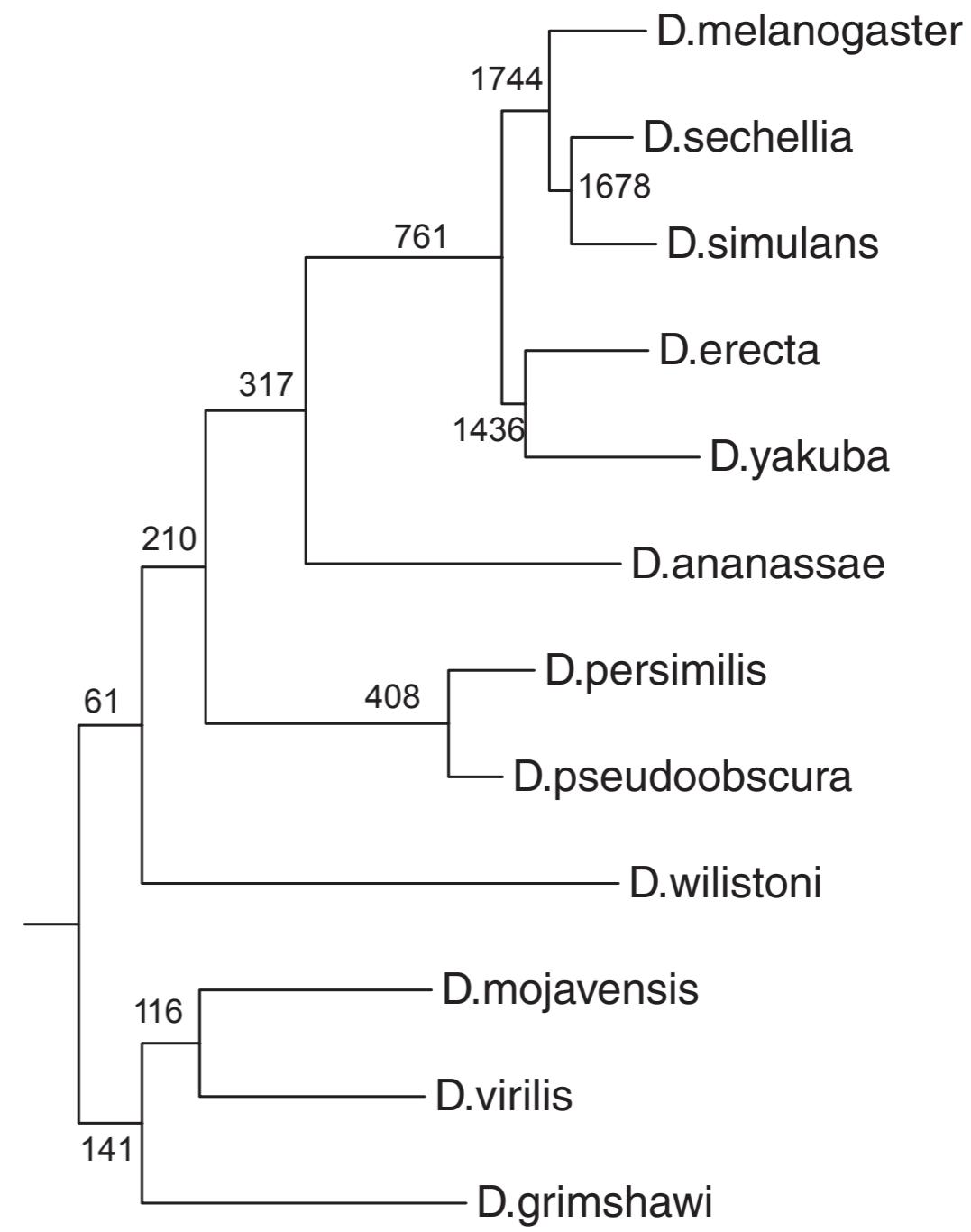
## Ecology and Evolution

Open Access

### Is RAD-seq suitable for phylogenetic inference? An *in silico* assessment and optimization

Marie Cariou, Laurent Duret & Sylvain Charlat

Université de Lyon, Université Lyon 1, CNRS, UMR 5558, Laboratoire de Biométrie et Biologie Evolutive, 43 bou  
Villeurbanne F-69622, France



# Phylogenetic relationships

---

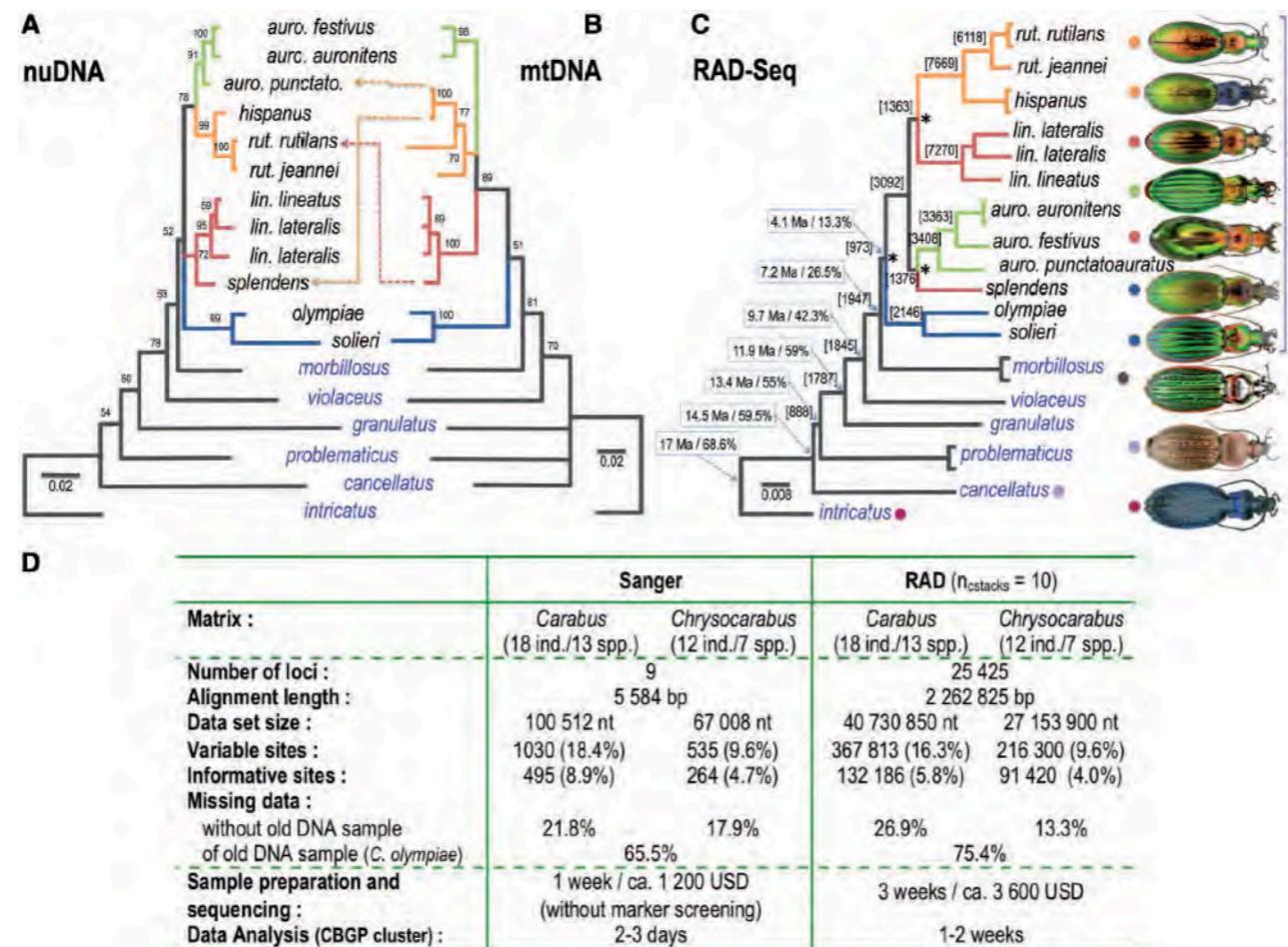
Species pair <i>D. melanogaster</i>	Node depth (My)	Orthologous tags	Retrieved orthologous tags (%)	In clusters including paralogs (%)
<i>D. sechellia</i>	5.4	2978	99	5
<i>D. simulans</i>	5.4	2892	99	4
<i>D. erecta</i>	12.6	2390	97	3
<i>D. yakuba</i>	12.8	2314	97	8
<i>D. ananassae</i>	44.2	916	68	9
<i>D. persimilis</i>	54.9	648	65	9
<i>D. pseudoobscura</i>	54.9	648	66	9
<i>D. wilistoni</i>	62.2	242	49	6
<i>D. grimshawi</i>	62.9	290	60	8
<i>D. virilis</i>	62.9	286	59	5
<i>D. mojavensis</i>	62.9	298	59	8

---

# Empirical Assessment of RAD Sequencing for Interspecific Phylogeny

Mol. Biol. Evol. 31(5):1272–1274

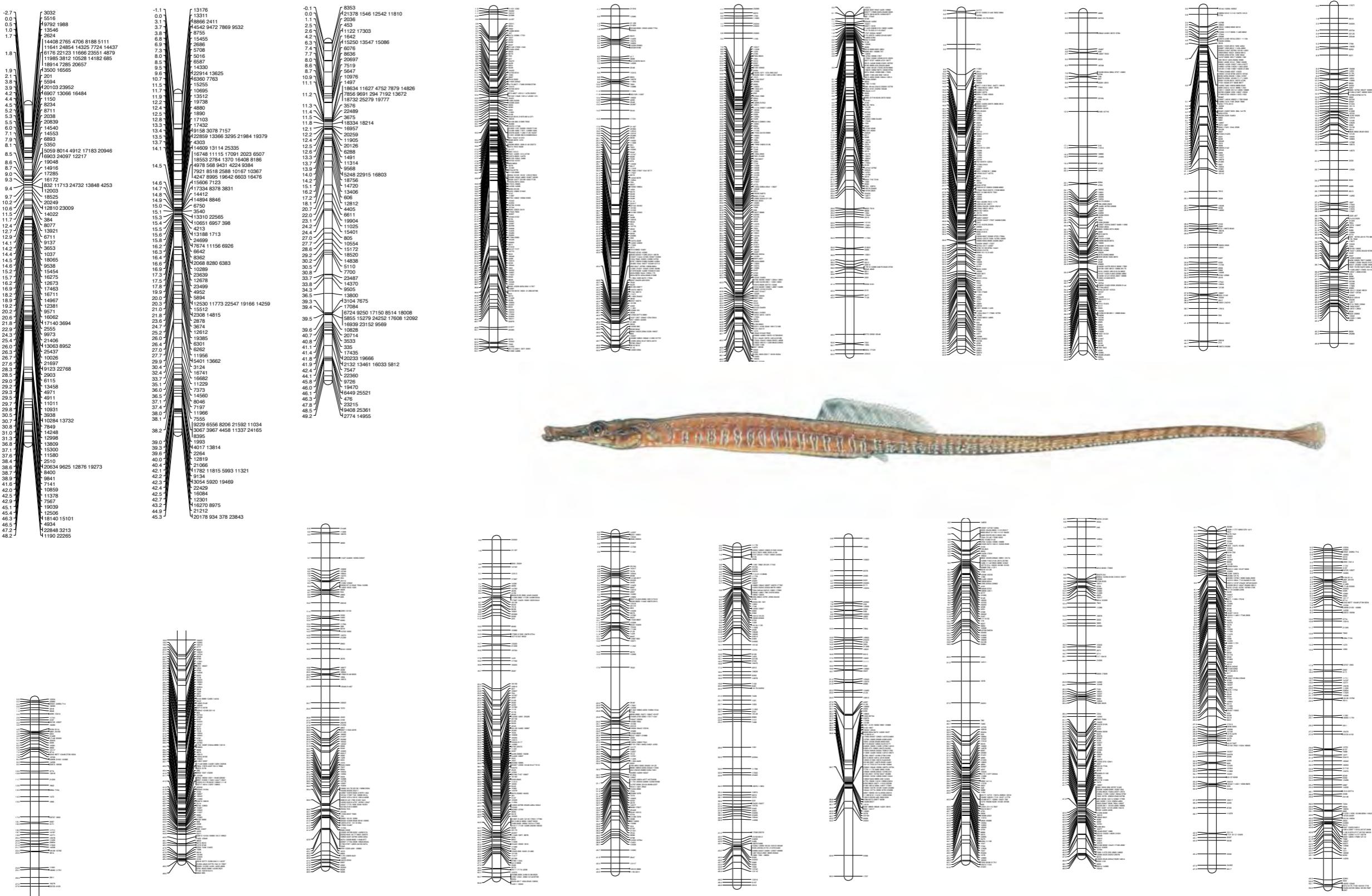
Astrid Cruaud,<sup>\*†,1</sup> Mathieu Gautier,<sup>†,1,2</sup> Maxime Galan,<sup>1</sup> Julien Foucaud,<sup>1</sup> Laure Sauné,<sup>1</sup> Gwenaëlle Genson,<sup>1</sup> Emeric Dubois,<sup>3</sup> Sabine Nidelet,<sup>3</sup> Thierry Deuve,<sup>4</sup> and Jean-Yves Rasplus<sup>1</sup>



## Abstract

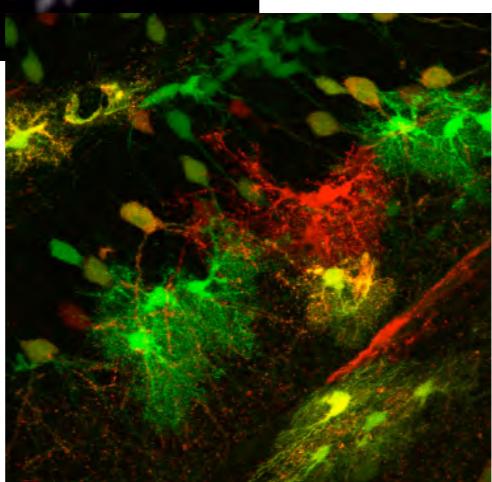
Next-generation sequencing opened up new possibilities in phylogenetics; however, choosing an appropriate method of sample preparation remains challenging. Here, we demonstrate that restriction-site-associated DNA sequencing (RAD-seq) generates useful data for phylogenomics. Analysis of our RAD library using current bioinformatic and phylogenetic tools produced 400× more sites than our Sanger approach (2,262,825 nt/species), fully resolving relationships between 18 species of ground beetles (divergences up to 17 My). This suggests that RAD-seq is promising to infer phylogeny of eukaryotic species, though potential biases need to be evaluated and new methodologies developed to take full advantage of such data.

# Improve genome assemblies with genetic maps

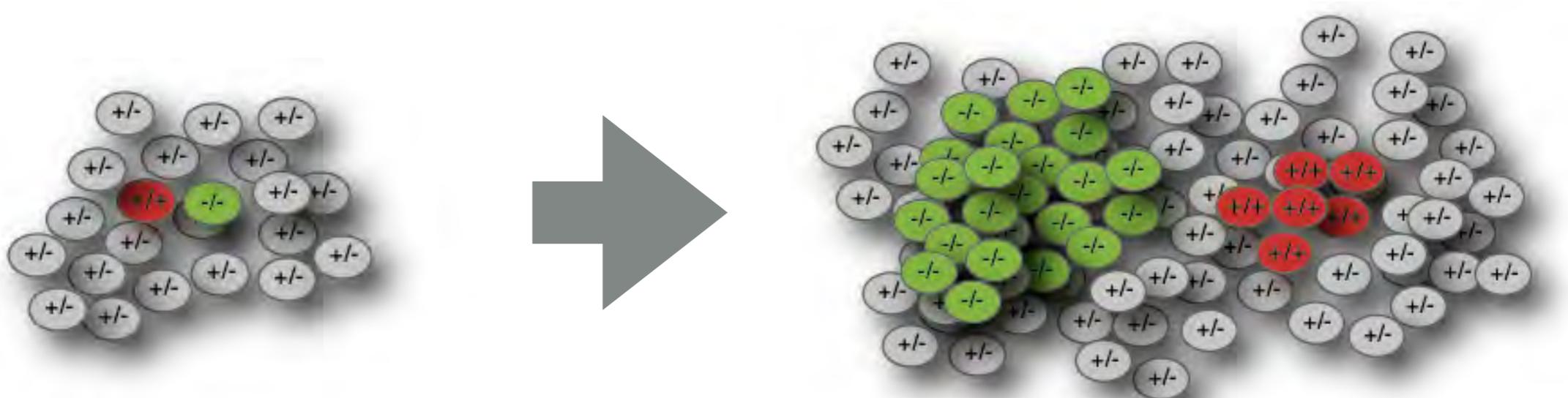
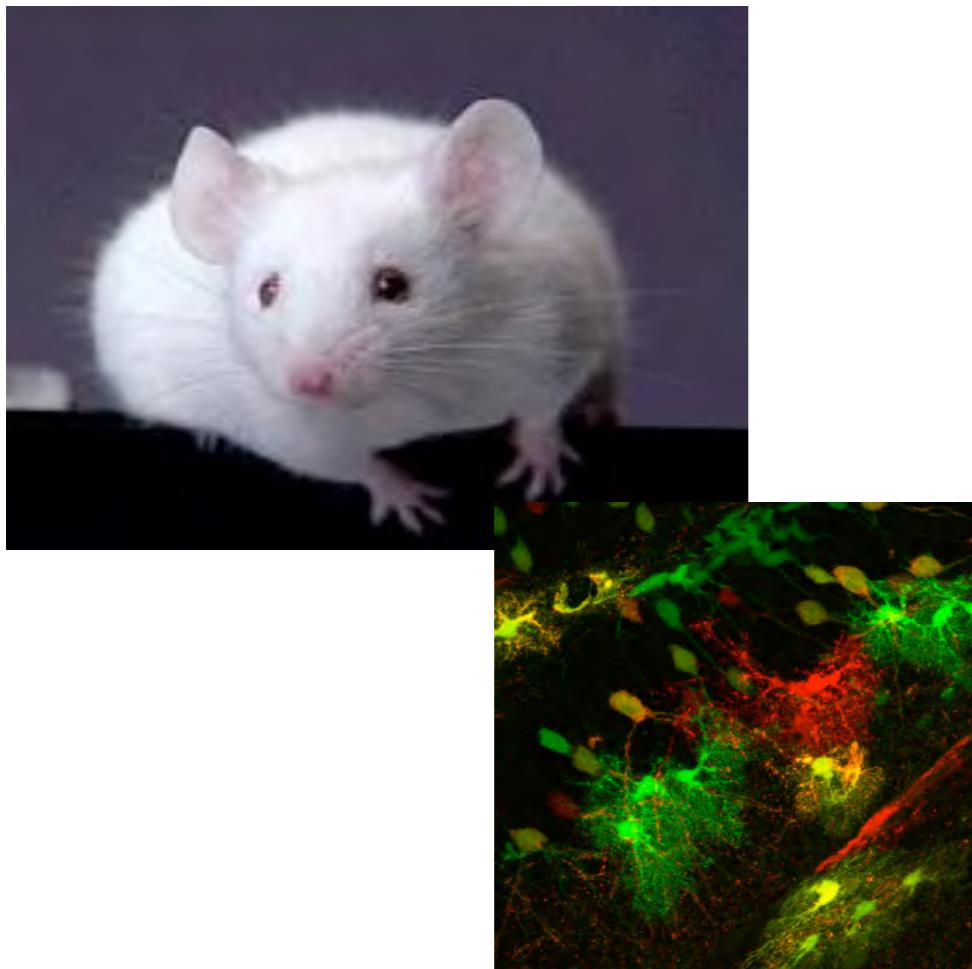


# Studying cancer as an evolutionary process

---



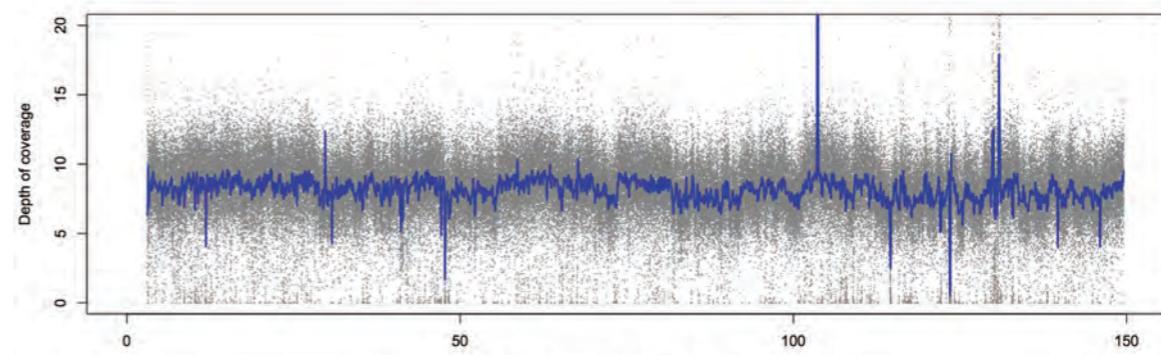
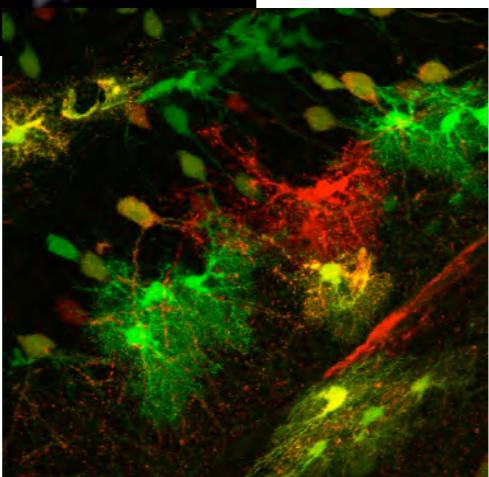
# Studying cancer as an evolutionary process



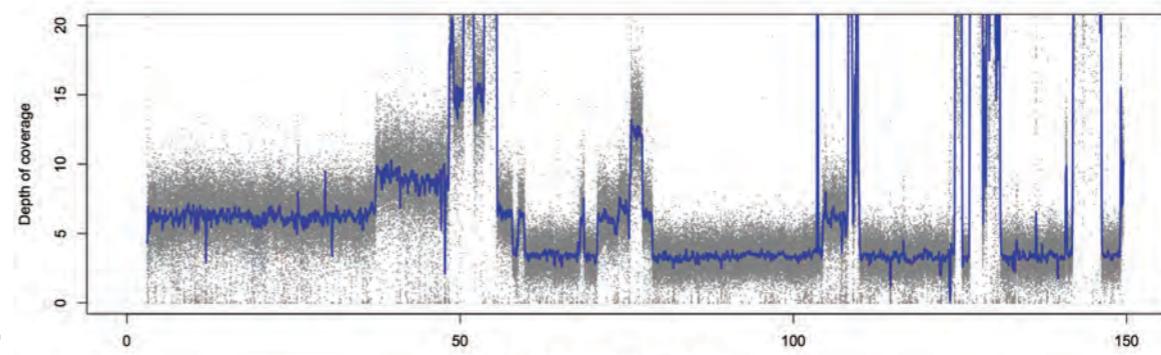
# Studying cancer as an evolutionary process



Wild  
Type



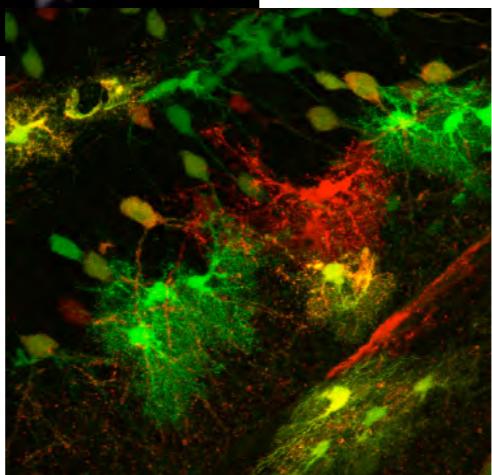
WGS



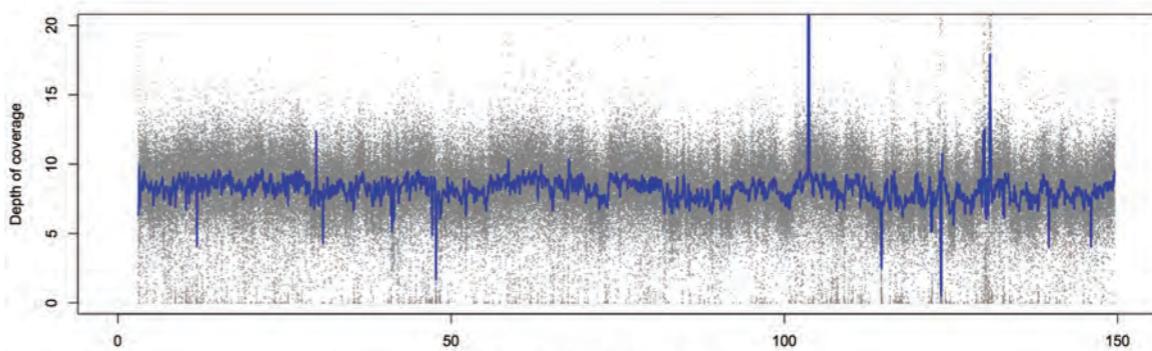
WGS

Mutant

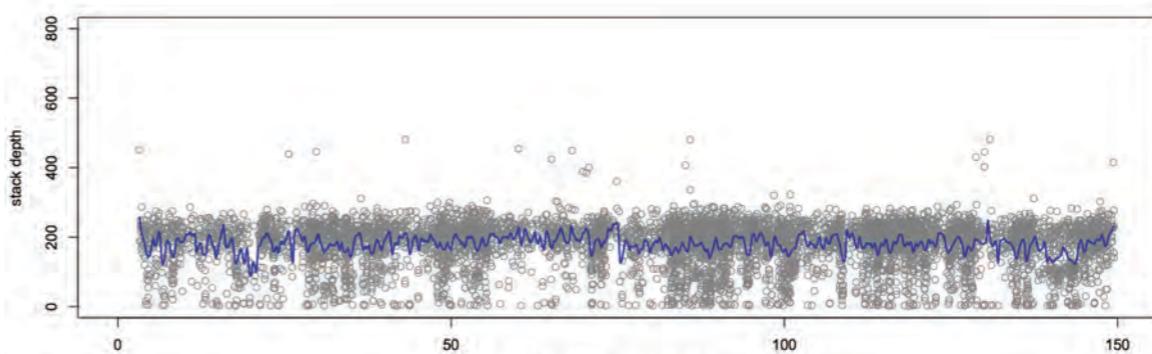
# Studying cancer as an evolutionary process



Wild  
Type

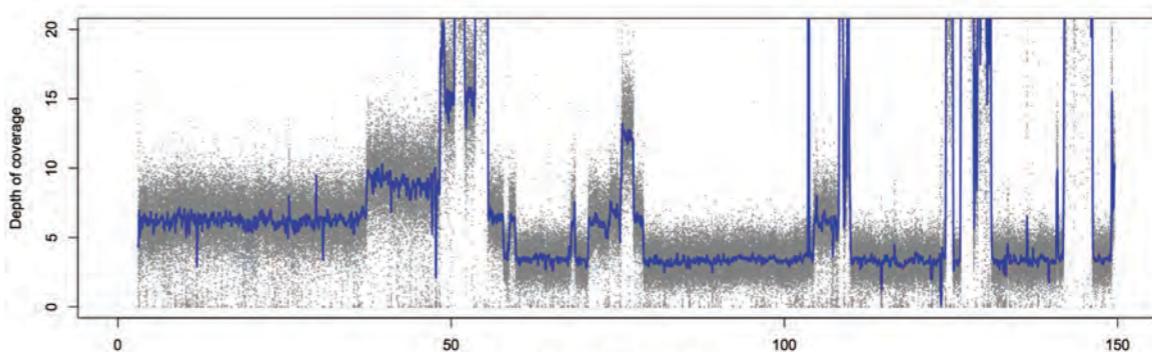


WGS

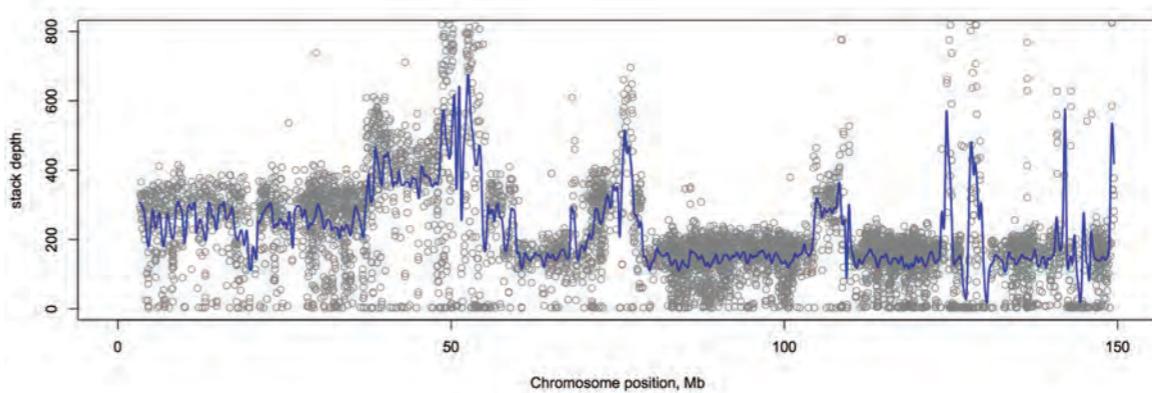


RAD

Mutant



WGS



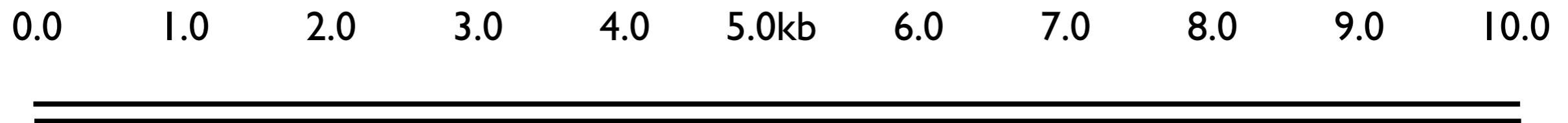
RAD

---

Let's go into a little more detail

# Restriction Enzyme (RE) digestion and first adaptor ligation

---



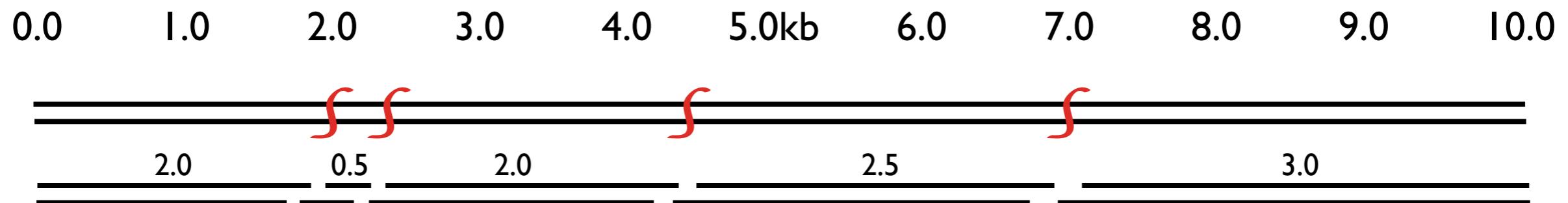
# Restriction Enzyme (RE) digestion and first adaptor ligation

---



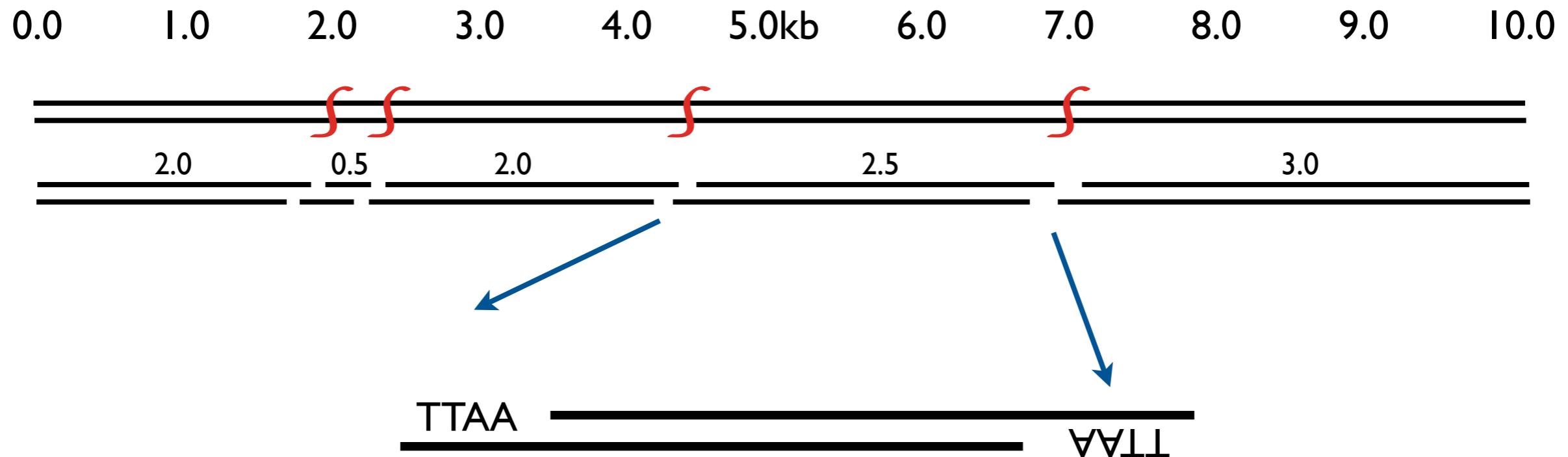
# Restriction Enzyme (RE) digestion and first adaptor ligation

---

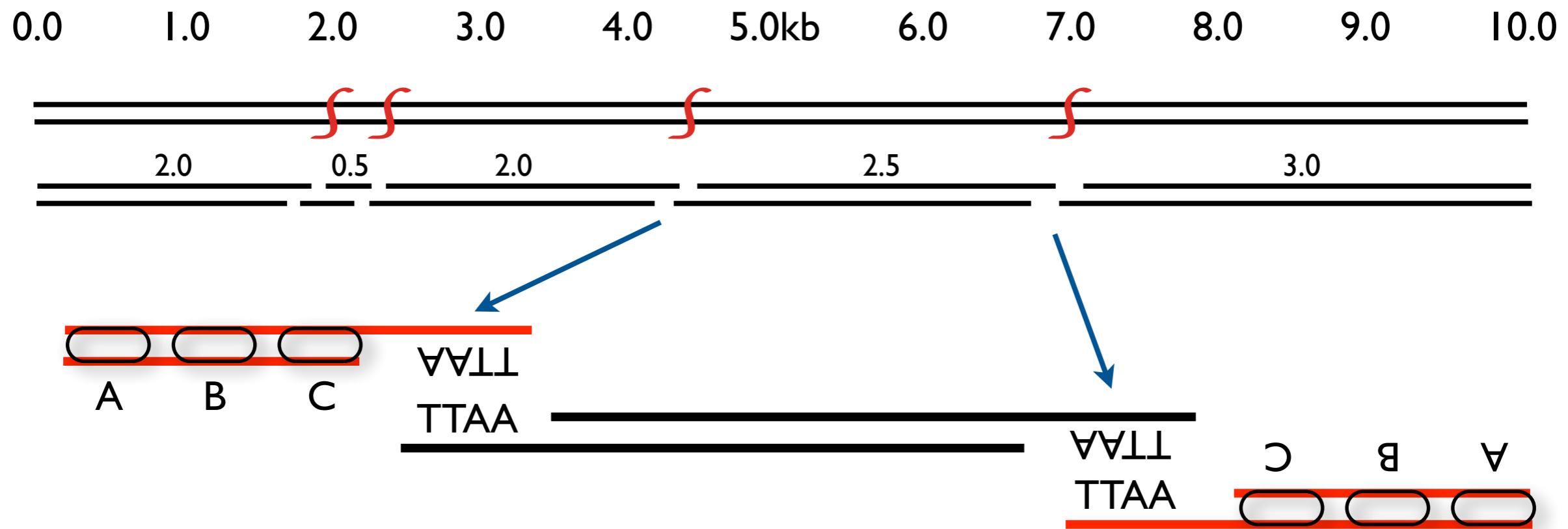


# Restriction Enzyme (RE) digestion and first adaptor ligation

---

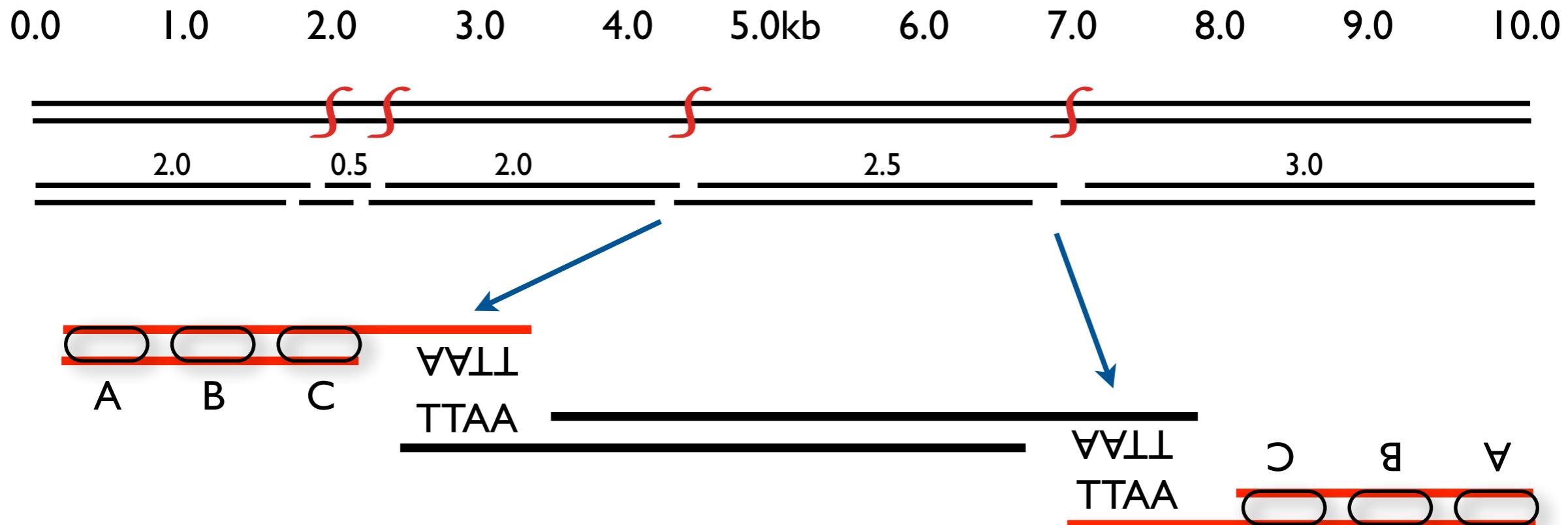


# Restriction Enzyme (RE) digestion and first adaptor ligation



A = Amplification primer  
B = Sequencing primer  
C = Barcode

# Restriction Enzyme (RE) digestion and first adaptor ligation

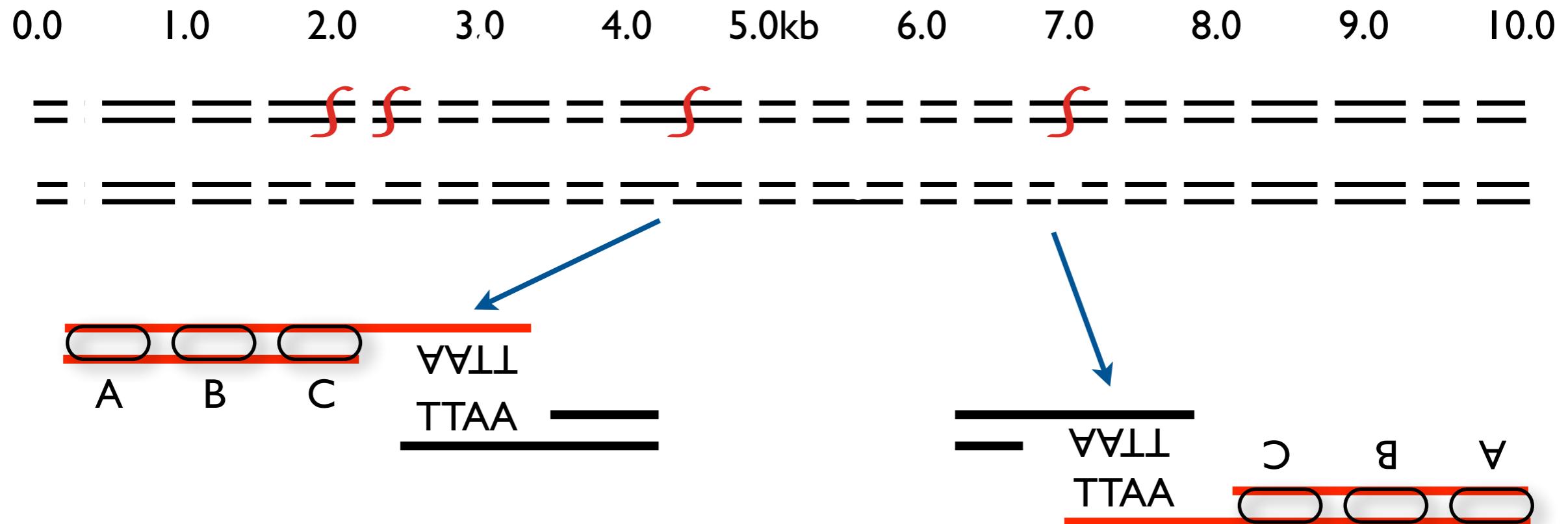


Note - there are now commonly two levels of barcodes used:

Sample Barcodes and  
Molecular Identification Barcodes (MIPs)

A = Amplification primer  
B = Sequencing primer  
C = Barcode

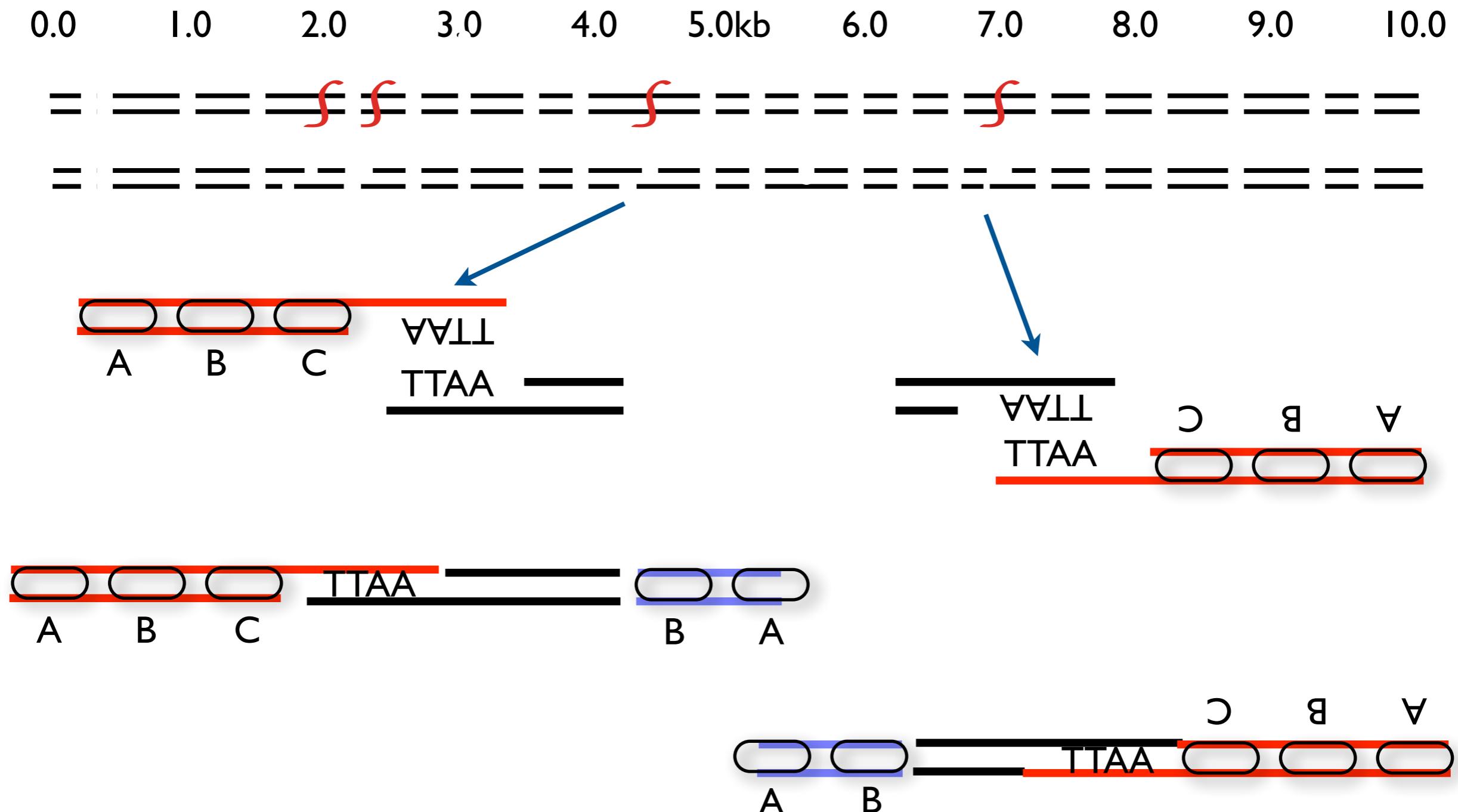
# Shearing and second adaptor ligation



\* Defining difference between original RAD and other approaches\*

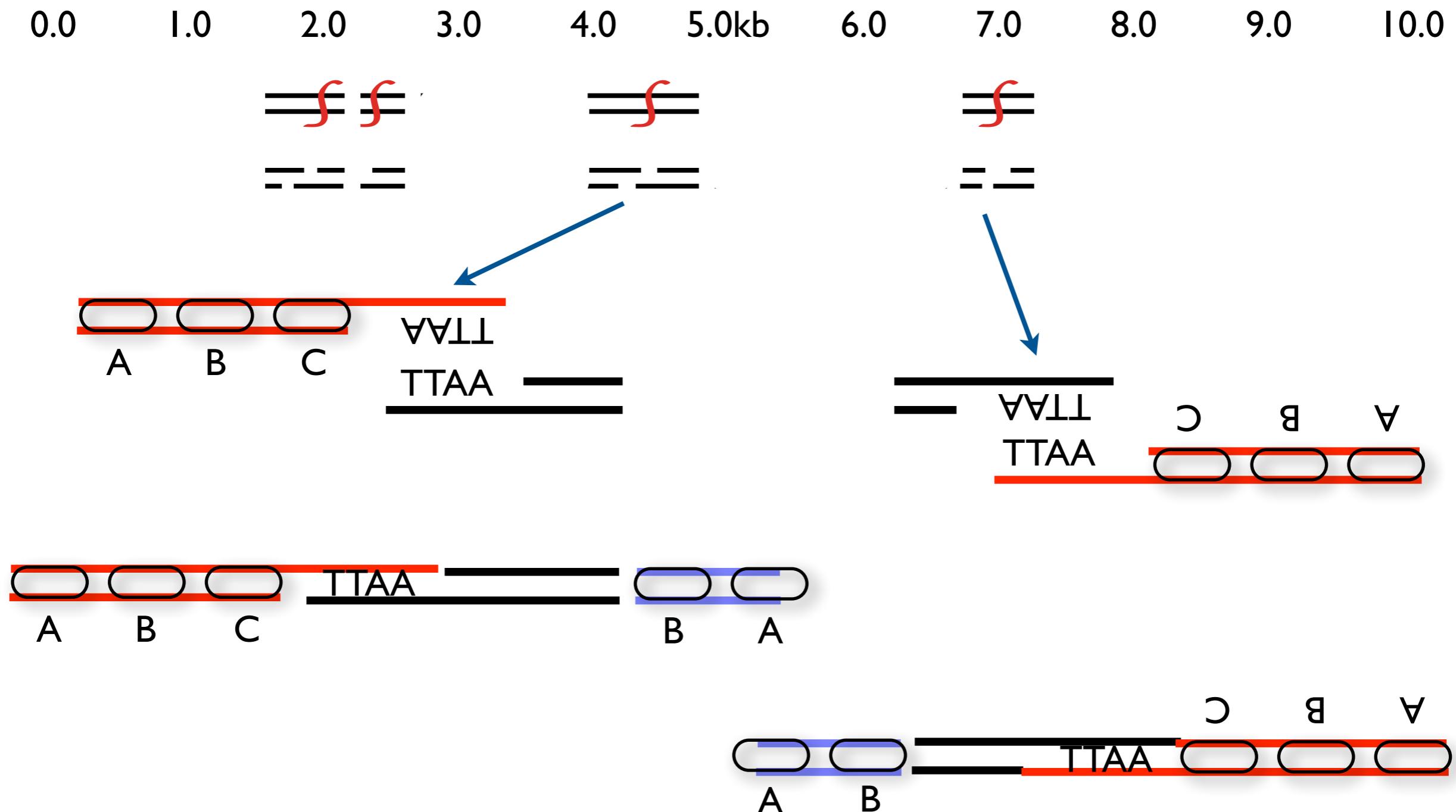
A = Amplification primer  
B = Sequencing primer  
C = Barcode

# Shearing and second adaptor ligation



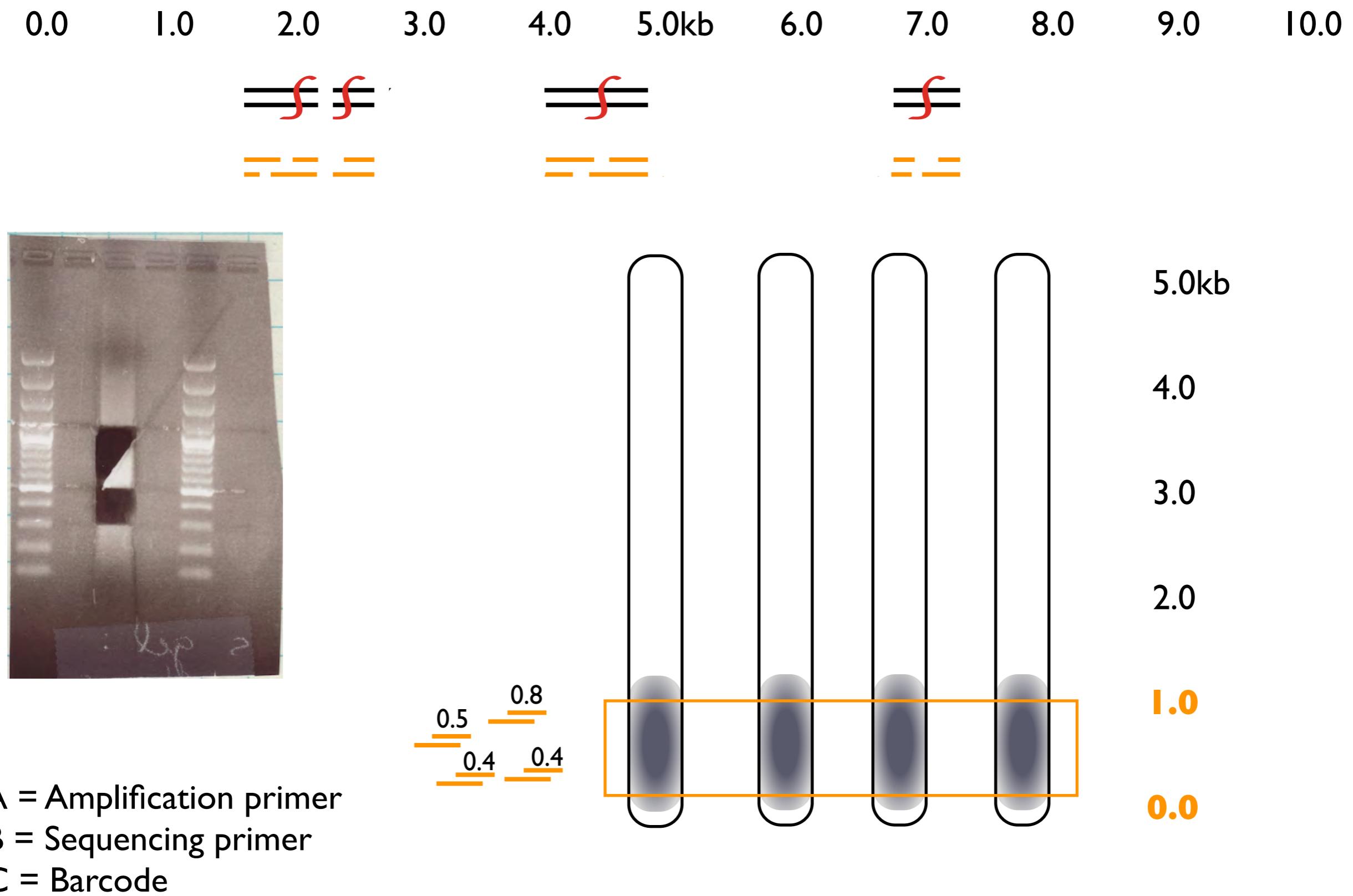
A = Amplification primer  
B = Sequencing primer  
C = Barcode

# Shearing and second adaptor ligation

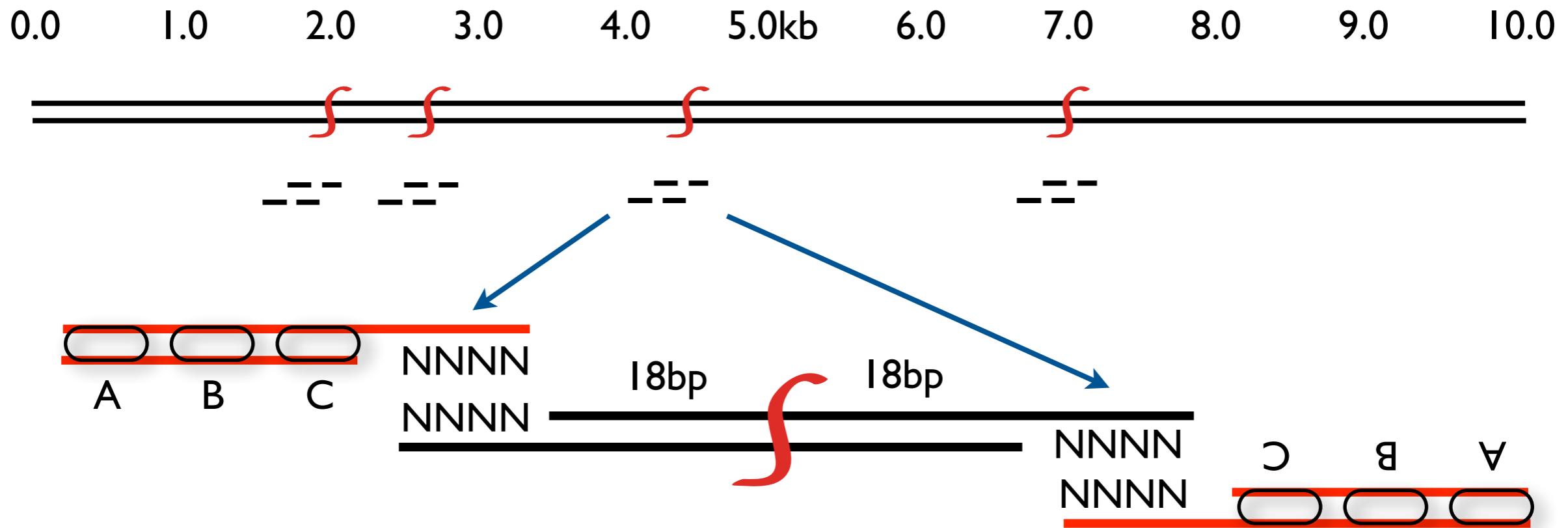


A = Amplification primer  
B = Sequencing primer  
C = Barcode

# Shearing makes consistent fragments for sequencing

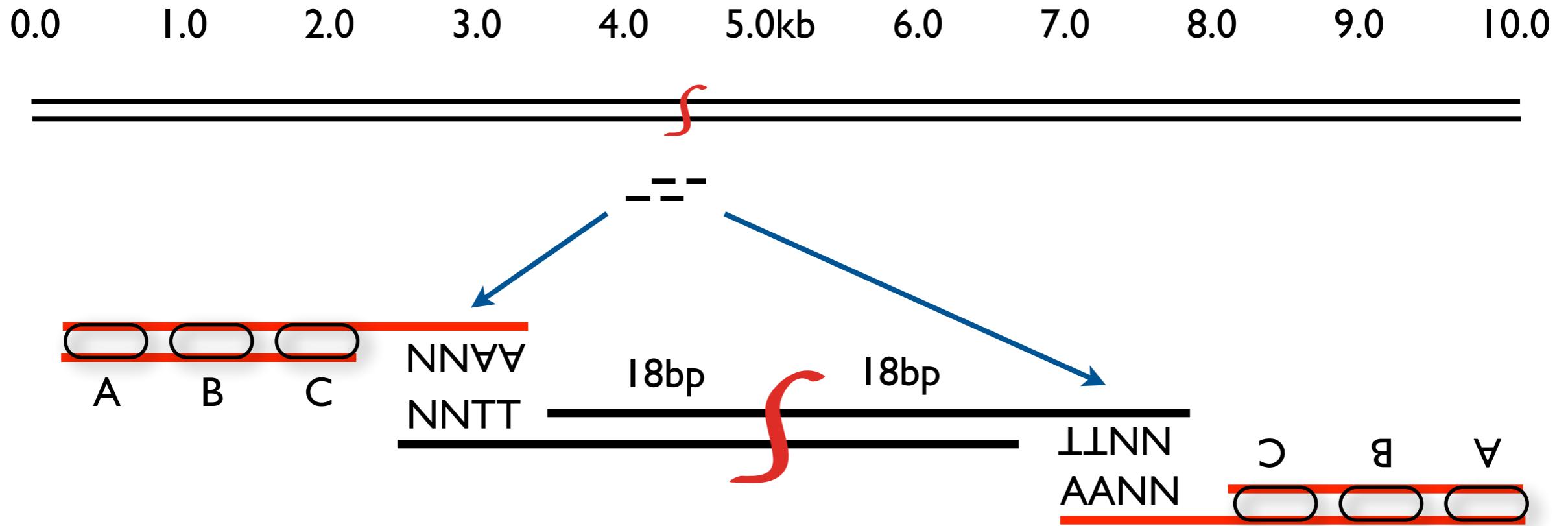


# 2bRAD - type 2b restriction enzyme



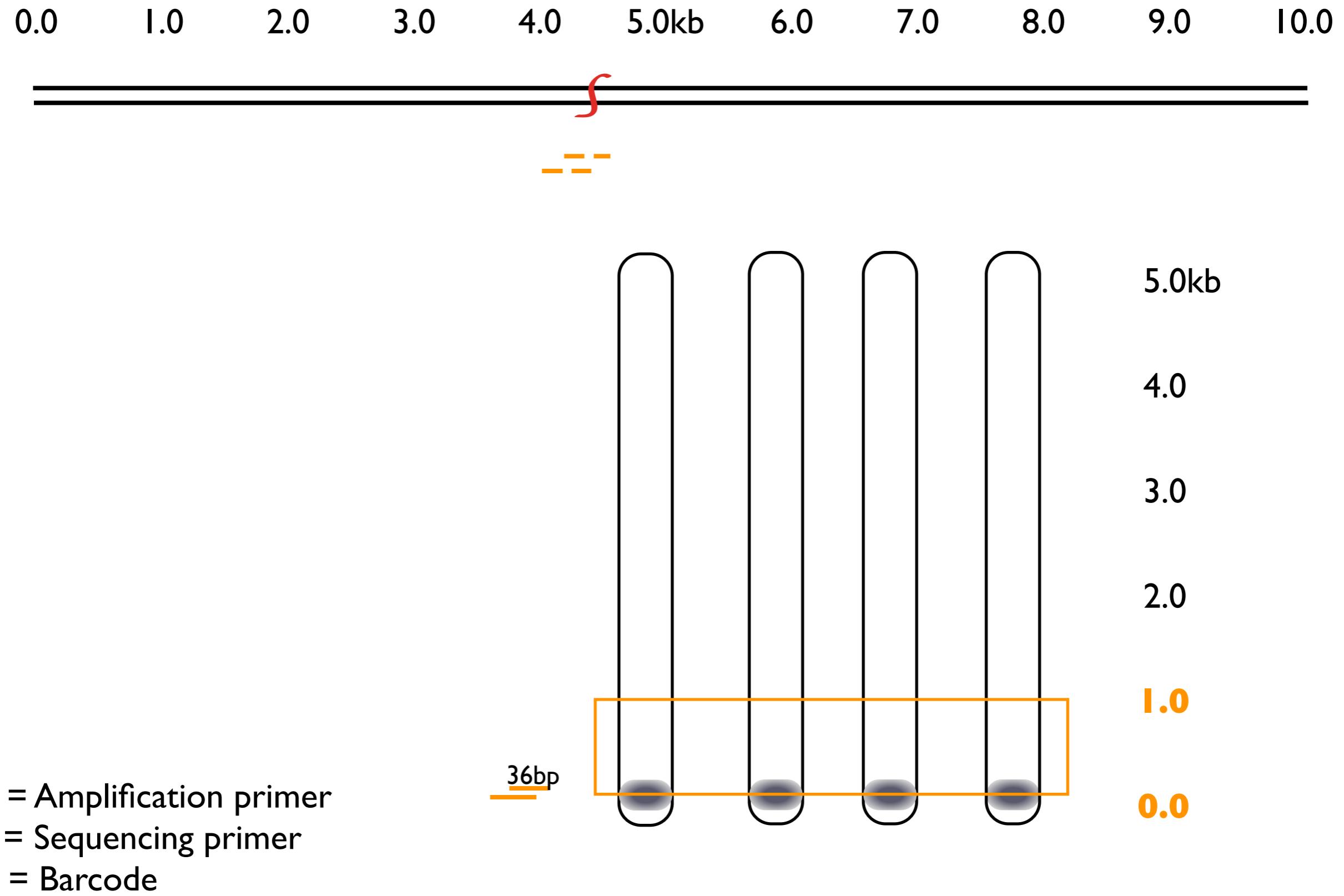
A = Amplification primer  
B = Sequencing primer  
C = Barcode

## 2bRAD - can scale number of markers easily

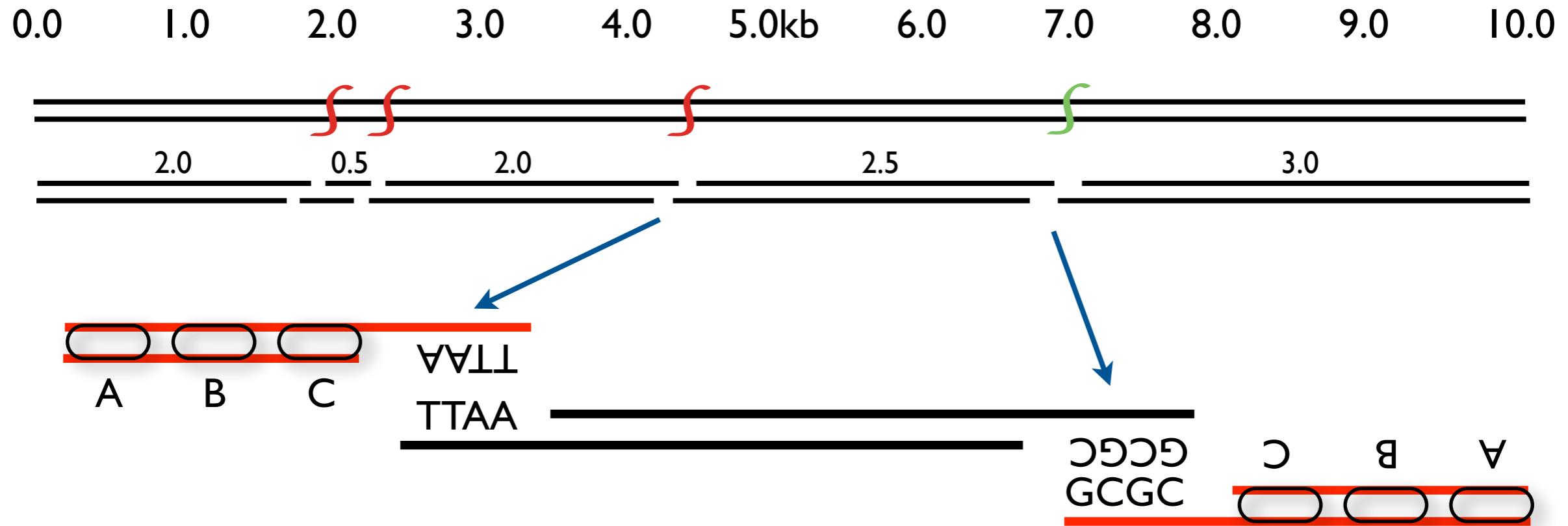


A = Amplification primer  
B = Sequencing primer  
C = Barcode

# 2bRAD - size selection is difficult

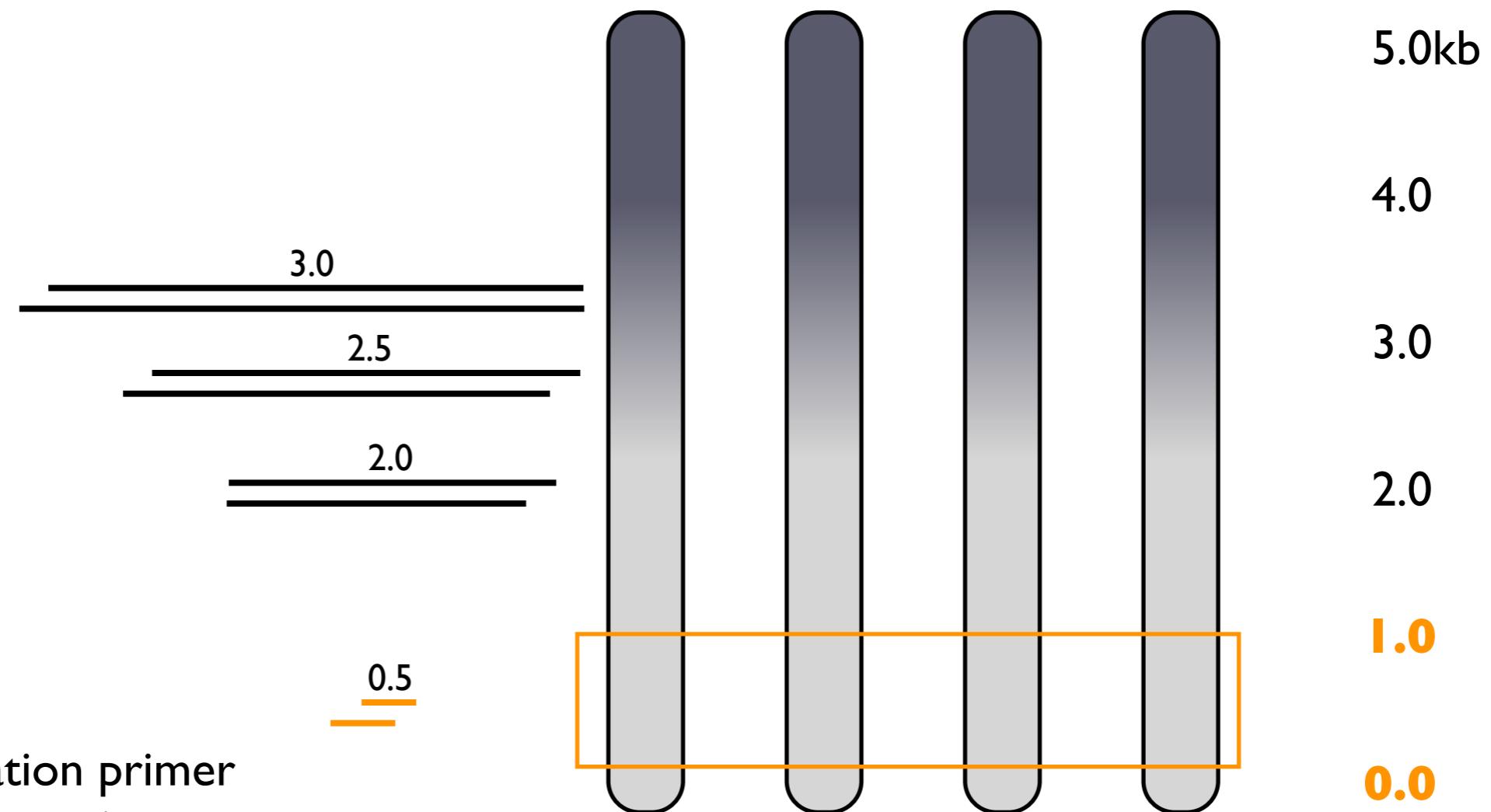
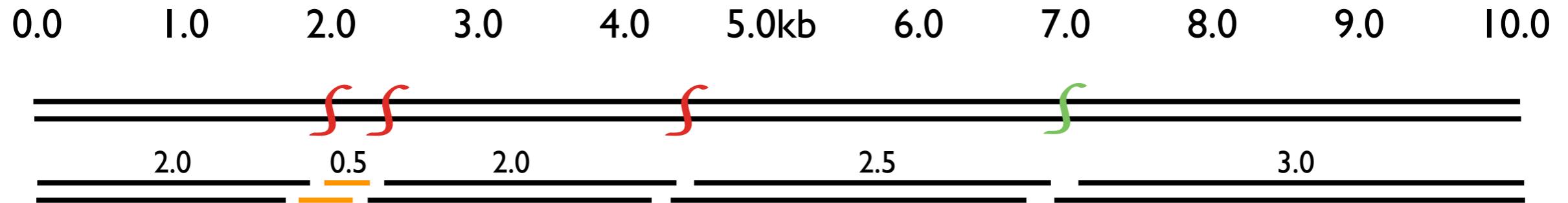


# Single (GBS) or Double Digest RAD (ddRAD)



A = Amplification primer  
B = Sequencing primer  
C = Barcode

# Size selection is more problematic without shearing

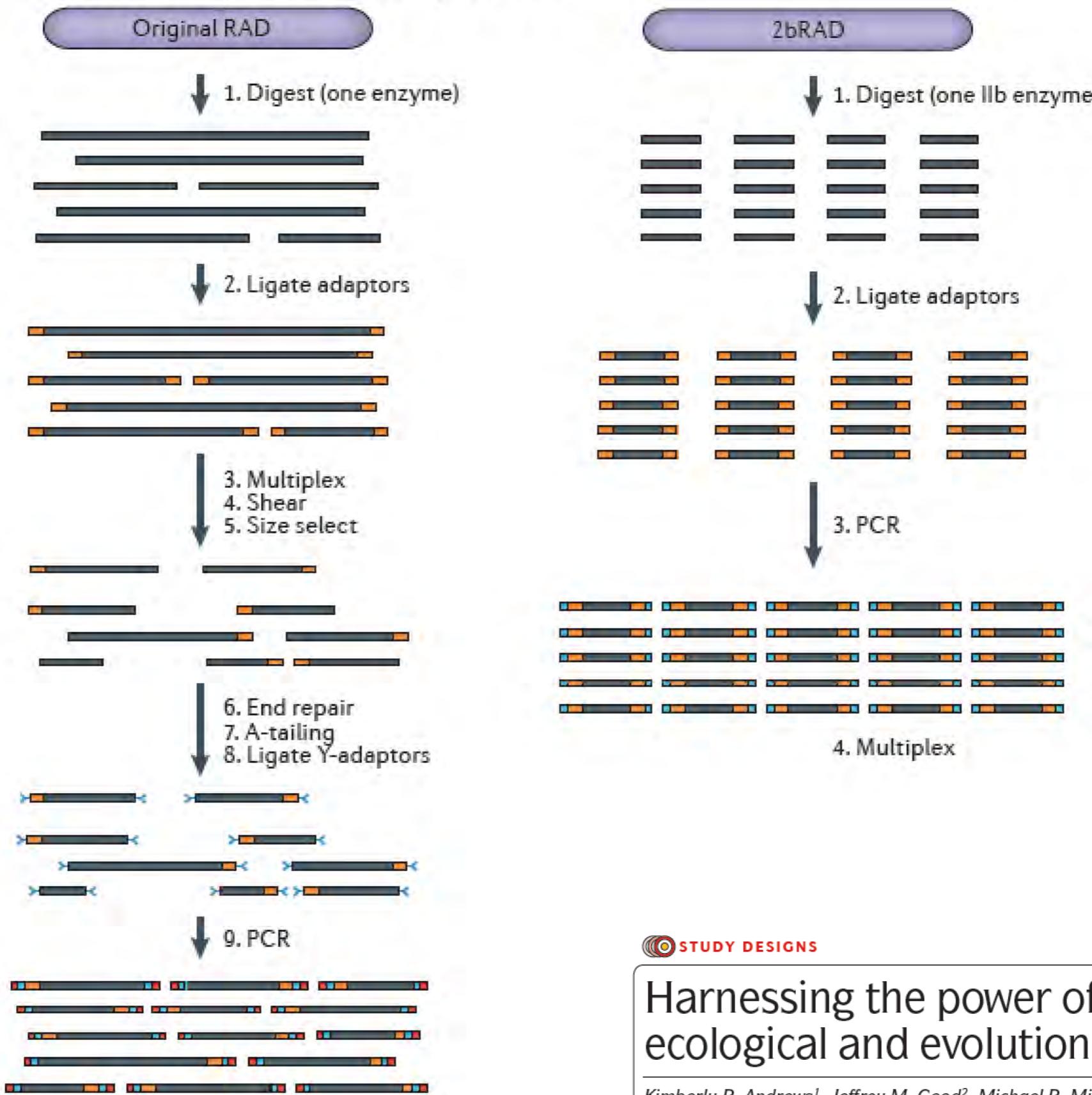


A = Amplification primer

B = Sequencing primer

C = Barcode

# RADseq with one enzyme digestion & shearing

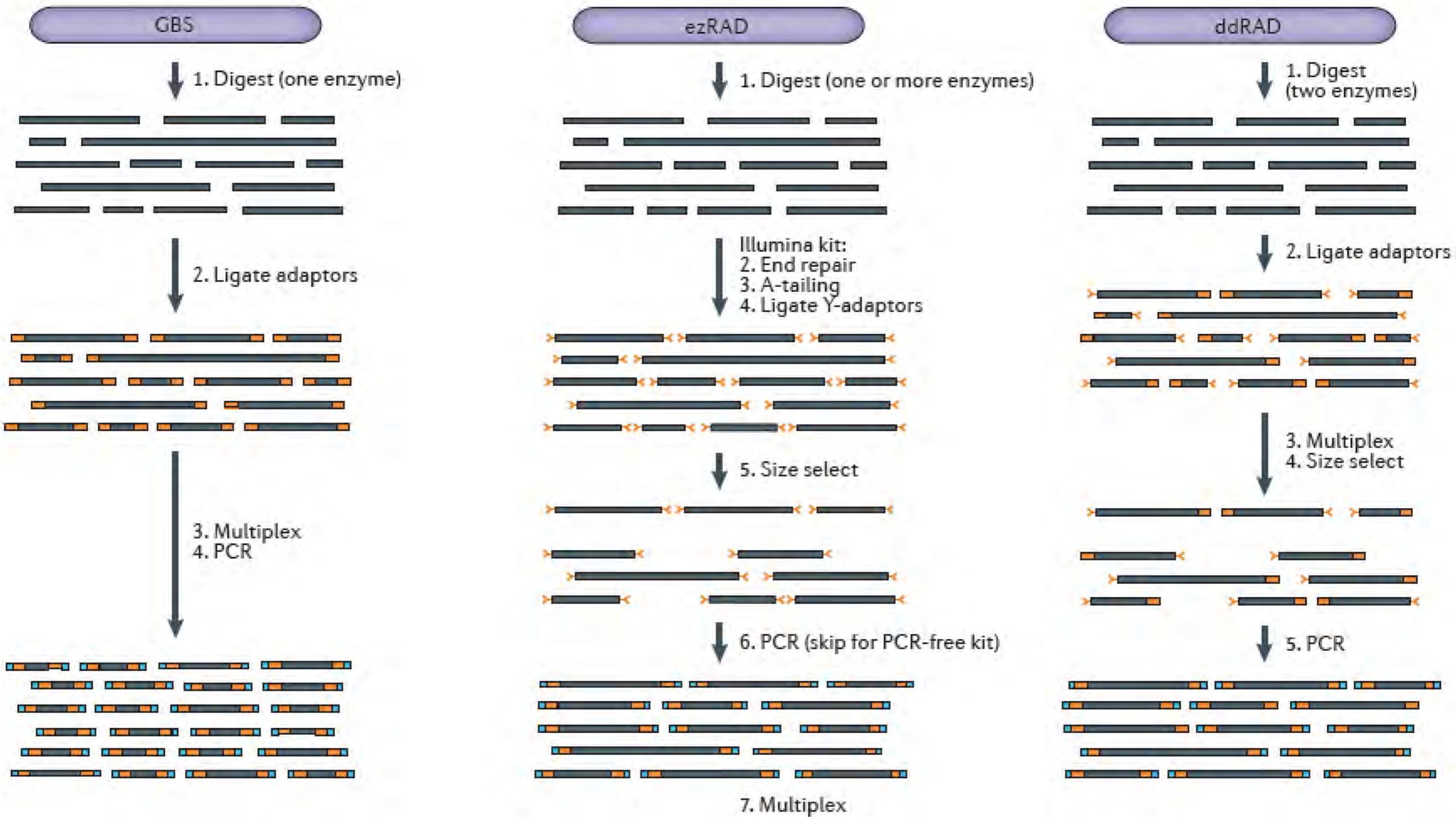


## STUDY DESIGNS

Harnessing the power of RADseq for ecological and evolutionary genomics

Kimberly R. Andrews<sup>1</sup>, Jeffrey M. Good<sup>2</sup>, Michael R. Miller<sup>3</sup>, Gordon Luikart<sup>4</sup> and Paul A. Hohenlohe<sup>5</sup>

# RADseq with one or two enzyme digestion



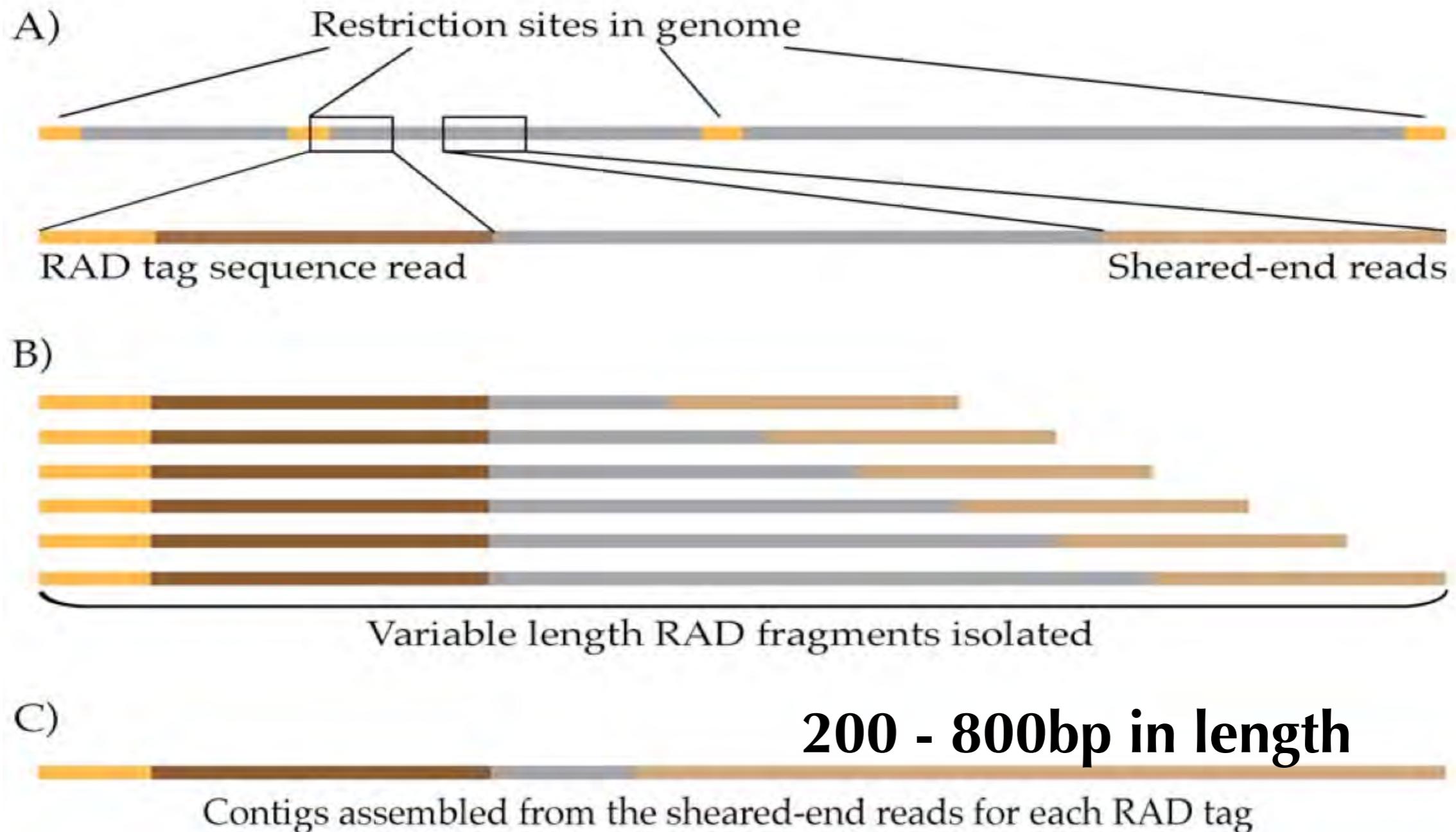
# Summary of plusses and minuses of RAD family

## Harnessing the power of RADseq for ecological and evolutionary genomics

Kimberly R. Andrews<sup>1</sup>, Jeffrey M. Good<sup>2</sup>, Michael R. Miller<sup>3</sup>, Gordon Luikart<sup>4</sup>  
and Paul A. Hohenlohe<sup>5</sup>

	Original RAD	2bRAD	GBS	ddRAD	ezRAD
Options for tailoring number of loci	Change restriction enzyme	Change restriction enzyme	Change restriction enzyme	Change restriction enzyme or size selection window	Change restriction enzyme or size selection window
Number of loci per 1 Mb of genome size*	30–500	50–1,000	5–40	0.3–200	10–800
Length of loci	≤1kb if building contigs; otherwise ≤300 bp <sup>‡</sup>	33–36 bp	<300 bp <sup>‡</sup>	≤300 bp <sup>‡</sup>	≤300 bp <sup>‡</sup>
Cost per barcoded or indexed sample	Low	Low	Low	Low	High
Effort per barcoded or indexed sample <sup>§</sup>	Medium	Low	Low	Low	High
Use of proprietary kit	No	No	No	No	Yes
Identification of PCR duplicates	With paired-end sequencing	No	With degenerate barcodes	With degenerate barcodes	No
Specialized equipment needed	Sonicator	None	None	Pippin Prep <sup>¶</sup>	Pippin Prep <sup>¶</sup>
Suitability for large or complex genomes <sup>¶</sup>	Good	Poor	Moderate	Good	Good
Suitability for de novo locus identification (no reference genome) <sup>#</sup>	Good	Poor	Moderate	Moderate	Moderate
Available from commercial companies	Yes	No	Yes	Yes	No

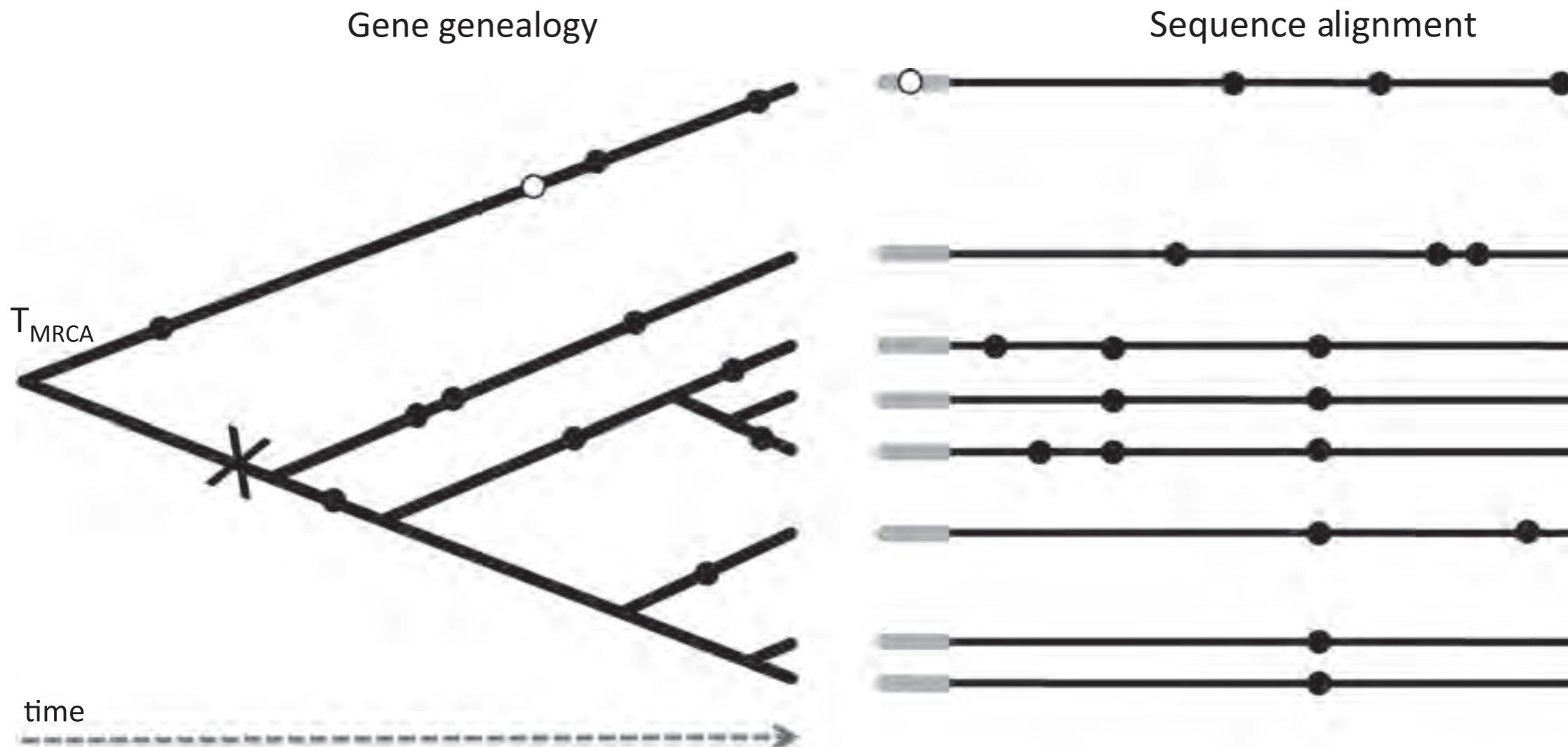
# Random shearing in original RAD - Local Paired End (PE) Assemblies



**RADseq underestimates diversity and introduces genealogical biases due to nonrandom haplotype sampling**

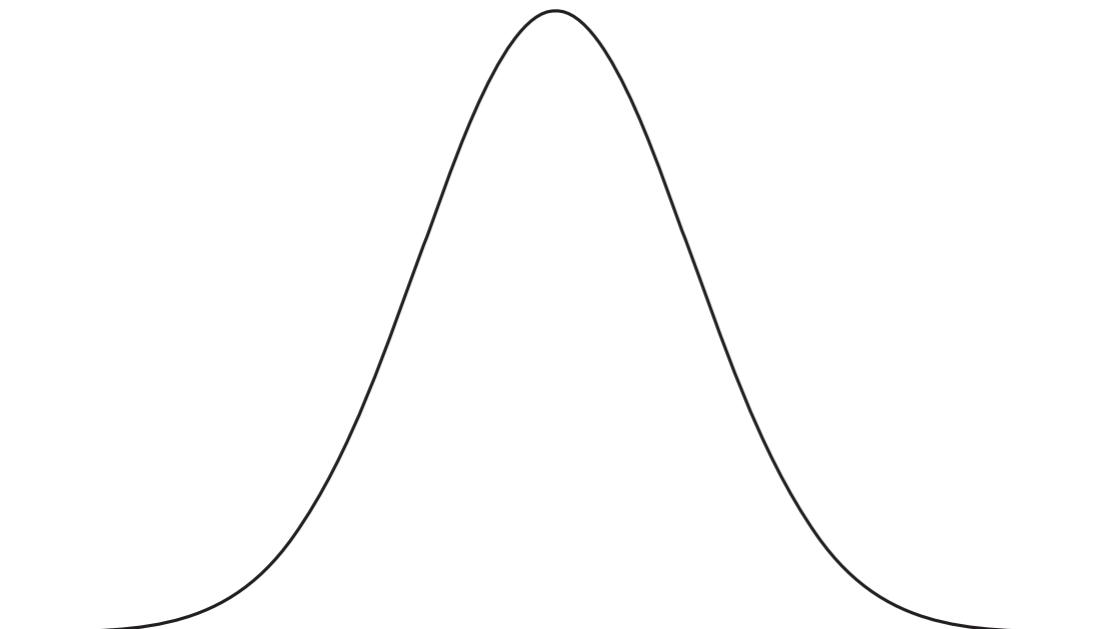
B. ARNOLD,<sup>1</sup> R. B. CORBETT-DETIG,<sup>1</sup> D. HARTL and K. BOMBLIES

*Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA 02138, USA*



# ‘Bias’ in RADseq

---



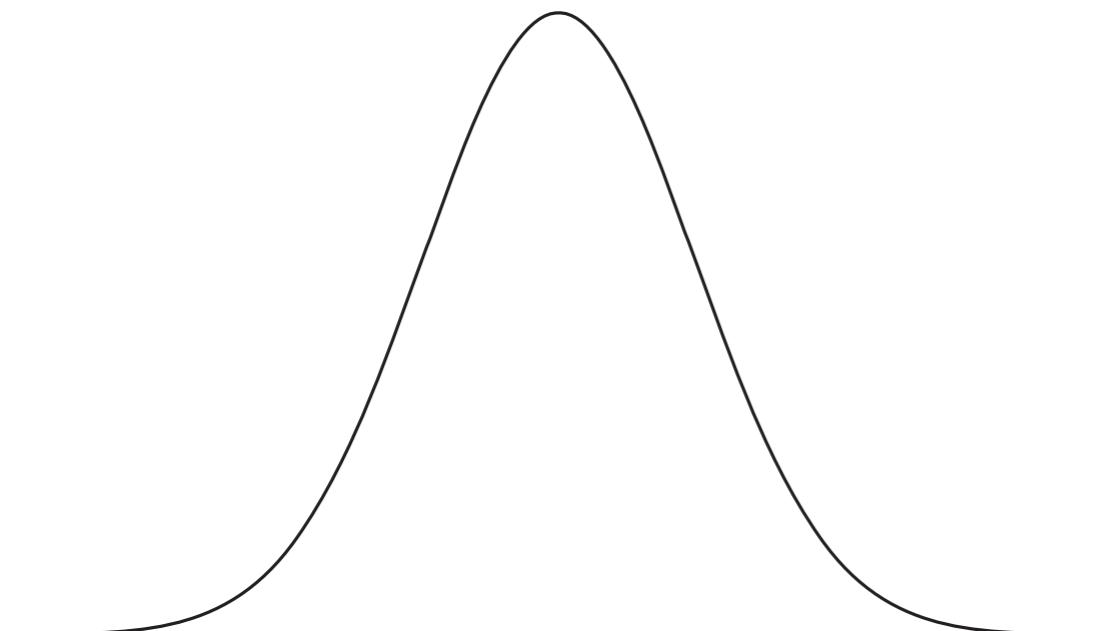
$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

$e = 2.7182\dots$

$\pi = 3.1415\dots$

# ‘Bias’ in RADseq

---



$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

# ‘Bias’ in RADseq is increased in some RAD protocols

Protocol	$\theta$ per bp	Mean	
		$\theta_{we}/\theta_{wa}$	$\pi_e/\pi_a$
Standard	0.0001	0.994	0.995
	0.001	0.987	0.982
	0.01	0.956	0.933
Double digest	0.0001	0.835	0.836
	0.001	0.858	0.851
	0.01	0.829	0.797

# ‘Bias’ in RADseq is increased in some RAD protocols

Protocol	$\theta$ per bp	Mean	
		$\theta_{we}/\theta_{wa}$	$\pi_e/\pi_a$
Standard	0.0001	0.994	0.995
	0.001	0.987	0.982
	0.01	0.956	0.933
Double digest	0.0001	0.835	0.836
	0.001	0.858	0.851
	0.01	0.829	0.797

# ‘Bias’ in RADseq is increased in some RAD protocols

Protocol	$\theta$ per bp	Mean	
		$\theta_{we}/\theta_{wa}$	$\pi_e/\pi_a$
Standard	0.0001	0.994	0.995
	0.001	0.987	0.982
	0.01	0.956	0.933
Double digest	0.0001	0.835	0.836
	0.001	0.858	0.851
	0.01	0.829	0.797

# Biological studies that benefit from whole genome approaches powered by RAD-seq

---

- Defining individuality, parentage and pedigrees
- Performing quantitative genetic studies in outbred populations
- Fine scale estimates of population structure
- Identifying the genetic basis of inbreeding depression
- Making management decisions for biological populations
- Genome Wide Association Studies (GWAS) studies of traits
- Building genetic maps to genetically enable non model organisms
- Estimating species and higher level phylogenetic relationships
- Population genomics - identifying the signatures of natural selection

# Biological studies that benefit from whole genome approaches powered by RAD-seq

---

- Defining individuality, parentage and pedigrees
  - Performing quantitative genetic studies in outbred populations
  - Fine scale evolution
  - Identifying transposons
  - Making maps
  - Genome Wide Association Studies
  - Building genomic resources for non-model organisms
  - Estimating species and higher level phylogenetic relationships
  - Population genomics - identifying the signatures of natural selection
- Any study where improving biological sample size would be beneficial

# Outline for today's lecture

---

RAD-seq for ecological and evolutionary genomics

## **Primer on Population Genomics**

Evolutionary genomics of stickleback fish

- Population genomics of rapid adaptation
- Using long read RAD-seq for coalescent analyses
- Genome Wide Association Studies using RAD-seq

Genomically enabling the Gulf pipefish

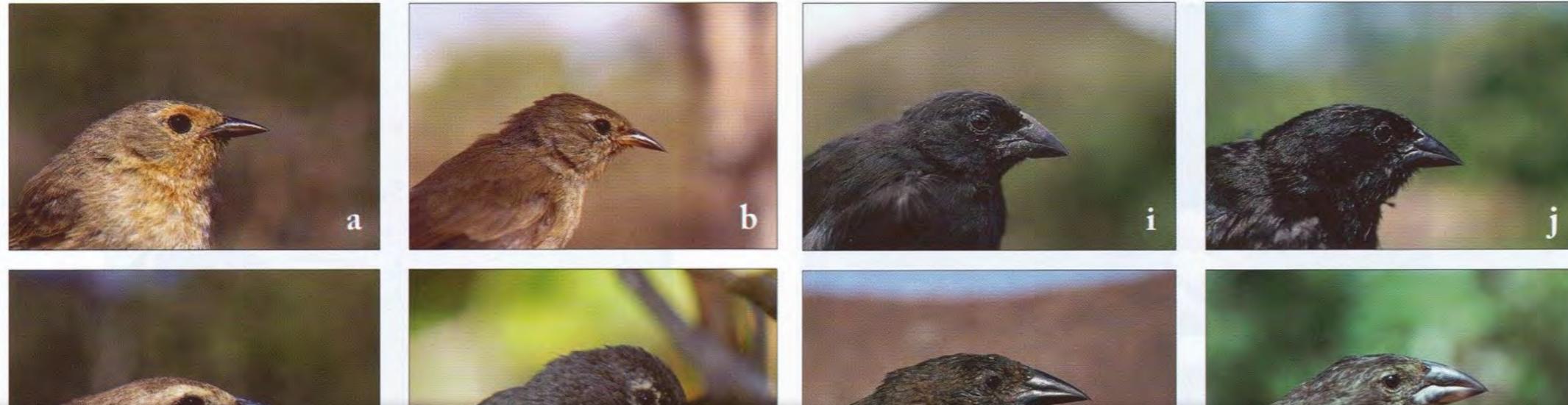
# How do organisms adapt to novel environments?

---



from Grant and Grant. 2007. How and why species multiply: The radiation of Darwin's finches. Princeton University Press

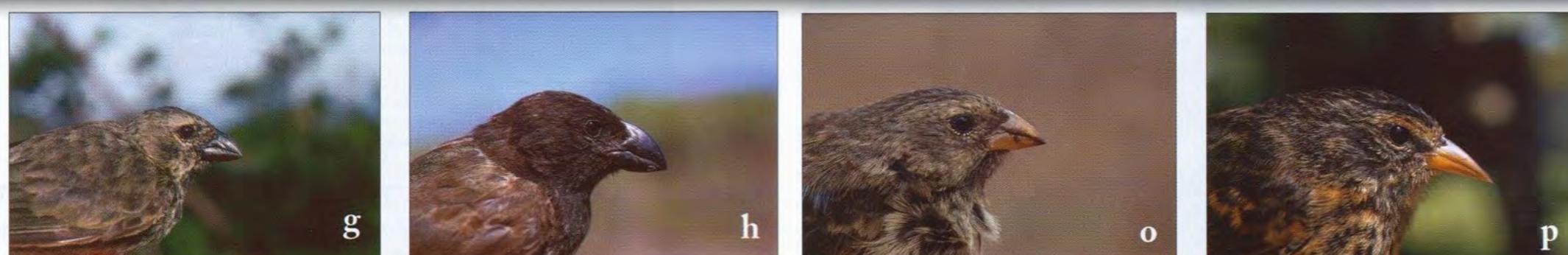
# How do organisms adapt to novel environments?



How is genetic diversity partitioned across individuals, populations and species?

What genomic regions are important for adaptation to novel environments?

How does genome architecture influence rapid evolution?



from Grant and Grant. 2007. How and why species multiply: The radiation of Darwin's finches. Princeton University Press

# Four fundamental processes in evolution

---

Origin of genetic variation

**mutation**

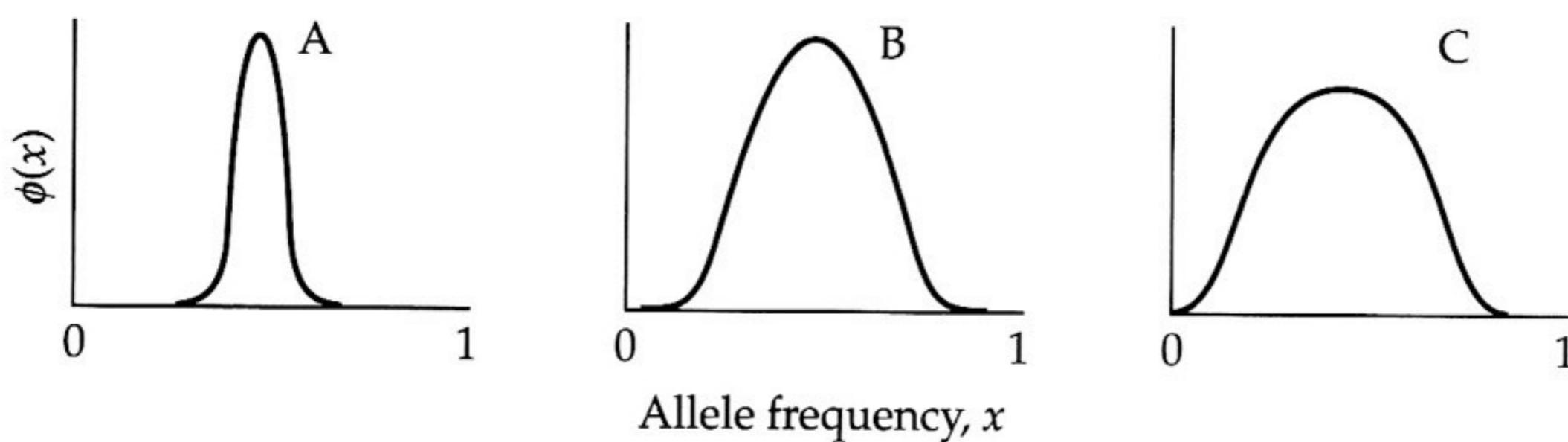
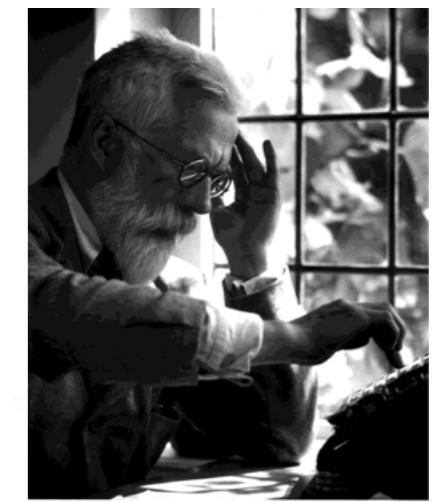
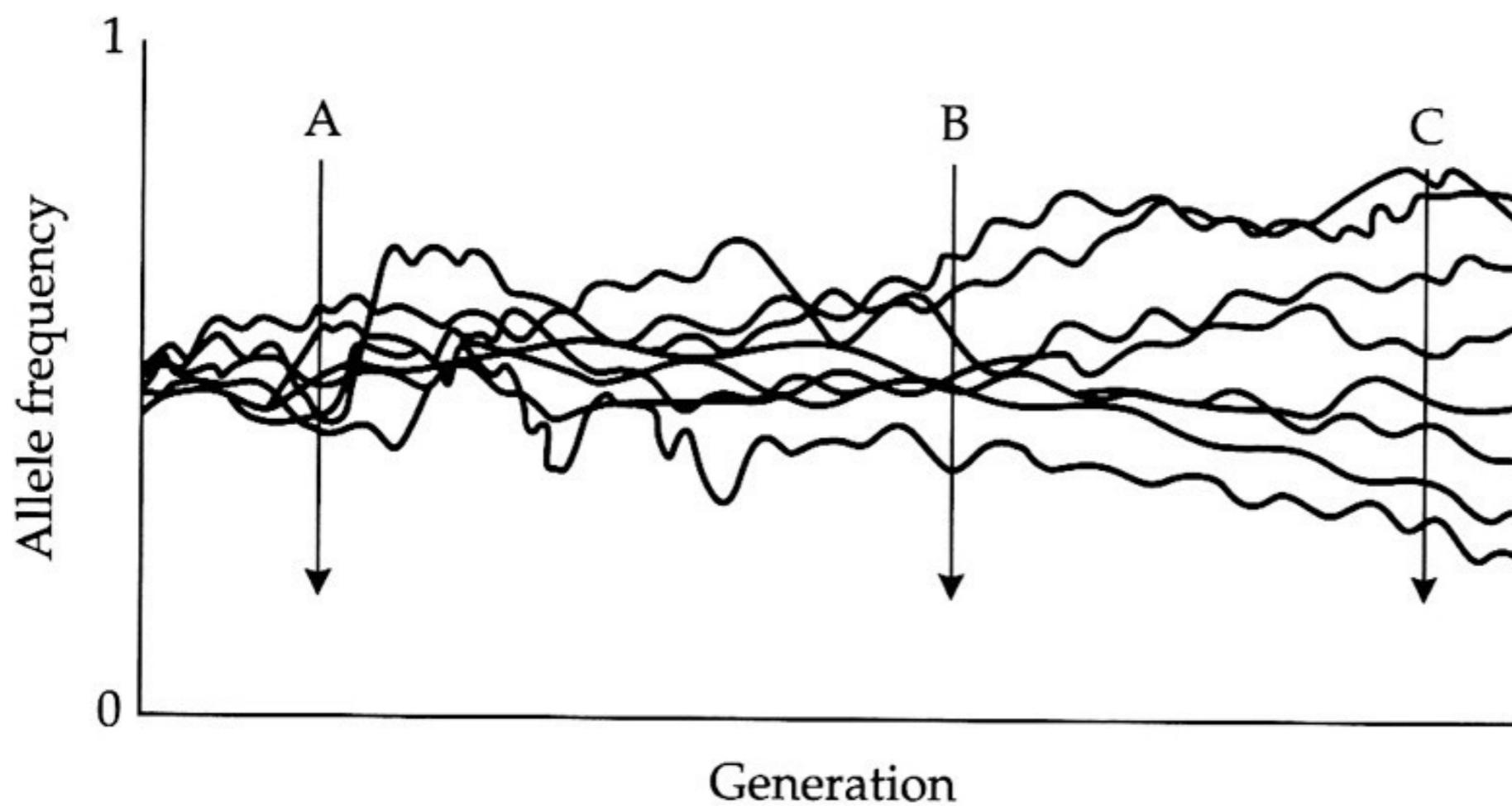
**migration**

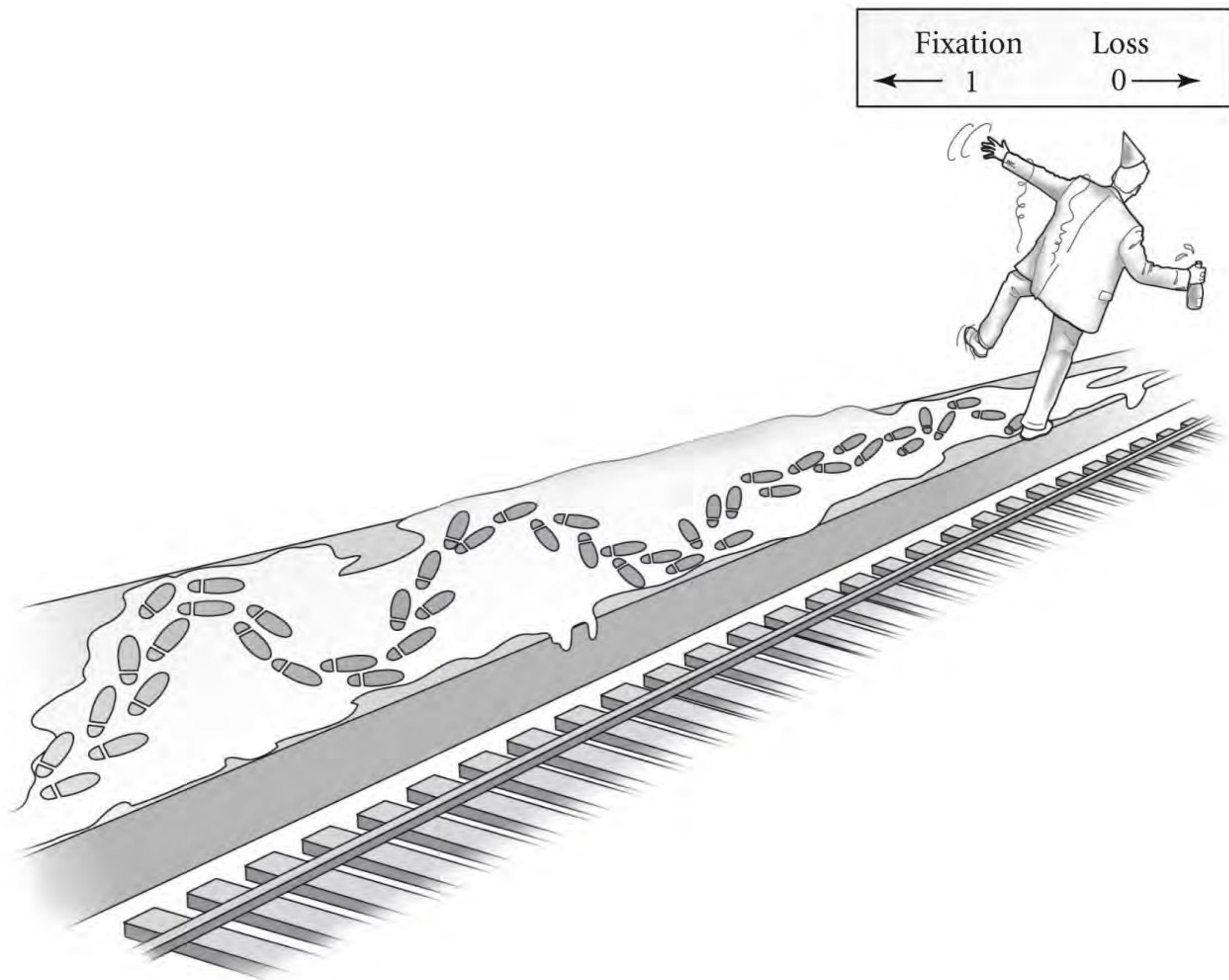
Sorting of variation

**genetic drift**

**natural selection**

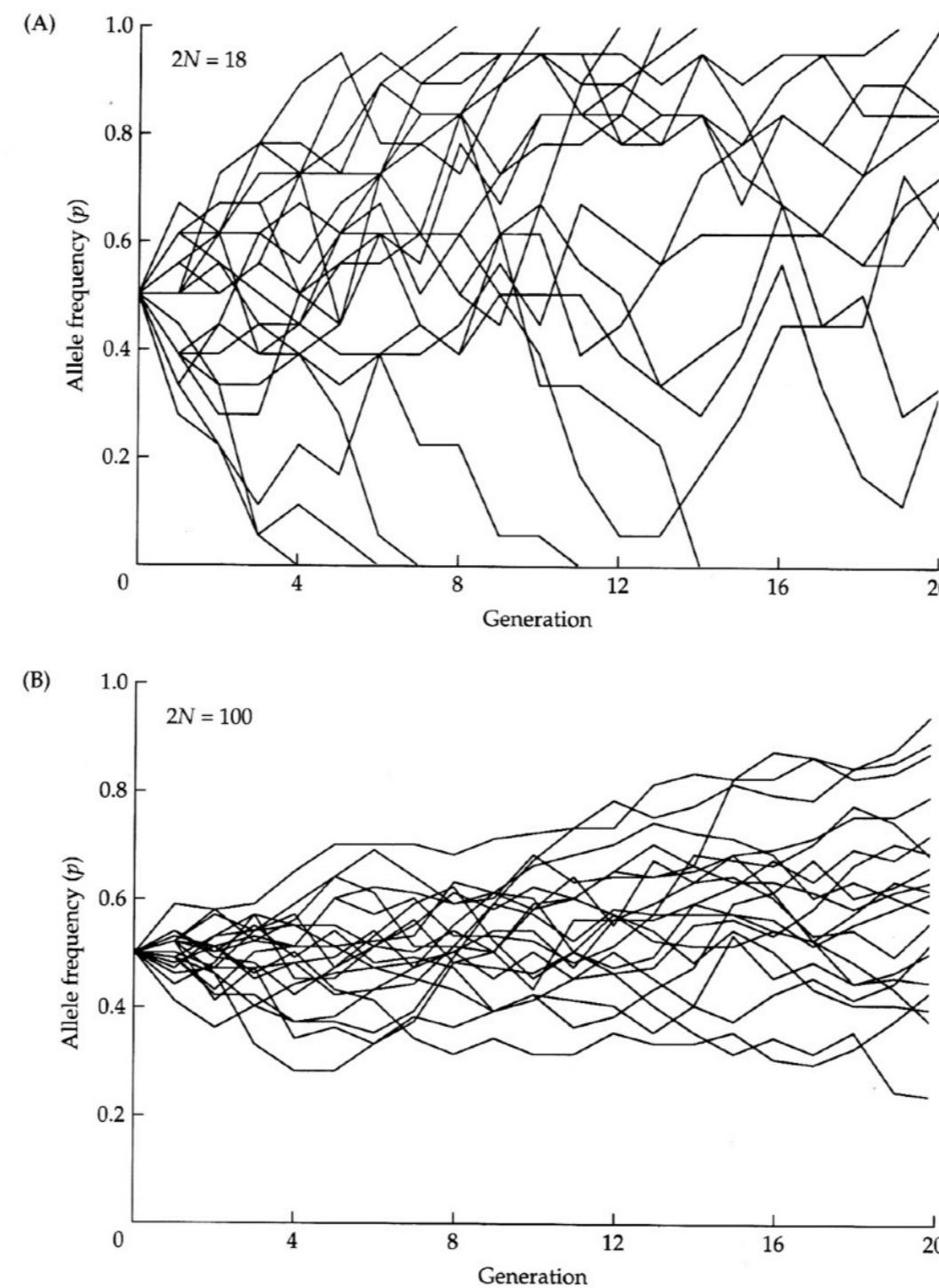
# Genetic drift is a null model





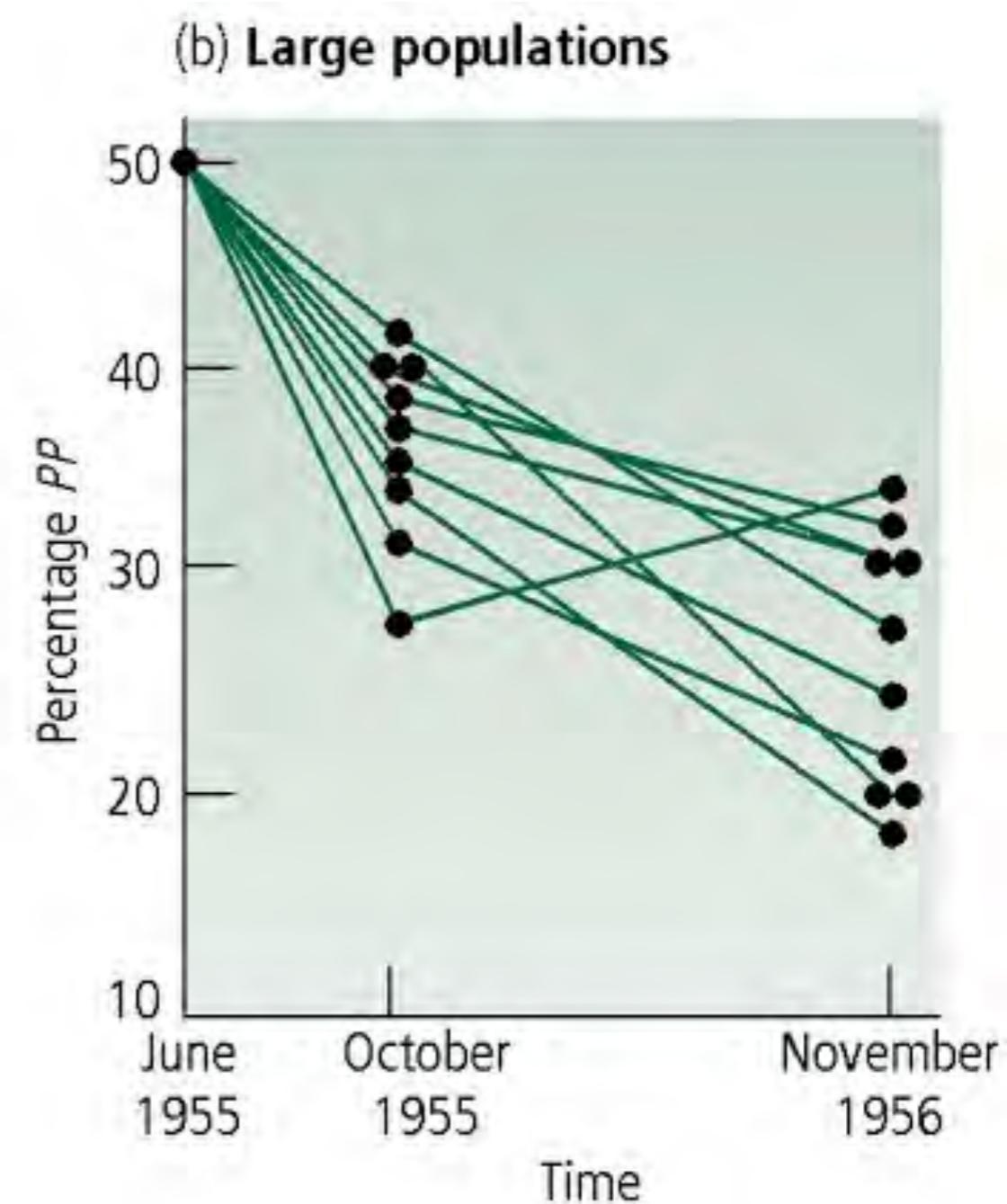
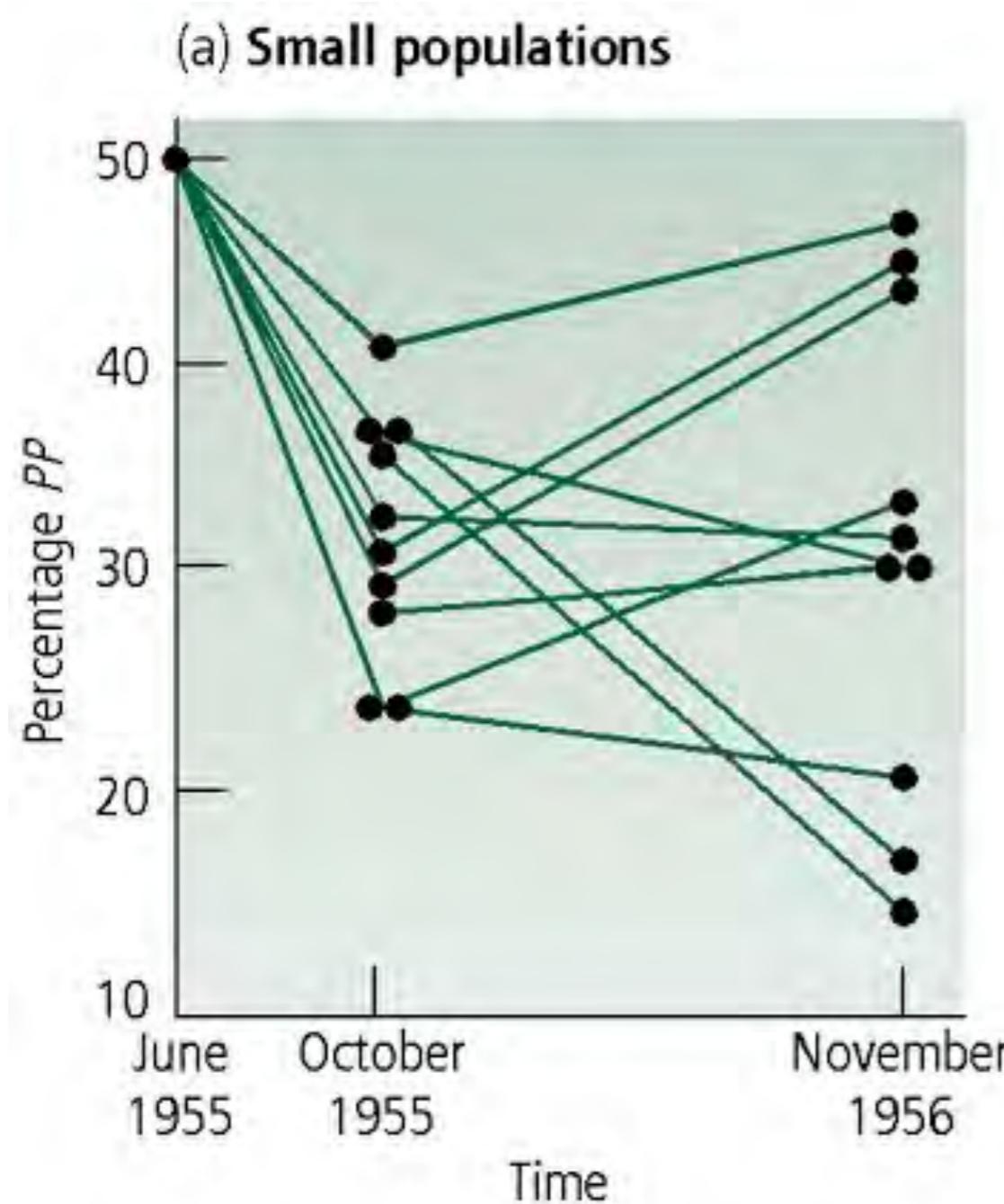
**EVOLUTION 2e, Figure 10.2**

# Size of population ( $N_e$ ) affects rate of spread (diffusion)



Affects the entire genome equally (on average)

Natural selection biases the allele frequency change, but drift is still occurring



\*\*\* The effects of selection can be genomically localized \*\*\*

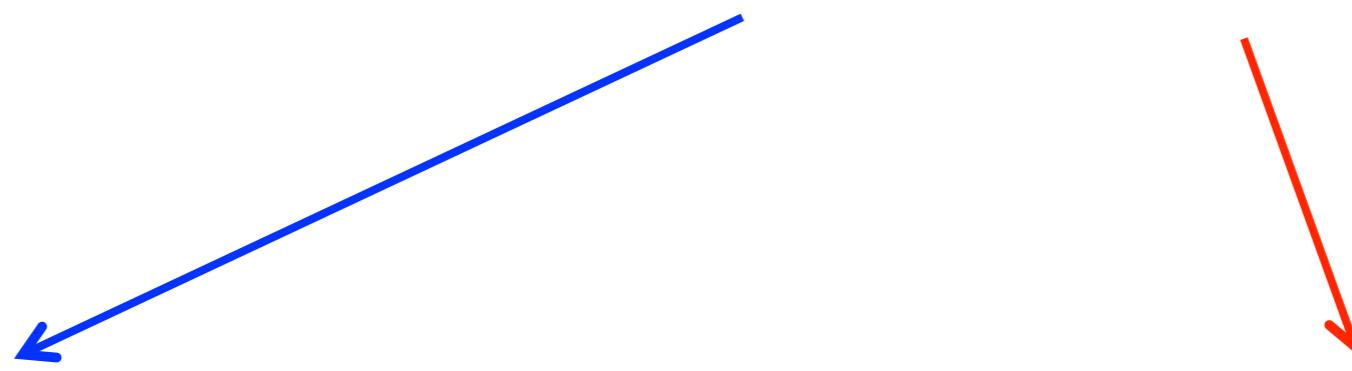
# Population genomics

---

Simultaneous genotyping of **neutral** and **adaptive** loci

Genome-wide background provides more precise estimates:

- Demographic processes (e.g.  $N_e$ )
- Phylogeography



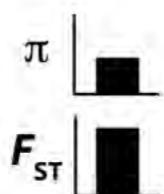
Outliers from background indicate:

- Selective sweeps
- Local adaptation

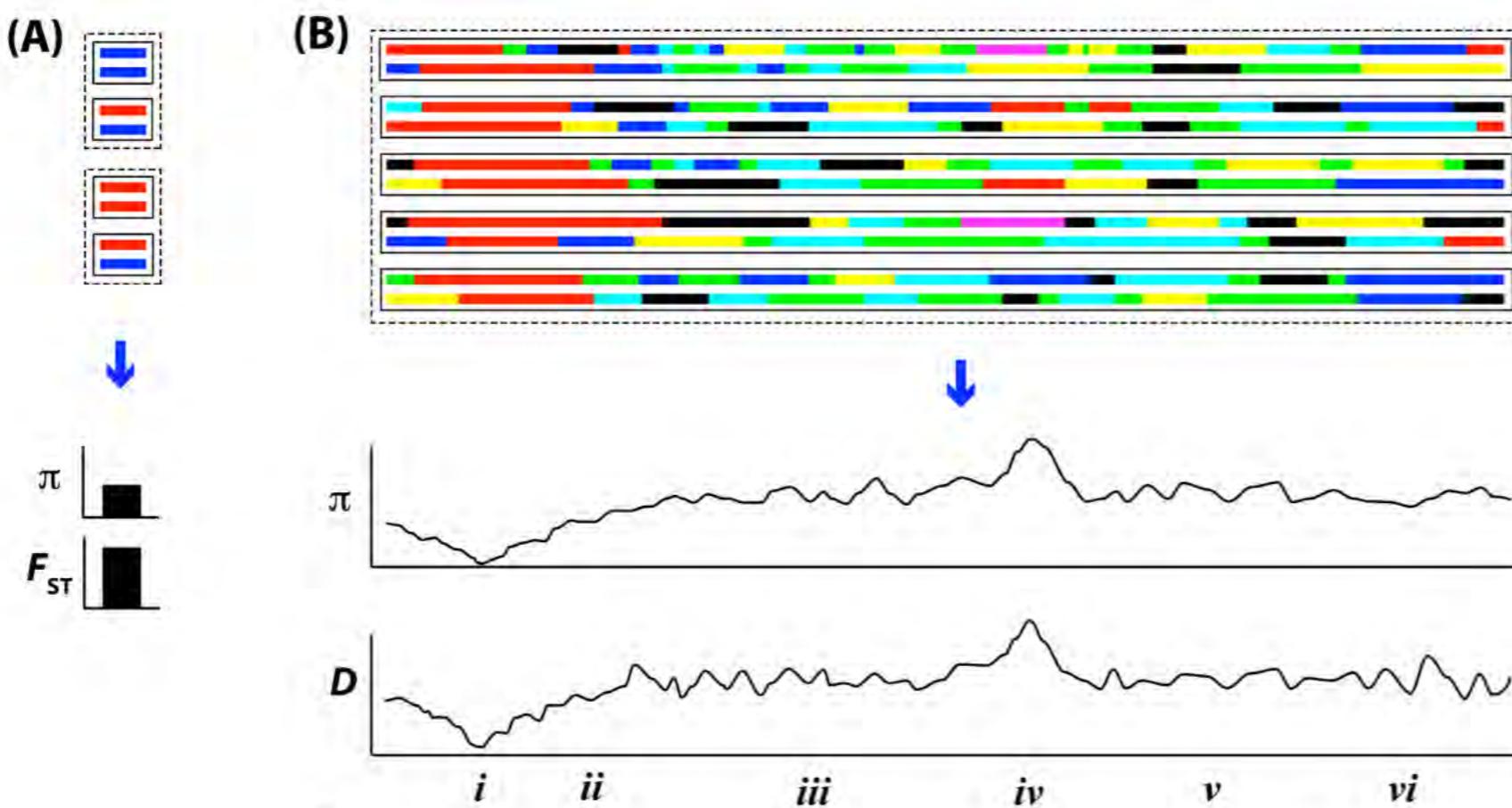


# Population genomics of ordered markers

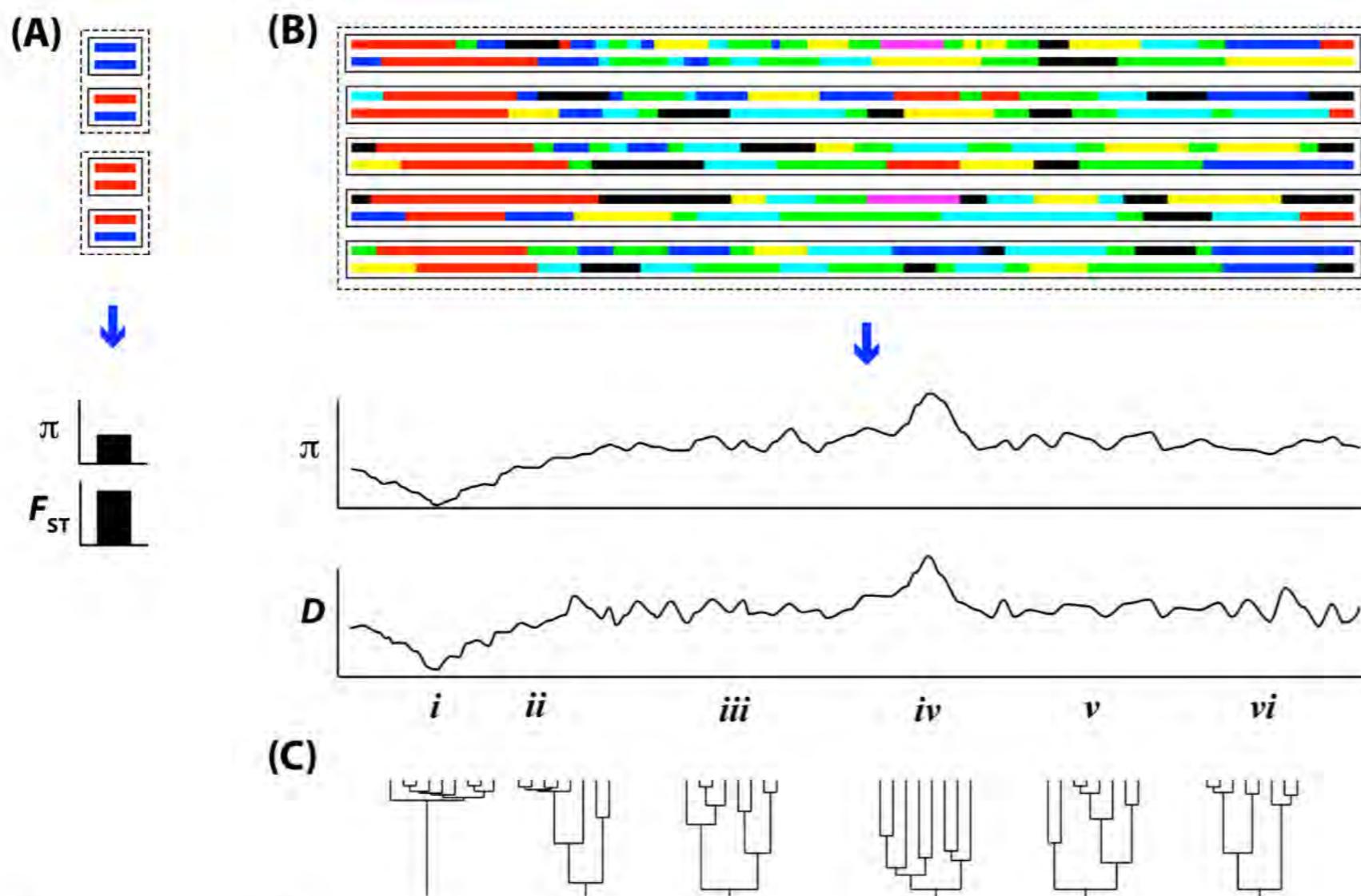
(A)



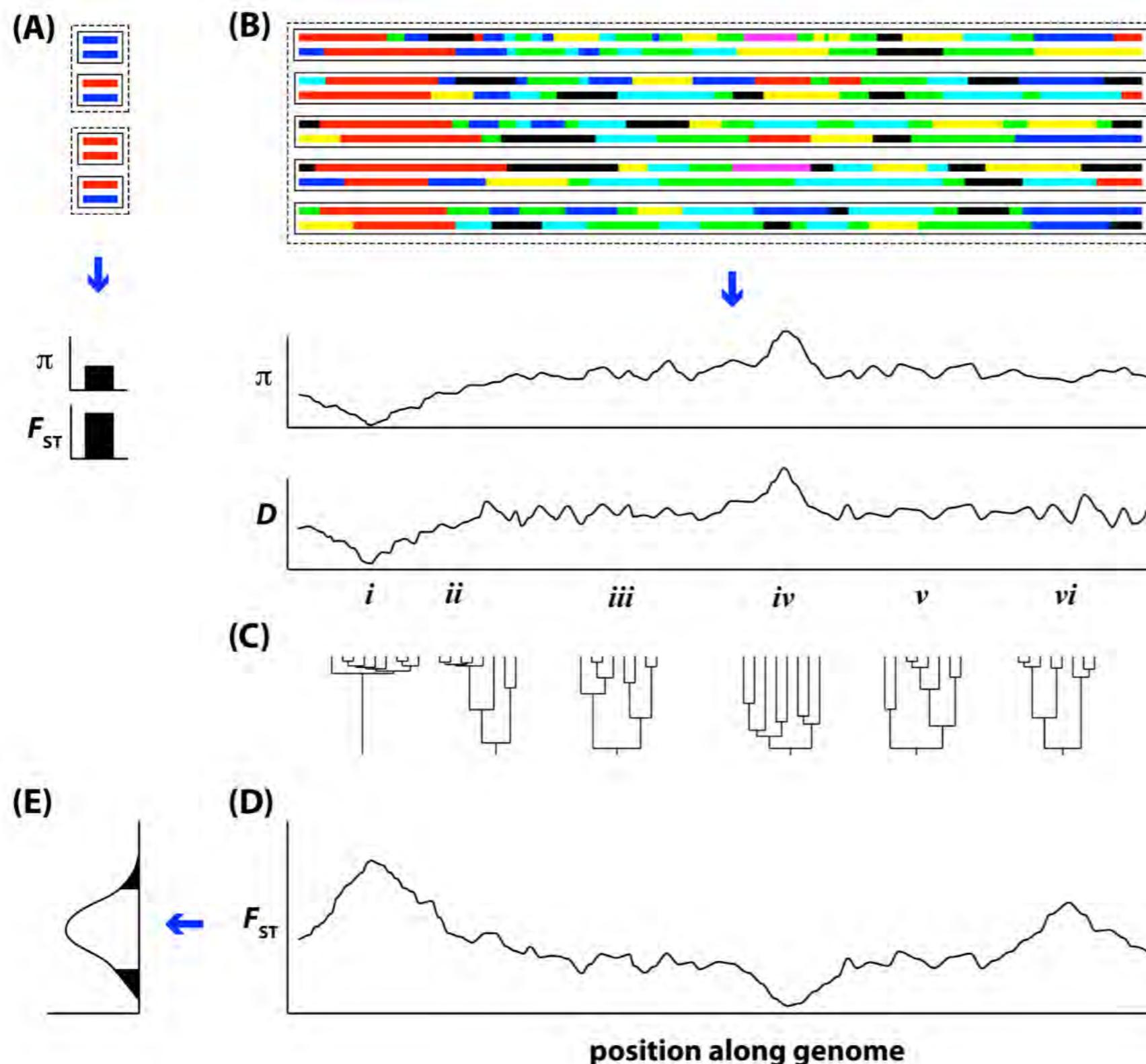
# Population genomics of ordered markers



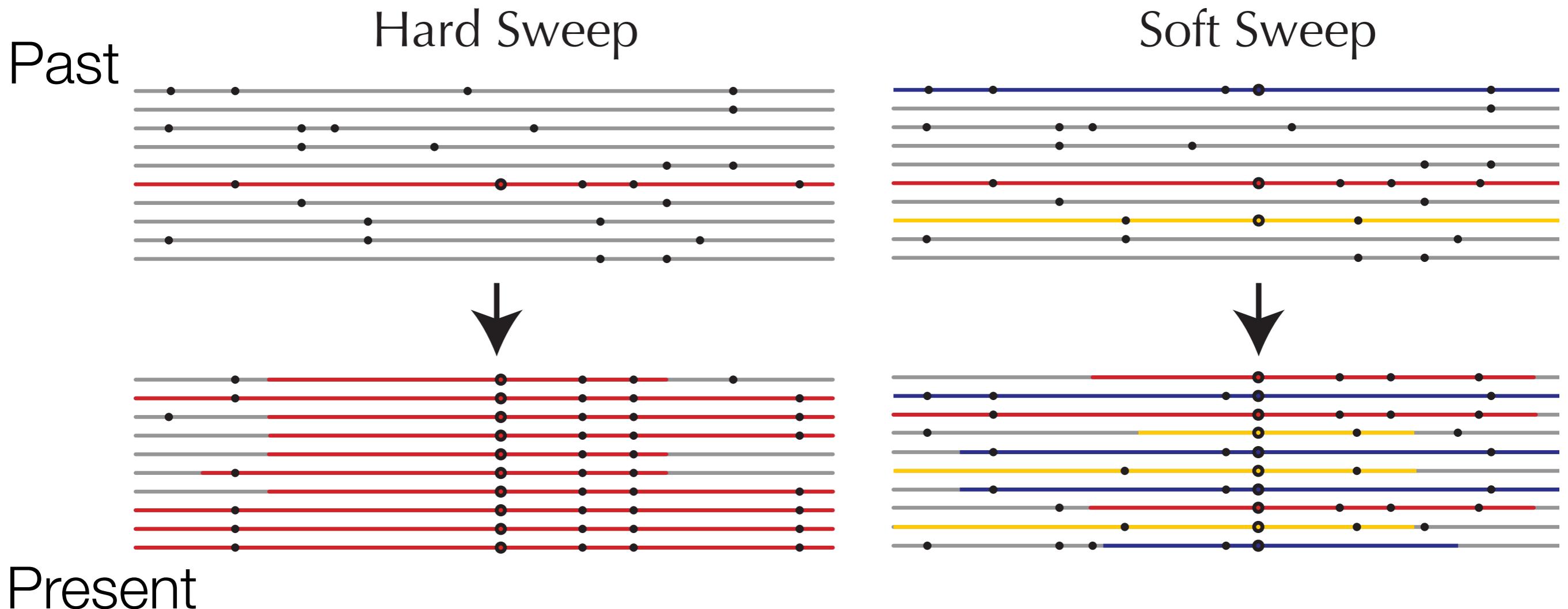
# Population genomics of ordered markers



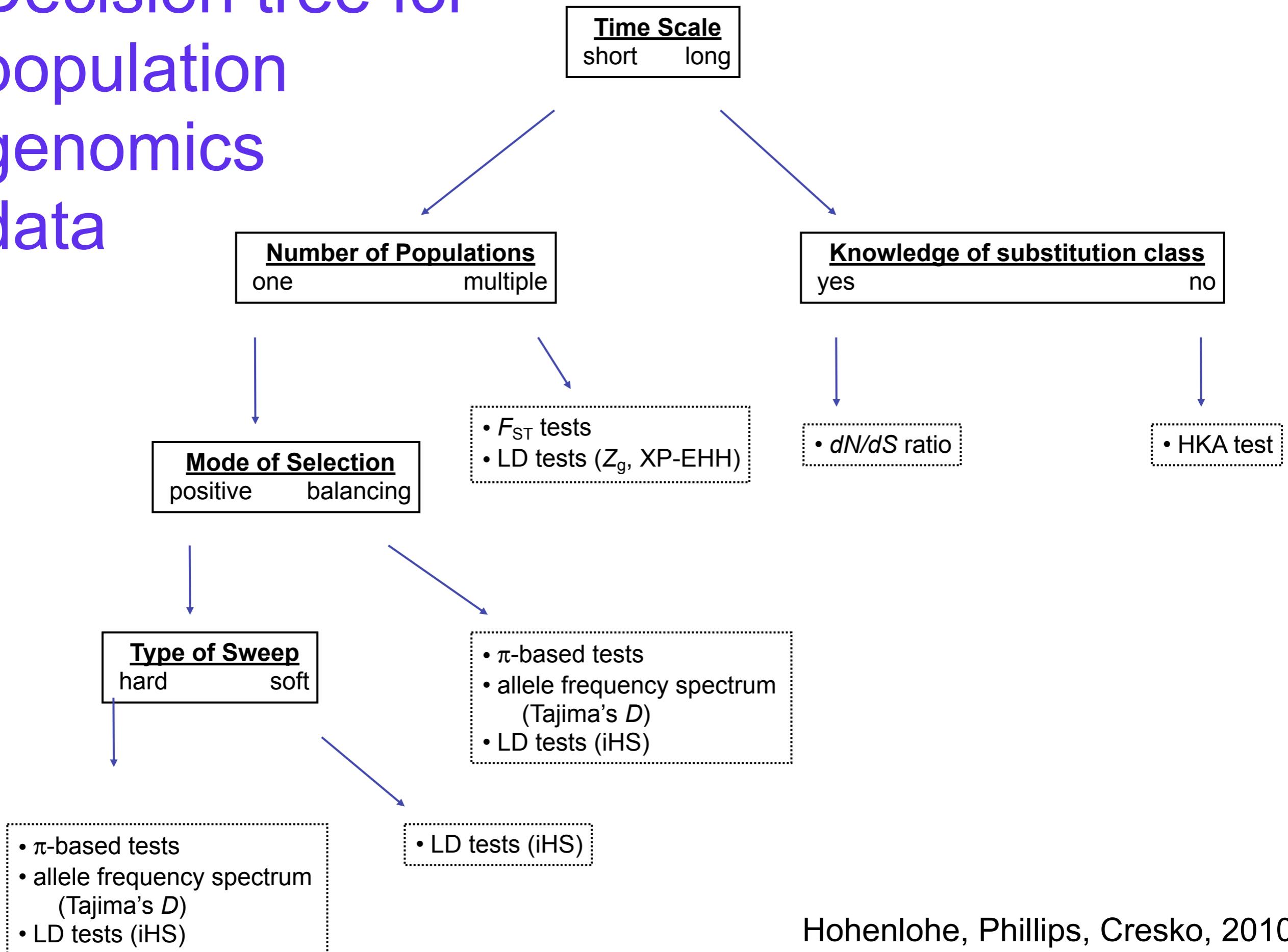
# Population genomics of ordered markers



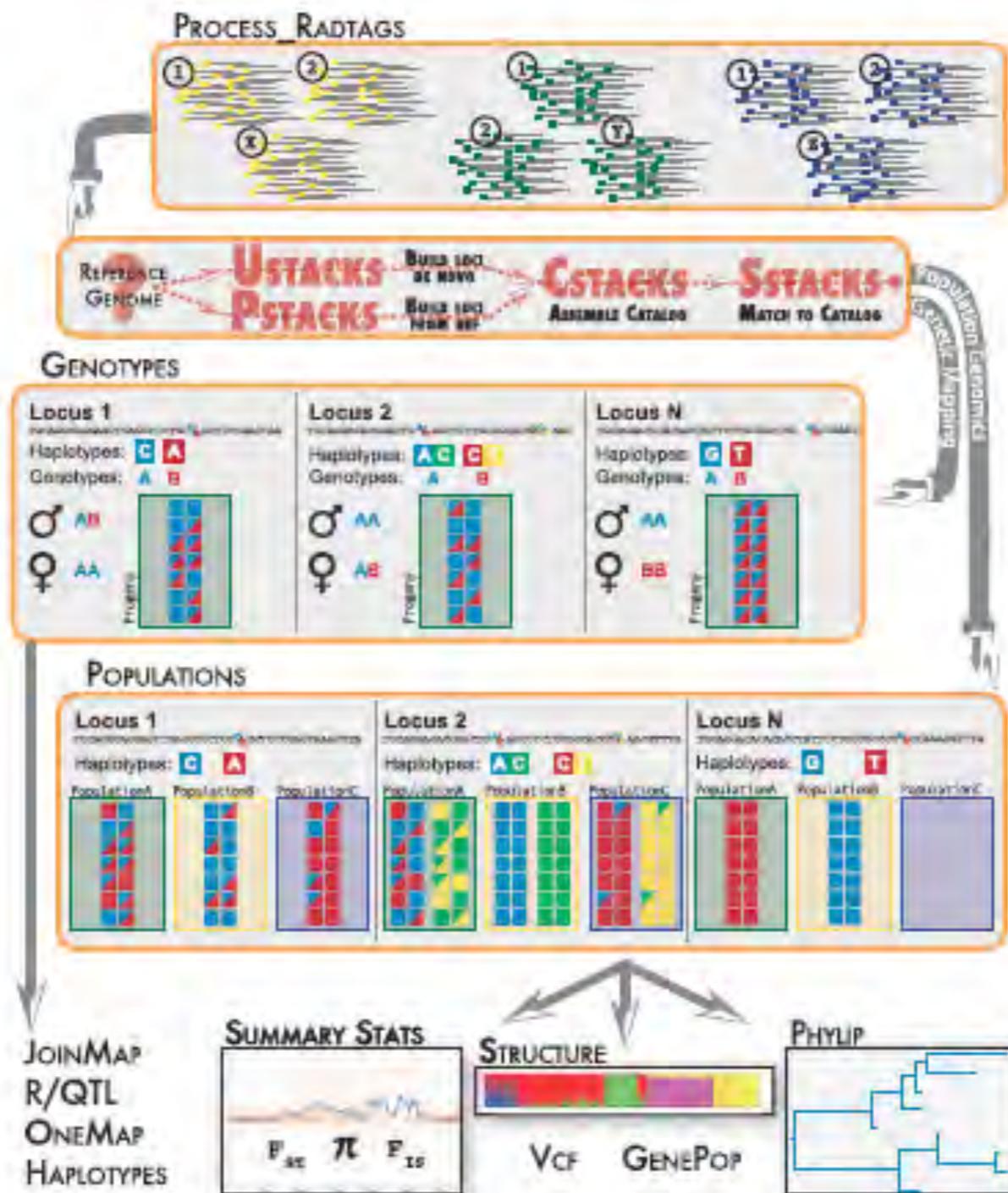
# Sweeps of directional selection across genomes



# Decision tree for population genomics data



# Stacks analysis pipeline for RAD-seq



## Stacks: Building and Genotyping Loci De Novo From Short-Read Sequences

Julian M. Catchen,\* Angel Amores,<sup>†</sup> Paul Hohenlohe,<sup>\*</sup> William Cresko,<sup>\*</sup> and John H. Postlethwait<sup>†,1</sup>

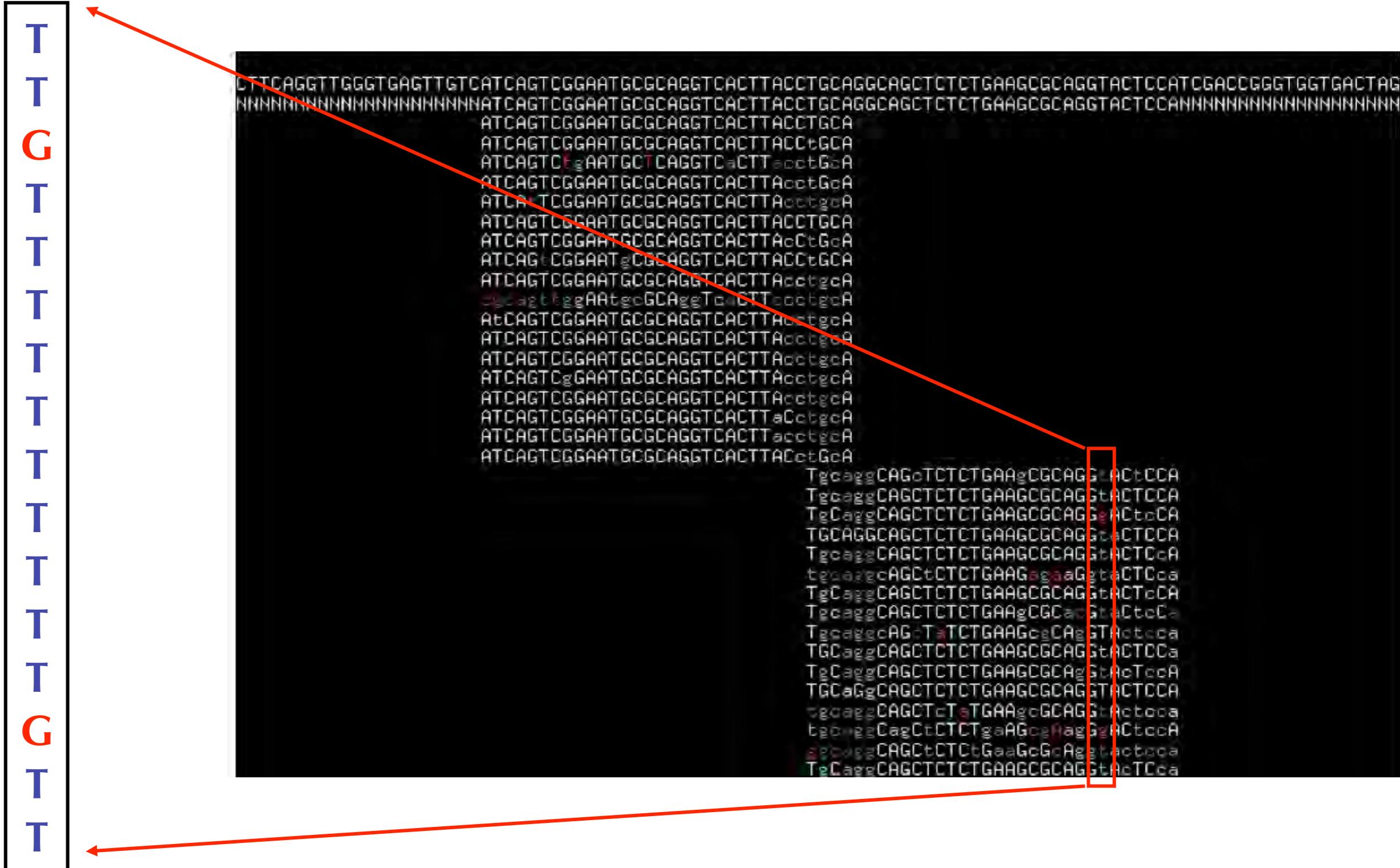
\*Center for Ecology and Evolutionary Biology and <sup>†</sup>Institute of Neuroscience, University of Oregon, Eugene, Oregon 97403

## Stacks: an analysis tool set for population genomics

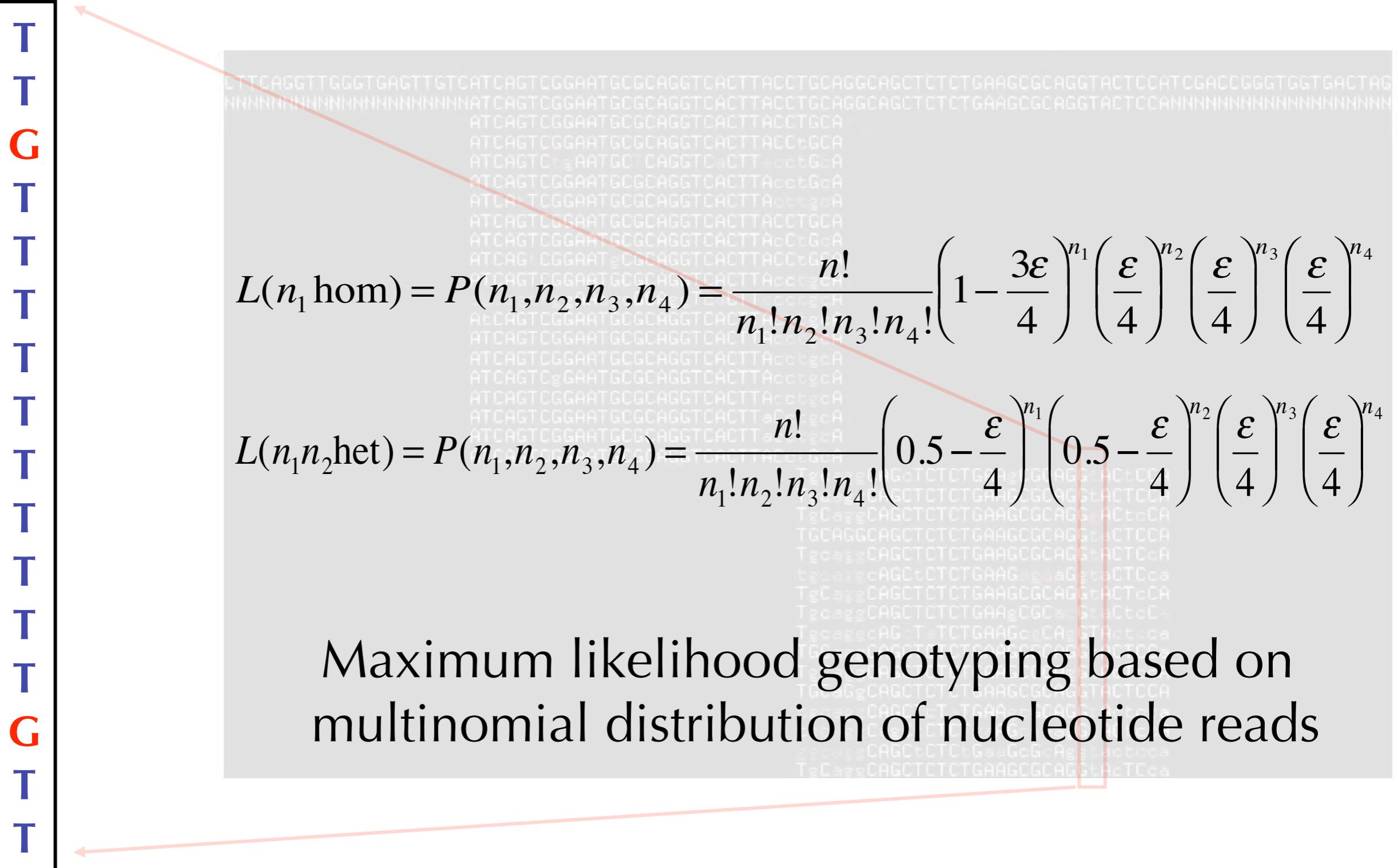
JULIAN CATCHEN,\* PAUL A. HOHENLOHE,\*<sup>†</sup> SUSAN BASSHAM,\* ANGEL AMORES<sup>‡</sup> and WILLIAM A. CRESKO\*

\*Institute of Ecology and Evolution, University of Oregon, Eugene, OR 97403-5289, USA, <sup>†</sup>Biological Sciences, University of Idaho, Moscow, ID 83844-3051, USA, <sup>‡</sup>Institute of Neuroscience, University of Oregon, Eugene, OR 97403-1254, USA

# Stacks

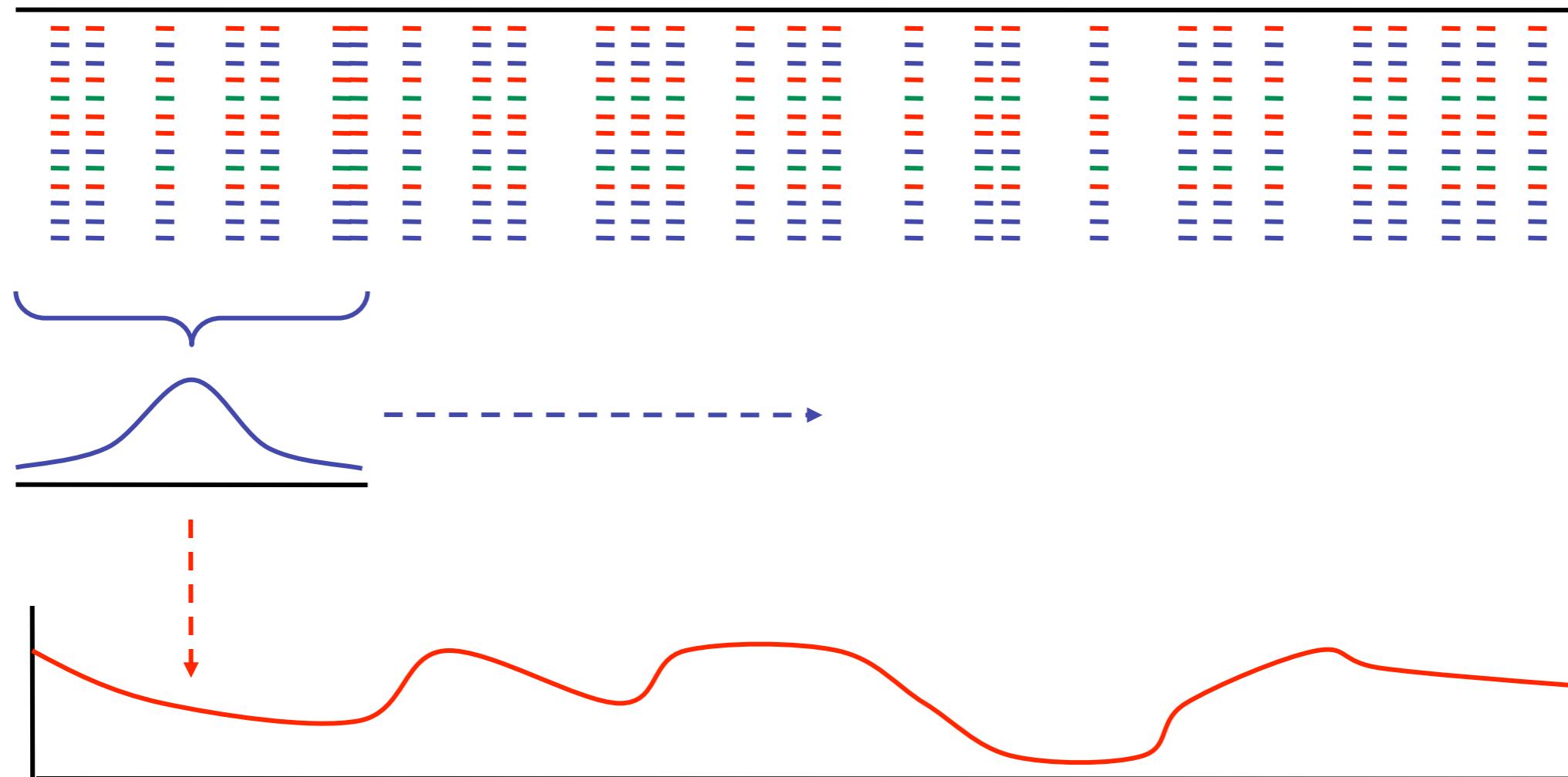


# Stacks



# Making statistics continuous across the genome

Kernel-smoothing average of summary statistics along genome



Bootstrap re-sampling to estimate significance of moving average

# Stacks

## 7.4.2 batch\_X.sumstats\_summary.tsv: Summary of summary statistics for each population

Column	Name	Description
1	Pop ID	Population ID as defined in the Population Map file.
2	Private	Number of private alleles in this population.
3	Number of Individuals	Mean number of individuals per locus in this population.
4	Variance	
5	Standard Error	
6	P	Mean frequency of the most frequent allele at each locus in this population.
7	Variance	
8	Standard Error	
9	Observed Heterozygosity	Mean observed heterozygosity in this population.
10	Variance	
11	Standard Error	
12	Observed Homozygosity	Mean observed homozygosity in this population.
13	Variance	
14	Standard Error	
15	Expected Heterozygosity	Mean expected heterozygosity in this population.
16	Variance	
17	Standard Error	
18	Expected Homozygosity	Mean expected homozygosity in this population.
19	Variance	
20	Standard Error	
21	$\Pi$	Mean value of $\pi$ in this population.
22	$\Pi$ Variance	
23	$\Pi$ Standard Error	
24	$F_{IS}$	Mean measure of $F_{IS}$ in this population.
25	$F_{IS}$ Variance	
26	$F_{IS}$ Standard Error	

# Stacks

## 7.4.4 batch\_X.hapstats.tsv: Haplotype-based summary statistics for each locus in each population

Column	Name	Description
1	Batch ID	The batch identifier for this data set.
2	Locus ID	Catalog locus identifier.
3	Chromosome	If aligned to a reference genome.
4	Basepair	If aligned to a reference genome.
5	Population ID	The ID supplied to the populations program, as written in the population map file.
6	N	Number of alleles/haplotypes present at this locus.
7	Haplotype count	
8	Gene Diversity	
9	Smoothed Gene Diversity	
10	Smoothed Gene Diversity P-value	
11	Haplotype Diversity	
12	Smoothed Haplotype Diversity	
13	Smoothed Haplotype Diversity P-value	
14	Haplotypes	A semicolon-separated list of haplotypes/haplotype counts in the population.



#### 7.4.3 batch\_X.fst\_Y-Z.tsv: $F_{ST}$ calculations for each pair of populations

Column	Name	Description
1	Batch ID	The batch identifier for this data set.
2	Locus ID	Catalog locus identifier.
3	Population ID 1	The ID supplied to the populations program, as written in the population map file.
4	Population ID 2	The ID supplied to the populations program, as written in the population map file.
5	Chromosome	If aligned to a reference genome.
6	Basepair	If aligned to a reference genome.
7	Column	The nucleotide site within the catalog locus, reported using a zero-based offset (first nucleotide is enumerated as 0).
8	Overall $\pi$	An estimate of nucleotide diversity across the two populations.
9	$F_{ST}$	A measure of population differentiation.
10	FET p-value	P-value describing if the $F_{ST}$ measure is statistically significant according to Fisher's Exact Test.
11	Odds Ratio	Fisher's Exact Test odds ratio.
12	CI High	Fisher's Exact Test confidence interval.
13	CI Low	Fisher's Exact Test confidence interval.
14	LOD Score	Logarithm of odds score.
15	Corrected $F_{ST}$	$F_{ST}$ with either the FET p-value, or a window-size or genome size Bonferroni correction.
16	Smoothed $F_{ST}$	A weighted average of $F_{ST}$ depending on the surrounding $3\sigma$ of sequence in both directions.
17	AMOVA $F_{ST}$	Analysis of Molecular Variance alternative $F_{ST}$ calculation. Derived from Weir, <a href="#">Genetic Data Analysis II</a> , chapter 5, "F Statistics," pp166-167.
18	Corrected AMOVA $F_{ST}$	AMOVA $F_{ST}$ with either the FET p-value, or a window-size or genome size Bonferroni correction.
19	Smoothed AMOVA $F_{ST}$	A weighted average of AMOVA $F_{ST}$ depending on the surrounding $3\sigma$ of sequence in both directions.
20	Smoothed AMOVA $F_{ST}$ P-value	If bootstrap resampling is enabled, a p-value ranking the significance of $F_{ST}$ within this pair of populations.
21	Window SNP Count	Number of SNPs found in the sliding window centered on this nucleotide position.

**Notes:** The preferred version of  $F_{ST}$  is the **AMOVA  $F_{ST}$**  in column 17, or the corrected version in column 18 if you have specified a correction to the **populations** program (option -f)

# Stacks

## 7.4.6 batch\_X.phistats\_Y-Z.tsv: Haplotype-based $F_{ST}$ calculations for each pair of populations

Column	Name	Description
1	Batch ID	The batch identifier for this data set.
2	Locus ID	Catalog locus identifier.
3	Population ID 1	The ID supplied to the populations program, as written in the population map file.
4	Population ID 2	The ID supplied to the populations program, as written in the population map file.
5	Chromosome	If aligned to a reference genome.
6	Basepair	If aligned to a reference genome.
7	$\Phi_{ST}$	
8	Smoothed $\Phi_{ST}$	
9	Smoothed $\Phi_{ST}$ P-value	
10	$F_{ST}'$	
11	Smoothed $F_{ST}'$	
12	Smoothed $F_{ST}'$ P-value	
13	$D_{EST}$	
14	Smoothed $D_{EST}$	
15	Smoothed $D_{EST}$ P-value	

# Outline for today's lecture

---

RAD-seq for ecological and evolutionary genomics

Primer on Population Genomics

## **Evolutionary genomics of stickleback fish**

- Population genomics of rapid adaptation
- Using haplo-RAD-seq for coalescent analyses
- Genome Wide Association Studies using RAD-seq

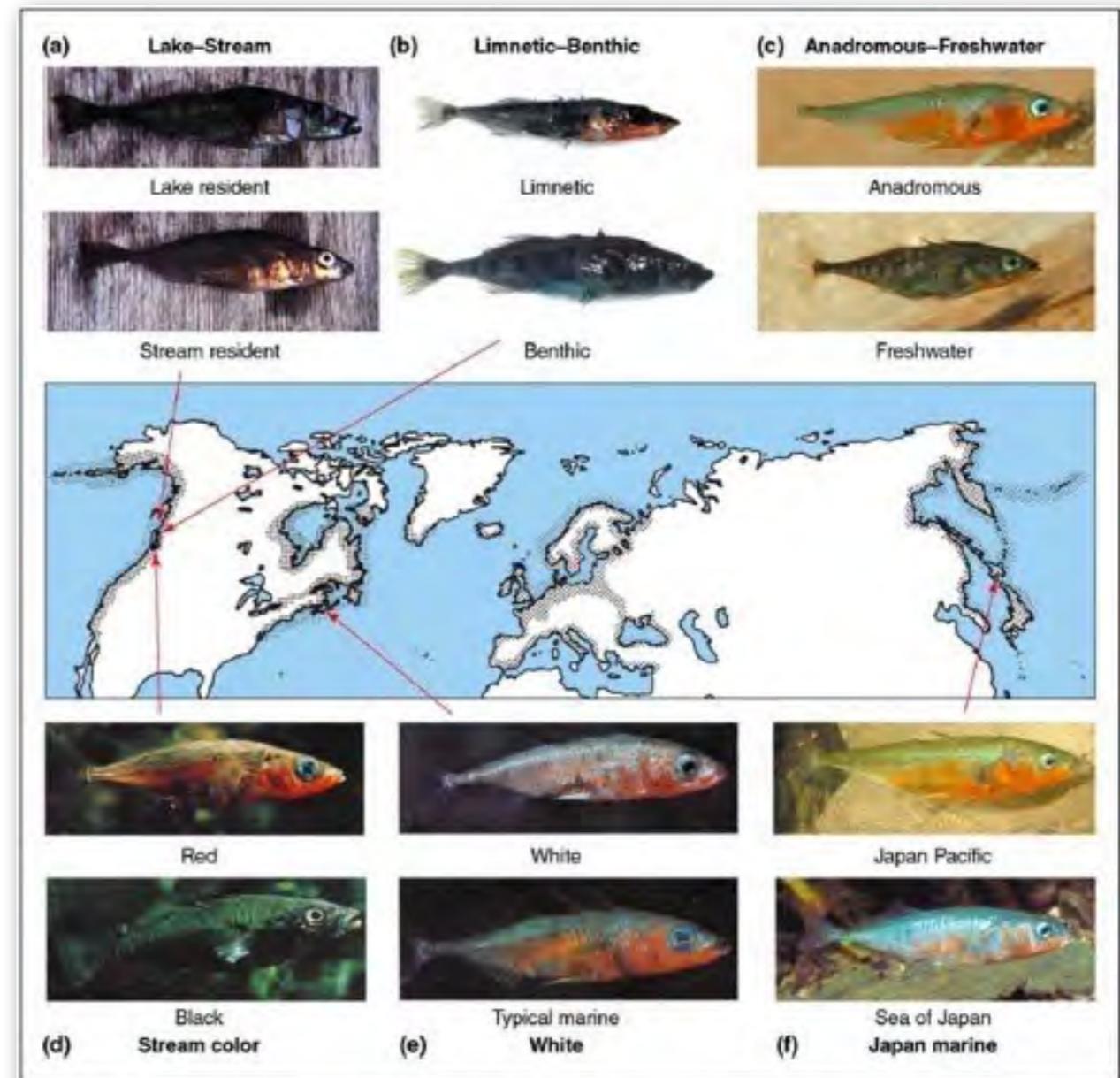
Genomically enabling the Gulf pipefish

# Threespine stickleback, *Gasterosteus aculeatus*

---



# Threespine stickleback, *Gasterosteus aculeatus*



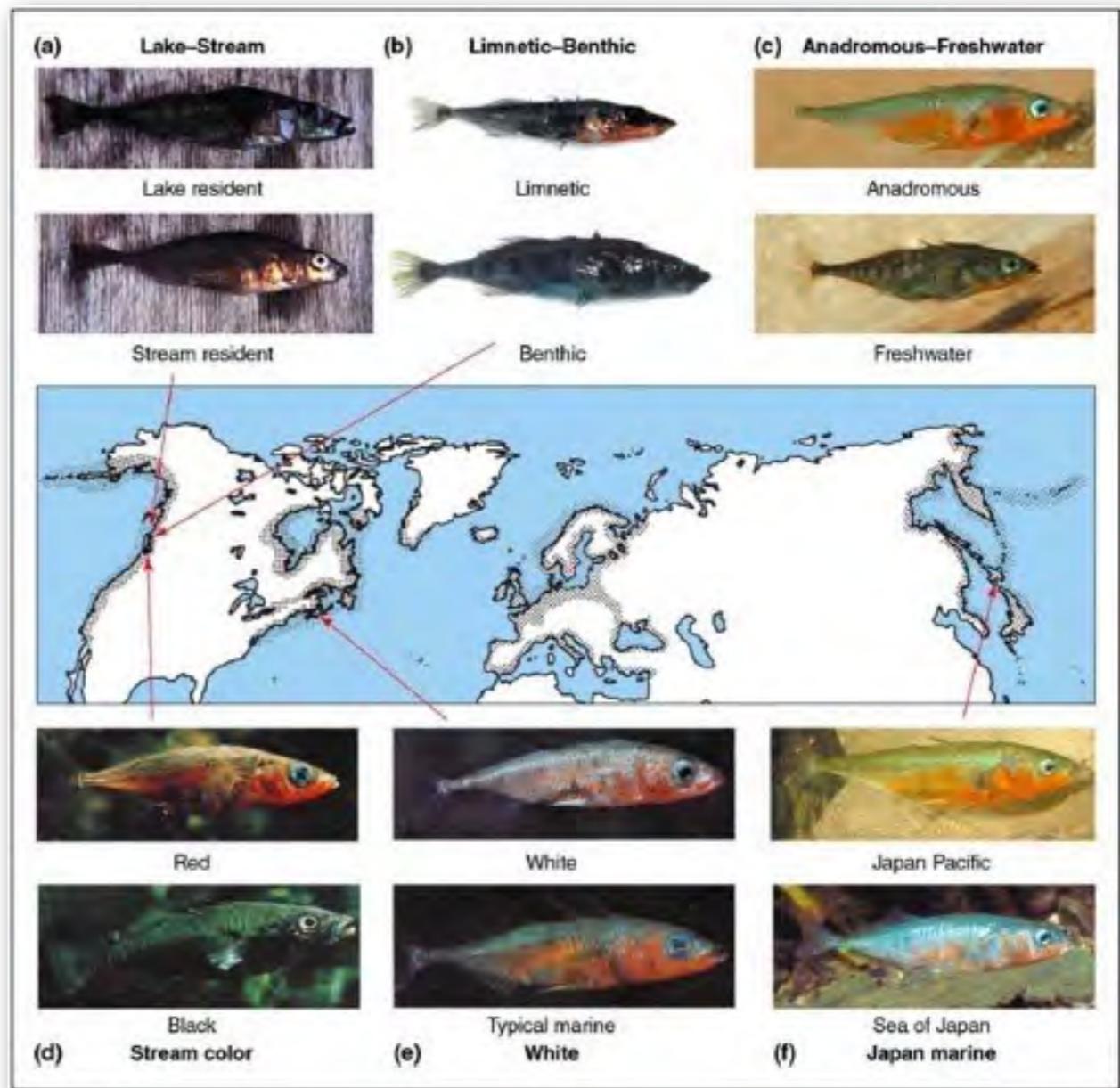
Rundle and McKinnon 2002

# Threespine stickleback, *Gasterosteus aculeatus*

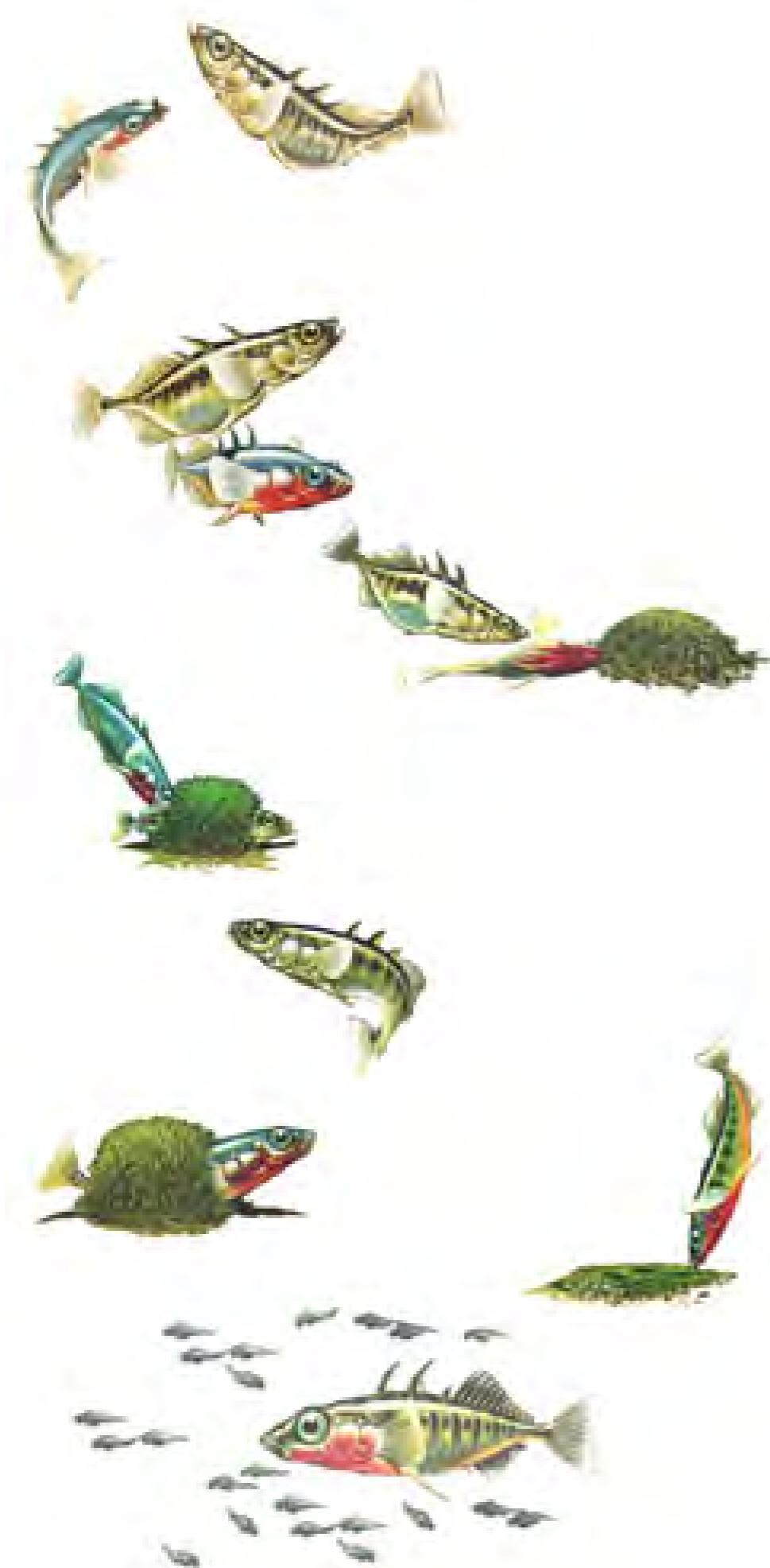
Pelvic  
Structure



Lateral  
Plates



Rundle and McKinnon 2002



# Stickleback phenotypes mapped in the lab so far....

---

Pelvic structure size and shape \*\*\* (*Eda*)

Lateral plate number \*\*\* (*Pitx1*)

Body coloration \*\*\* (*KitL*)

Opercle bone shape

Pelvic spine length

Body shape

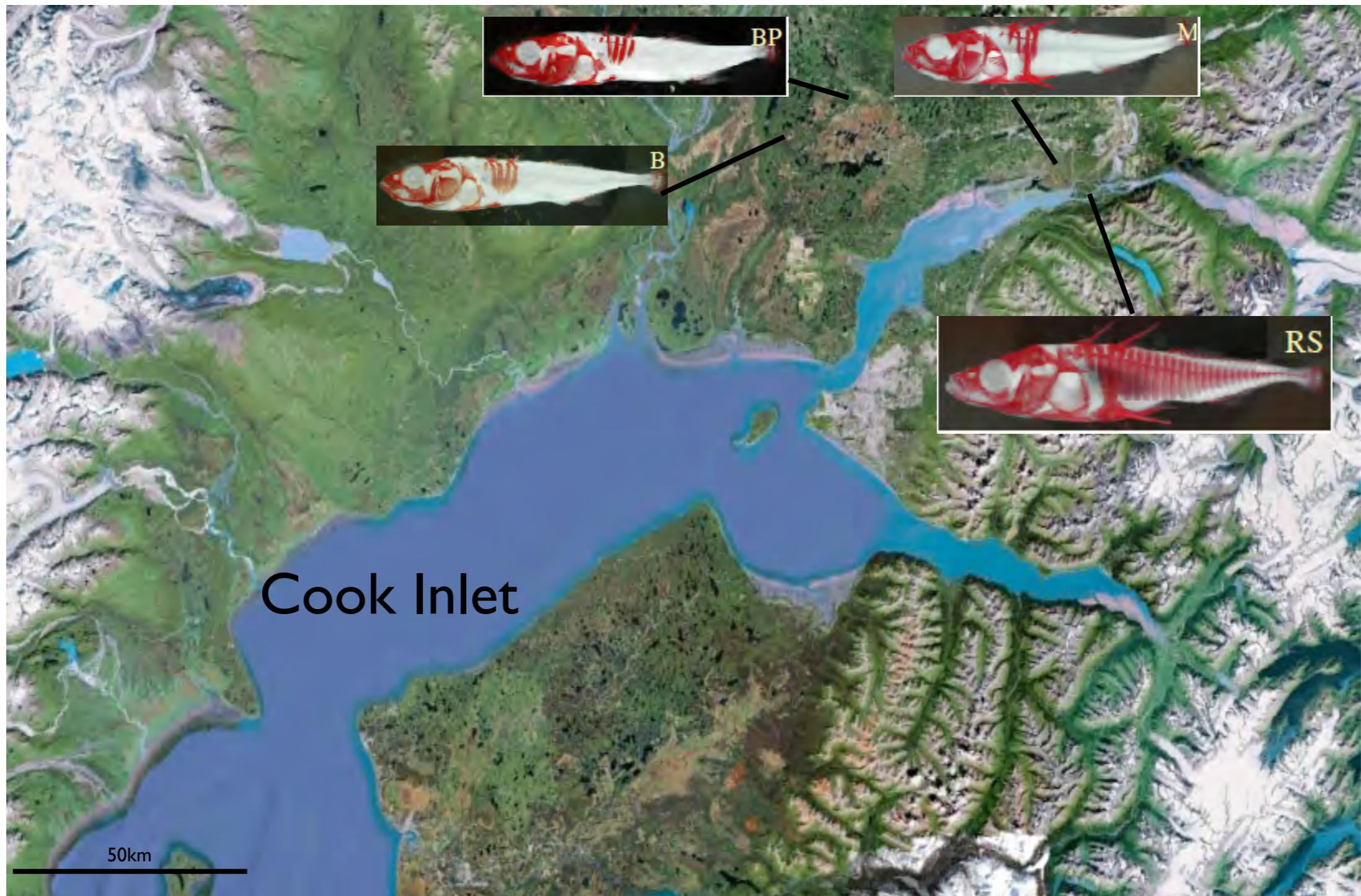
Courtship behavior

Gill raker size

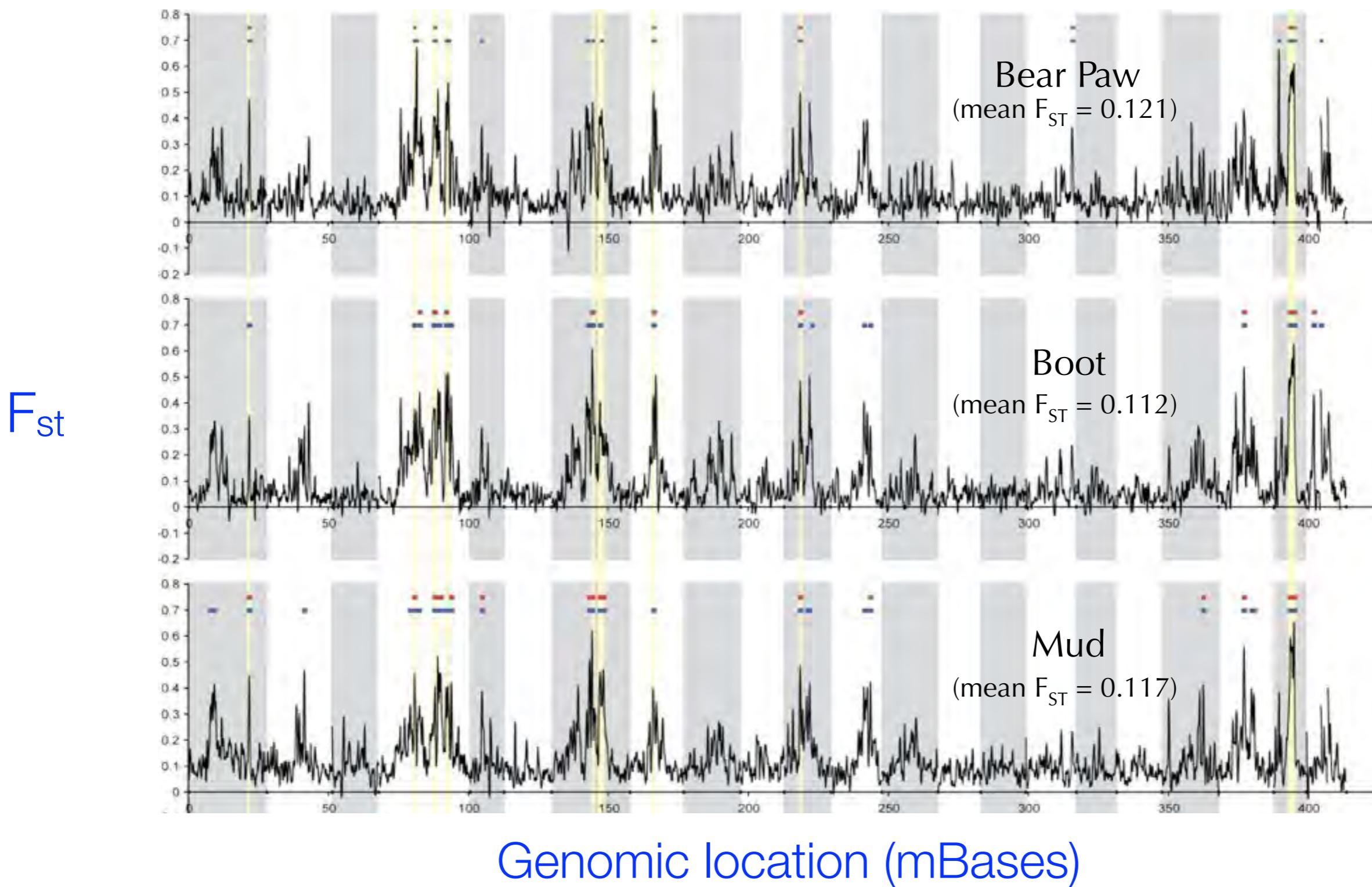
Dorsal spine length

- 
- A trend of large effect loci identified in the laboratory
  - Similar genomic regions and sometimes alleles mapped in independent populations
  - A question is whether population genomics studies can provide complementary and more complete information.

# Signatures of natural selection in 13,000 years

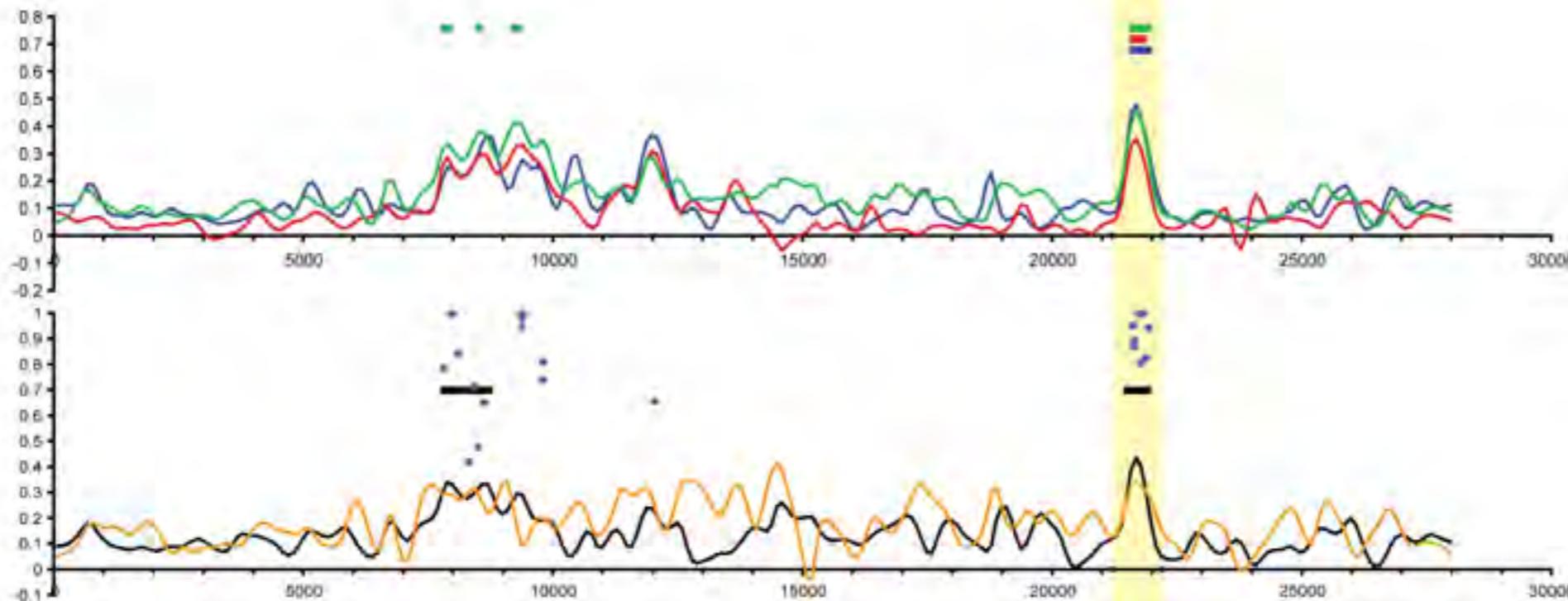


# Signatures of natural selection in 13,000 years



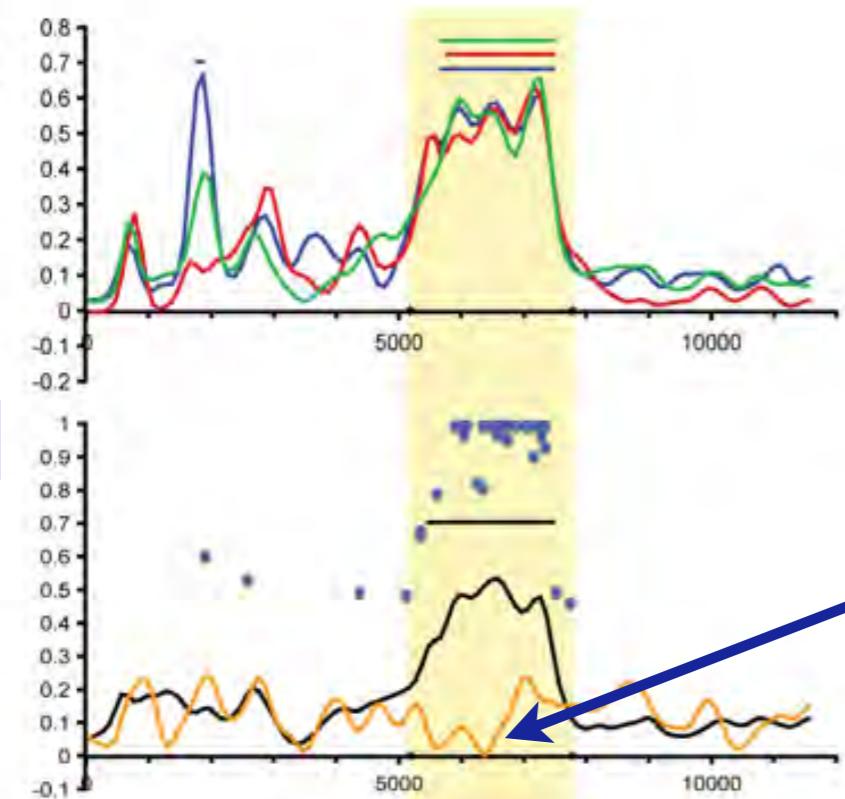
# Numerous novel regions identified

LGI



Different alleles

LGXXI



More often the same alleles

Emily  
Lescak



Julian  
Catchen



Susan  
Bassham



Mary  
Sherbick



Frank  
von Hippel

# Evolution of stickleback in 50 years on earthquake-uplifted islands

Emily A. Lescak<sup>a,b</sup>, Susan L. Bassham<sup>c</sup>, Julian Catchen<sup>c,d</sup>, Ofer Gelson<sup>b,1</sup>, Mary L. Sherbick<sup>b</sup>, Frank A. von Hippel<sup>b</sup>, and William A. Cresko<sup>c,2</sup>

<sup>a</sup>School of Fisheries and Ocean Sciences, University of Alaska Fairbanks, Fairbanks, AK 99775; <sup>b</sup>Department of Biological Sciences, University of Alaska Anchorage, Anchorage, AK 99508; <sup>c</sup>Institute of Ecology and Evolution, University of Oregon, Eugene, OR 97403; and <sup>d</sup>Department of Ecology and Evolution, University of Illinois at Urbana-Champaign, Urbana, IL 61801

Edited by John C. Avise, University of California, Irvine, CA, and approved November 9, 2015 (received for review June 19, 2015)



# Middleton Island - 50 year old populations

1955



2008

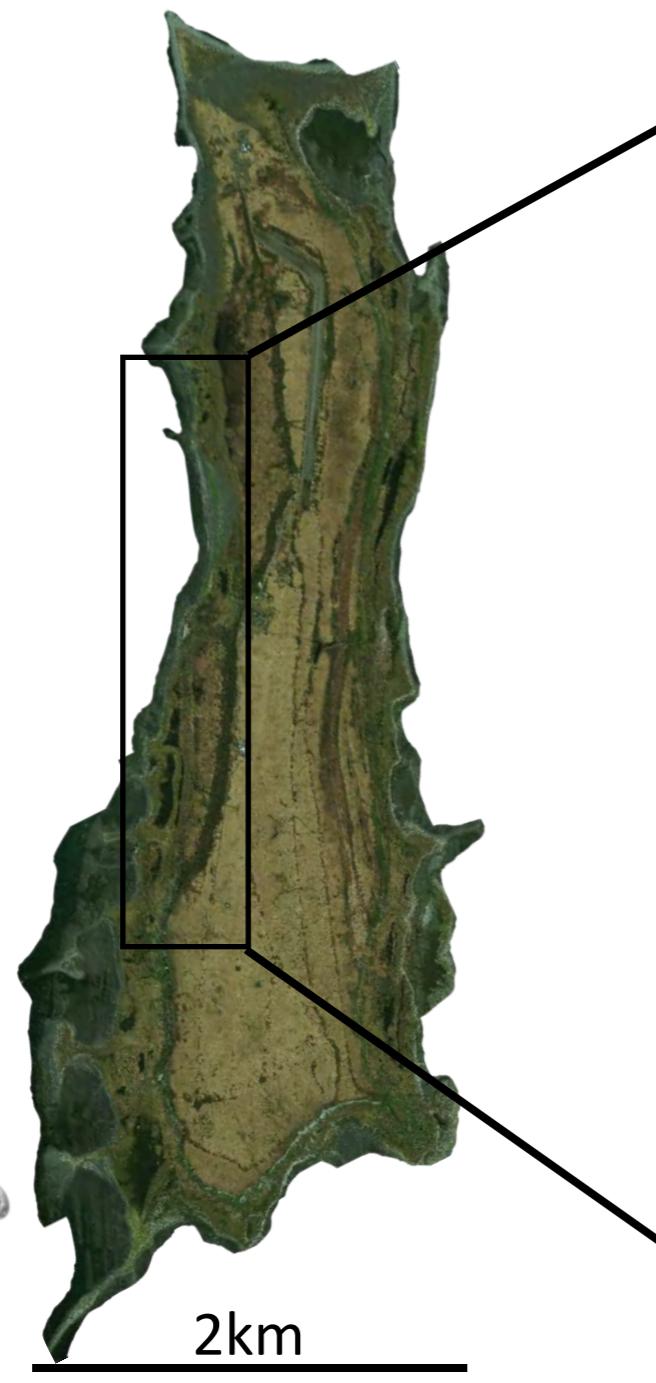


# Middleton Island - 50 year old populations

1955



2008

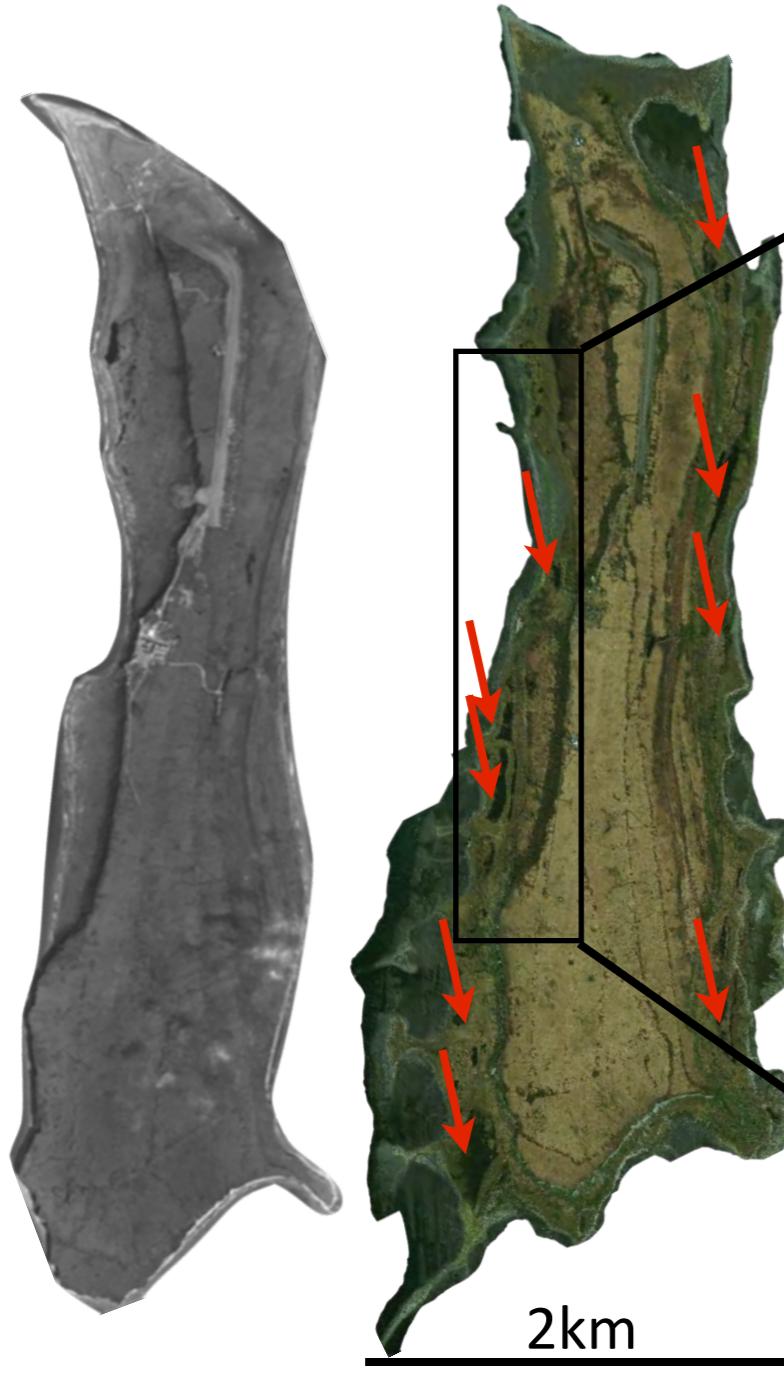


# Middleton Island - 50 year old locations

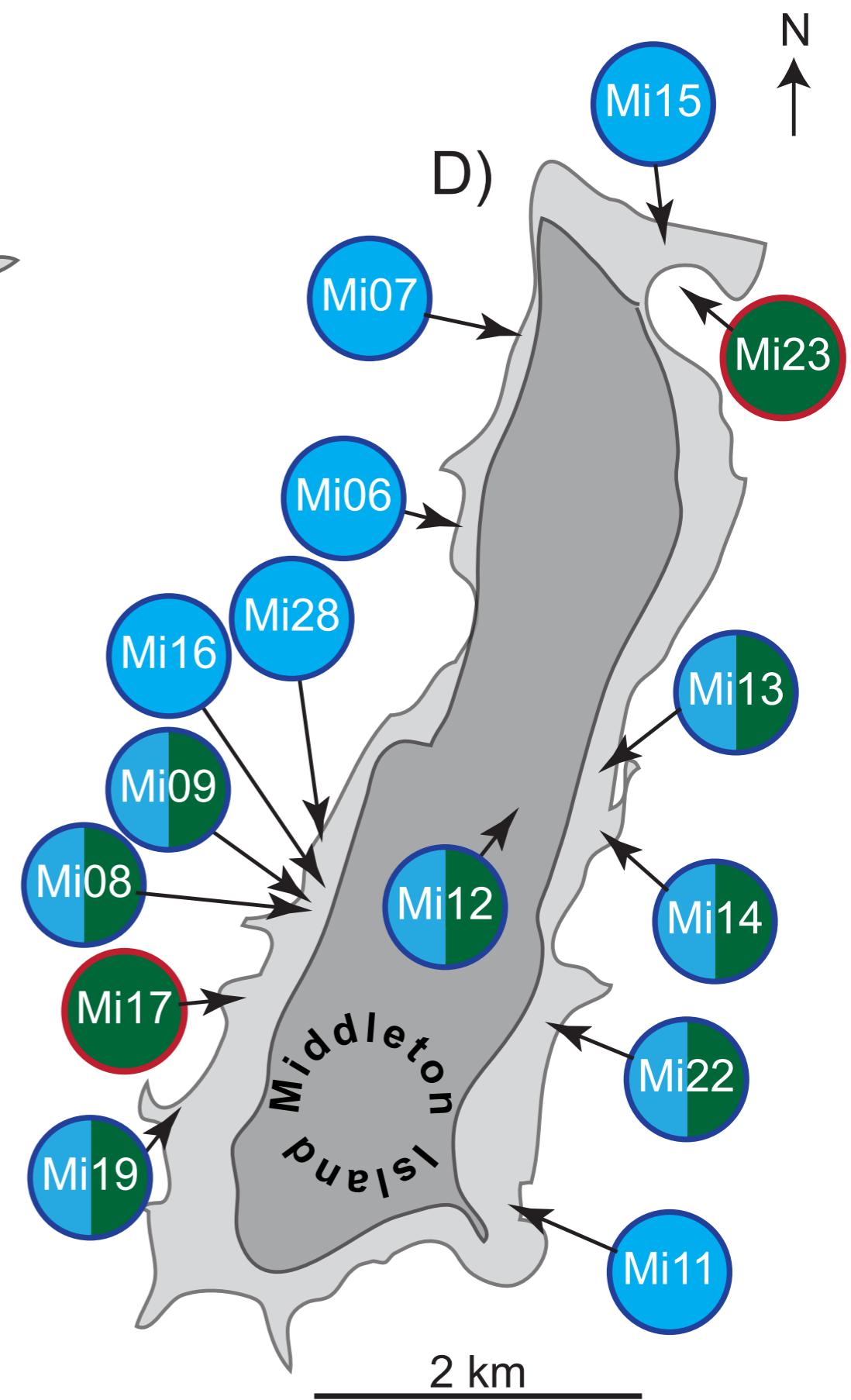
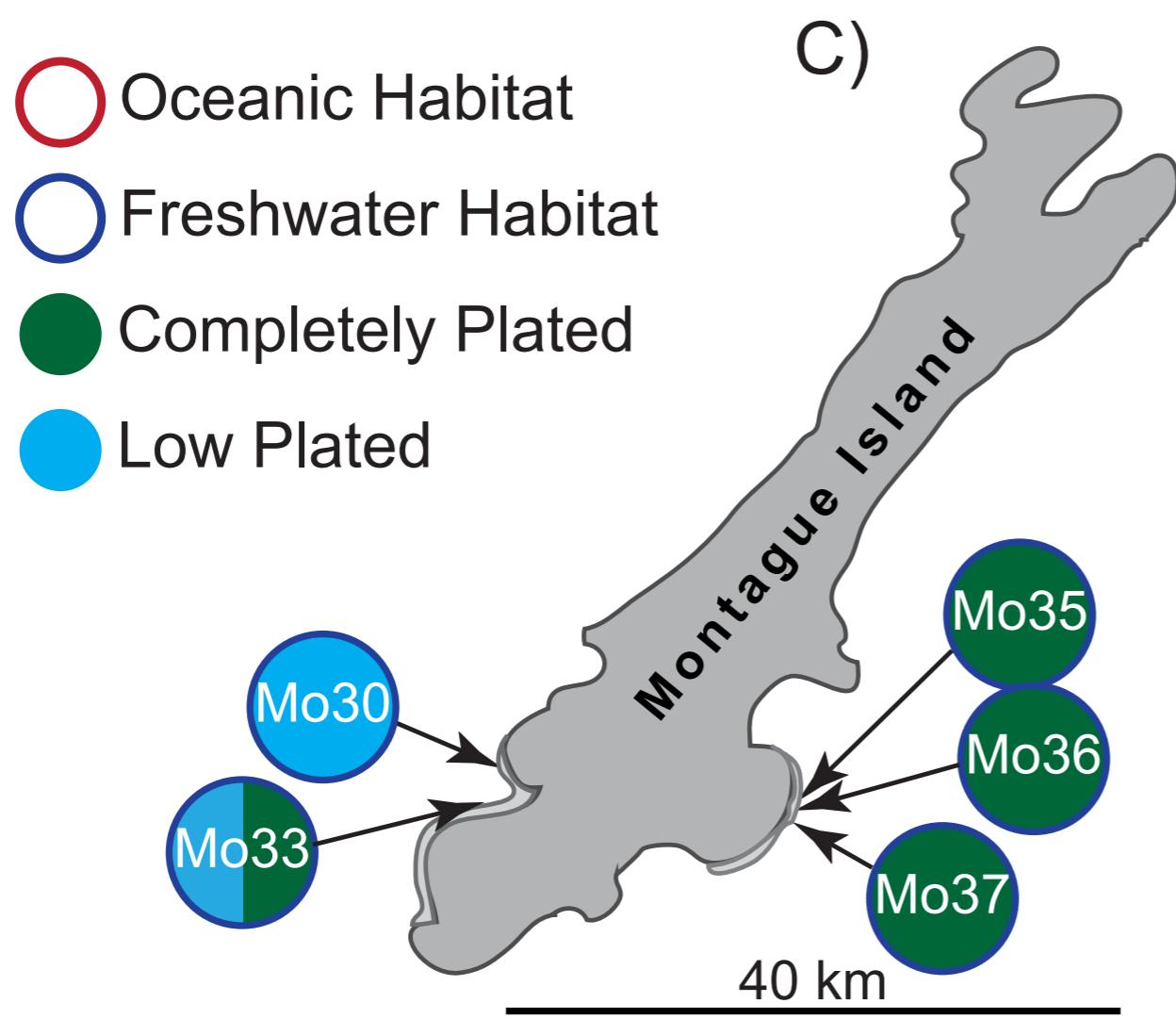
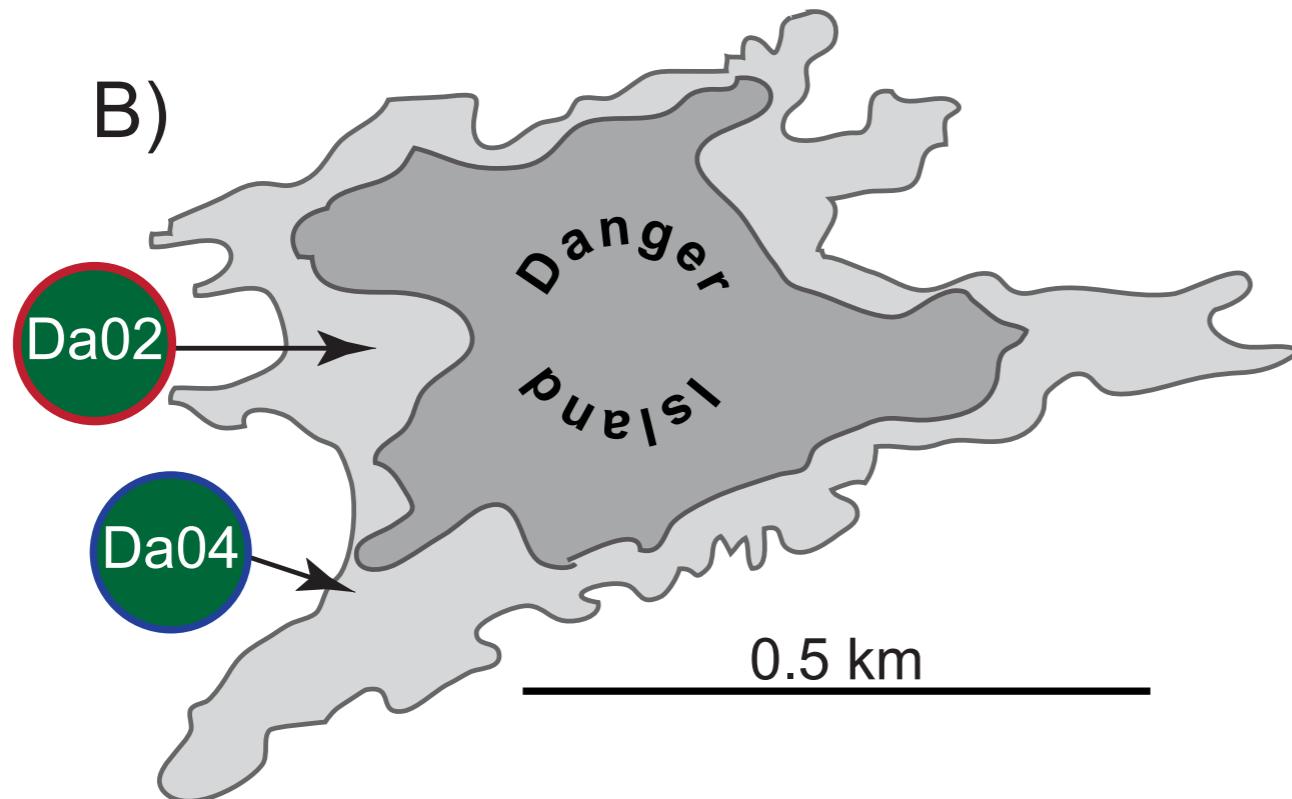
1955

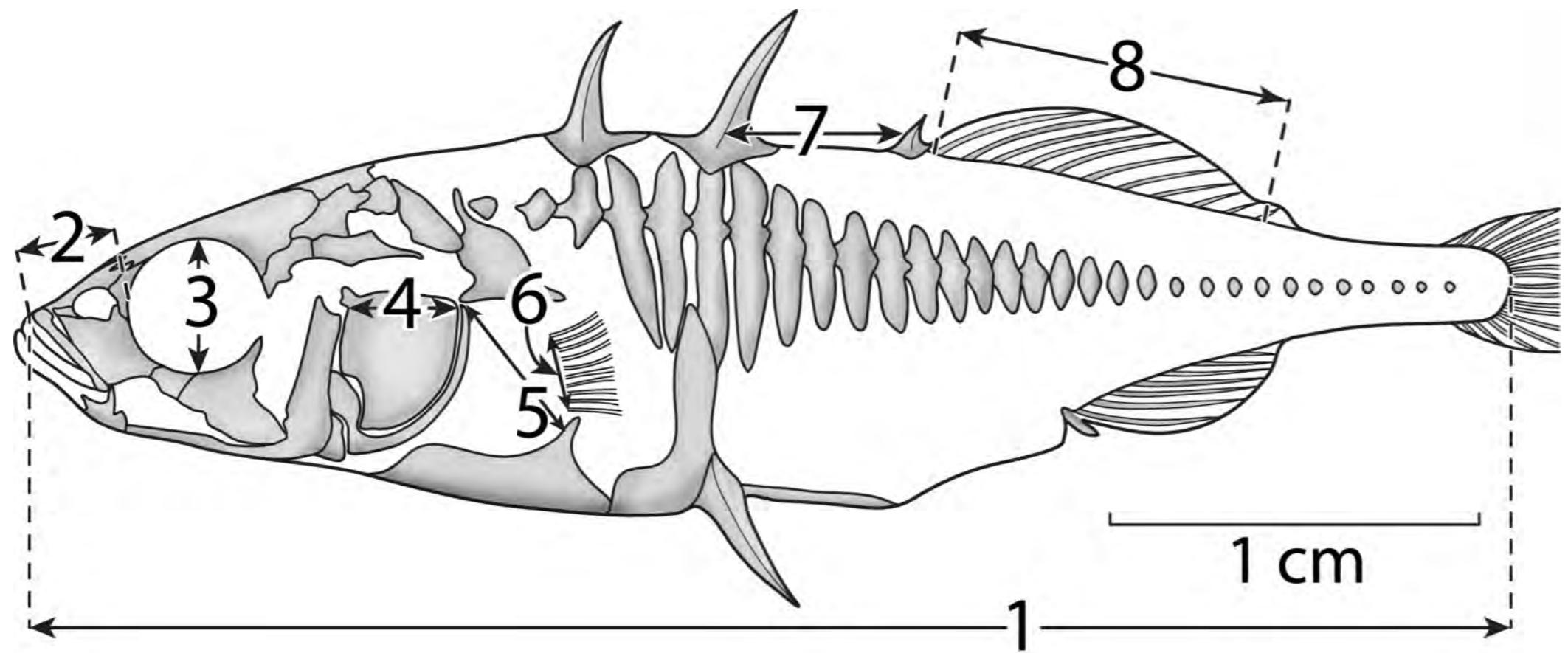
1955

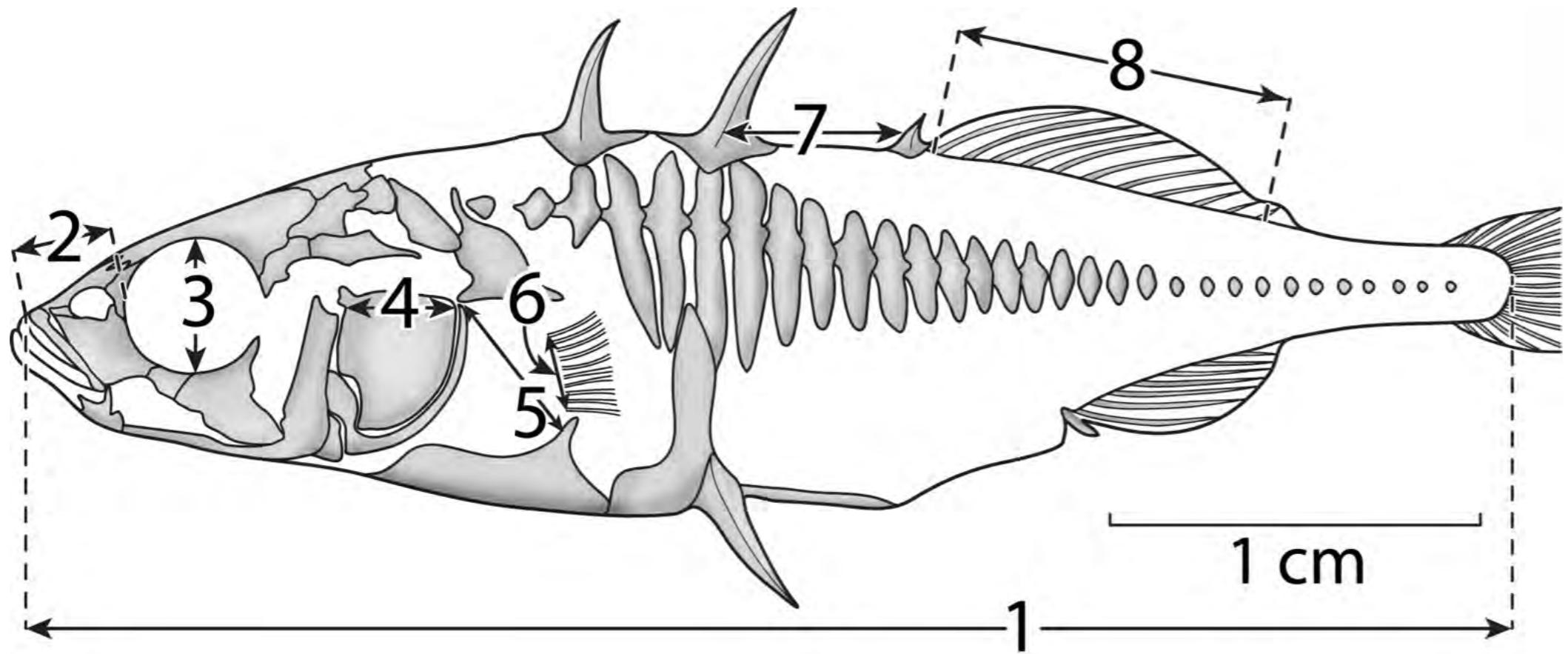
2008







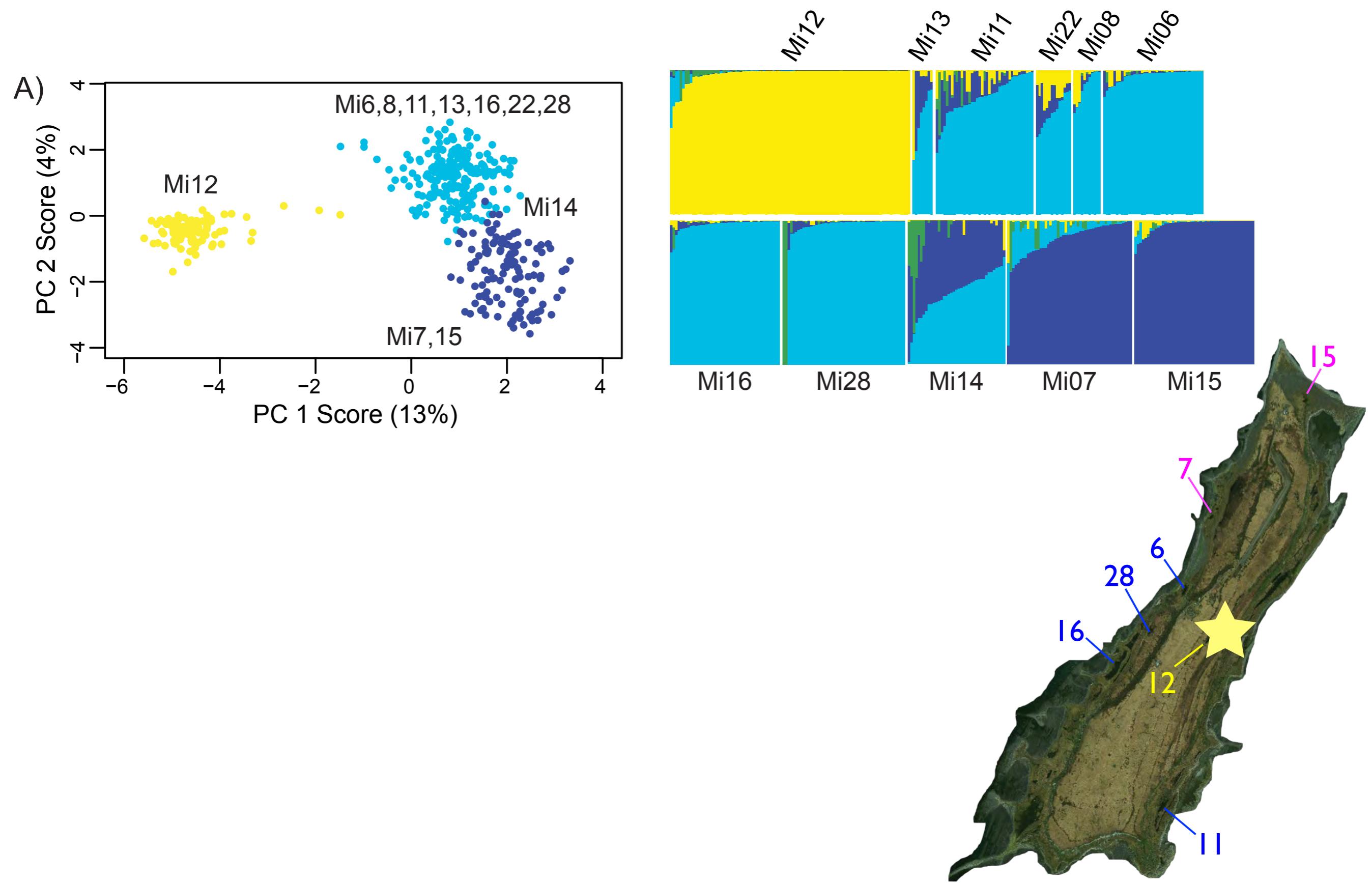




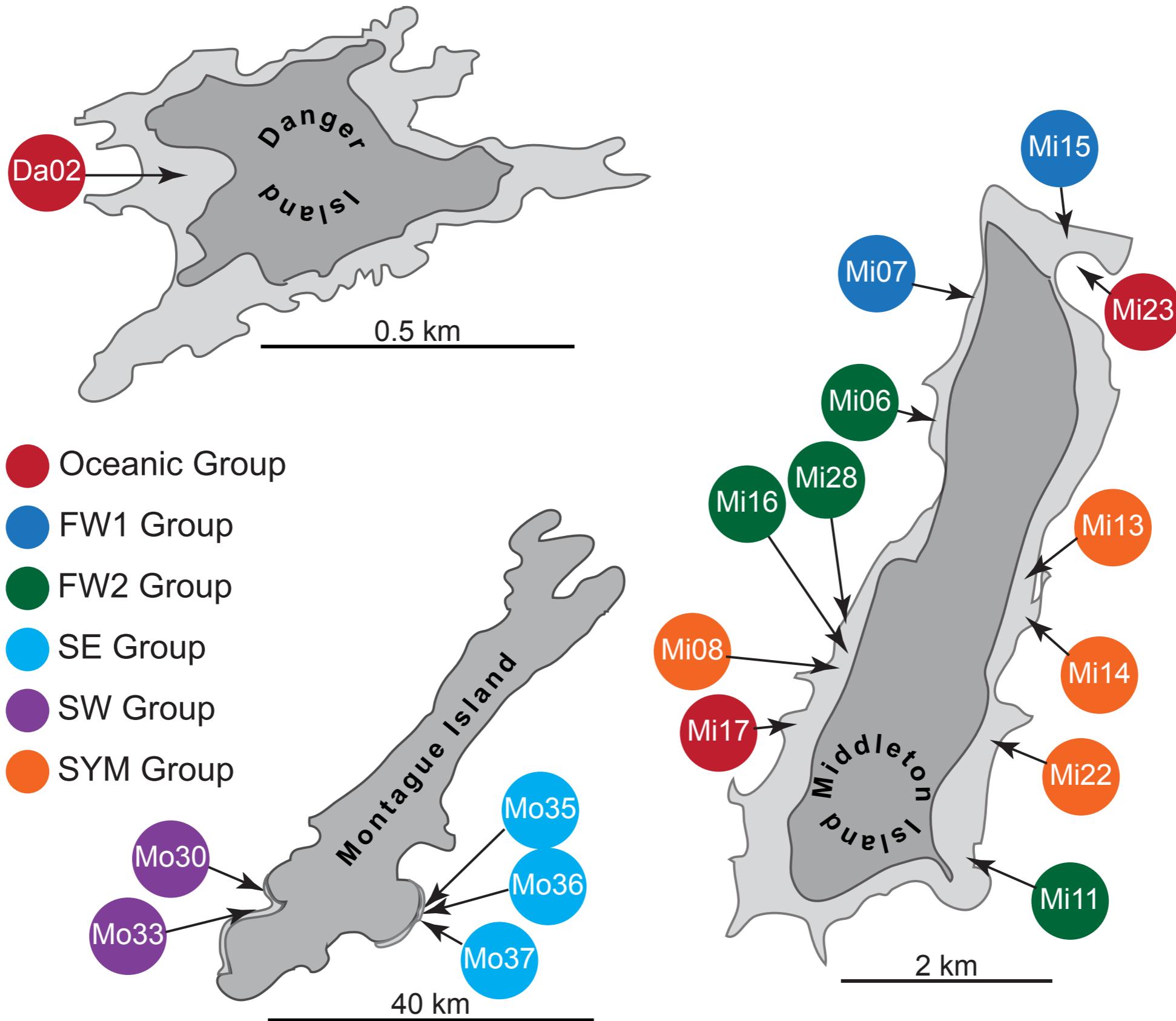
## RAD-seq analysis

110,000 SNPs per individual  
>1000 Individuals  
20 million genotypes

# Structure analysis shows independent evolution even among populations on a single island



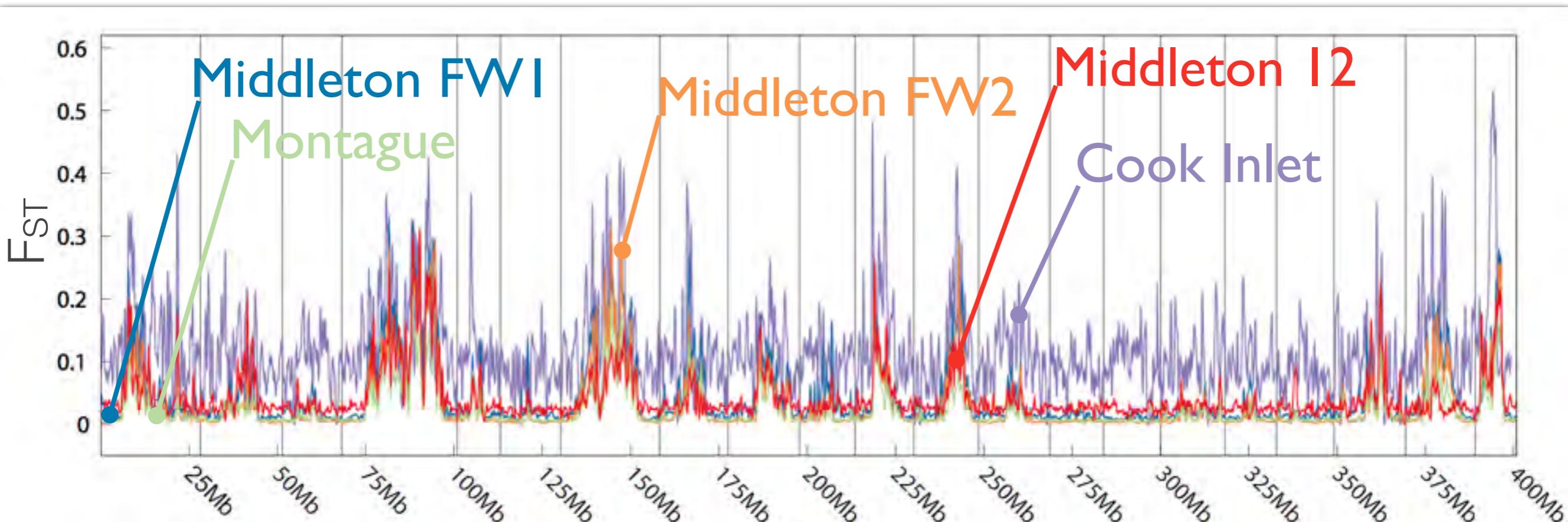
# At least six independent evolutionary events in freshwater in the last 50 years



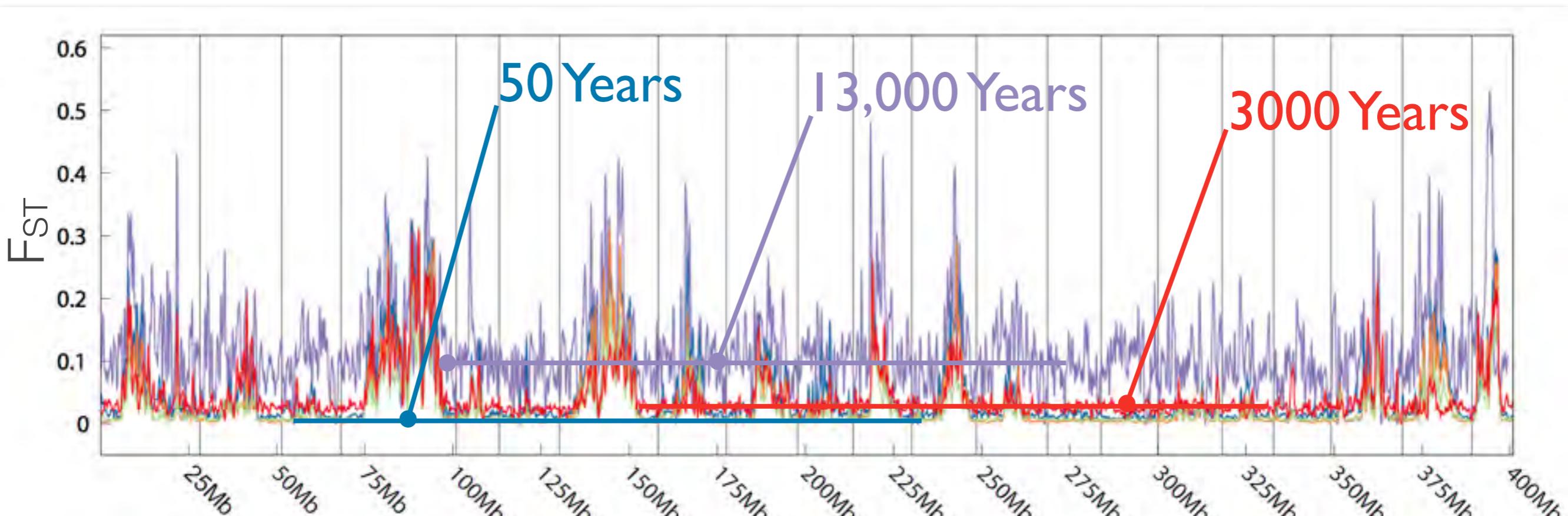
How much of the genome is differentiated?

---

How similar are the genomic patterns of differentiation?

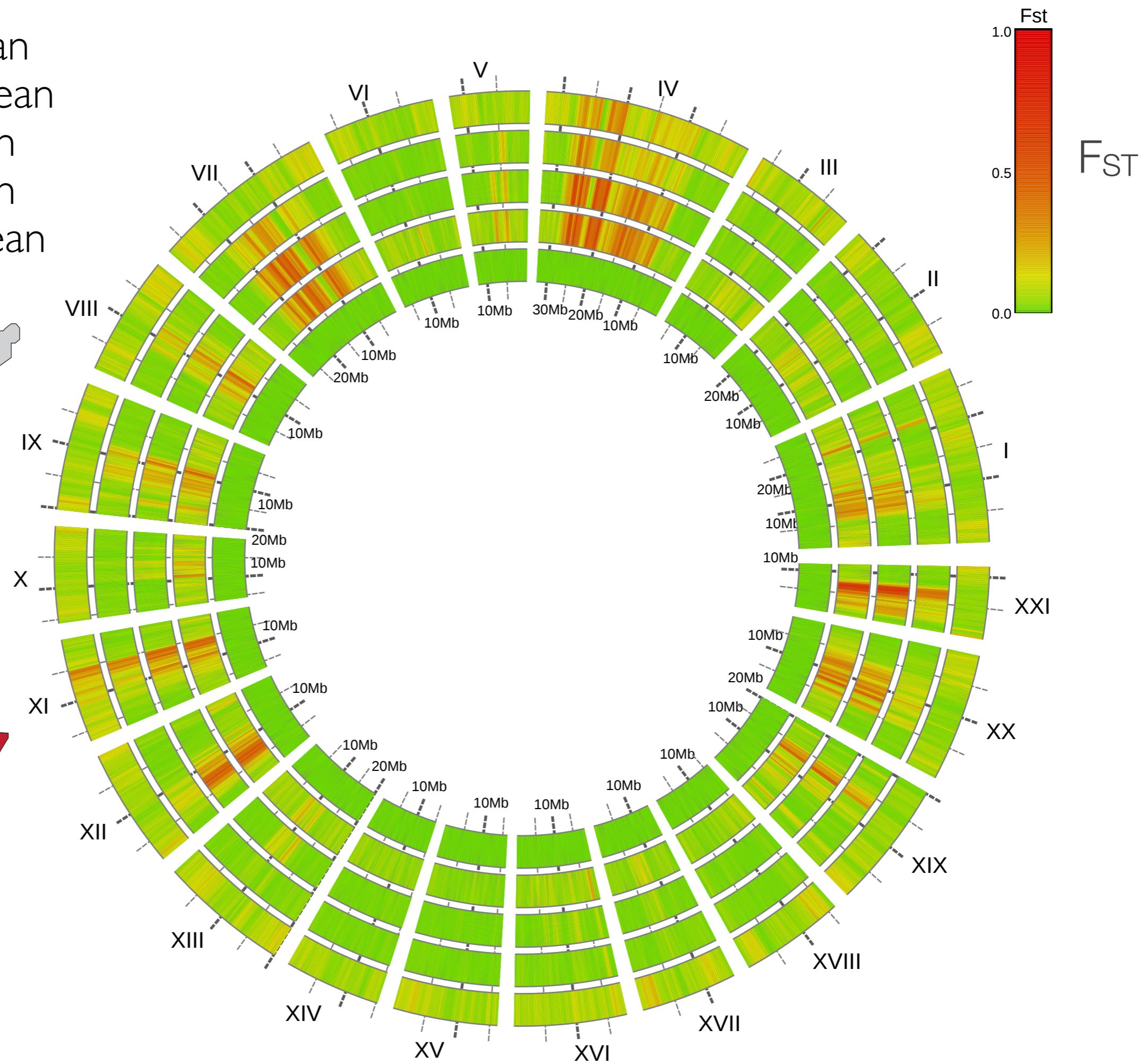
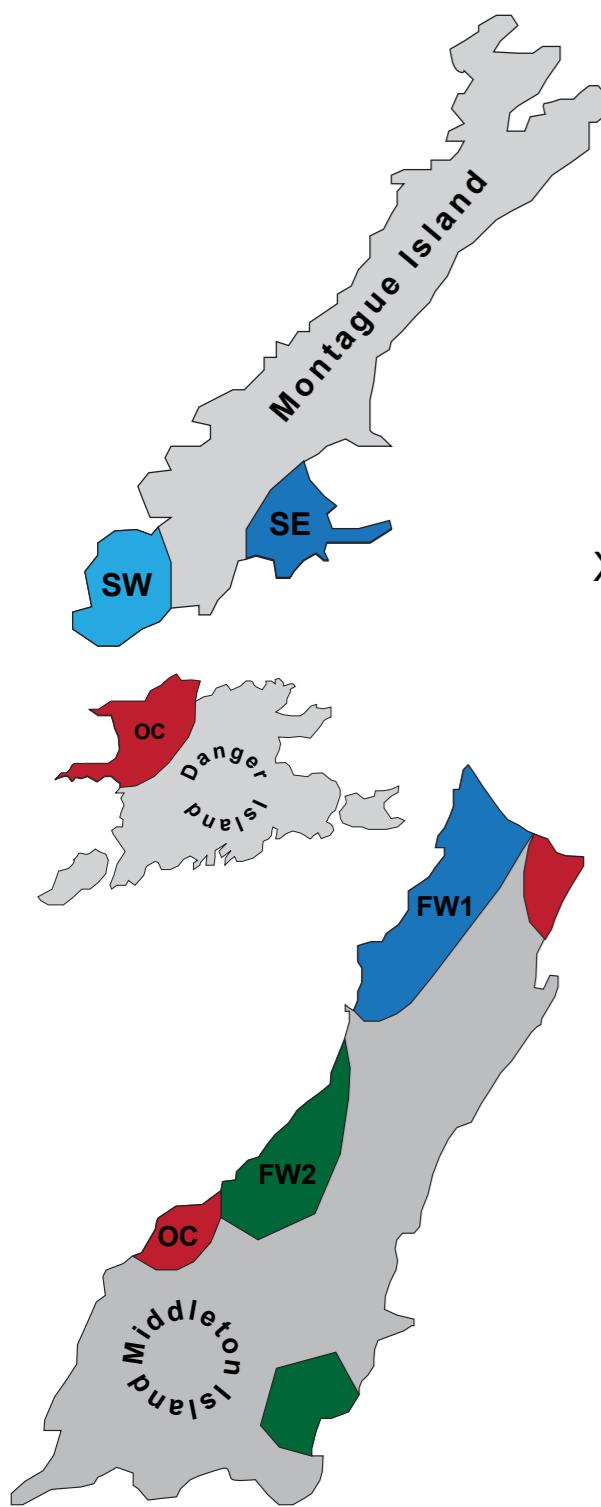


Freshwater Populations (grouped) vs. Marine

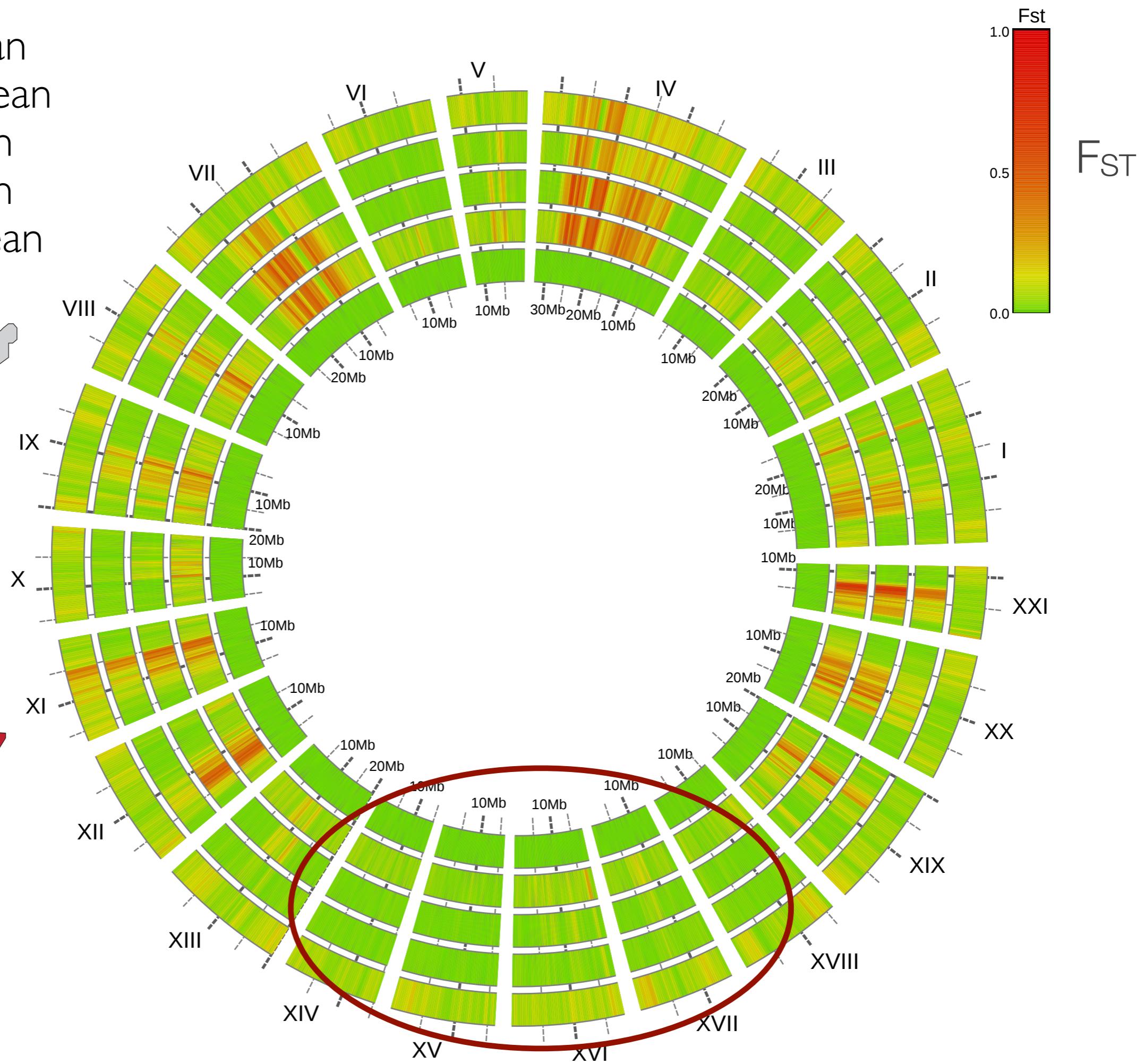
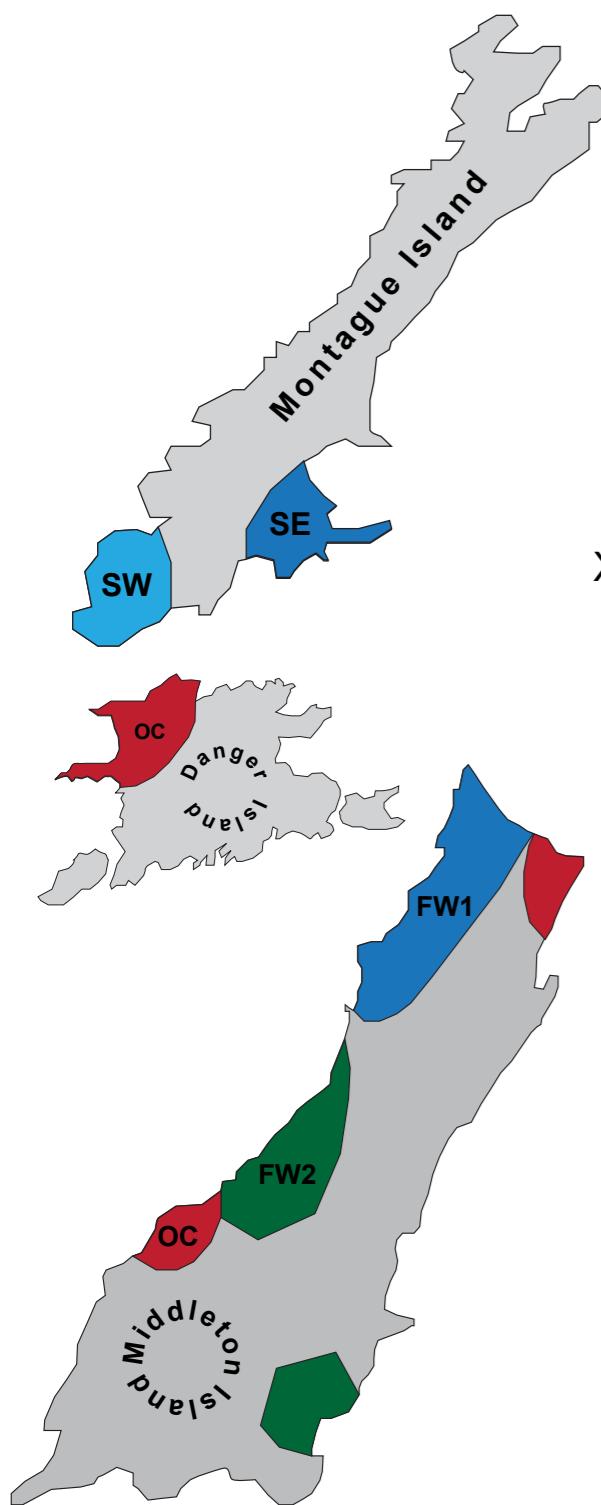


Freshwater Populations (grouped) vs. Marine

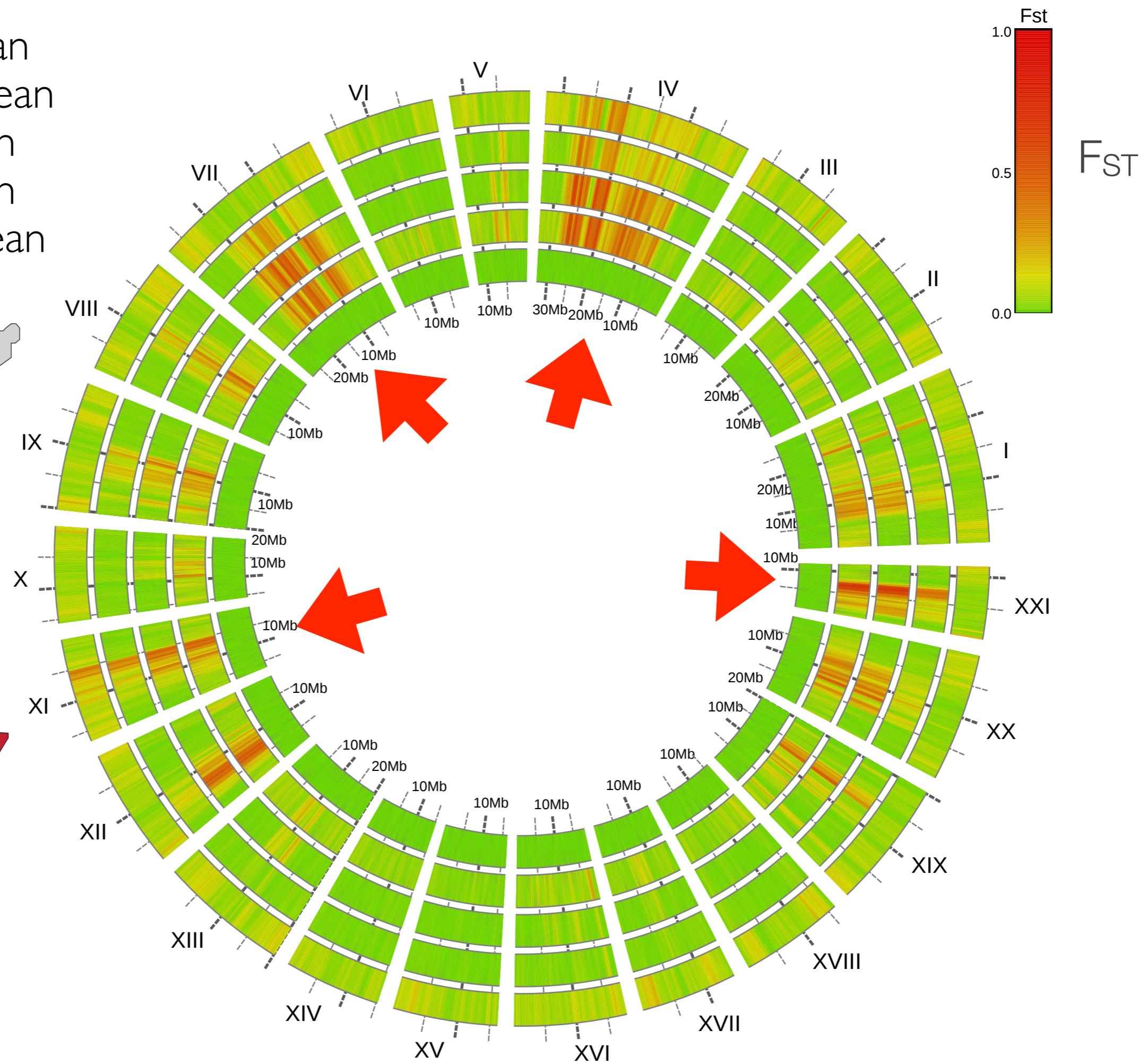
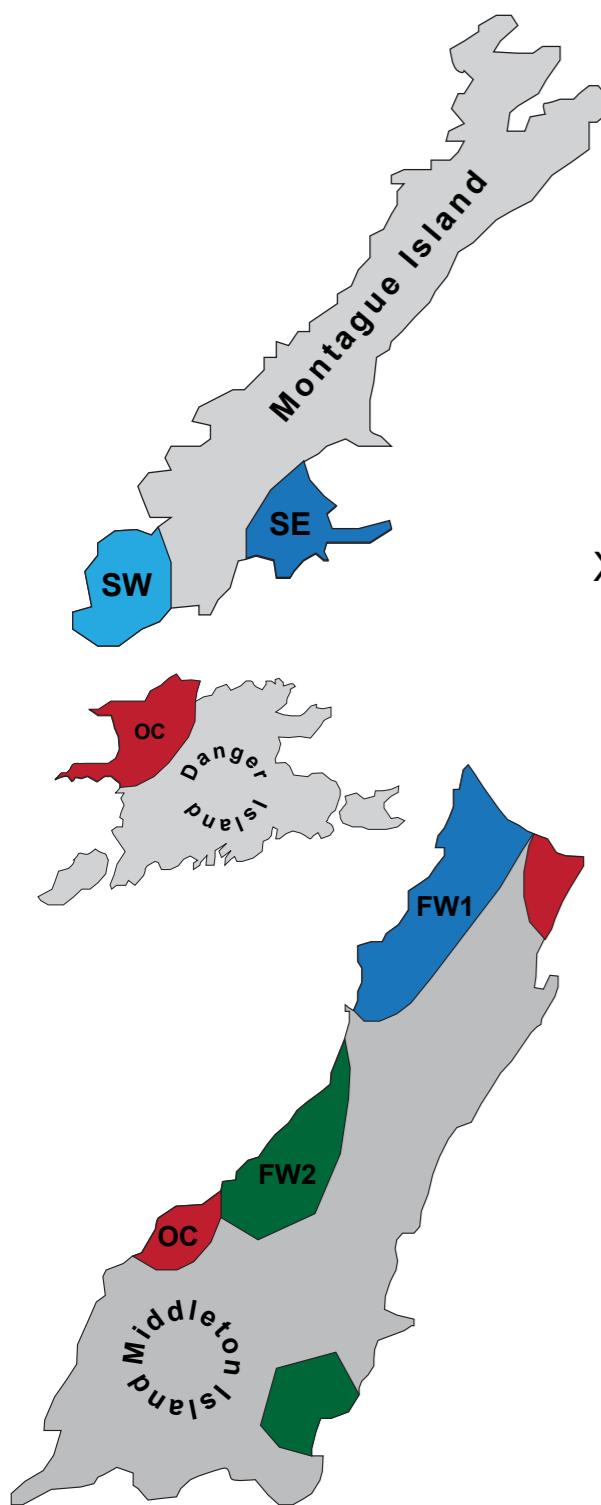
MoSe vs. Ocean  
MoSW vs. Ocean  
FW2 vs. Ocean  
FW1 vs. Ocean  
Ocean vs. Ocean



MoSe vs. Ocean  
MoSW vs. Ocean  
FW2 vs. Ocean  
FW1 vs. Ocean  
Ocean vs. Ocean



MoSe vs. Ocean  
MoSW vs. Ocean  
FW2 vs. Ocean  
FW1 vs. Ocean  
Ocean vs. Ocean



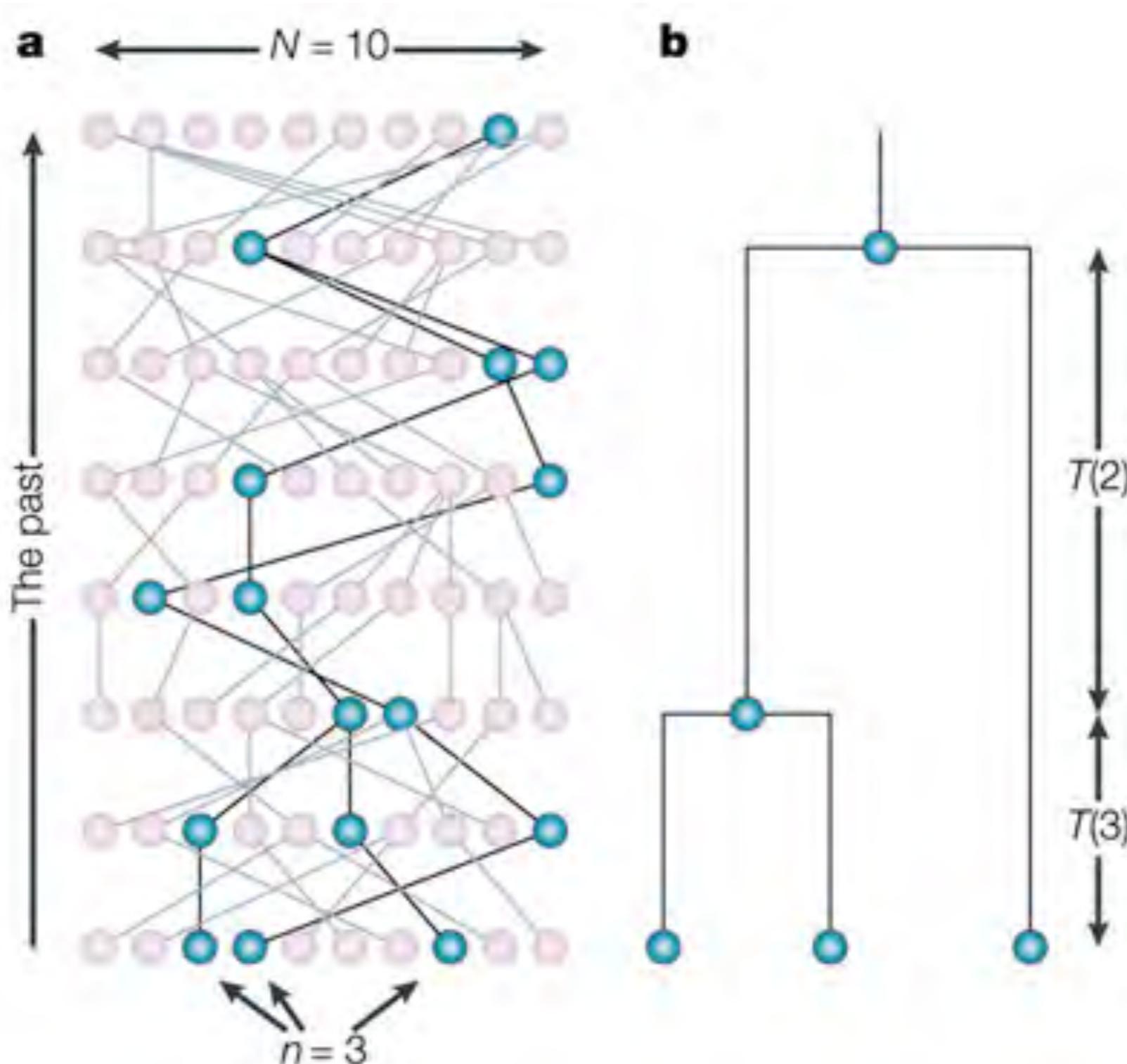
# Coalescent analyses using RADseq

---



Thom Nelson

# The coalescent in population genetics



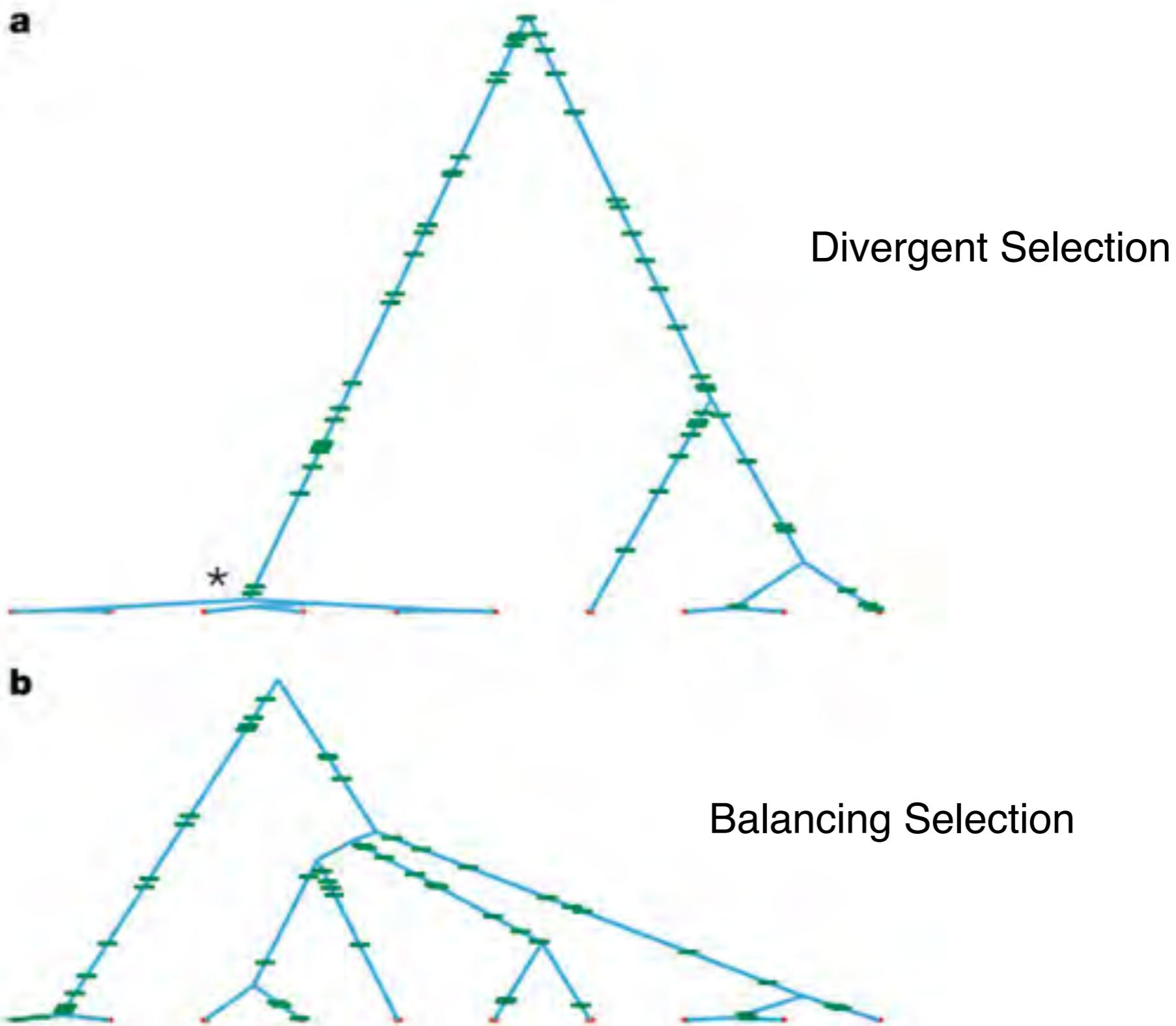
# Neutral coalescent expectations

---

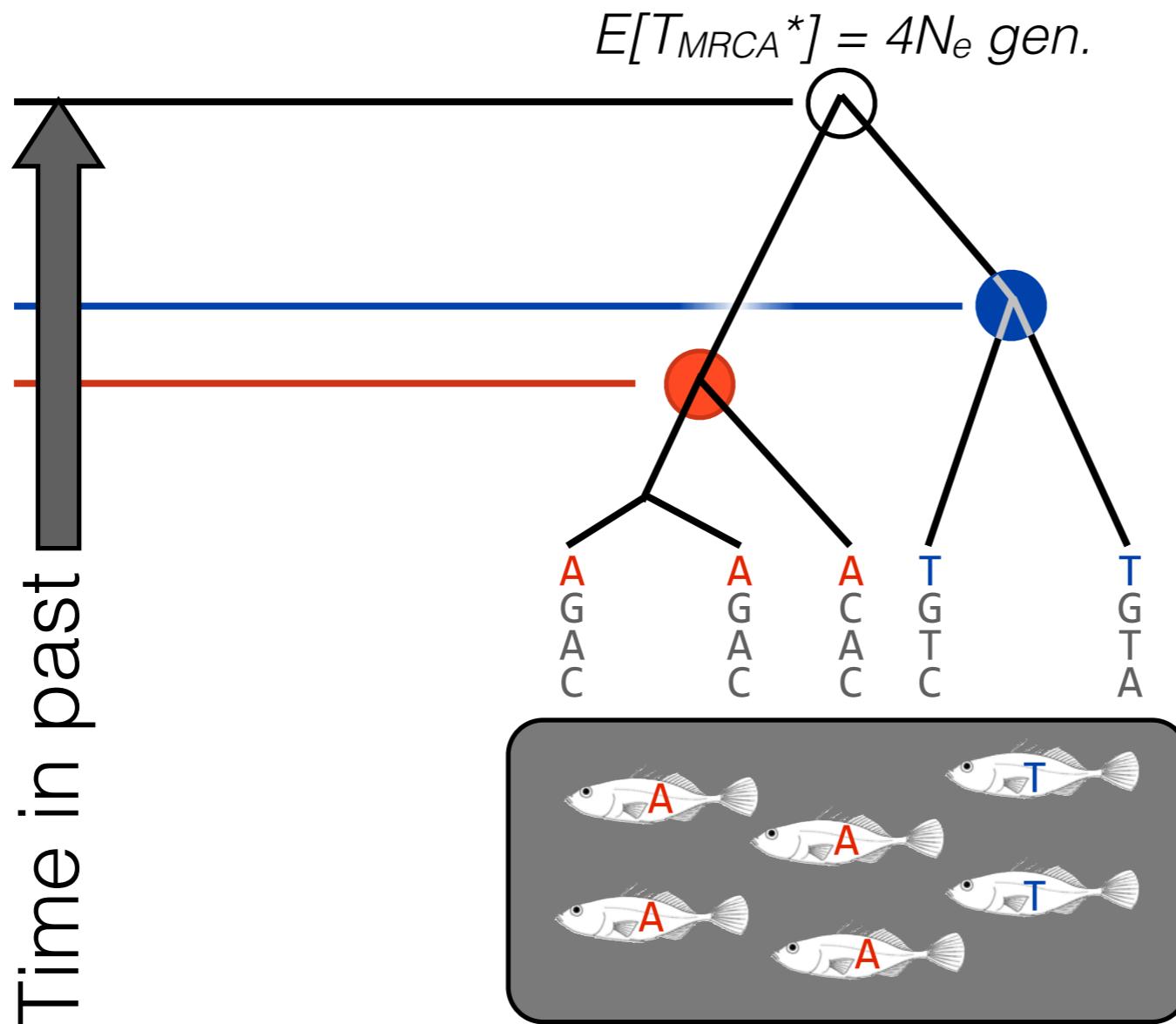


# Natural selection and the coalescent

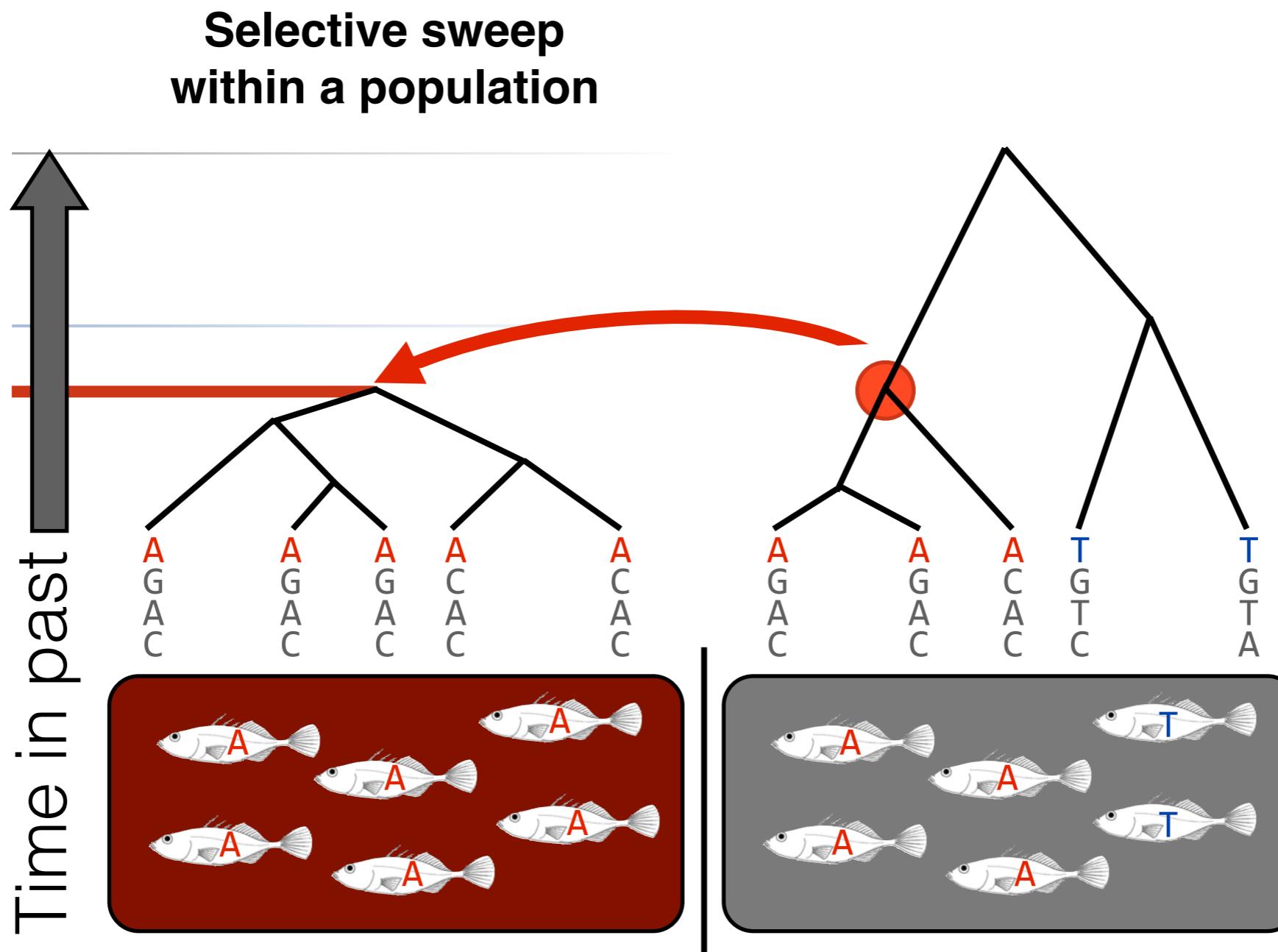
---



# Coalescent analyses in stickleback



# Coalescent analyses in stickleback

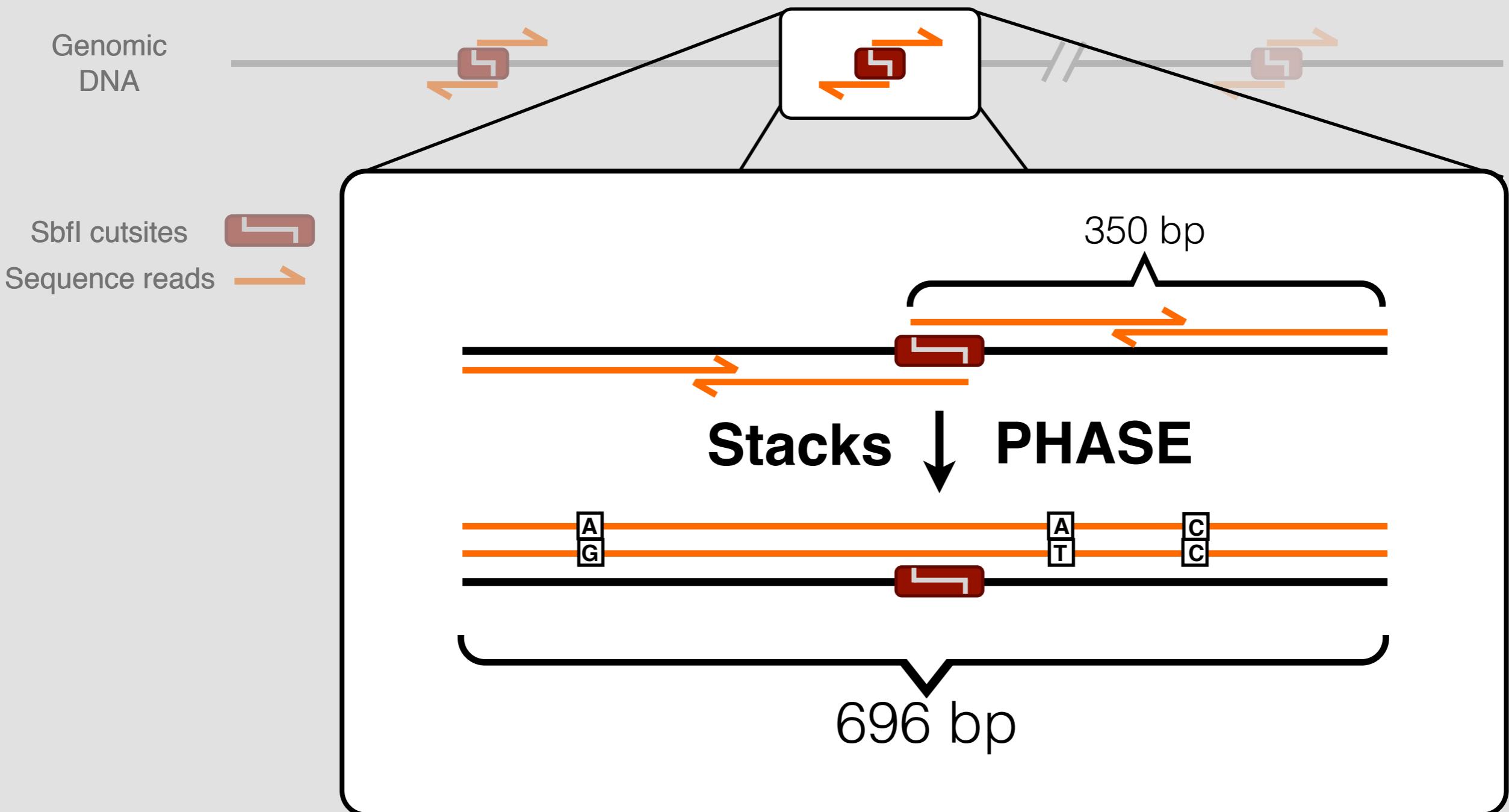


# Coalescent analyses with RAD-seq

---

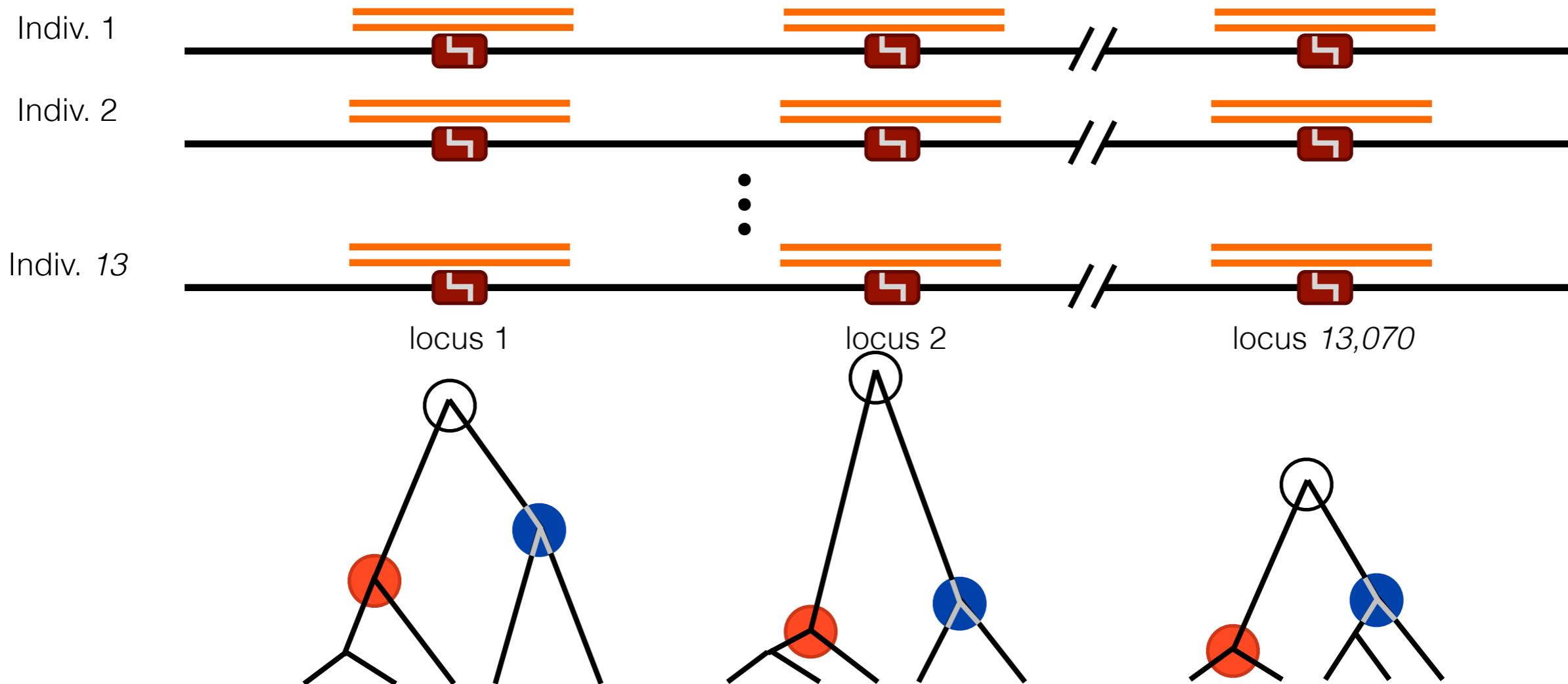


# Coalescent analyses with RAD-seq

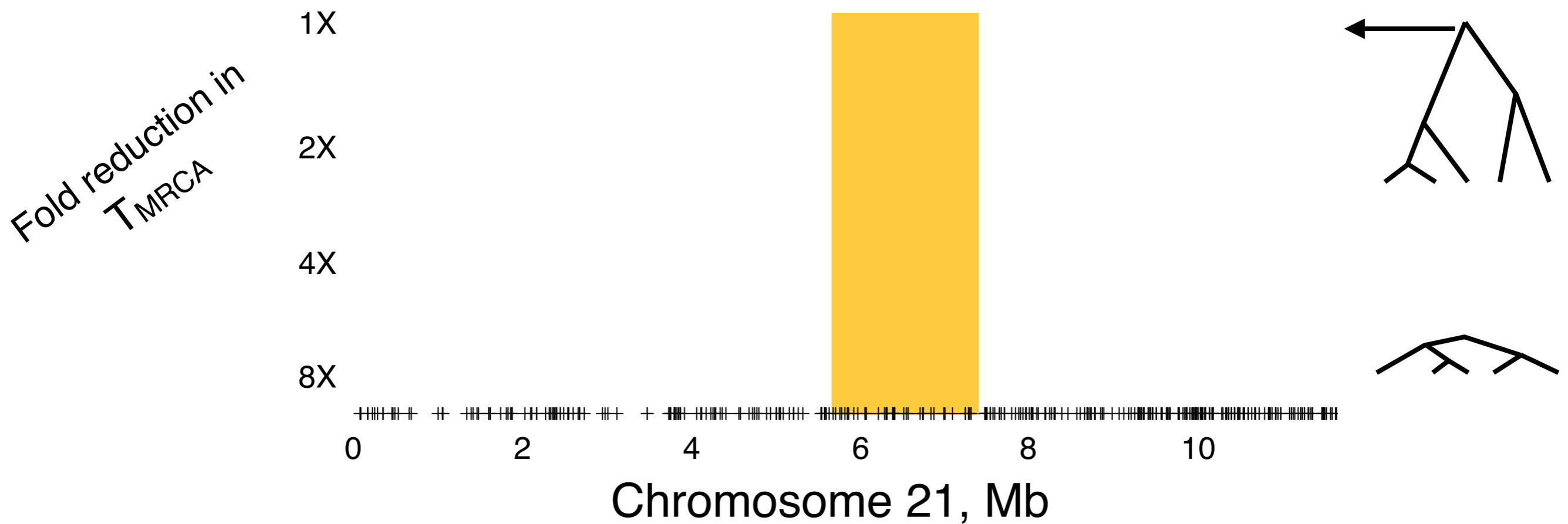


# Coalescent analyses with RAD-seq

---

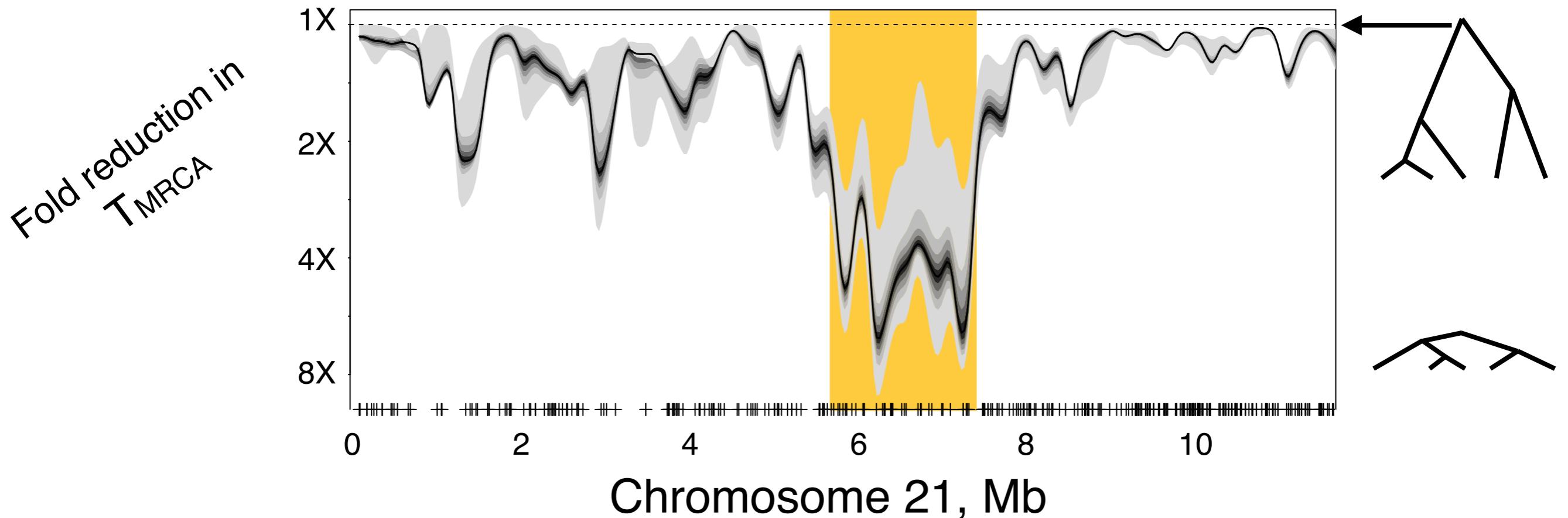


# Selection in one population reduces coalescence time



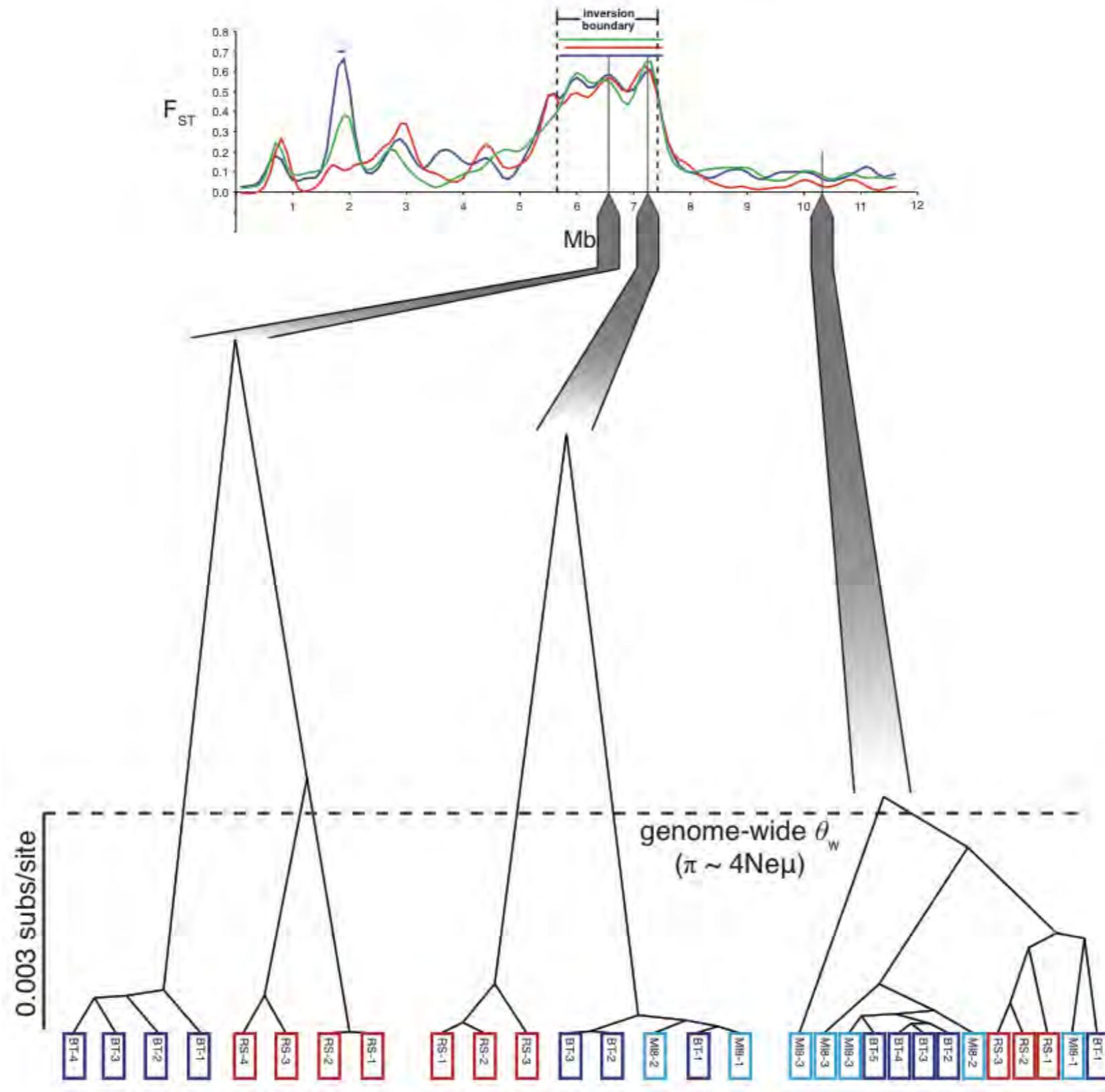
$$\log_2 \left( \frac{T_{\text{MRCA}} \text{FW}}{T_{\text{MRCA}} \text{ALL}} \right)$$

# Selection in one population reduces coalescence time



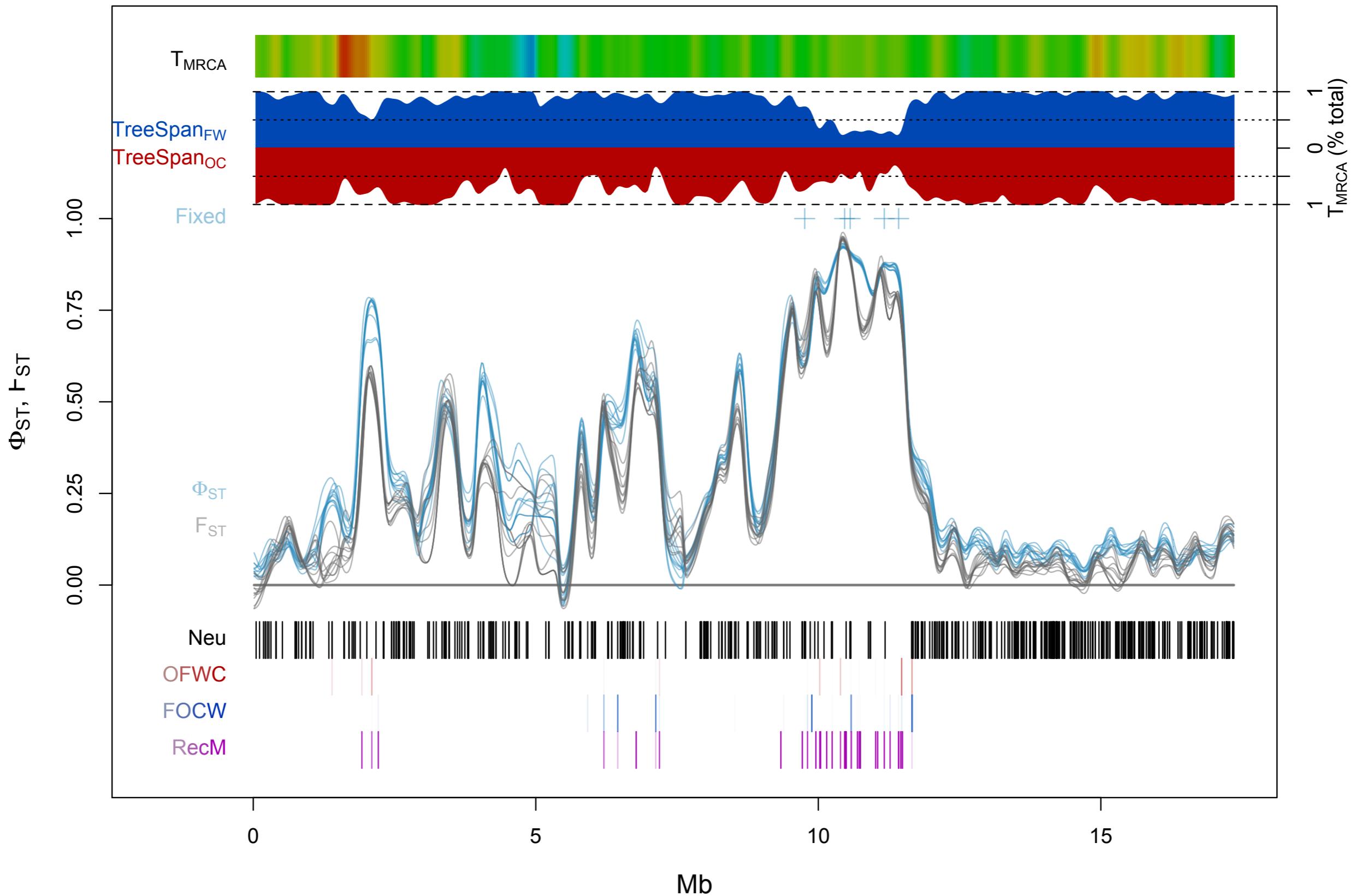
$$\log_2 \left( \frac{T_{\text{MRCA}} \text{FW}}{T_{\text{MRCA}} \text{ALL}} \right)$$

However, across populations the coalescence time can increases significantly in a genomic region



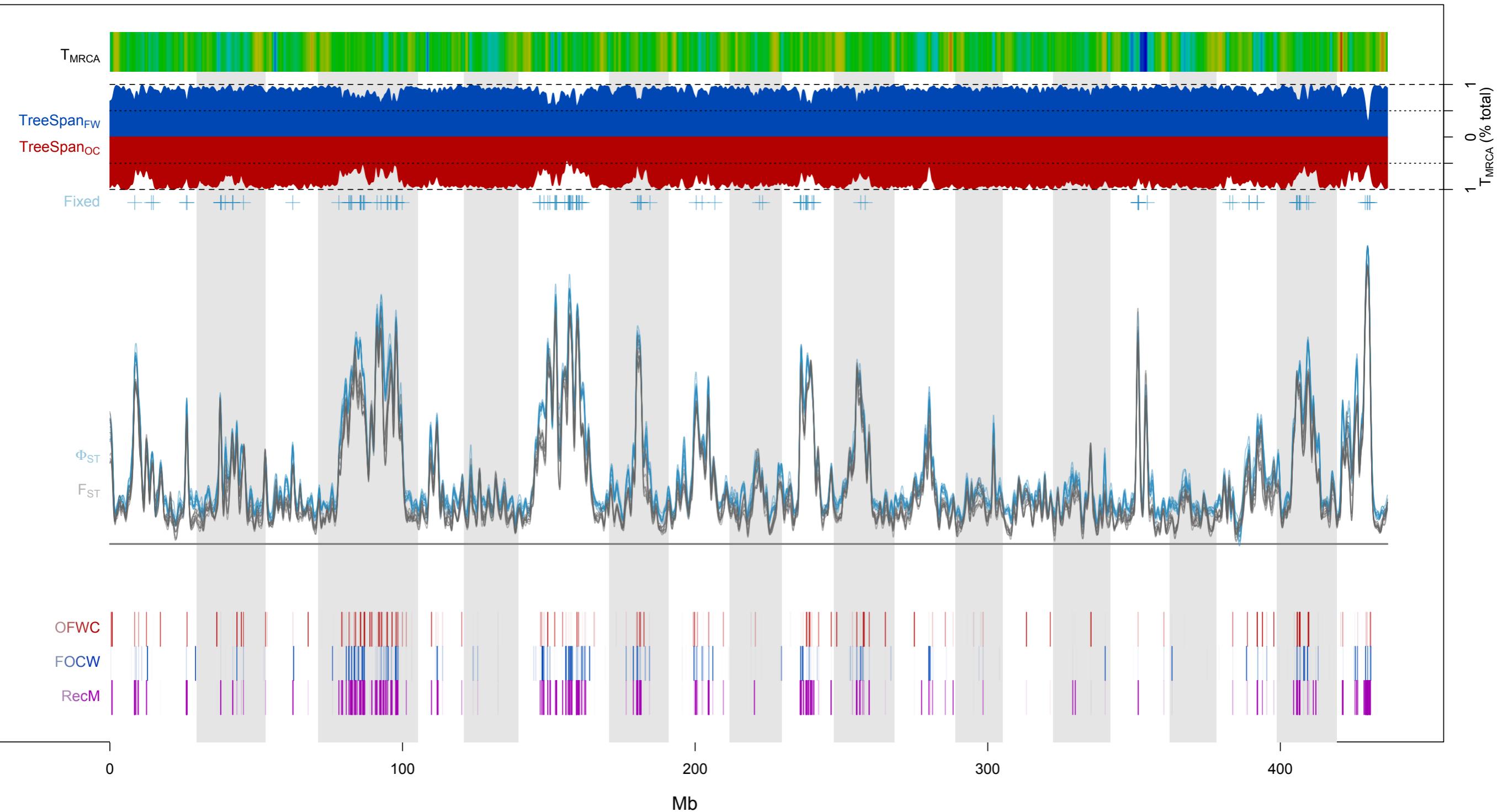
# Relative proportion of $T_{\text{mrca}}$ across different habitats

## LG21

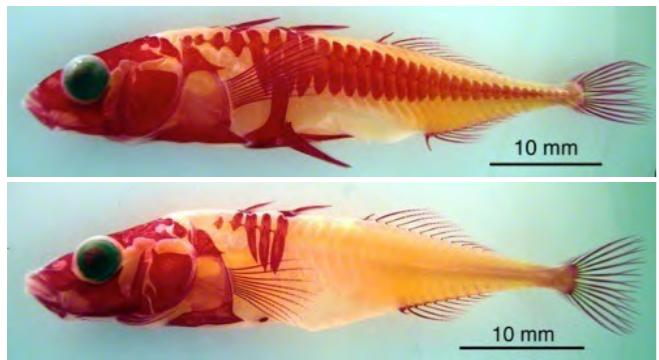


# Relative proportion of $T_{mrca}$ across different habitats

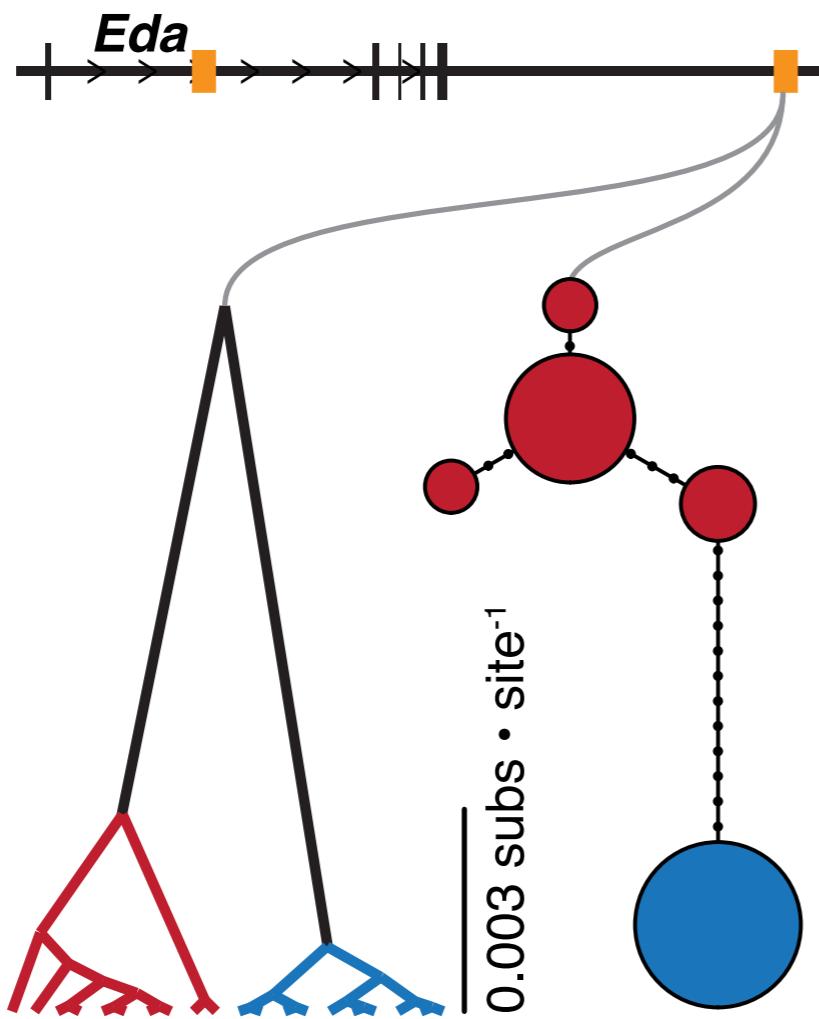
## Entire genome



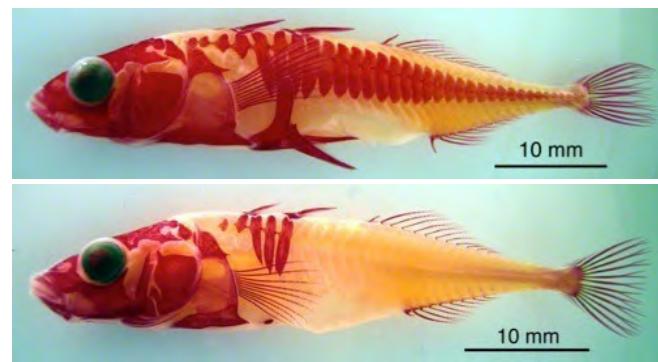
# Increased absolute divergence between habitats



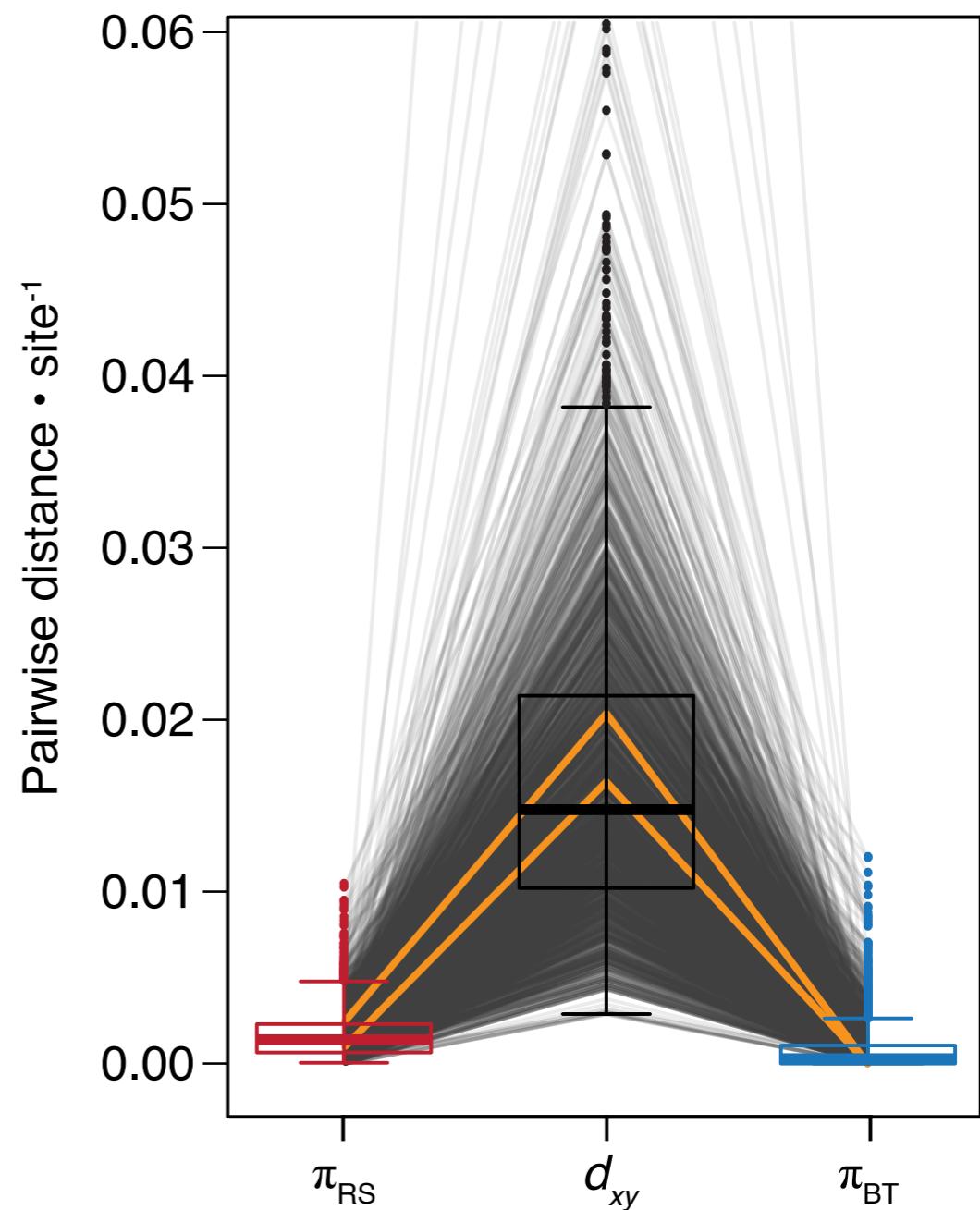
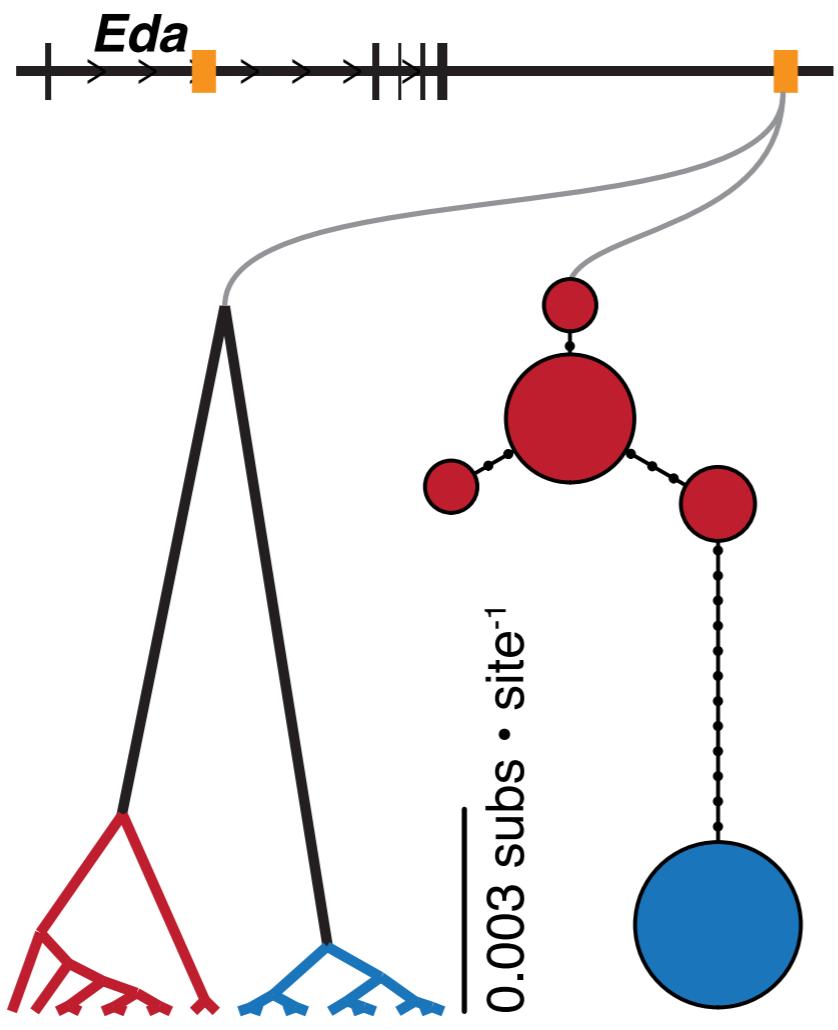
Stickleback chr IV: 12.80 - 12.82 Mb



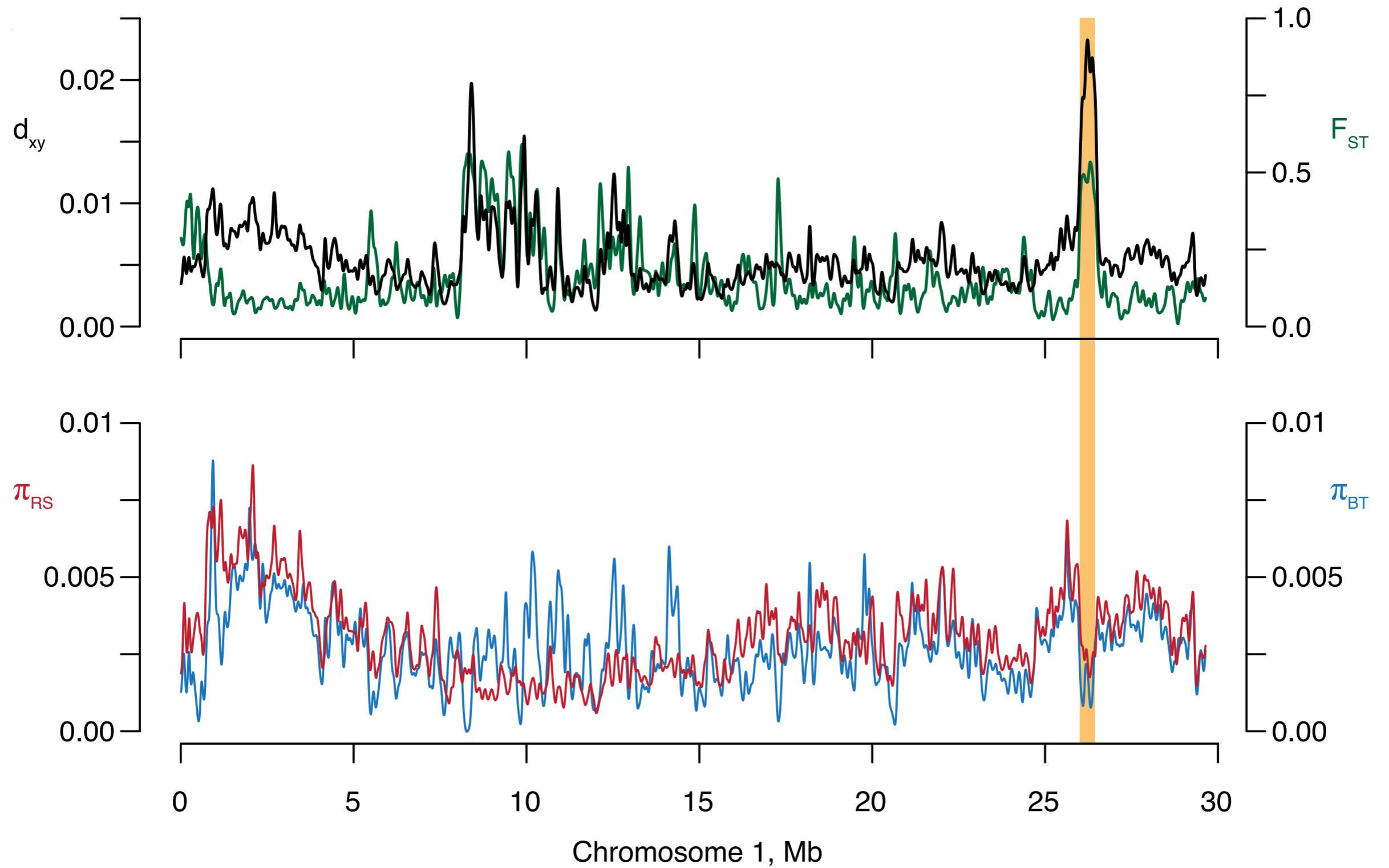
# Increased absolute divergence between habitats



Stickleback chr IV: 12.80 - 12.82 Mb



# Absolute divergence co-localizes with inversions



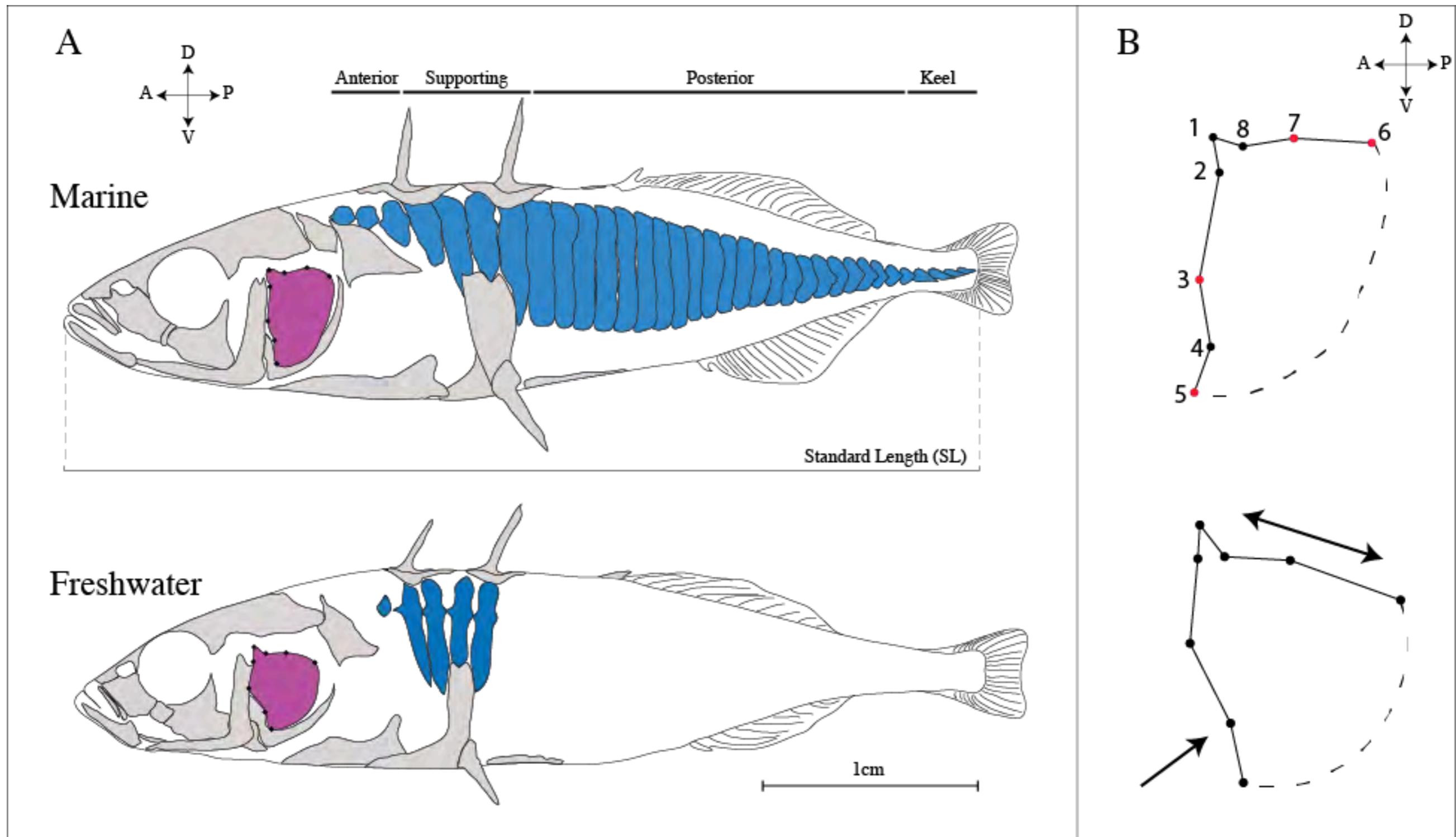
How do the genomic patterns of divergence link to phenotypic diversification?

---

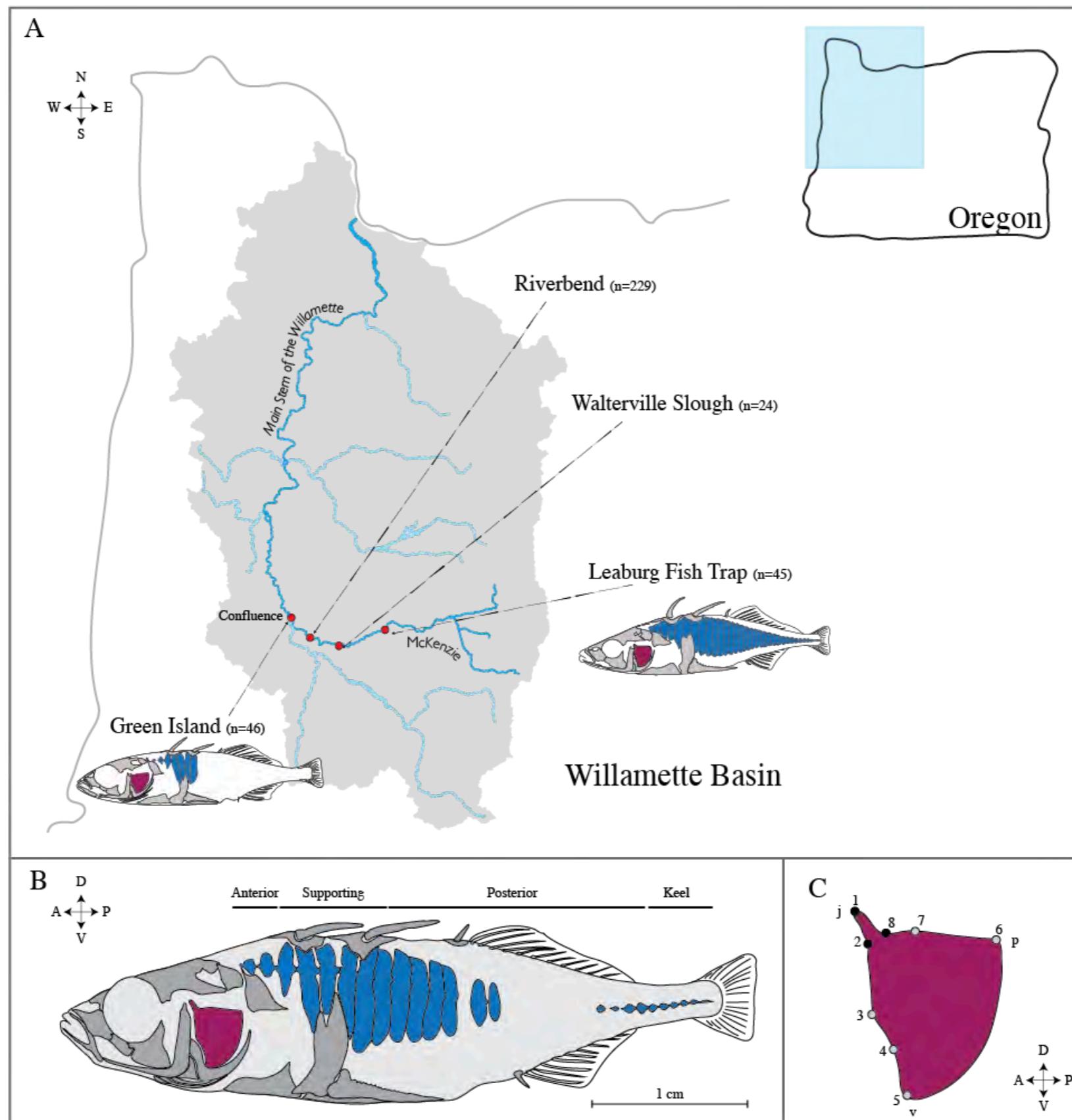


Kristin Alligood

# Lateral plate and opercle shape co-vary in the wild

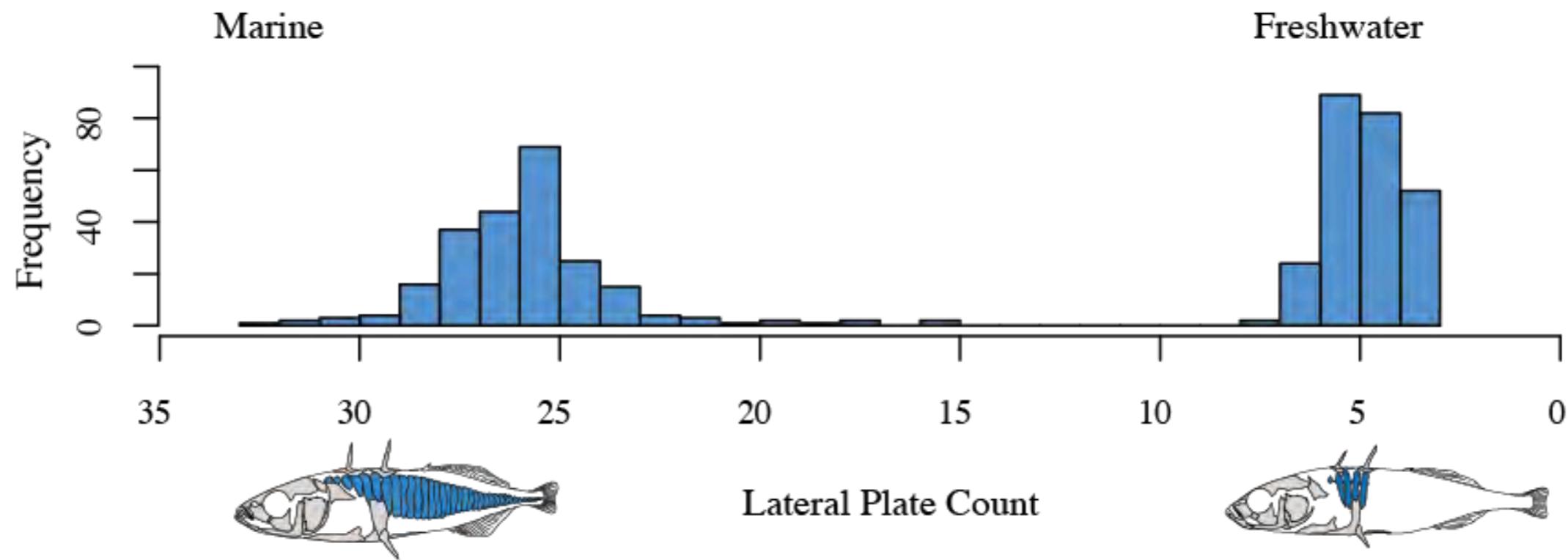


# An interesting stickleback hybrid population in Oregon

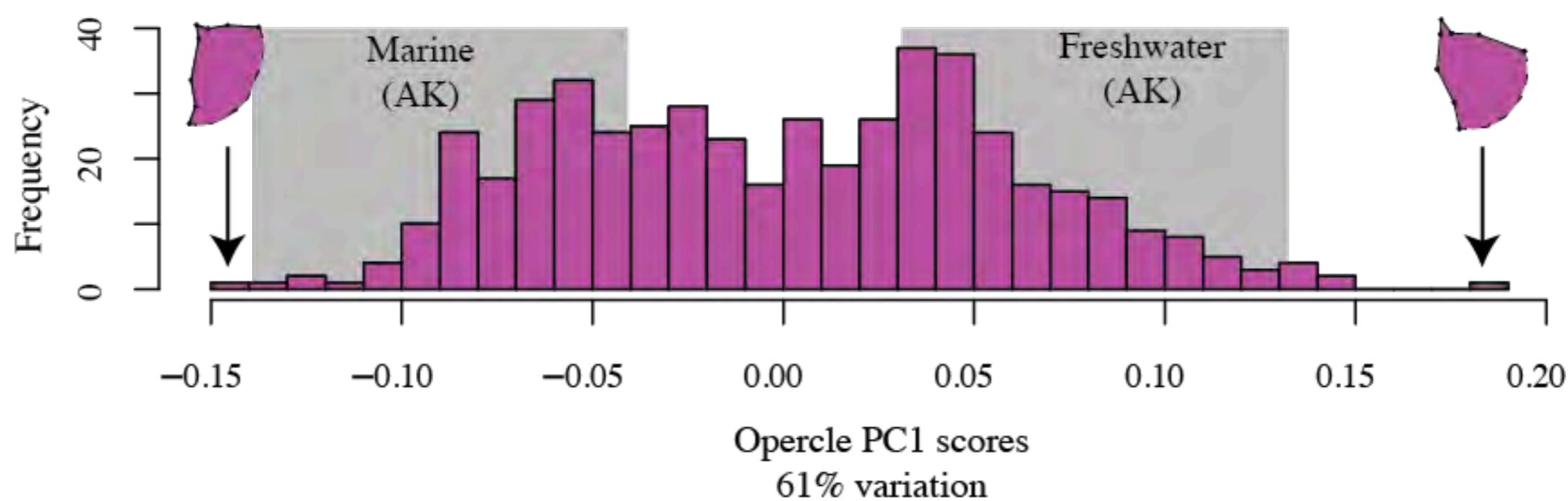


# An interesting stickleback hybrid population in Oregon

A

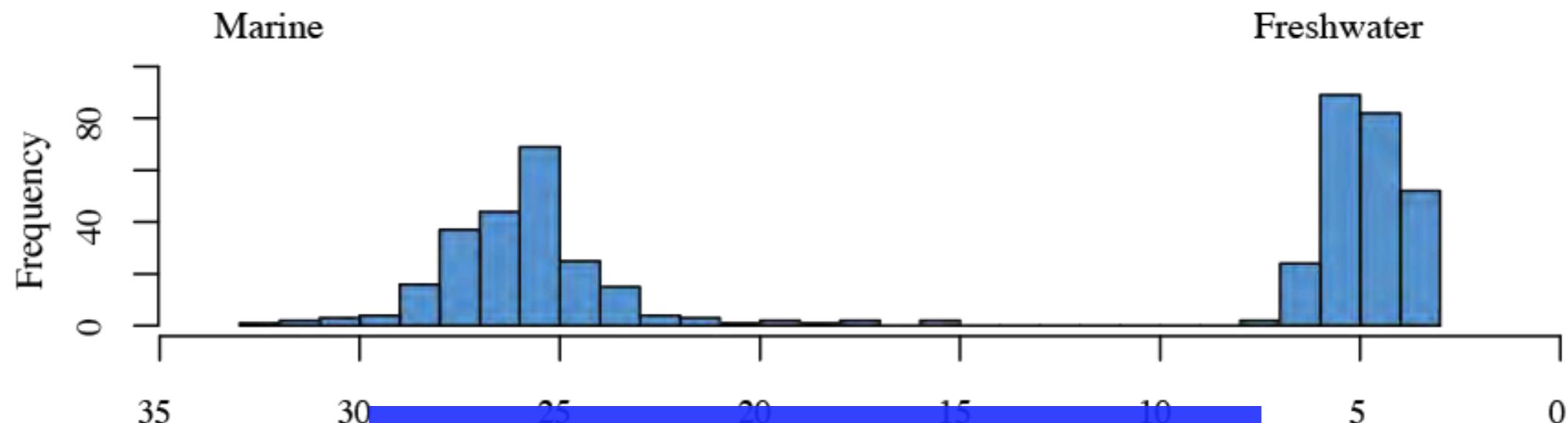


B



# An interesting stickleback hybrid population in Oregon

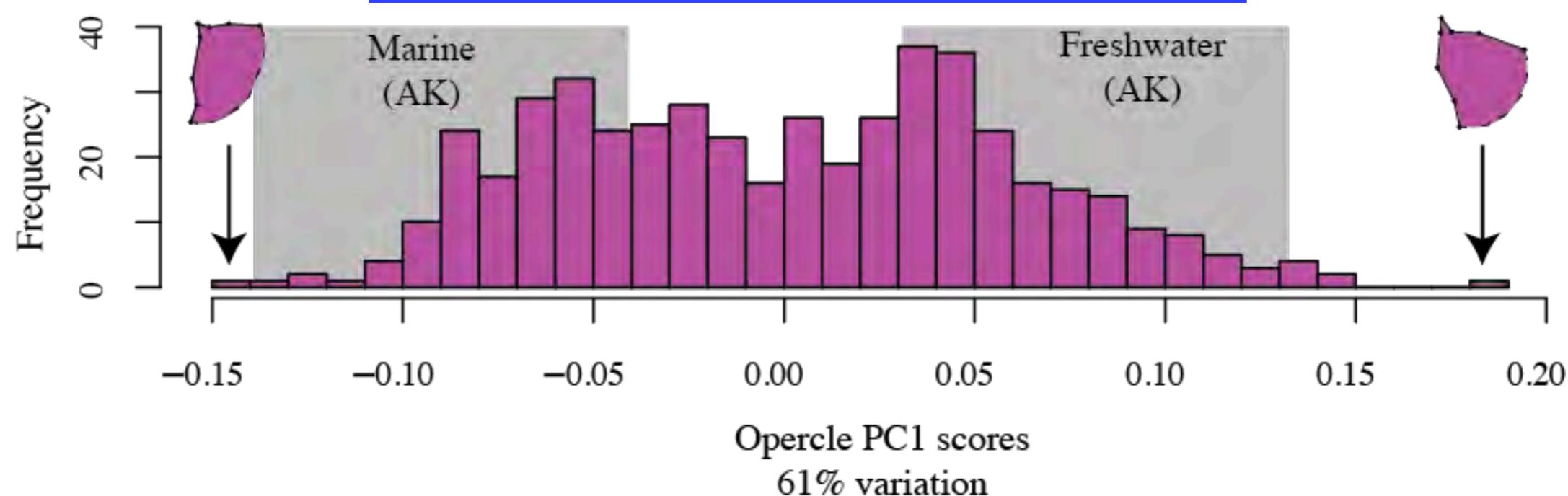
A



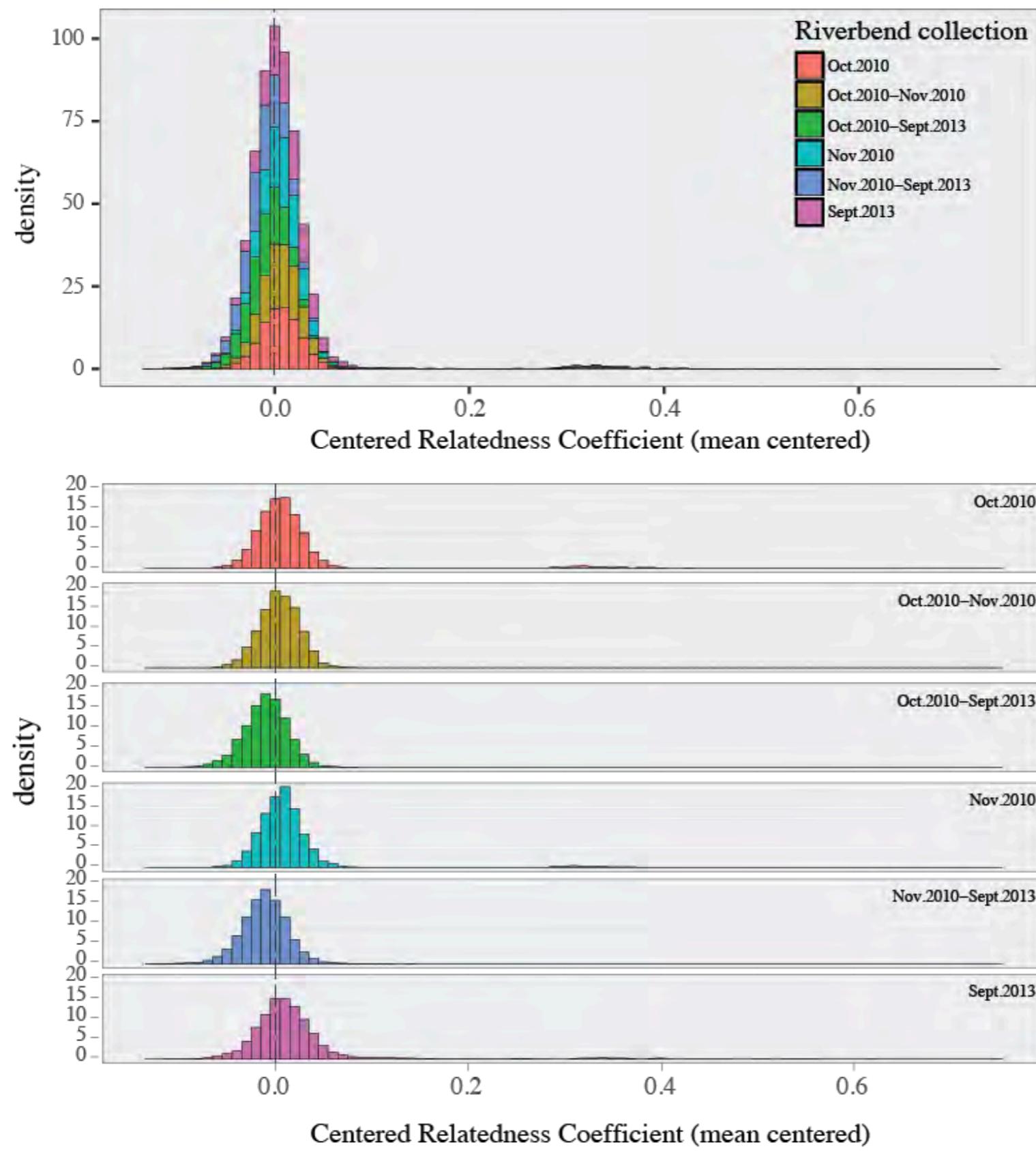
Genome Wide Association:  
GEMMA

Zhou and Stephens 2012

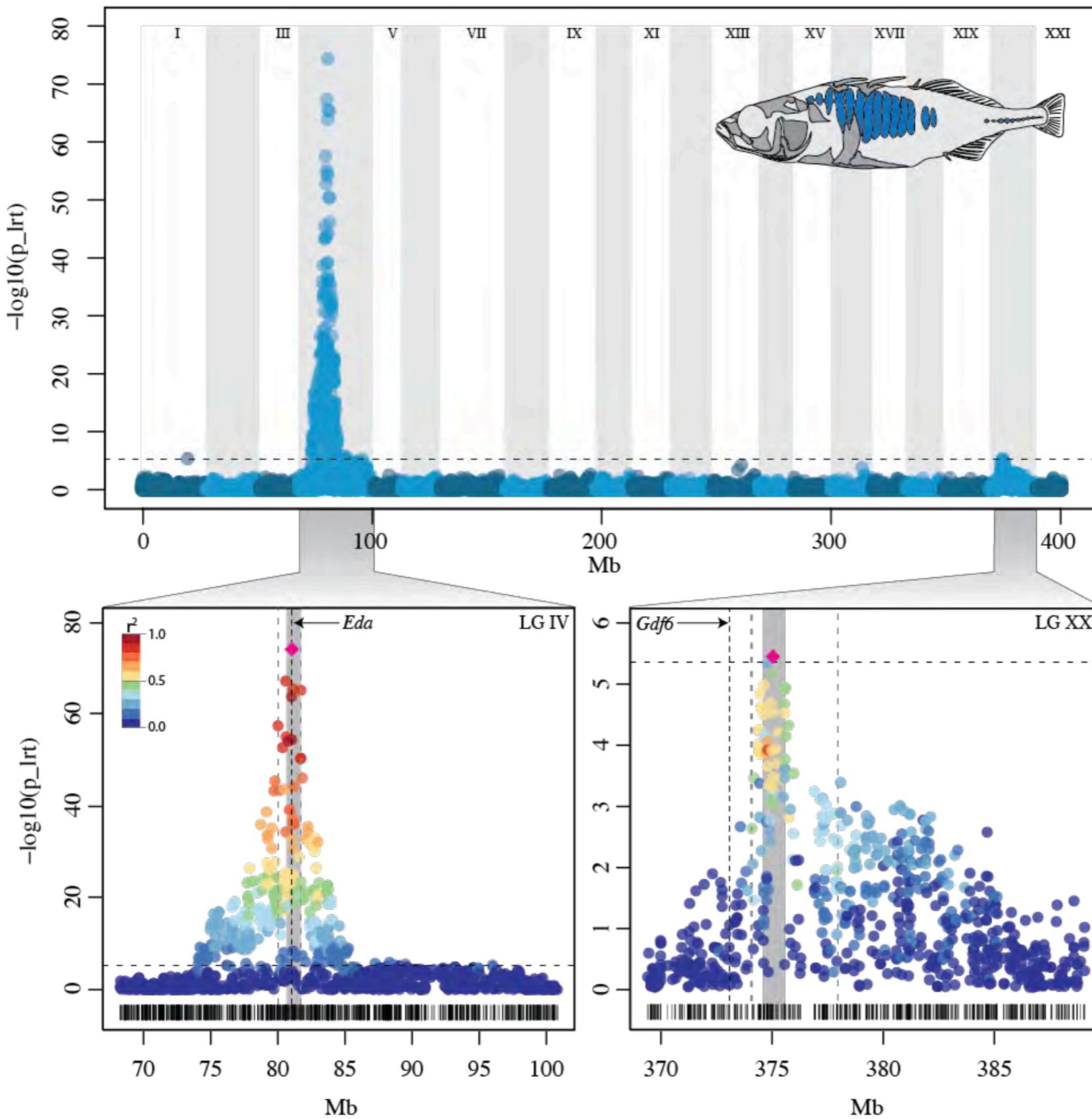
B



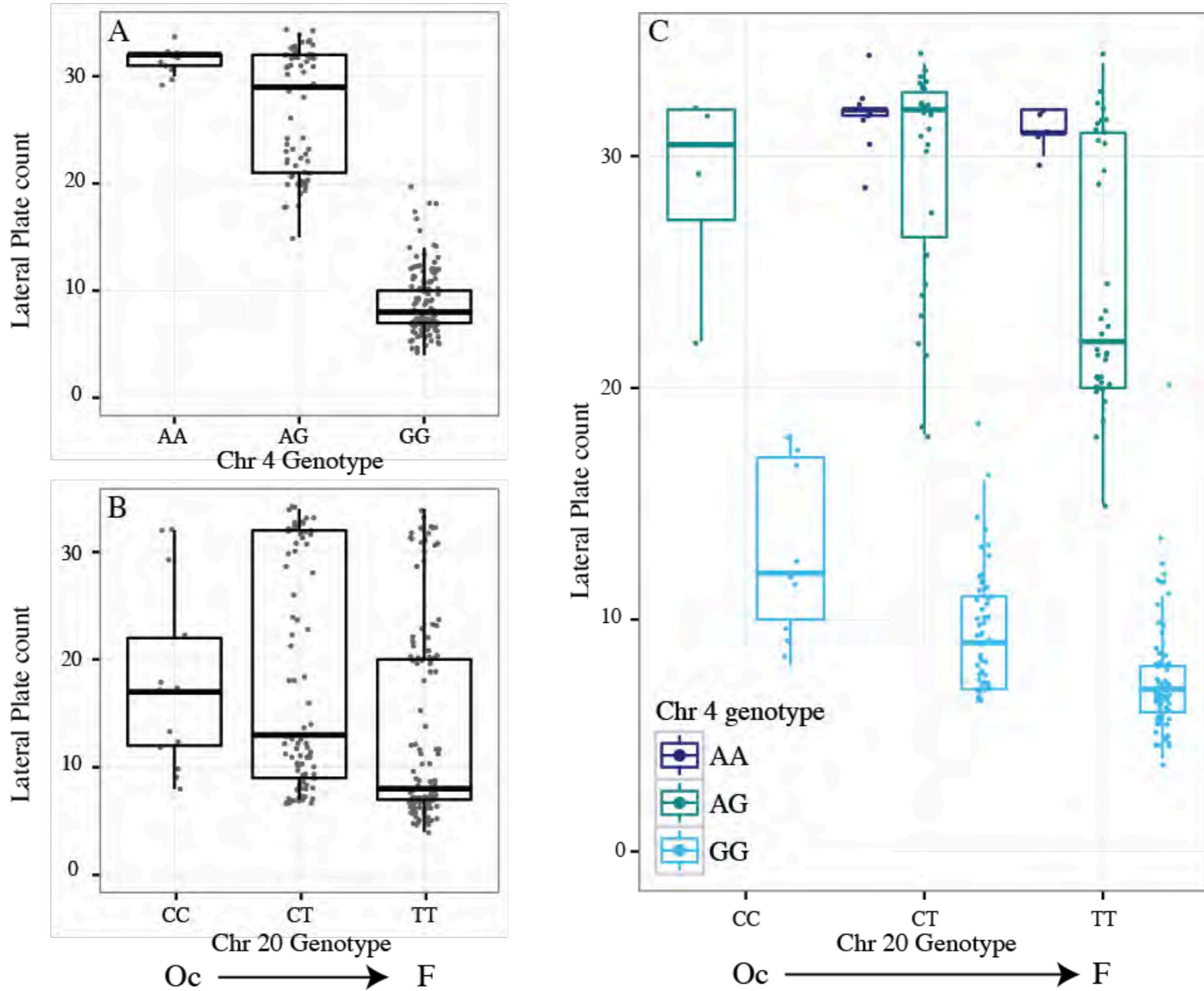
# Pairwise relatedness among individuals



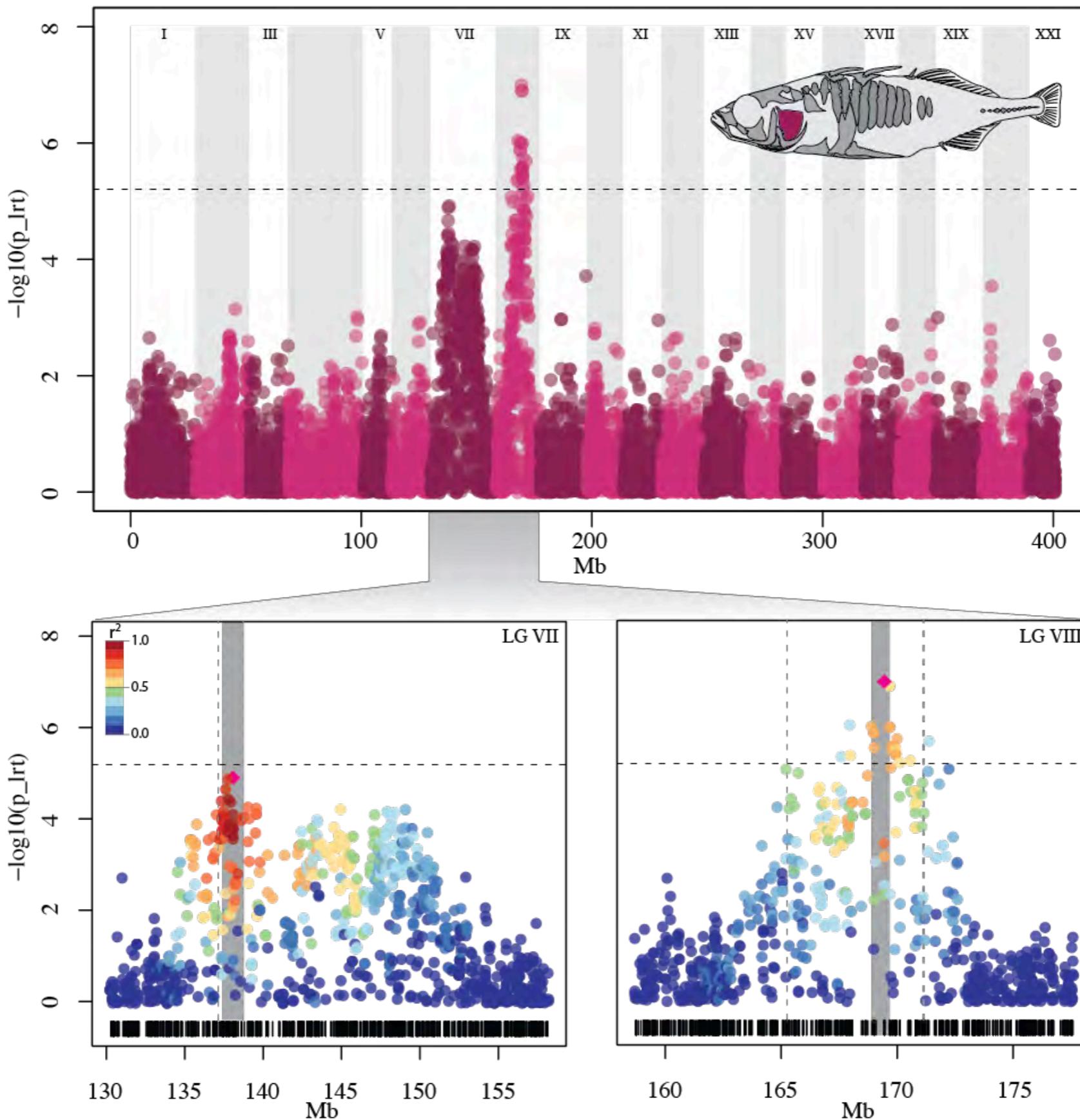
# Lateral plate count: A good hit to a novel locus on LGXX



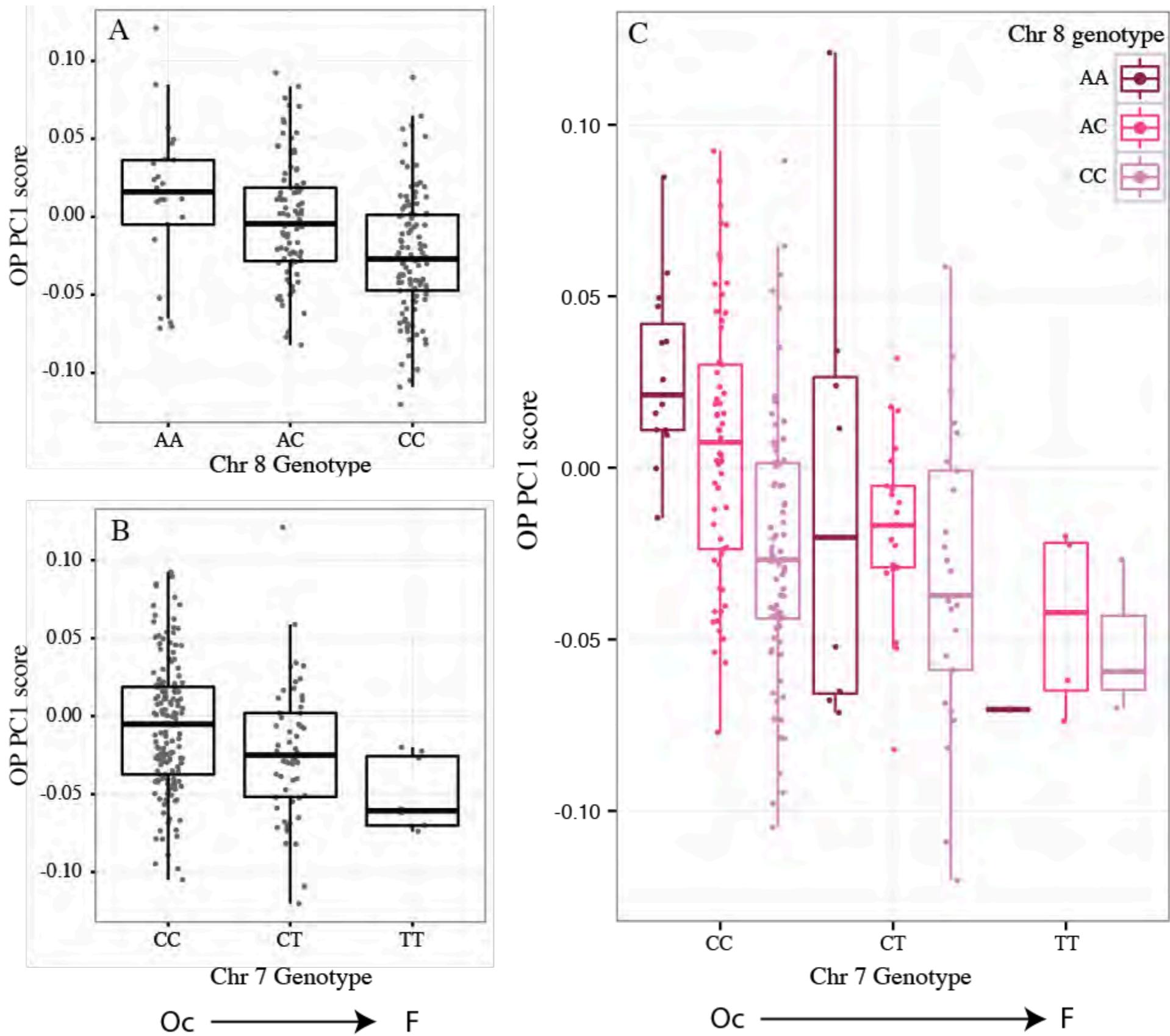
# Lateral plate count: A good hit to a novel locus on LGXX



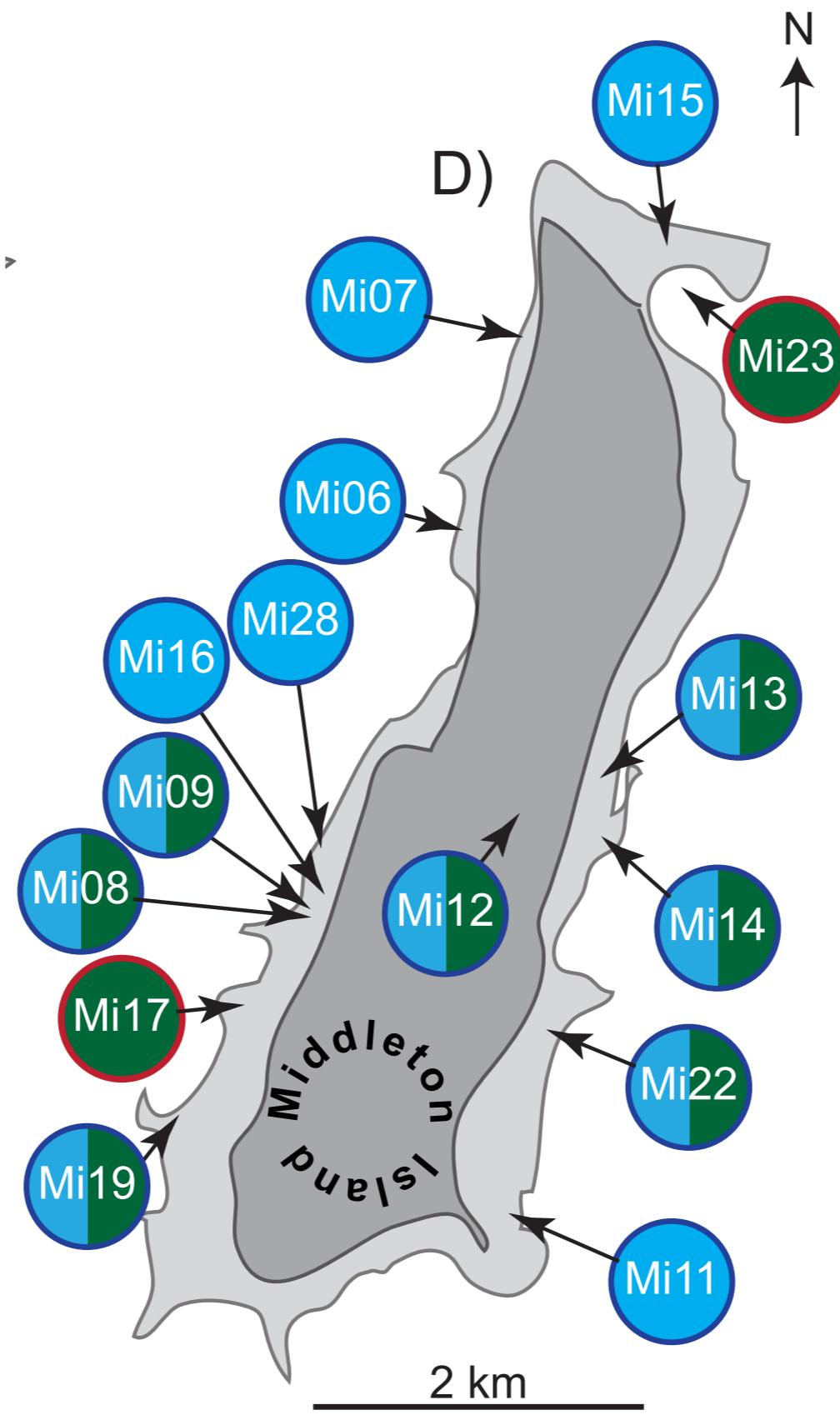
# Opercle shape: Good hits on LGVII and LGVIII



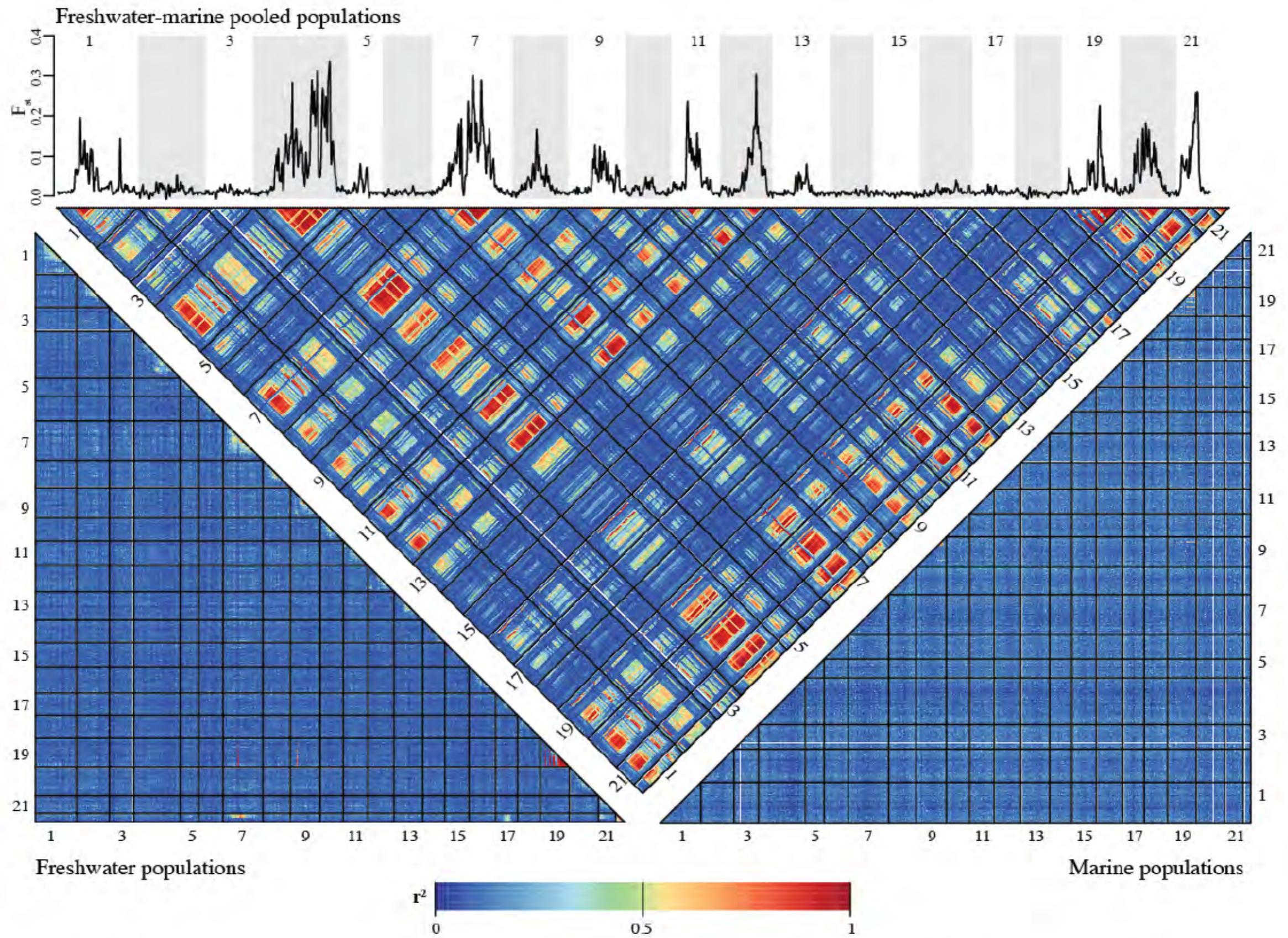
# Opercle shape: Good hits on LGVII and LGVIII



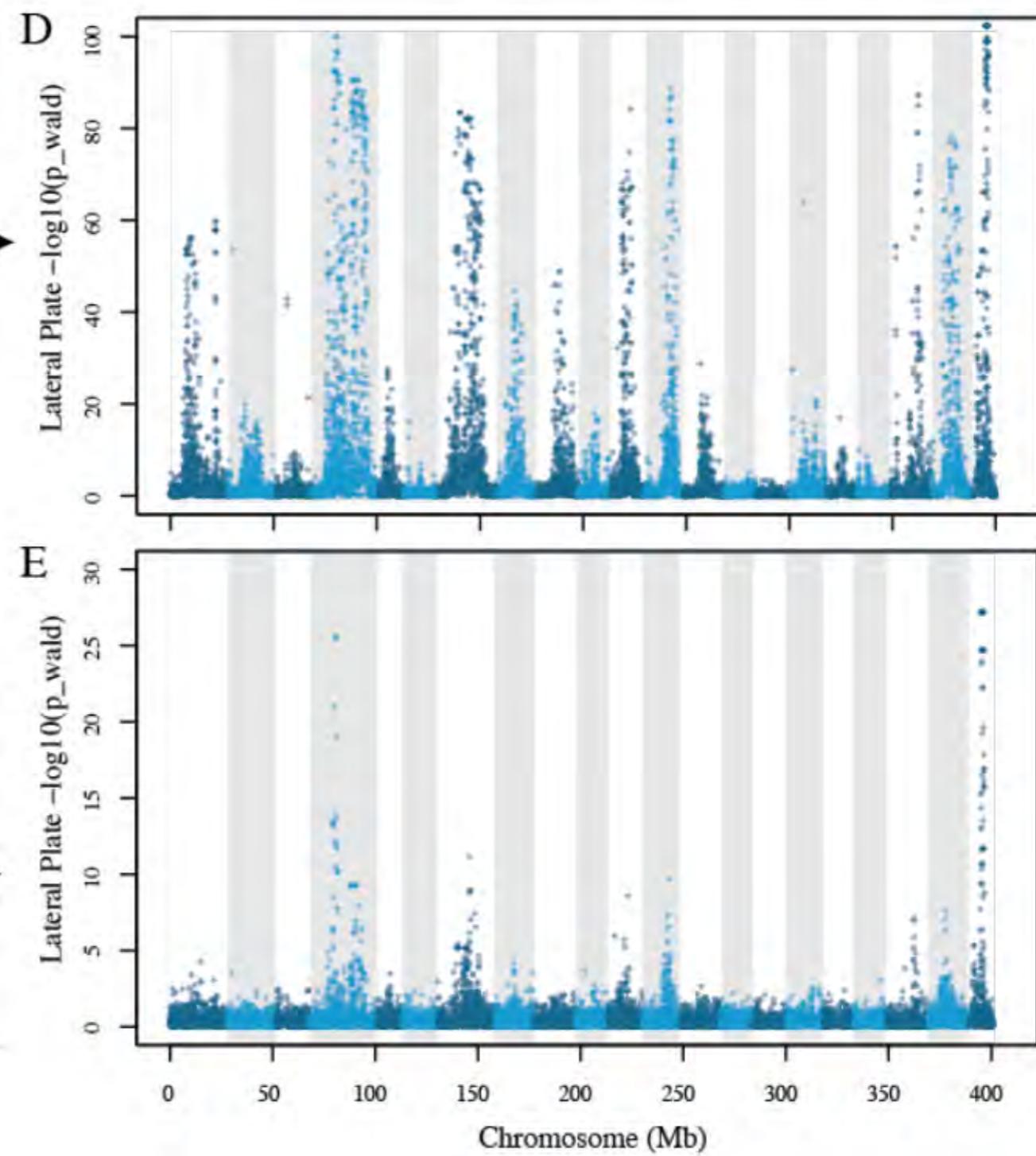
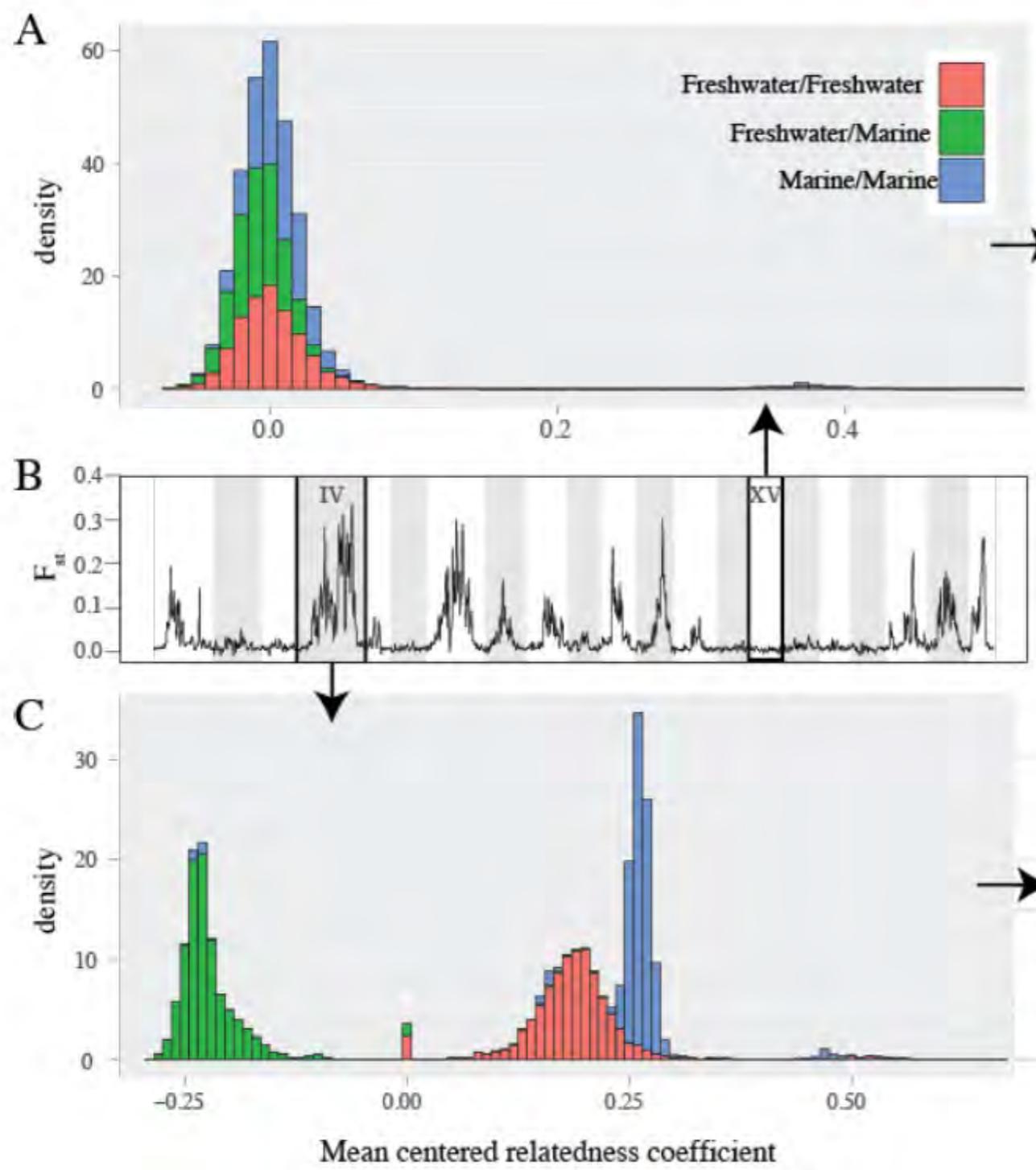
# Can't we just pool together young divergent populations and do GWAS?



# Linkage disequilibrium created by natural selection

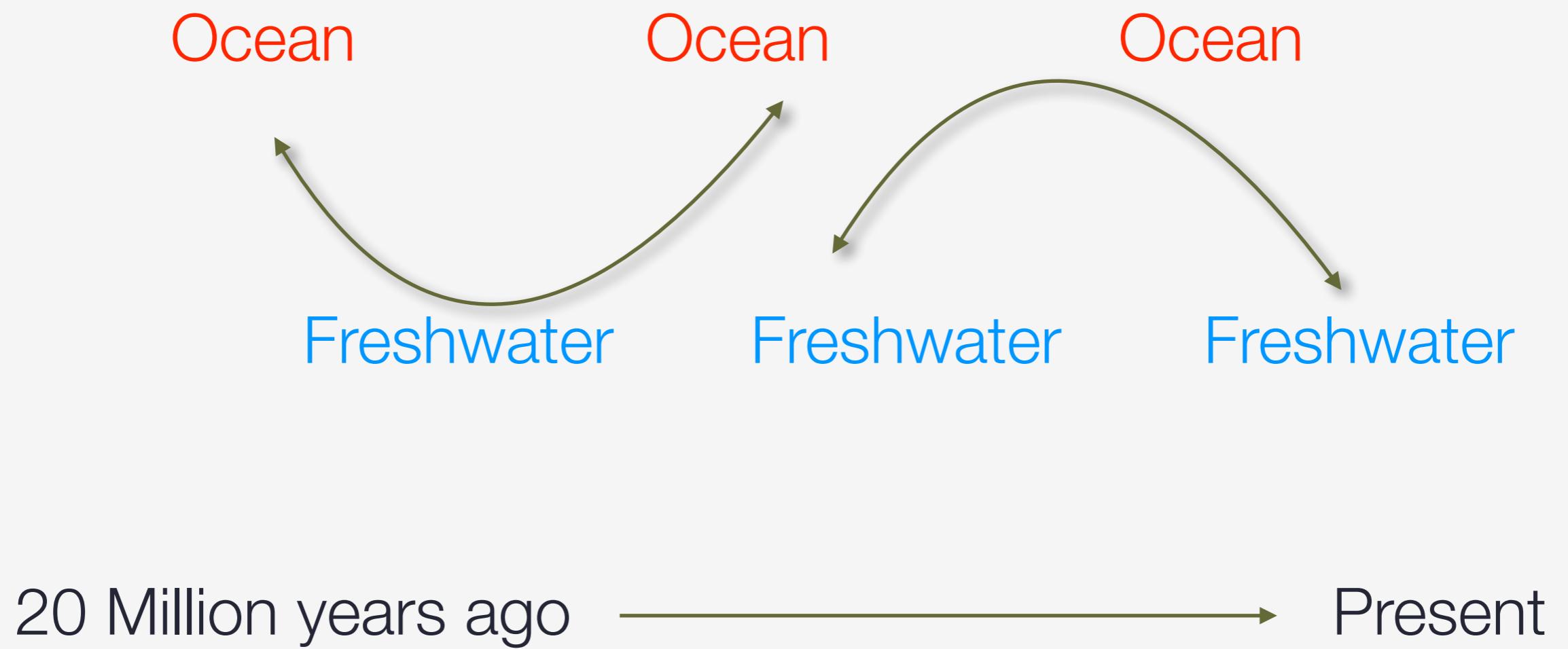


# This sort of population structure is difficult to control



- Previous work has shown that the freshwater genomes evolve in 13,000 years.
- These new Middleton Island data shows that the phenotype can appear in as little as 50 years.
- Much of the divergence involves soft sweeps.
- This could represent thousands of haplotypes reassembling, but the genome appears chunkier.
- Many haplotypes are *habitat specific*, quite *ancient* and coincide with *structural variation*
- Diverging phenotypes map to these same genomic regions

Hypothesis: Old genomic architecture variation is a product of the metapopulation structure of stickleback, and this architecture strongly influences subsequent rapid evolution.



# Outline for today's lecture

---

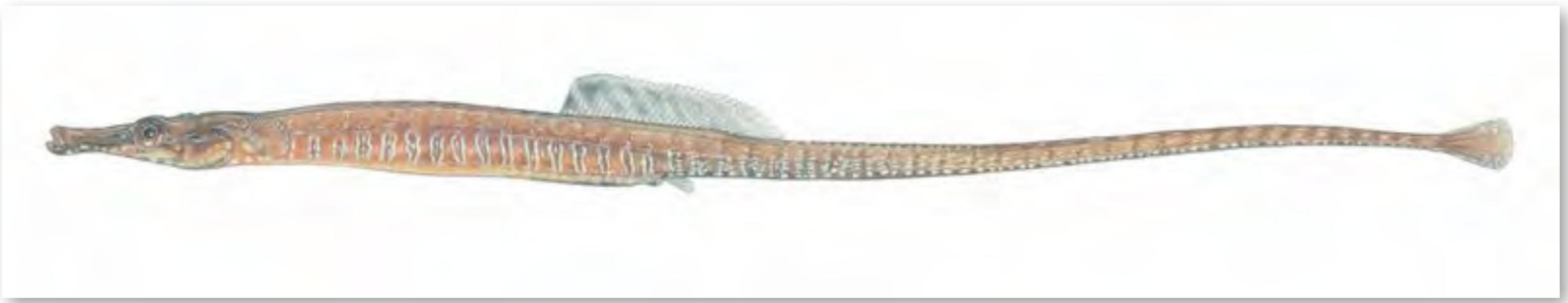
RAD-seq for ecological and evolutionary genomics

Primer on Population Genomics

Evolutionary genomics of stickleback fish

- Population genomics of rapid adaptation
- Using long read RAD-seq for coalescent analyses
- Genome Wide Association Studies using RAD-seq

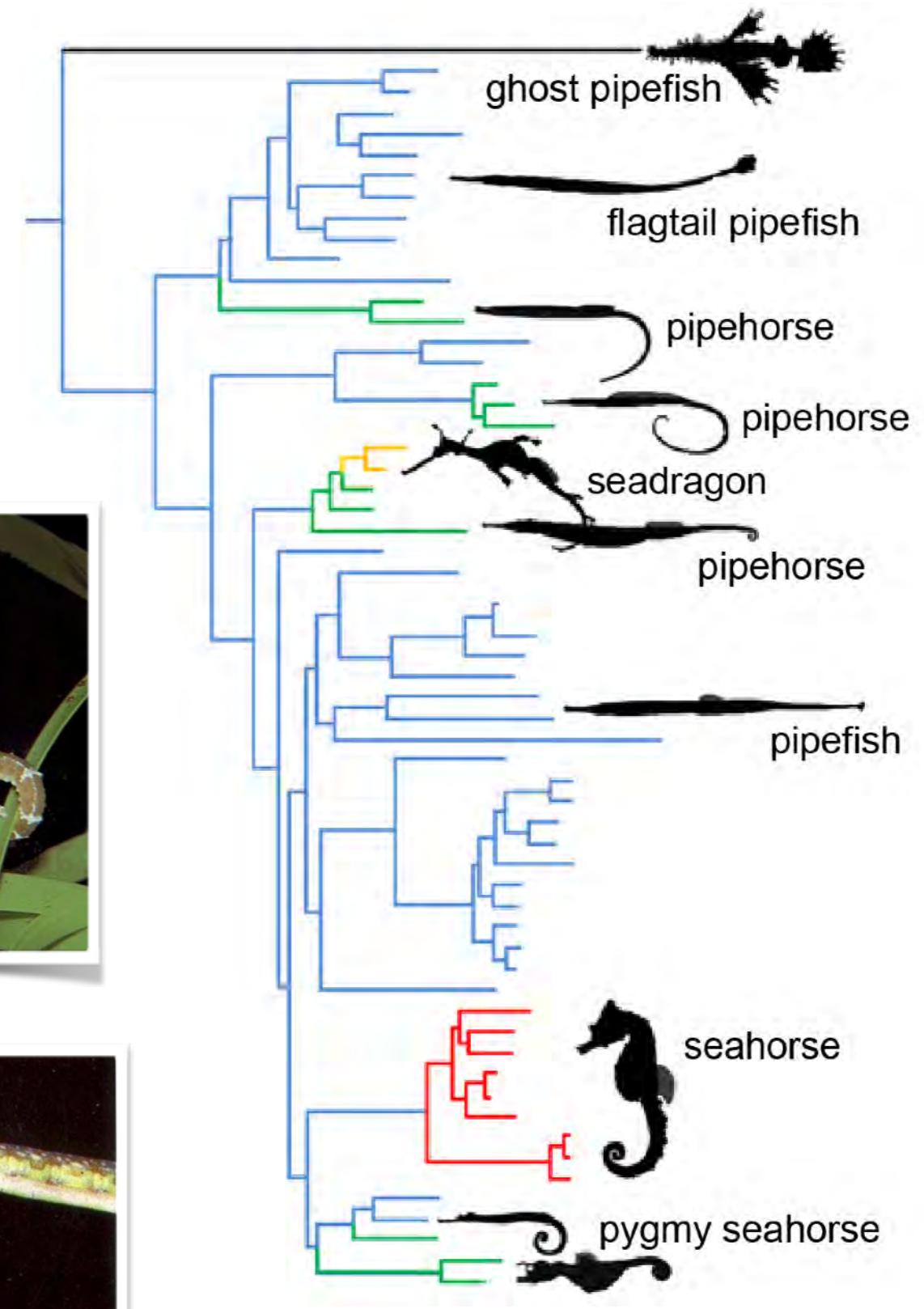
**Genomically enabling the Gulf pipefish**



What if you don't have a genome sequence?

A case study of RAD-seq and genome assembly

# Seahorses, sea dragons and pipefishes



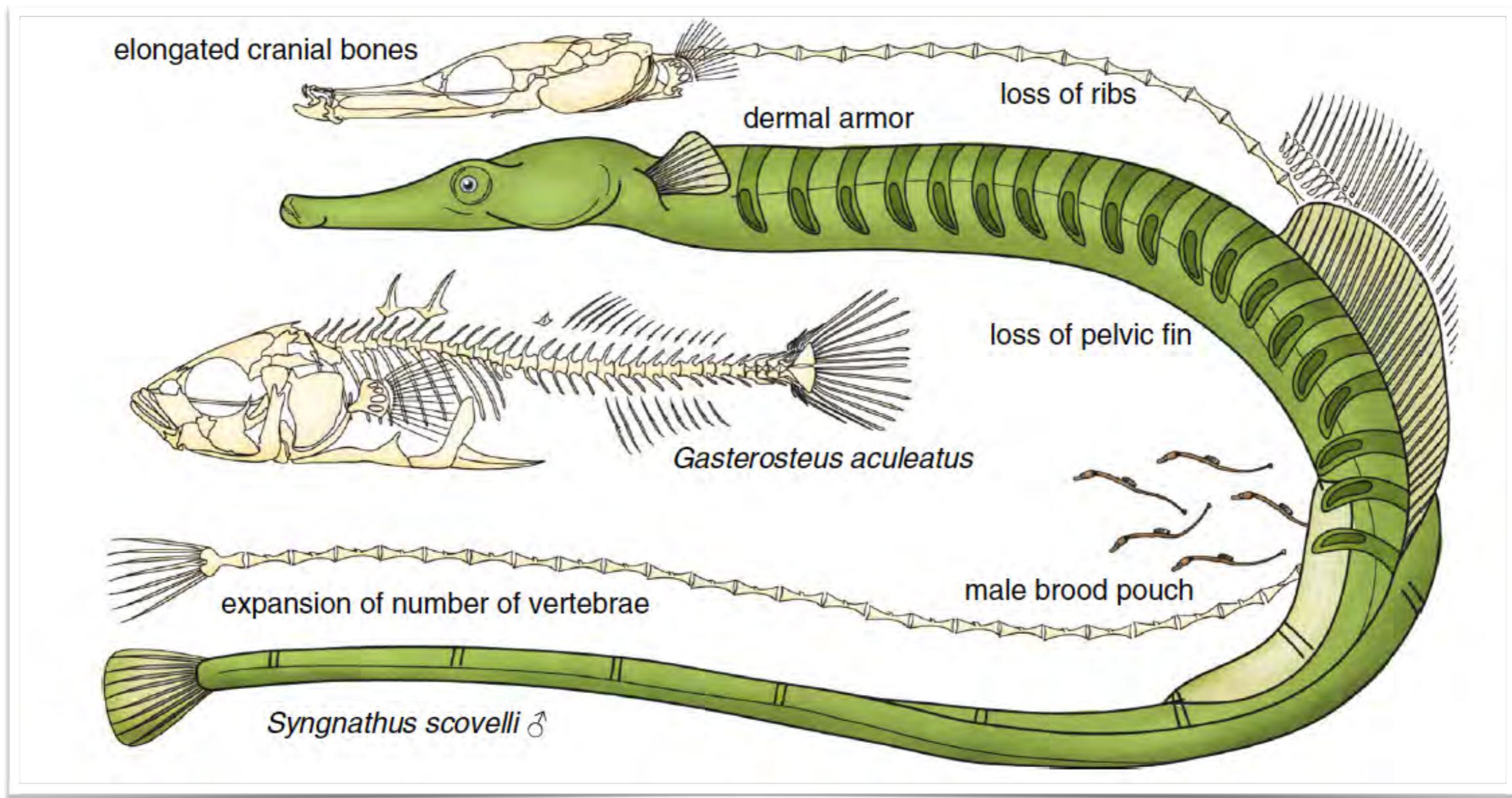
RESEARCH

Open Access



# The genome of the Gulf pipefish enables understanding of evolutionary innovations

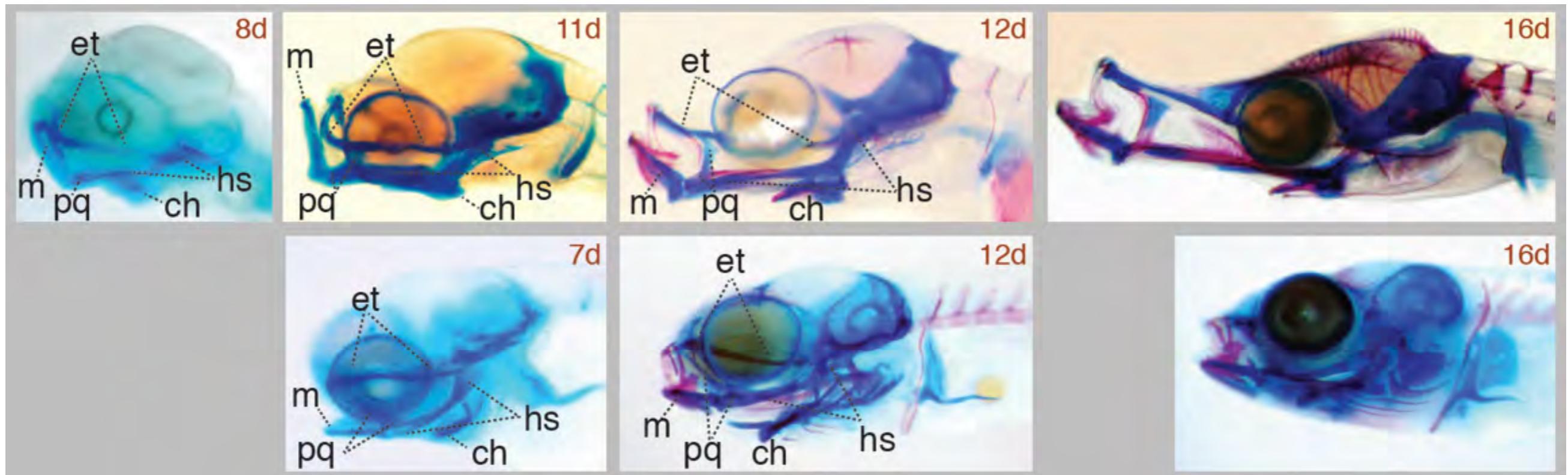
C. M. Small<sup>1†</sup>, S. Bassham<sup>1†</sup>, J. Catchen<sup>1,2†</sup>, A. Amores<sup>3</sup>, A. M. Fuiten<sup>1</sup>, R. S. Brown<sup>1,4</sup>, A. G. Jones<sup>5</sup> and W. A. Cresko<sup>1\*</sup>



# We're really interested in head development



Pipefish



Stickleback

## *How did we genomically enable pipefish*

---

- 1) A high quality transcriptome
- 2) Very dense RAD genetic map
- 3) Deep shotgun sequencing of the genome
- 4) Order contigs against the RAD reference map

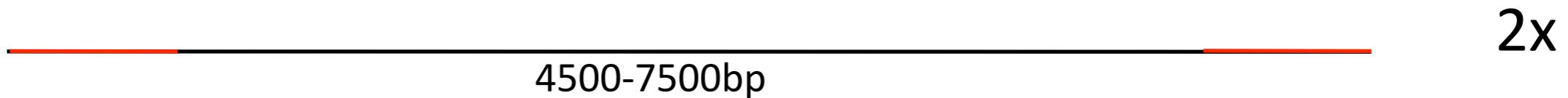
# Illumina genomic libraries for pipefish genome

---

paired end 101bp



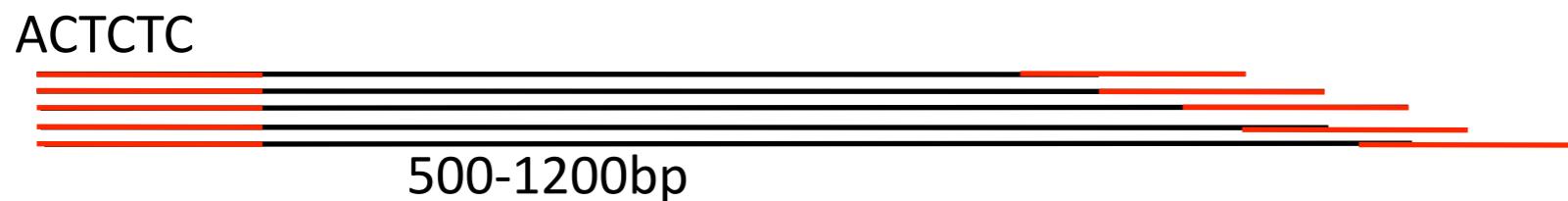
mate pair



overlapping



paired end RAD



15-25x of  
3% of the  
genome

# Nearly the whole genome is covered

---

**Table 1** Scaffold-level assembly statistics for the Gulf pipefish genome

Genome	Scaffolds (n)	Longest scaffold	Scaffold N50	Contig N50	Assembly length	Gaps in assembly (%)	CEGs complete (%)
Gulf pipefish ( <i>Syngnathus scovelli</i> )	2104	6.7 Mb	640.4 kb	32.2 kb	307.0 Mb	6.6	98.8
African turquoise killifish ( <i>Nothobranchius furzeri</i> )	29,054	0.7 Mb	119.7 kb	8.7 kb	1010.9 Mb	7.7	94.8
Blind cave fish ( <i>Astyanax mexicanus</i> )	10,542	9.8 Mb	1775.3 kb	14.7 kb	1191.1 Mb	19.1	87.9
Spotted gar ( <i>Lepisosteus oculatus</i> )	2105	21.3 Mb	6928.1 kb	68.3 kb	945.8 Mb	8.1	90.7

# Created a genetic map

---

Generated an F1 family of 103 individuals

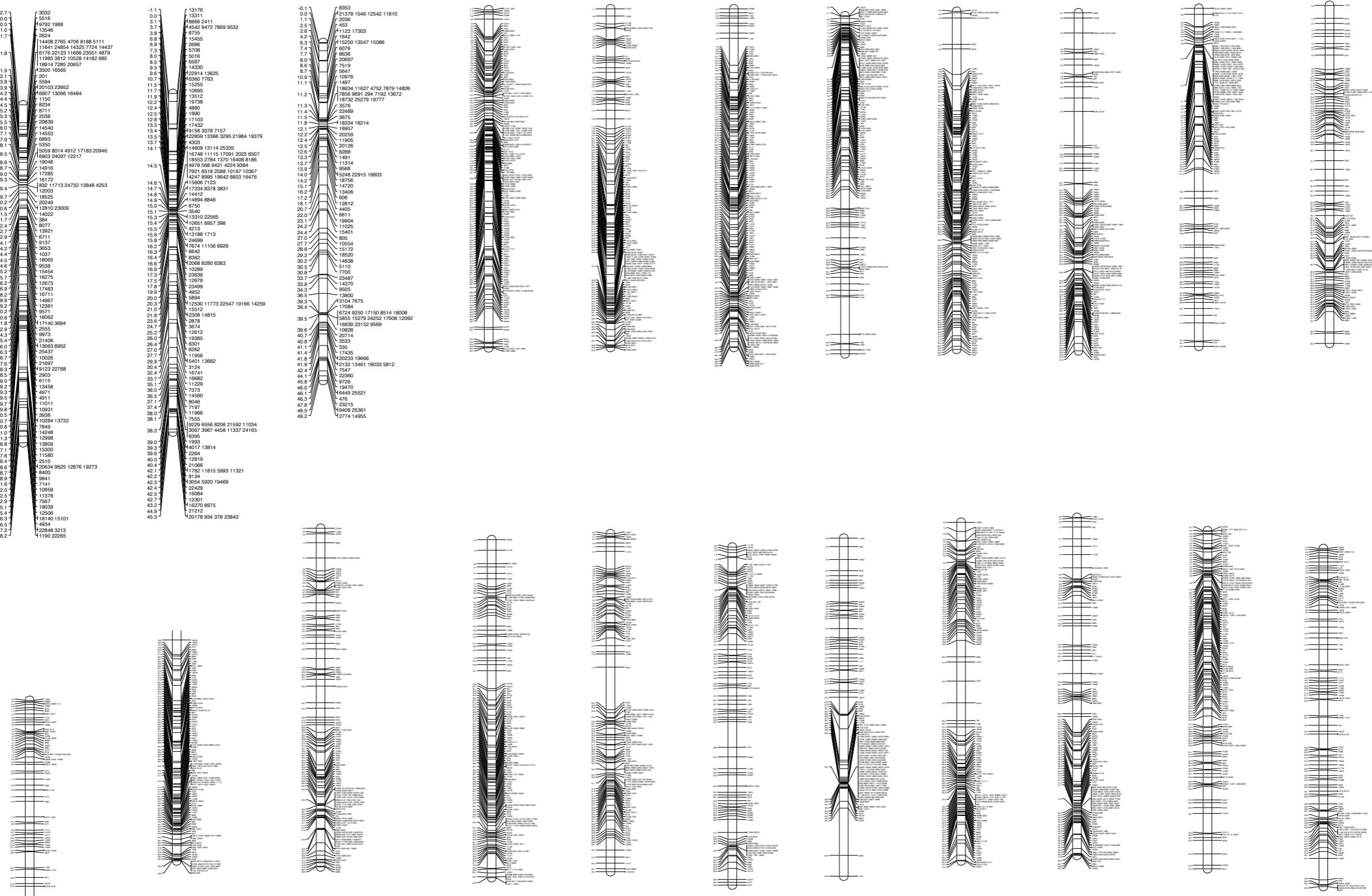
RAD sequenced the parents and offspring

Analyzed the data using *Stacks*

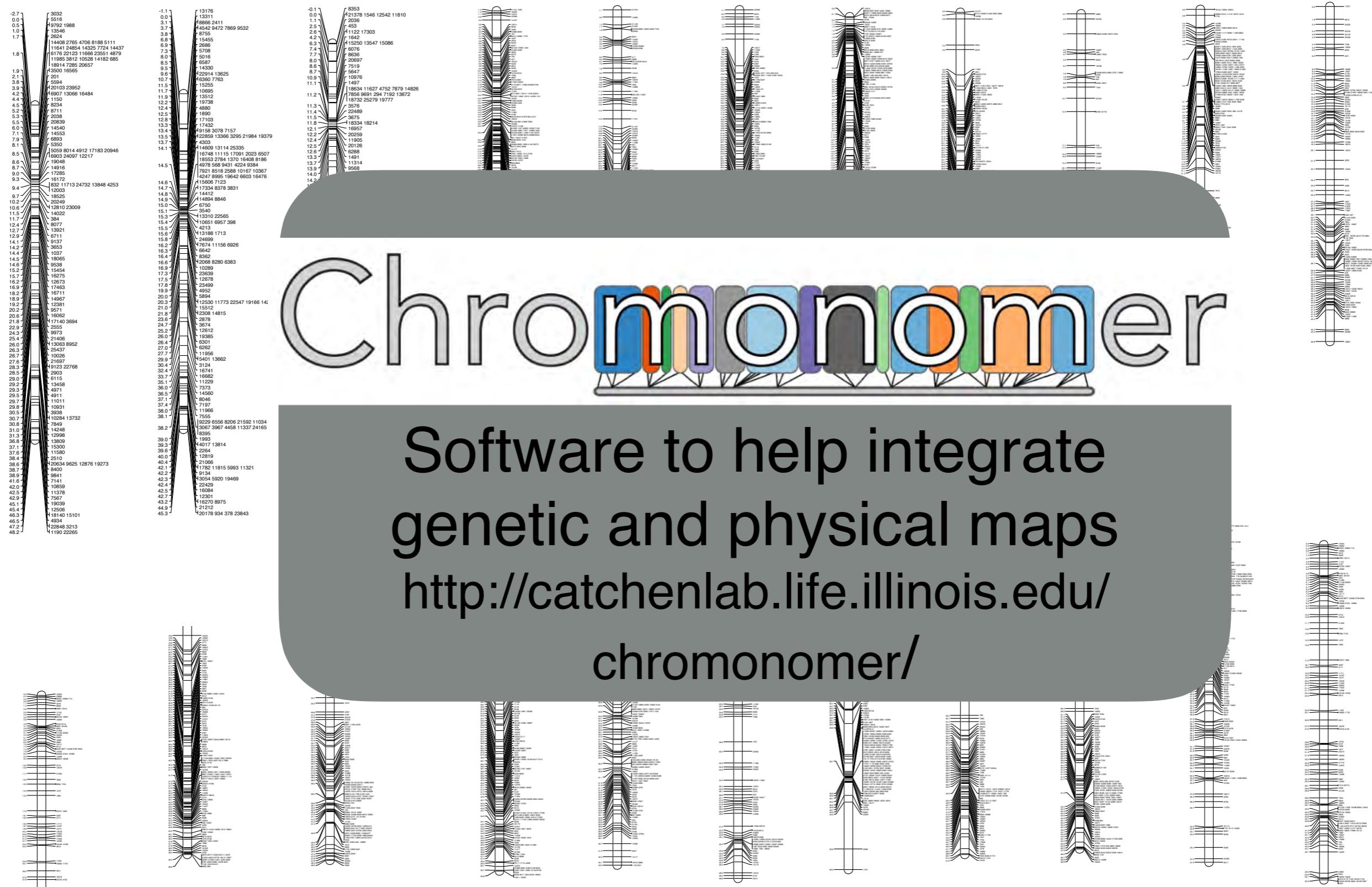
Output to JoinMap format

Created Linkage map

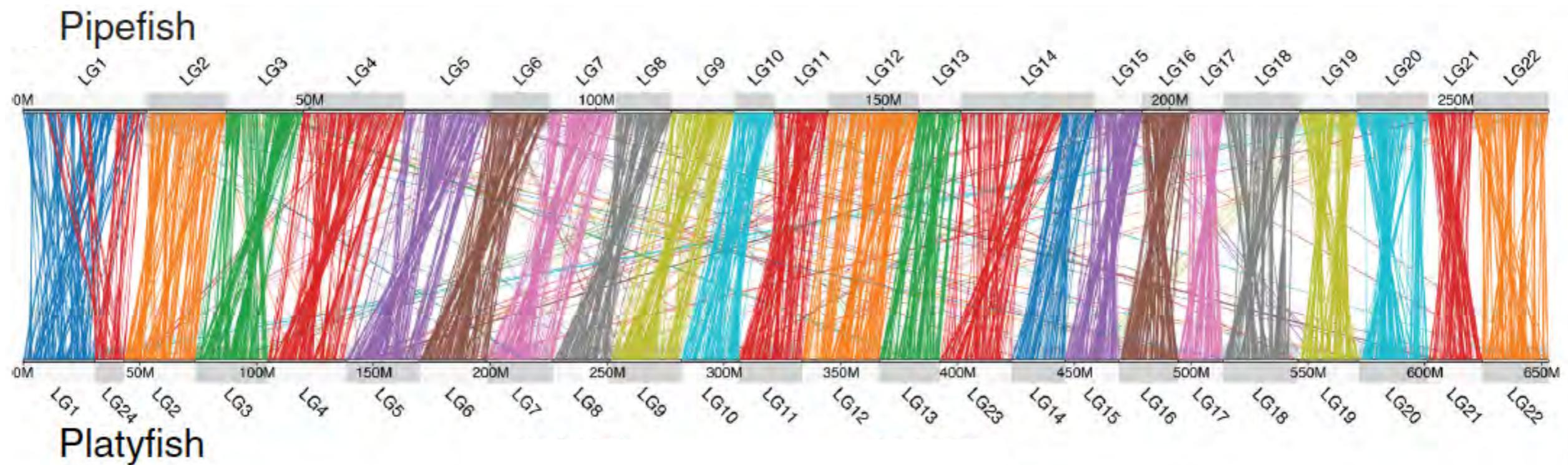
# 22 LGs; 6000 segregating SNPs; 30,000 RAD sites



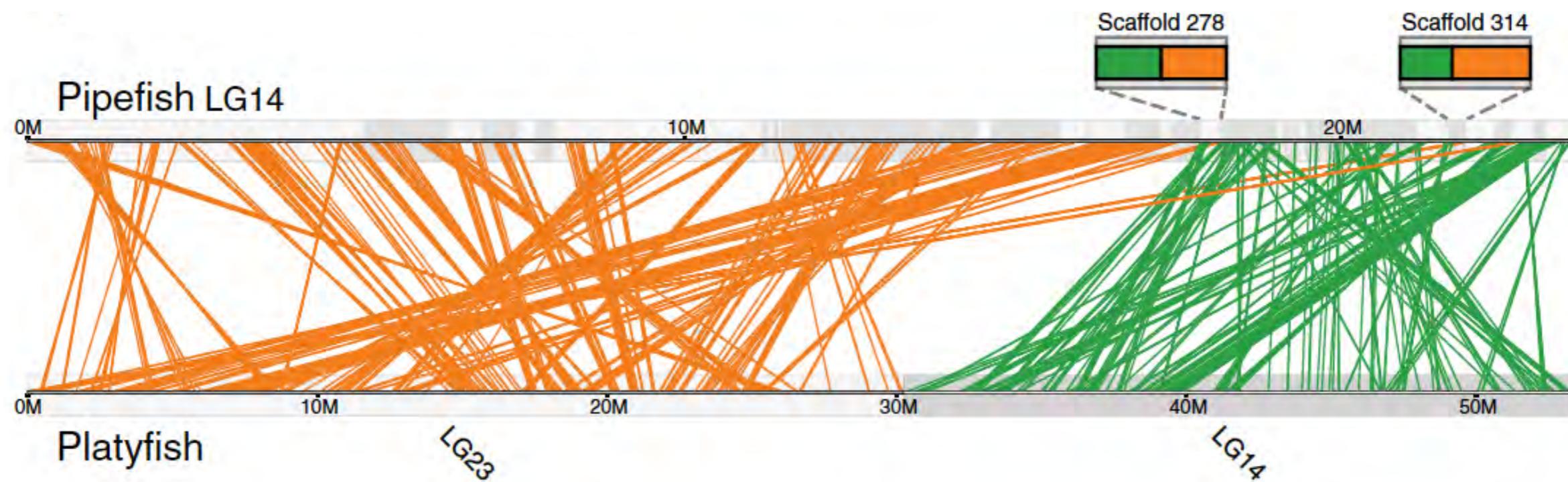
# The pipefish genetic map and genome together



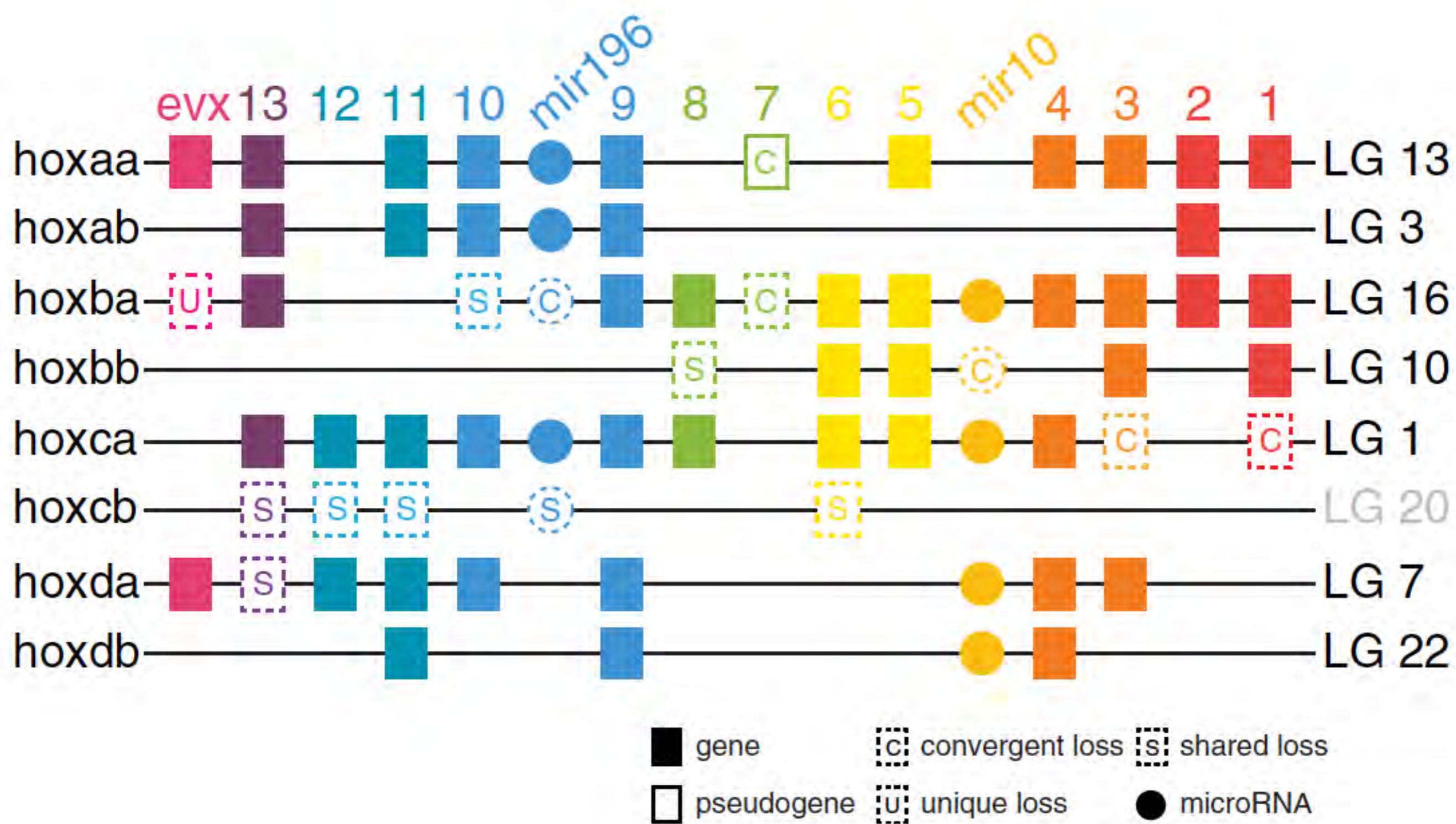
# Two instances of chromosome fusion in pipefish



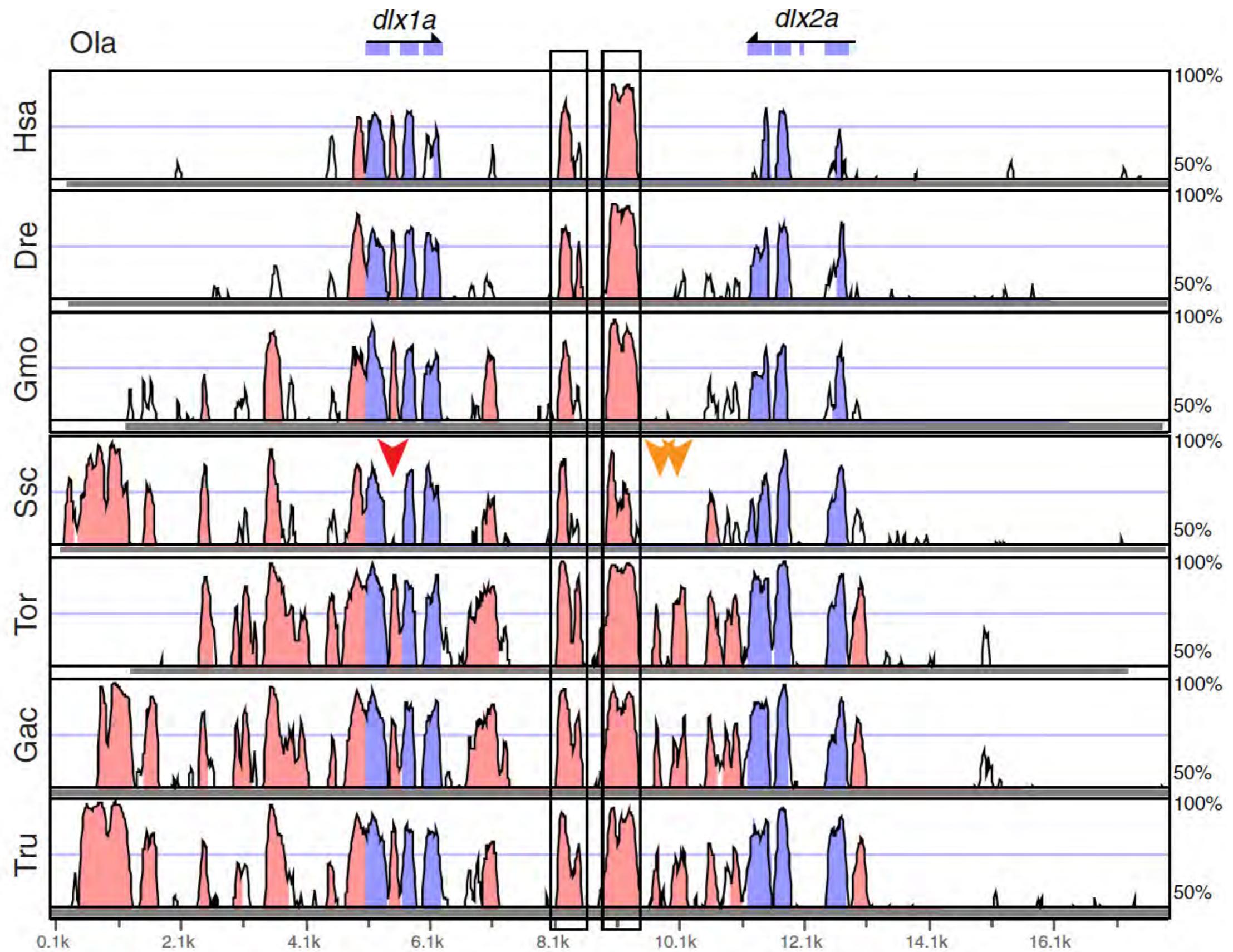
# Two instances of chromosome fusion in pipefish



# Key losses in *Hox* gene clusters



# Key losses in conserved non-coding elements



# Disruption of core hind fin patterning program

---

No evidence of *Tbx4* in the assembly



# Disruption of core hind fin patterning program

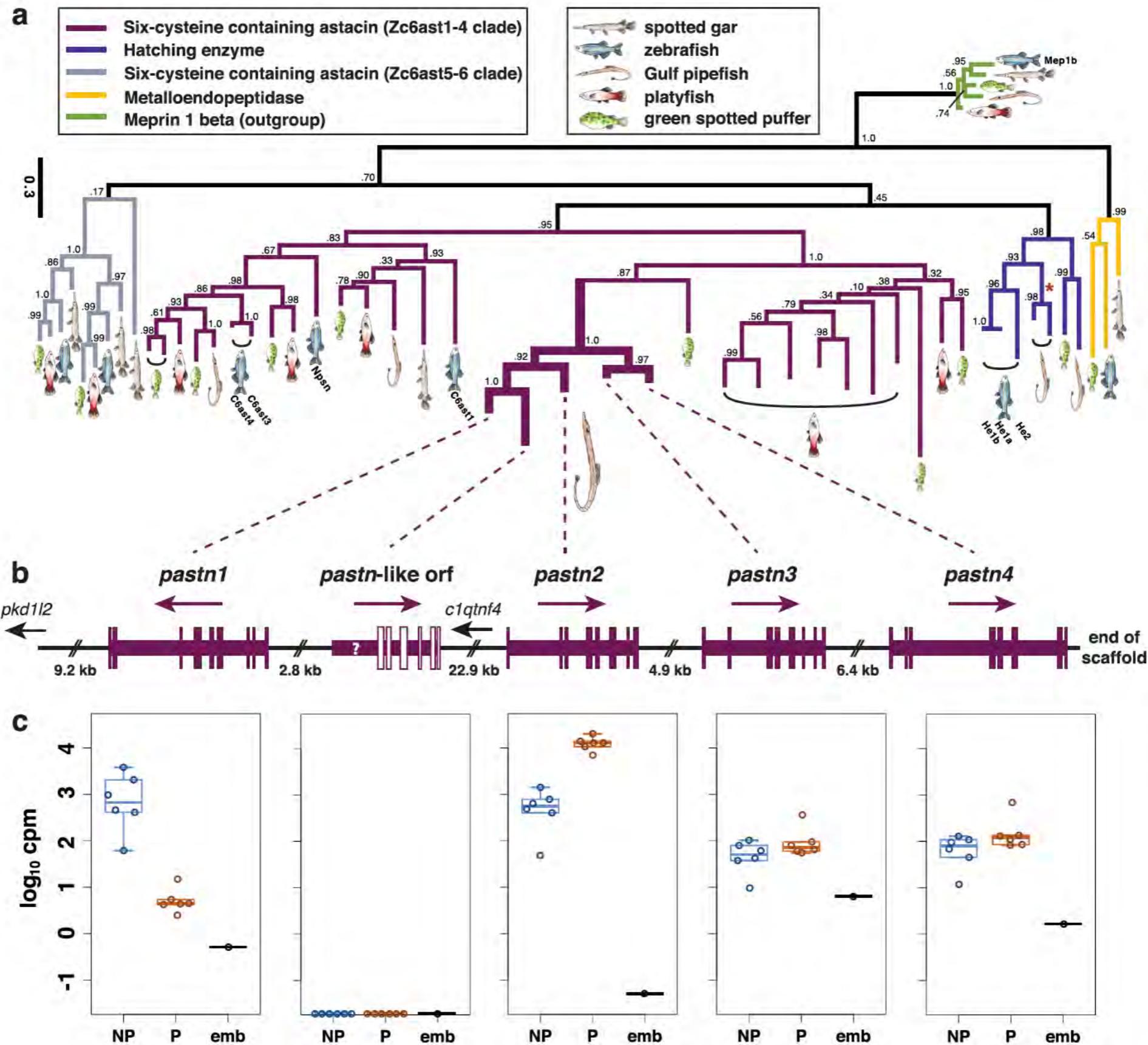
## Disruption in coding sequence of *Pitx1*

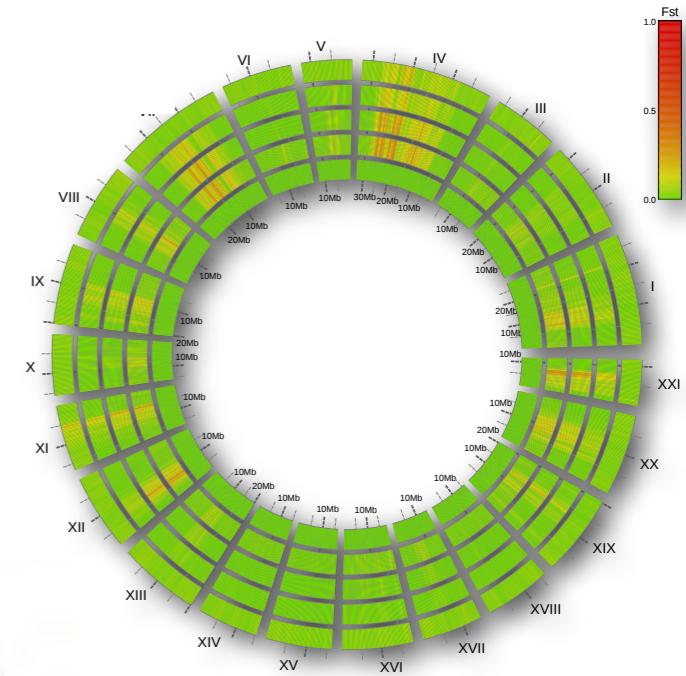


**Pitx1**

WRKRERNQQQLDLCGGYVPQFSGLVQPYED-MYA-----	AGYS-YNNWAAKSL-APAPLSTKSFTP-FNSM--SPLSSQ--SMF-SA-PSSIISSTM-----PSSMGPAG	human
WRKRERNQQMDLCNGYVPQFNGLMQSYDE-MYA-----	GYH-YNNWATKSL-TPAPLSTKGFTF-FNSM--SPLPSQ--SMF-SA-PSTIISMMN-----SSTMGHSGV	coelacanth
WRKRERNQQMDLCNSYLPQFSGLVQPYDD-MYP-----	AYT-YNNWTNKGL-APAPLSTKNFTF-FNSM--SPLTSQ--SMF-SA-PNSIISMMN-----PPTMAHTAV	gar
WRKRERNQQMDLCNSYLPQFSGLMQPYDD-VYP-----	TYT-YNNWTNKGL-TPAPLSTKNFTF-FNSM--SPLTSQ--SMF-SA-PSSIISMSM-----ASGMGHSAV	zebrafish
WRKRERNQQMDLCSGYLPQFSGLVQPYED-MYP-----	PYT-YNNWTNKGL-GPAPLSTKSFTP-FNSM--SPLTSQ--SVF-SA-PSSIISSTM-----ASGMAHSAV	cavefish
WRKRERNQQMDLCNSYLPQFSGLMQPYDD-MYP-----	ERNQQMDLCNSYLPQFSGLMQPYDD-MYP-----AYT-YNNWTNKGL-APAPLSTKNFTF-FNSM--STLTSQ--SMF-SA-PNSIISMSM-----SSGMGHSAV	ghost pipefish
WRKRERNQQMDLCKNAYLPQFSGLMQPYDDPMYP	<b>AAAAAAA</b> AYT-YNNWPNKSLHGP----KNFPF-FNSM--SPLTSQPVTMFSSS-PAPITTMHSVQ <b>AAAAAA</b> AAAHAHGGM	<b>Gulf pipefish</b>
WRKRERNQQMDLCNTYLPQFSGLMQPYED-MYP-----	ERNQQMDLCNTYLPQFSGLMQPYED-MYP-----AYT-YNNWTNKGLHGPTAPLAAKNFPFFNSM--SPLASQ--SVFSSSPTSISGMSMQH <b>AAAAAA</b> SAGMAHSGV	<b>messmate pipefish</b>
WRKRERNQQMDLCNSYLPQFSGLMQPYDD-MYP-----	WRKRERNQQMDLCNSYLPQFSGLMQPYDD-MYP-----AYT-YNNWTNKGL-APAPLSTKNFTF-FNSM--SPLTSQ--SMF-SA-PNSIISSTM-----ASGMGHSAV	medaka
WRKRERNQQMDLCNSYLPQFSGLMQPYDD-MYP-----	WRKRERNQQMDLCNSYLPQFSGLMQPYDD-MYP-----AYT-YNNWTNKGL-TPAPLSTKNFTF-FNSM--SPLTSQ--SMF-SA-PNSIISMM-----ASGMGHSAV	tilapia
WRKRERNQQMDLCNSYLPQFSGLMQPYDD-MYP-----	WRKRERNQQMDLCNSYLPQFSGLMQPYDD-MYP-----AYT-YNNWTNKGL-APAPLSTKNFTF-FNSM--SPLTSQ--SMF-SA-PNSIISMM-----ASGMGPSAV	pufferfish
WRKRERNQQMDLCSSYLPQFSGLMQPYED-MYP-----	WRKRERNQQMDLCSSYLPQFSGLMQPYED-MYP-----TYS-YNNWPNKGL-APAPLSSKNFTF-FNSM--SPLTSQ--SMF-SA-PNSIISMM-----APGMGPPAA	stickleback
WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----	WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----GYS-YNNWAAKGL-TSASLSTKSFTP-FNSMNVNPLSSQ--SMF-SP-PNSIISMSM-----SSMVPSAV	human
WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----	WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----SYS-YNNWAAKGL-TSASLSTKSFTP-FNSMNVNPLSSQ--TMF-SP-PNSIISMSM-----SSMVPS-V	coelacanth
WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----	WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----SYT-YNNWAAKGL-TSASLSTKSFTP-FNSMNVNPLSSQ--TMF-SP-PNSIISMSM-----SSMVPSAV	gar
WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----	WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----SYT-YNNWAAKGL-TSASLSTKSFTP-FNSMNVNPLSSQ--TMF-SP-PNSIISMSM-----SSMVPSAV	zebrafish
WRKRERNQQAELCKGGFGAQFNGLVQPYED-MYA-----	WRKRERNQQAELCKGGFGAQFNGLVQPYED-MYA-----SYPPYNNWAAKSL-APASLSAKSFTP-FNSVNVSPLSSQ--AVF-SP-PTSISMSV-----SSGMVPT--	cavefish
WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----	WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----GYT-YNNWAAKGL-TSASLSTKSFTP-FNSMNVNPLSSQ--AMF-SP-PNSIISM-----TSGMVPSAV	<b>Gulf pipefish</b>
WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----	WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----SYT-YNNWAAKGL-TSASLSTKSFTP-FNSMNVNPLSSQ--TMF-SP-PNSIISM-----TSSMVPAAV	medaka
WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----	WRKRERNQQAELCKNGFGPQFNGLMQPYDD-MYP-----SYT-YNNWAAKGL-TSASLSTKSFTP-FNSMNVNPLSSQ--TMF-SP-PNSIISM-----TSSMVPSAV	tilapia
WRKRERNQQAELCKNGFGPQFNGLMQPYED-MYP-----	WRKRERNQQAELCKNGFGPQFNGLMQPYED-MYP-----SYT-YNNWAAKGL-TPASLSTKSFTP-FNSMNVNPLSSQ--TMF-SP-SPAPNSIISM-----TSGMVPSAV	pufferfish
WRKRERNQQAELCKNGFGPQFNGLMQPYED-MYP-----	WRKRERNQQAELCKNGFGPQFNGLMQPYED-MYP-----SYT-YNNWAAKGL-TSASLSTKSFTP-FNSMNVNPLSSQ--TMF-SP-SNSIISM-----TSSMVPSAV	stickleback
WRKRERSQQAELCKGSFAAPLGLLVPPYEE-VYP-----	WRKRERSQQAELCKGSFAAPLGLLVPPYEE-VYP-----GYS-YGNWPPKAL-A-PPLAAKTFPFCAFNSVNVLGPLASQ--PVF-SP-PSSIAASMV-----PSAAAAPGT	human
WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYS-----	WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYS-----GYS-YNNWATKGL-ATSPLSAKSFTP-FNSMNVSPLSSQ--PMF-SP-PSSIASMTM-----PSSMVPSAV	coelacanth
WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYS-----	WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYS-----GYS-YNNWATKSL-ATSPLSAKSFTP-FNSMNVSPLSSQ--PMF-SP-PSSIISMMN-----ASSMVPSAV	gar
WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYS-----	WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYS-----GYS-YNNWATKSL-ASSPLSAKSFTP-FNSMNVSPLSSQ--PMF-SP-PSSIIPSMNM-----ASSMVPSAV	zebrafish
WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYS-----	WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYS-----GYS-YNNWATKSL-ASSPLSAKSFTP-FNSMNVSPLSSQ--PMF-SP-PSSIIPSMNM-----ASSMVPSAV	cavefish
WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYA-----	WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYA-----GYS-YNNWASKSL-AGGQLSAKSFTP-FNSMNVSPLSSQ--PMF-SP-PSSMPSMNM-----ASGMVPSAV	<b>Gulf pipefish</b>
WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYT-----	WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYT-----GYS-YNNWATKSL-ASSPLSTKSFTP-FNSMNVSPLSTQ--PMF-SP-PSSIIPSMNM-----ASSMVPSAV	medaka
WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYT-----	WRKRERNQQAELCKNGFGAQFNGLMQPYDD-MYT-----GYS-YNNWATKSLAASSPLSAKSFTP-FNSMNVSPLSSQ--PMF-SP-PSSIIPSMNM-----ASSMVPTAV	tilapia
WRKRERNQQAELCKNGFGAQFNGLMQPYDD-VYS-----	WRKRERNQQAELCKNGFGAQFNGLMQPYDD-VYS-----GYS-YNNWAAKGL-ASSPLSAKSFTP-FNSMNVSPLSSQ--SMF-SP-PSSLPSMNM-----ASSMVPSAV	pufferfish
WRKRERNQQAELCKNGFGSQFNGLMQPYDD-MYT-----	WRKRERNQQAELCKNGFGSQFNGLMQPYDD-MYT-----GYS-YNNWATKSL-ASSPLSAKSFTP-FNSMNVSPLSSQ--PMF-SP-PSSIIPSMNM-----ASSMVPSAV	stickleback

# Expansion of male pregnancy specific gene family





Ecological & evolutionary genomic  
analyses using RAD-seq -  
what have we learned?

# Overall Conclusions

---

RAD-seq can be a tool for enabling new research in models & nonmodels

- SNP identification and genotyping
- documenting patterns of genetic variation
- identifying the molecular genetic basis of important phenotypic variation
- assessing how ecological processes structure this genetic variation in genomes
- analytical and computational approaches are challenging but manageable

Not your father's (or your mother's) genome assembly

- a mixture of data types can be efficiently combined
- a genetic map is extremely useful for pulling it all together
- having a tiled genome of smaller contigs is often good enough

*Open Source Genomics* provides a suite of breakthrough technologies

- the molecular approaches are not as daunting as they first appear
- analytical and computational approaches are challenging
  - **New software tools can help, but knowledge of Unix and Scripting is essential**
  - **Also essential to be comfortable with classical and modern statistics**

# Acknowledgments

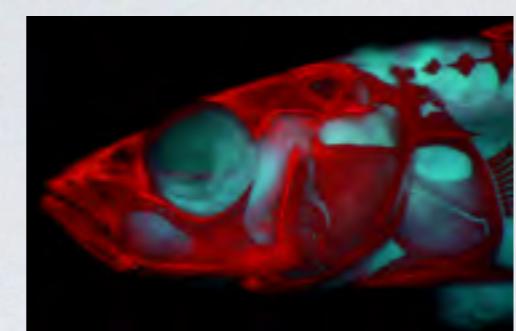
---



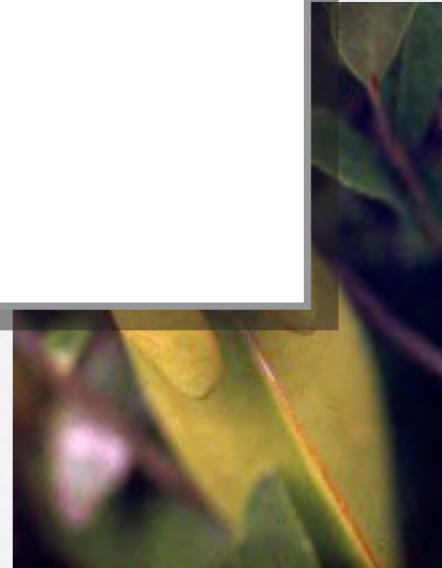
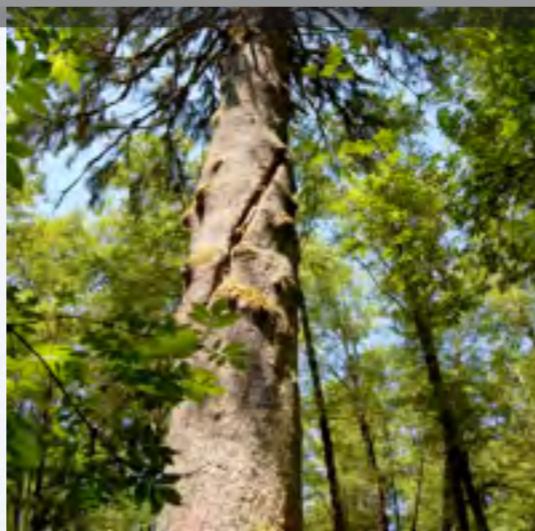
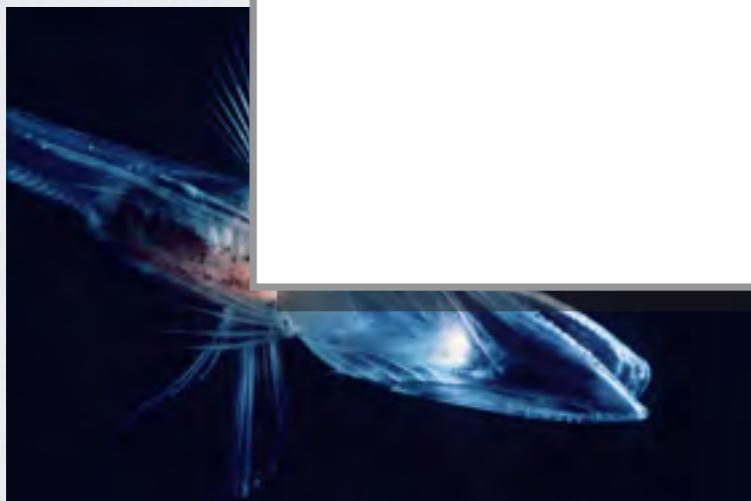
- *Past and present lab members **Paul Hohenlohe, Thom Nelson, Joe Dunham, Nicole Nishimura & Mark Currey***
- *Collaborators **Eric Johnson, Patrick Phillips, Chuck Kimmel, John Postlethwait***
- *Funding from NSF & NIH, as well as Keck & Murdock Foundations*







# Lab bench considerations for RADseq studies



# Statistical considerations in RAD-seq

T T G T T T T T T T T T T T T T T G T T

T  
T  
G  
T  
T  
T  
T  
T  
T  
T  
G  
T  
T

The reads are 14 T and 2 G:

GT heterozygote?

GG homozygote with error?

AA homozygote with lots of error?

Needed a rigorous method to call genotypes

T  
T  
G  
T  
T  
T  
T  
T  
T  
T  
G  
T

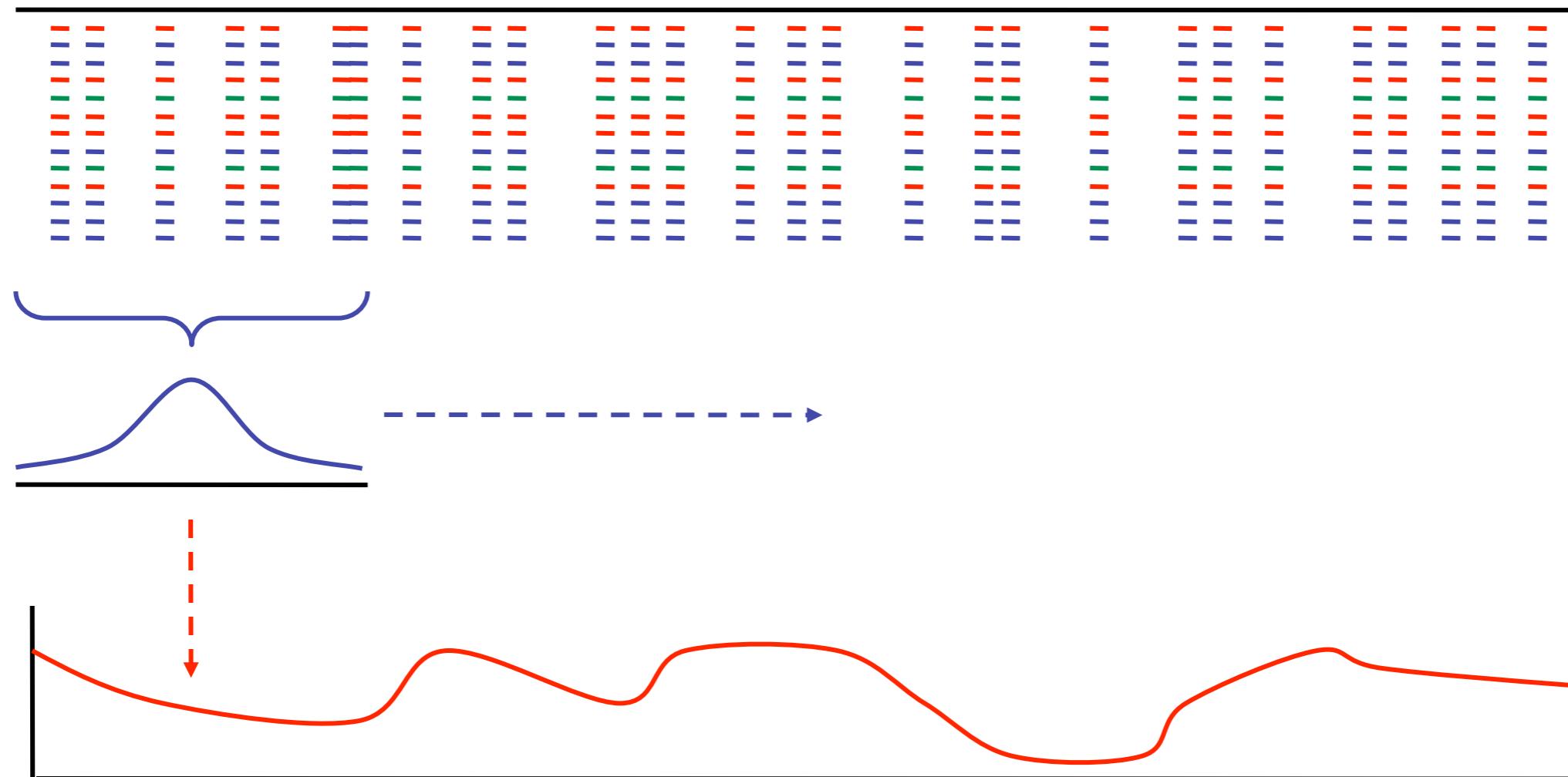
$$L(n_1 \text{ hom}) = P(n_1, n_2, n_3, n_4) = \frac{n!}{n_1! n_2! n_3! n_4!} \left(1 - \frac{3\epsilon}{4}\right)^{n_1} \left(\frac{\epsilon}{4}\right)^{n_2} \left(\frac{\epsilon}{4}\right)^{n_3} \left(\frac{\epsilon}{4}\right)^{n_4}$$

$$L(n_1 n_2 \text{ het}) = P(n_1, n_2, n_3, n_4) = \frac{n!}{n_1! n_2! n_3! n_4!} \left(0.5 - \frac{\epsilon}{4}\right)^{n_1} \left(0.5 - \frac{\epsilon}{4}\right)^{n_2} \left(\frac{\epsilon}{4}\right)^{n_3} \left(\frac{\epsilon}{4}\right)^{n_4}$$

Maximum likelihood genotyping based on multinomial distribution of nucleotide reads

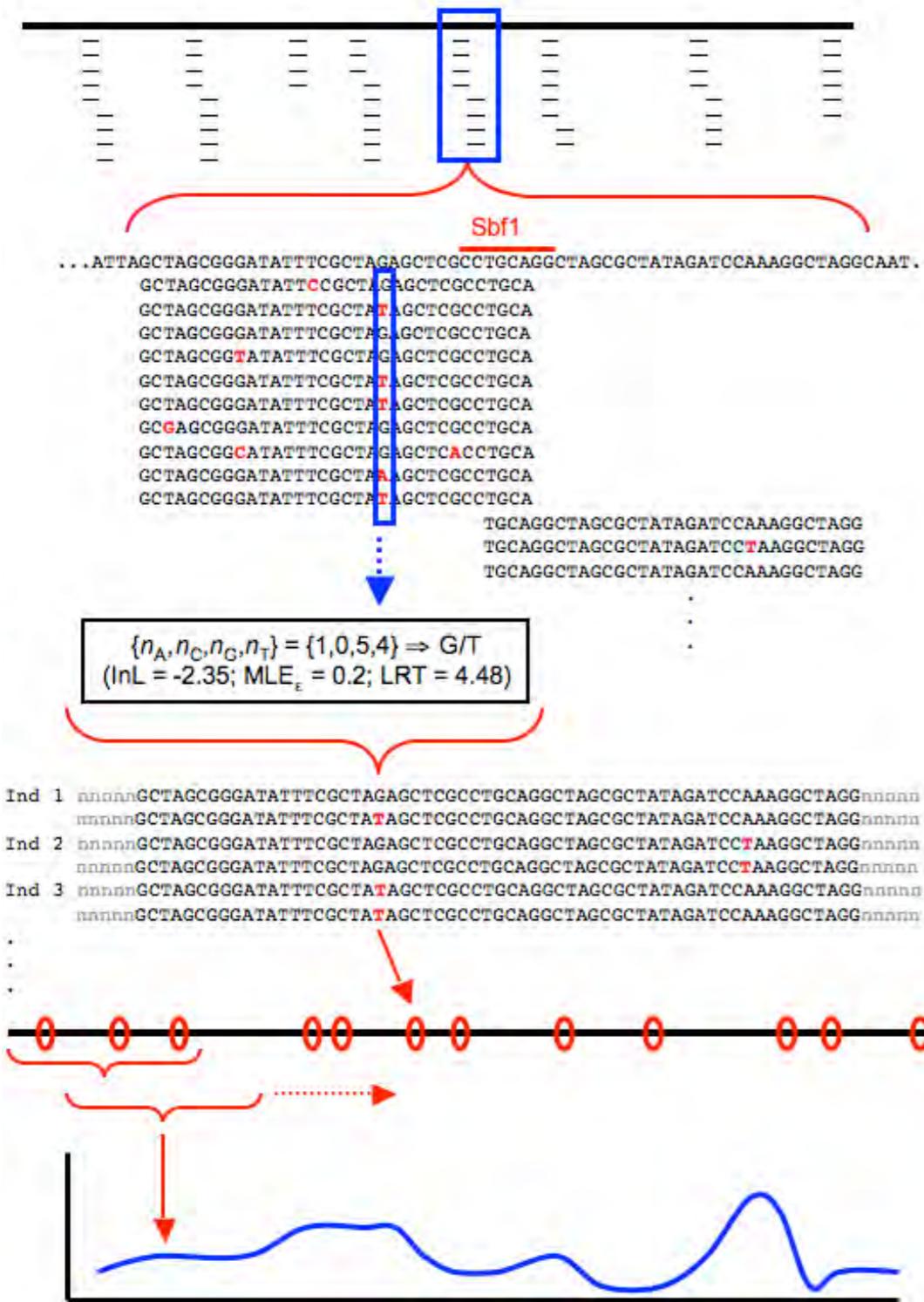
# Making statistics continuous across the genome

Kernel-smoothing average of summary statistics along genome



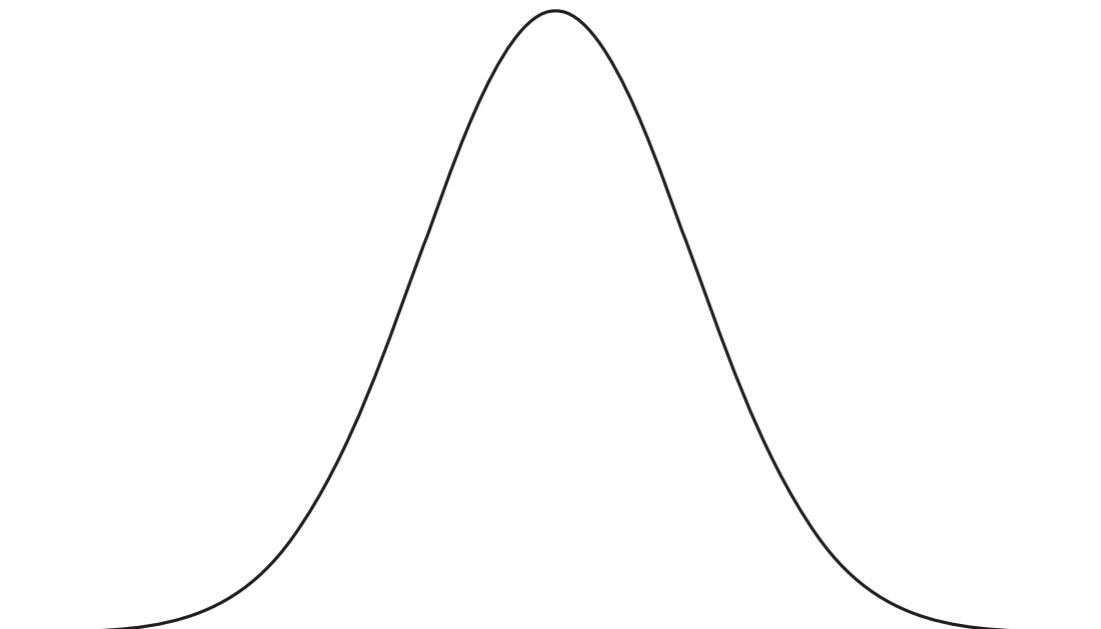
Bootstrap re-sampling to estimate significance of moving average

# Overall pipeline



# ‘Bias’ in RAD-sequencing

---



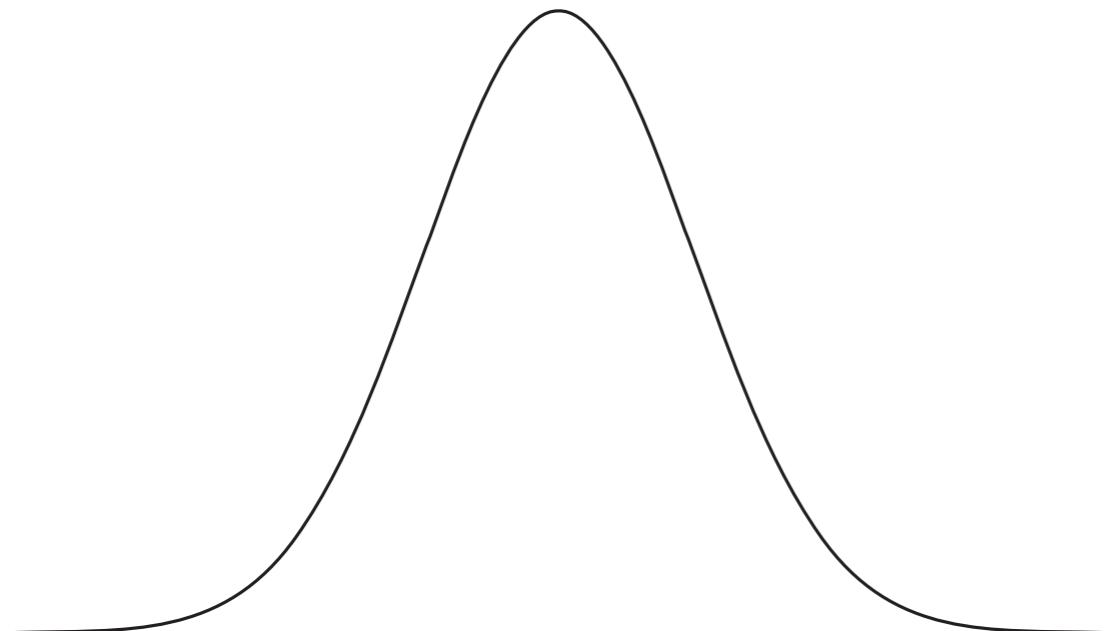
$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

$e = 2.7182\dots$

$\pi = 3.1415\dots$

# ‘Bias’ in RAD-sequencing

---



$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

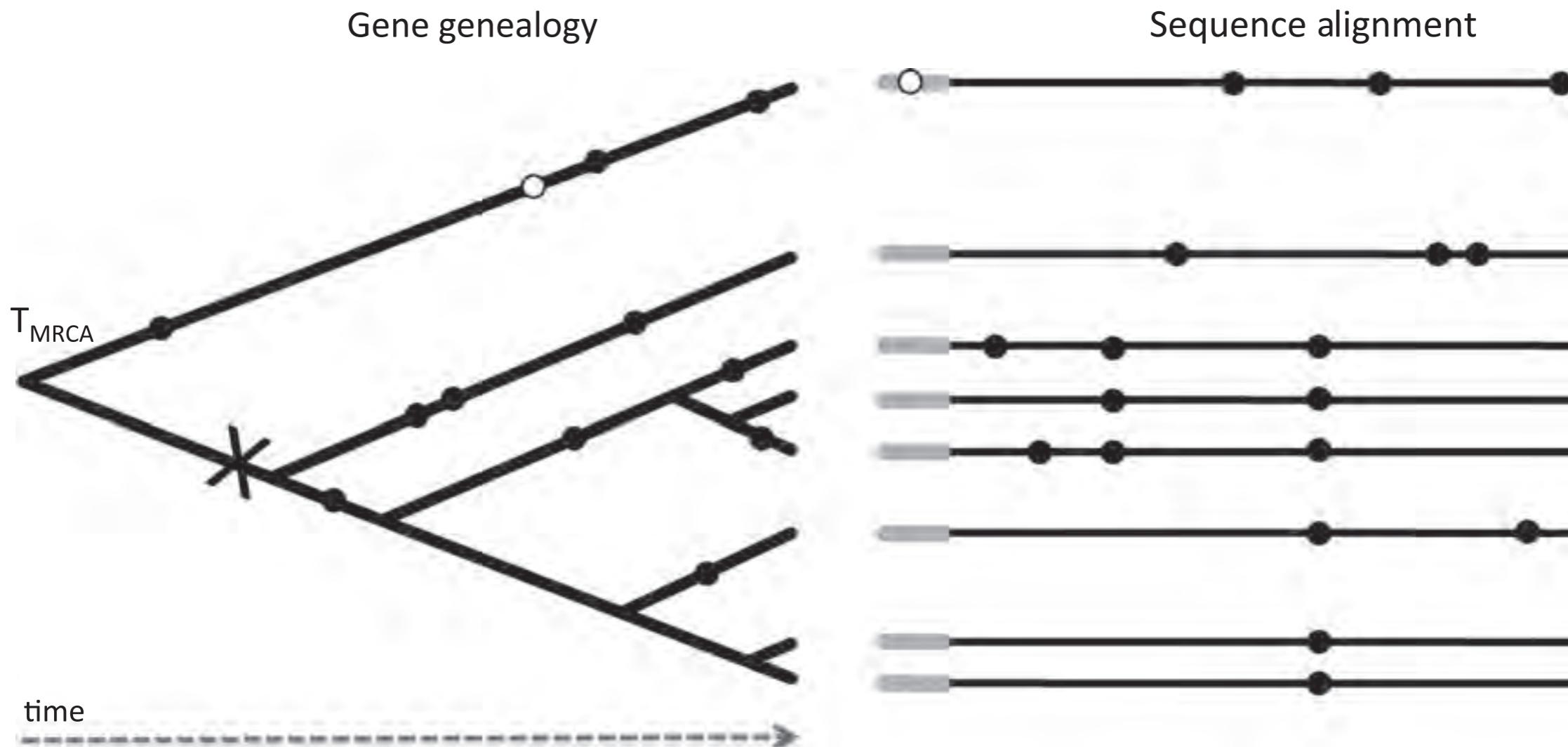
$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

# Bias in RAD-sequencing

**RADseq underestimates diversity and introduces genealogical biases due to nonrandom haplotype sampling**

B. ARNOLD,<sup>1</sup> R. B. CORBETT-DETIG,<sup>1</sup> D. HARTL and K. BOMBLIES

*Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA 02138, USA*



# Bias in RAD-sequencing summary

Protocol	$\theta$ per bp	Mean	
		$\theta_{we}/\theta_{wa}$	$\pi_e/\pi_a$
Standard	0.0001	0.994	0.995
	0.001	0.987	0.982
	0.01	0.956	0.933
Double digest	0.0001	0.835	0.836
	0.001	0.858	0.851
	0.01	0.829	0.797

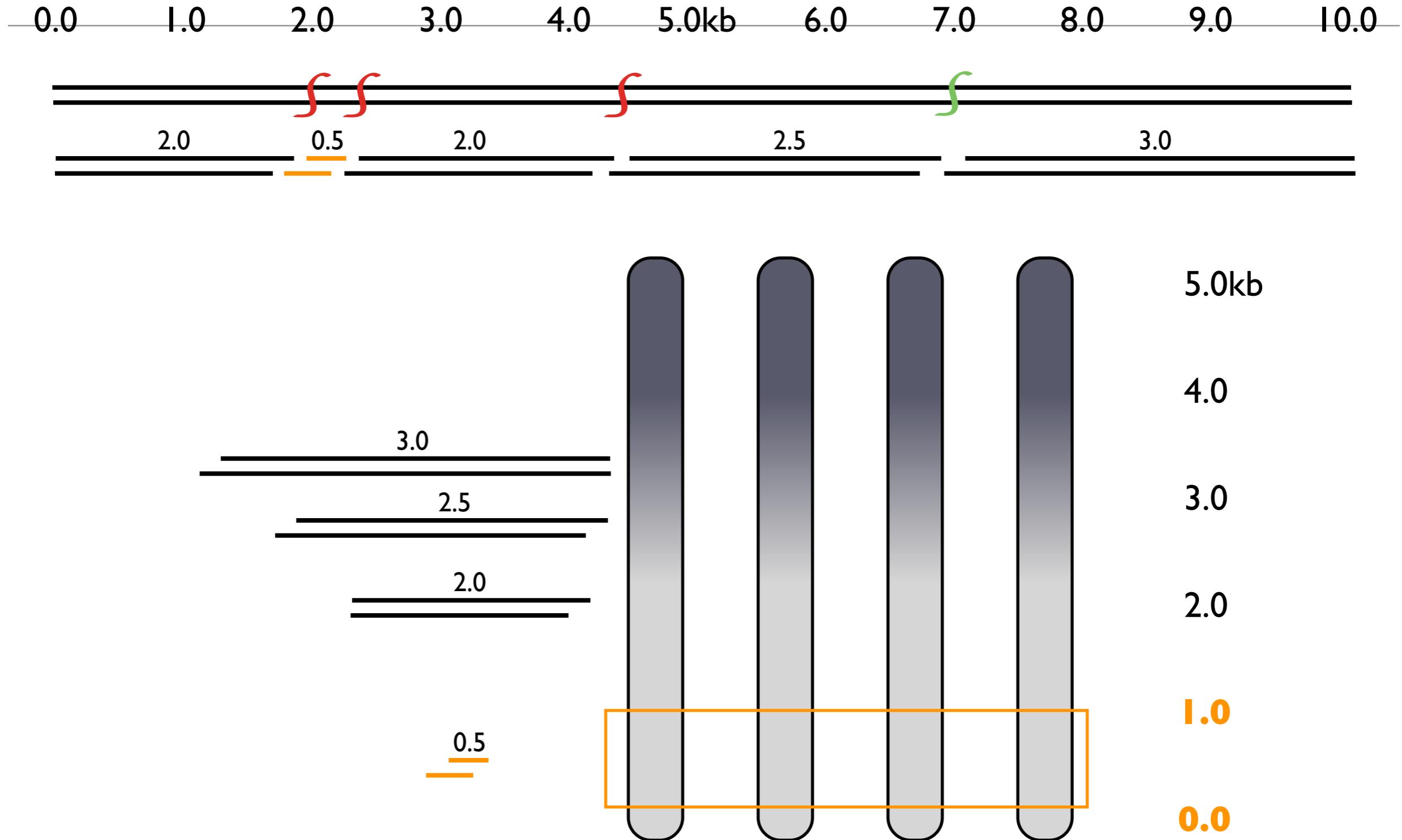
# Bias in RAD-sequencing summary

Protocol	$\theta$ per bp	Mean	
		$\theta_{we}/\theta_{wa}$	$\pi_e/\pi_a$
Standard	0.0001	0.994	0.995
	0.001	0.987	0.982
	0.01	0.956	0.933
Double digest	0.0001	0.835	0.836
	0.001	0.858	0.851
	0.01	0.829	0.797

# Bias in RAD-sequencing summary

Protocol	$\theta$ per bp	Mean	
		$\theta_{we}/\theta_{wa}$	$\pi_e/\pi_a$
Standard	0.0001	0.994	0.995
	0.001	0.987	0.982
	0.01	0.956	0.933
Double digest	0.0001	0.835	0.836
	0.001	0.858	0.851
	0.01	0.829	0.797

# Why is ddRAD so much more biased?



# Experimental design considerations for RAD

---

*Tradeoffs:*

**Number** of sites versus **Depth** of sequencing per site versus **Number of samples**

# Experimental design considerations for RAD

---

*Tradeoffs:*

**Number** of sites versus **Depth** of sequencing per site versus **Number of samples**

raw reads / samples / sites = coverage at each RAD locus

1,000,000 / 100 / 1,000 = 10x coverage

25 to 50x average coverage per RAD locus is a good goal

# Experimental design considerations for RAD

---

Tradeoffs:

**Number** of sites versus **Depth** of sequencing per site versus **Number of samples**

How many tags do I need?

Things to consider

Choice of enzyme and genome size       $(0.25)^n \times \text{genome size} = \text{expected } \# \text{ sites}$

Genomes are biased:

expect 112,300 six-cutter sites in stickleback (460 Mb)	actual <b>EcoRI</b> sites = 90,000
expect 7000 eight-cutter sites in stickleback	actual <b>SbfI</b> sites = 22,800
expect 32,900 six-cutter sites in <i>C. remanei</i> (135 Mb)	actual <b>EcoRI</b> sites = 73,200

# Experimental design considerations for RAD

---

*Tradeoffs:*

**Number** of sites versus **Depth** of sequencing per site versus **Number of samples**

How many tags do I need?

Things to consider

Choice of enzyme and genome size

Polymorphism and read length

Nucleotide polymorphism rate = 0.01 to 0.001 for most vertebrates

Stickleback populations: 0.01 to 0.02. At least 1 SNP every 100 bp, on average

# Experimental design considerations for RAD

*Tradeoffs:*

**Number** of sites versus **Depth** of sequencing per site versus **Number of samples**

How many samples should be multiplexed?

*Things to consider*

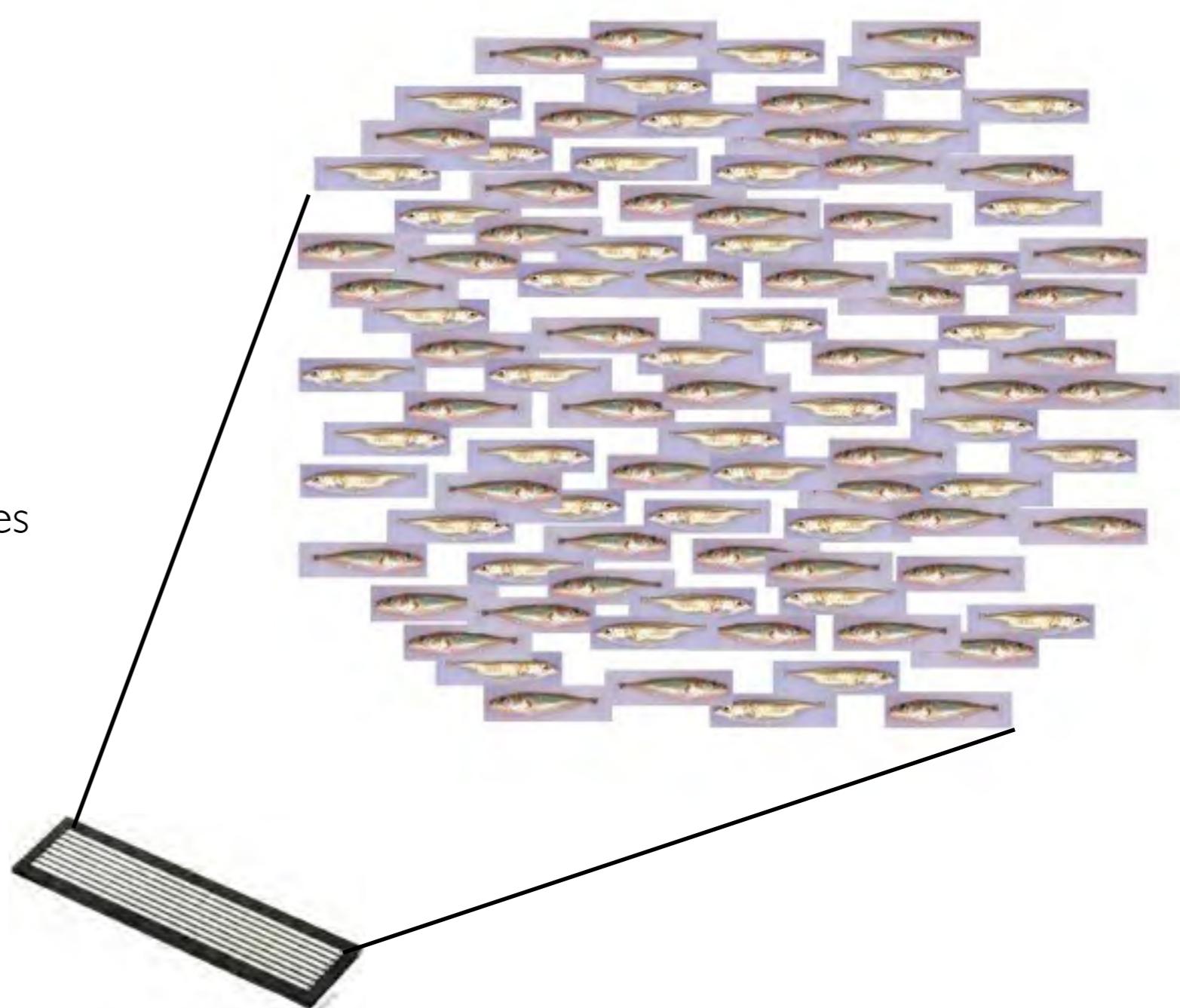
Barcoded adapters

5 to 8nt barcodes

Variable length barcodes

Combinatorial barcodes (PE)

Barcode distance - two mismatches



# Molecular considerations in library building

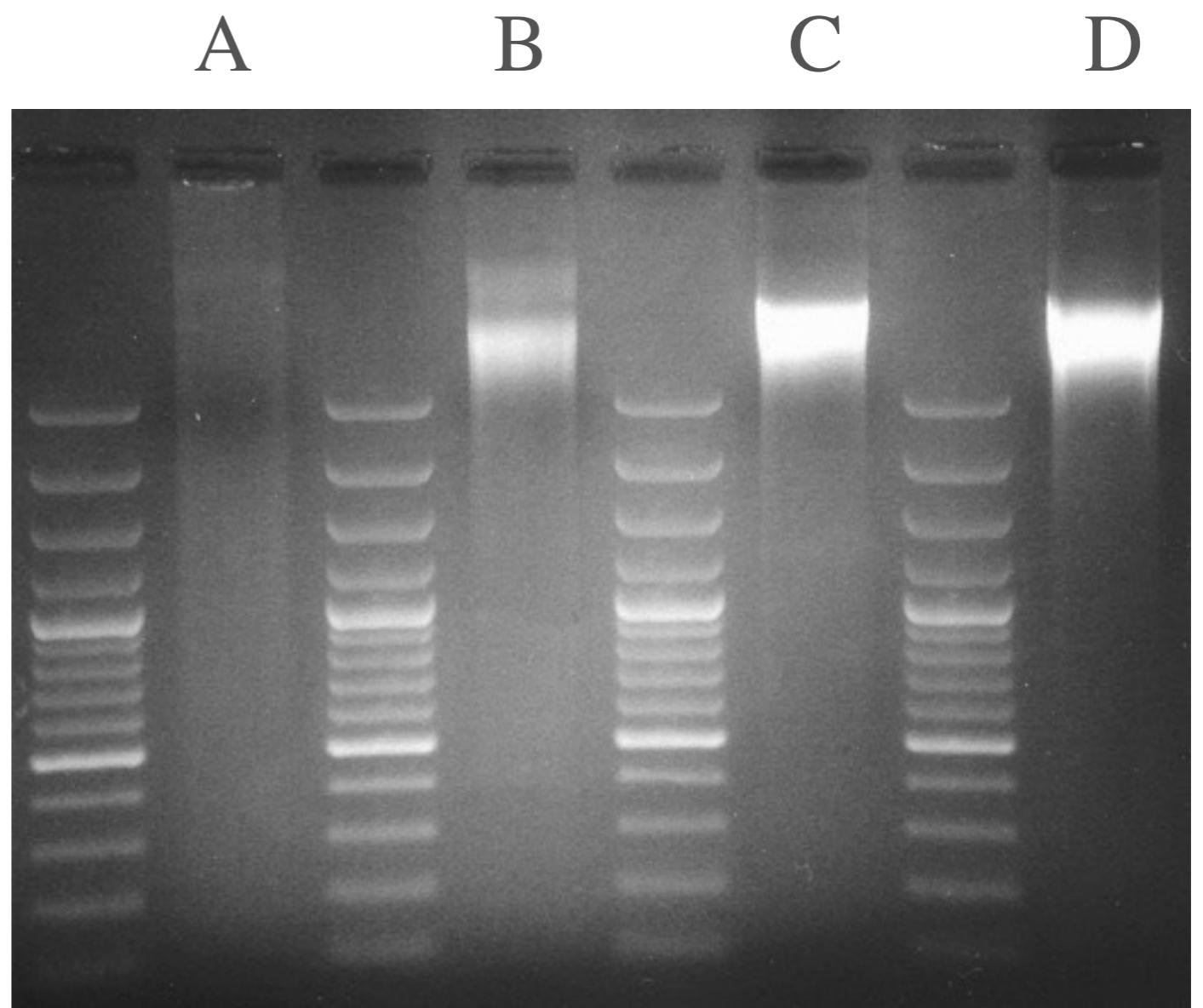
---

How many samples should be multiplexed?

Things to consider

## DNA Quality

Multiplex only like samples to help equalize representation of poor quality samples



# Molecular considerations in library building

---

How many samples should be multiplexed?

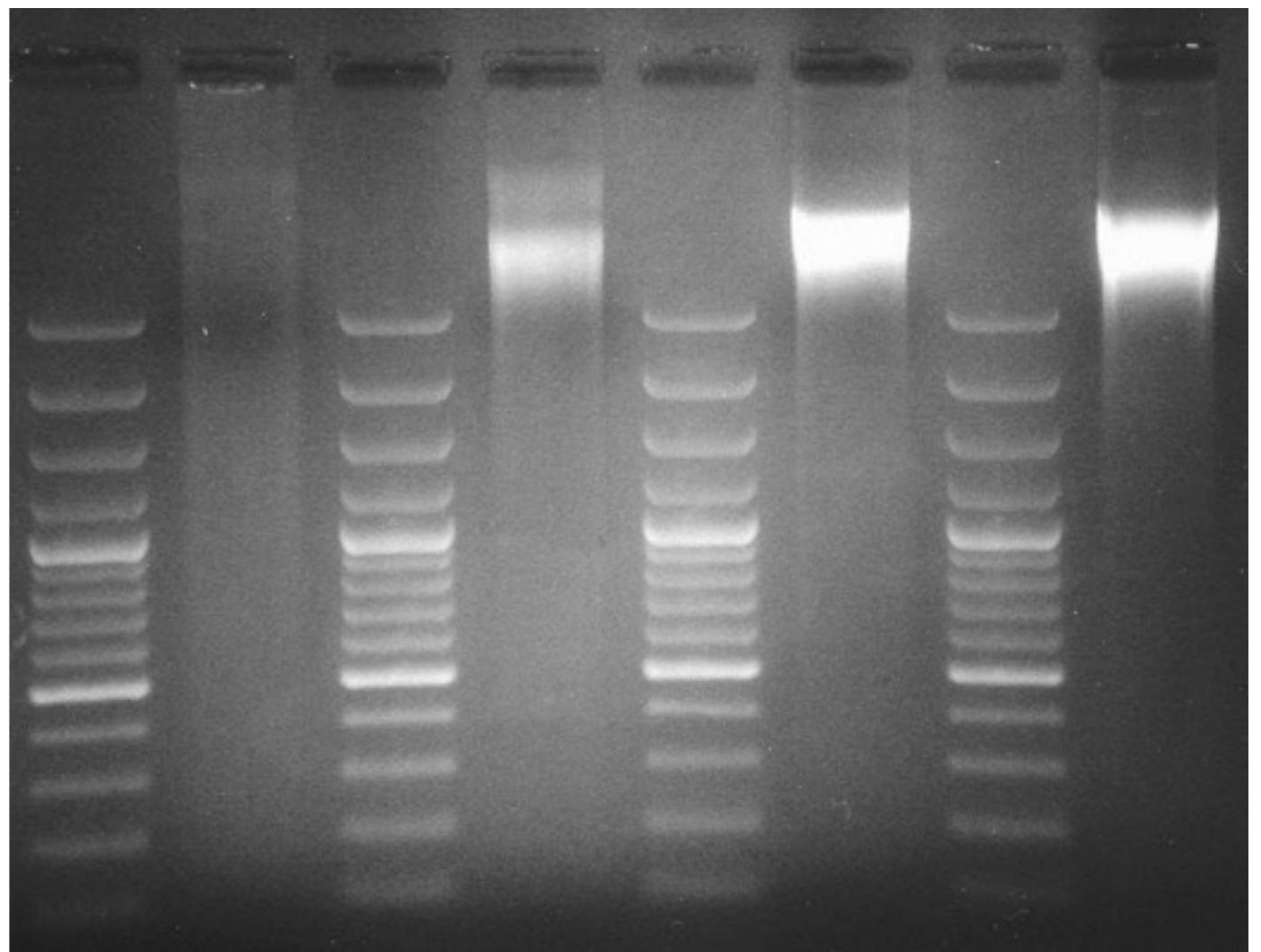
Things to consider

DNA Quality

[Diversify barcodes](#)

Illumina cluster calling is  
confused by repetition in first  
4 bases - can offset barcodes

CGATA      GTACA      TAGCC      ACTGC



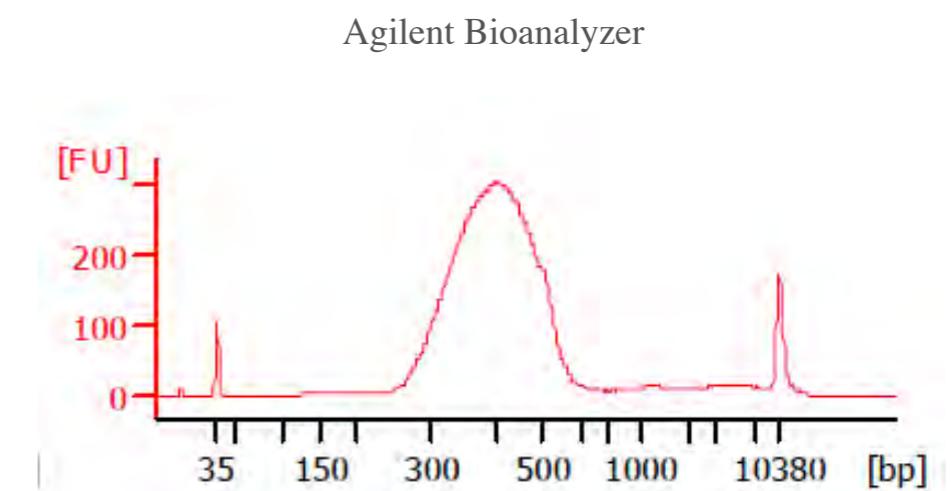
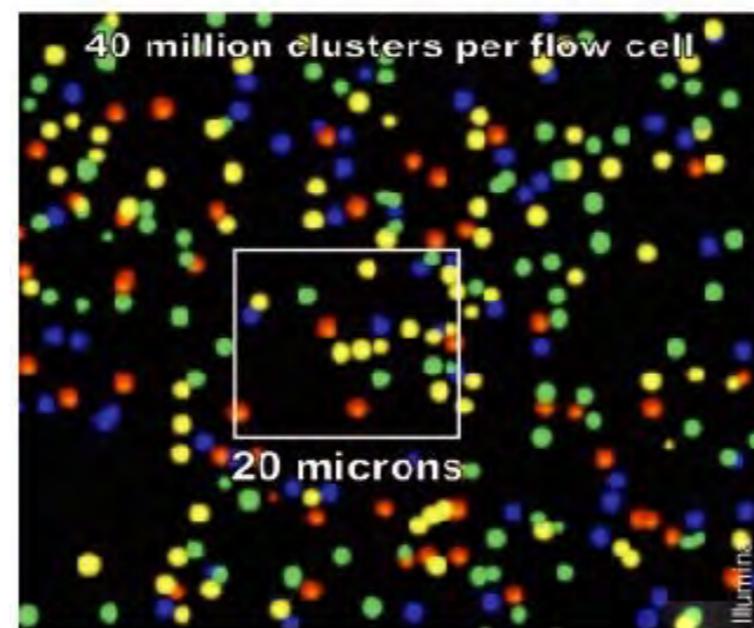
# Molecular considerations in library building

How can I get the best depth of coverage?

Things to consider

Fragment size

Smaller/tighter is better



# Molecular considerations in library building

How can I get the best depth of coverage?

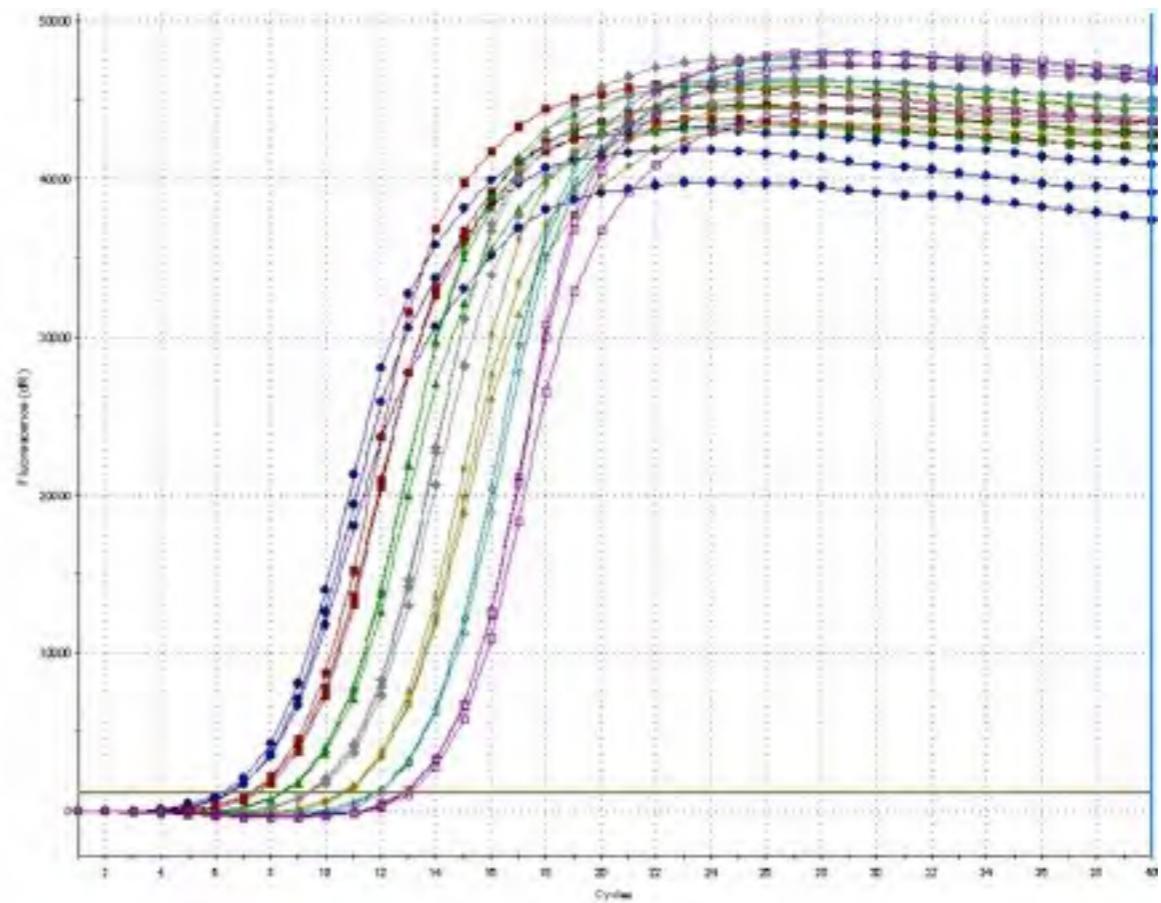
Things to consider

Fragment size

Library quality

qPCR

qPCR control should be similar to measured sample:



# Molecular considerations in library building

How can I get the best depth of coverage?

Things to consider

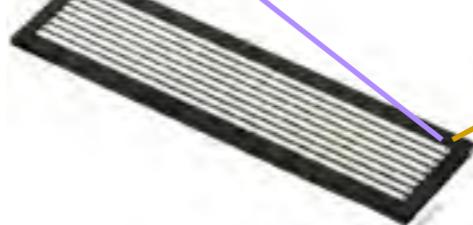
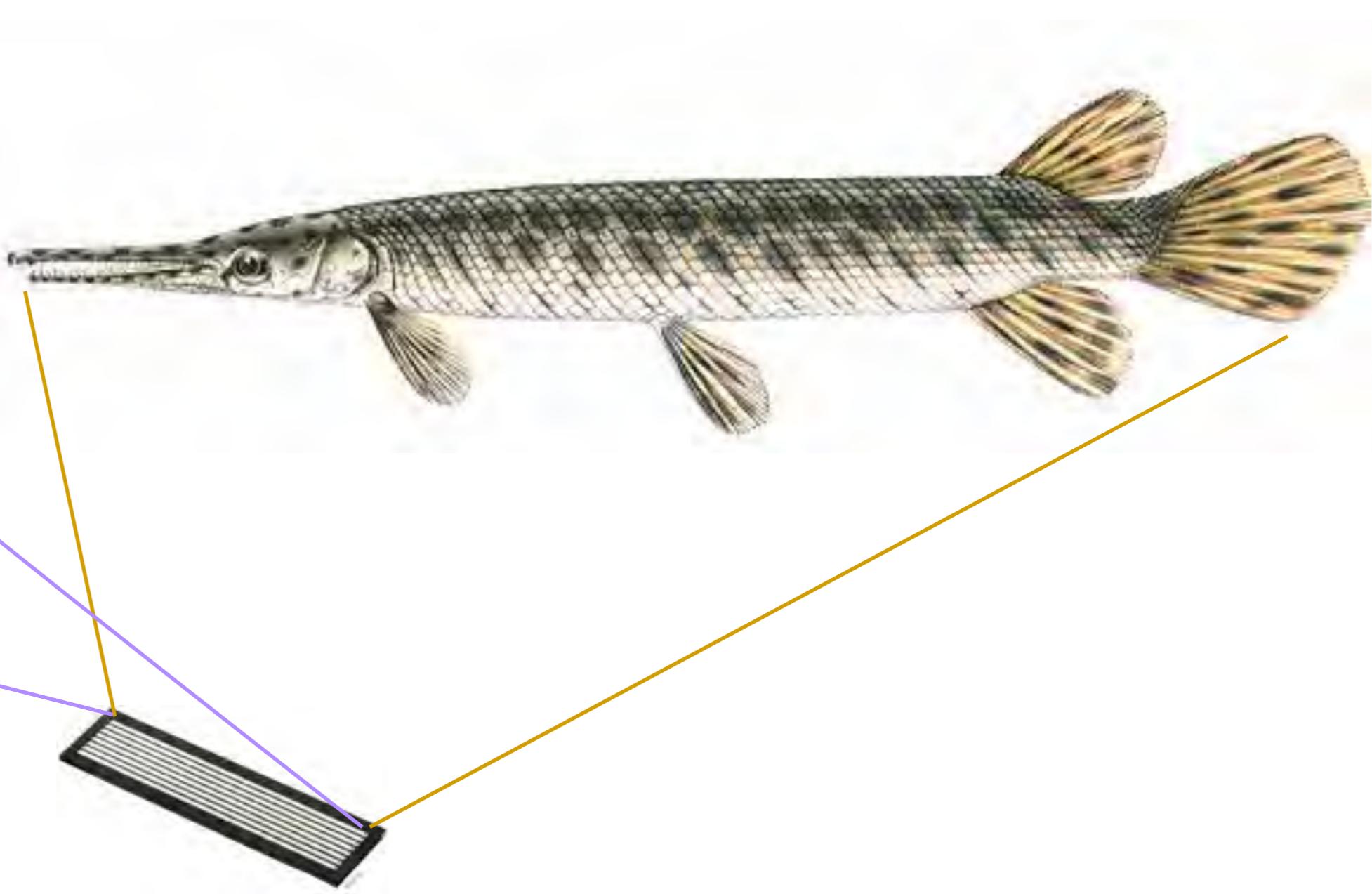
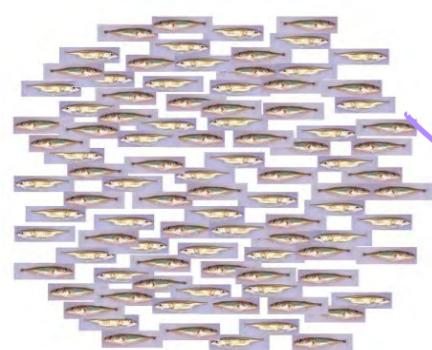
Fragment size

[Library quality](#)

qPCR

Pilot Experiment:

Spike or split a lane



# The pipefish genetic map is closed; 22 LGs 6000 segregating SNPs; 30,000 RAD sites

