

Ecological and evolutionary genomic analyses using RAD-seq

2018 Workshop on Genomics
Český Krumlov

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



Outline for today's lecture

RAD-seq for genomic questions

Primer on population genomics

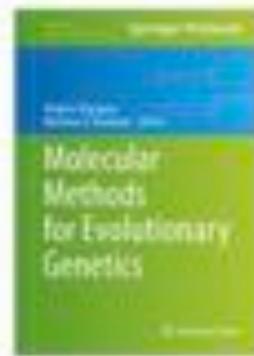
Evolutionary genomics of stickleback fish

- 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)



[Molecular Methods for Evolutionary Genetics](#) pp 157-178 | [Cite as](#)

SNP Discovery and Genotyping for Evolutionary Genetics Using RAD Sequencing

Authors

[Authors and affiliations](#)

Paul D. Etter, Susan Bassham, Paul A. Hohenlohe, Eric A. Johnson, William A. Cresko [✉](#)

STUDY DESIGNS

Harnessing the power of RADseq for ecological and evolutionary genomics

Kimberly R. Andrews¹, Jeffrey M. Good², Michael R. Miller³, Gordon Luikart⁴ and Paul A. Hohenlohe⁵

nature
REVIEWS GENETICS

nature
protocols

Altmetric: 54

[More detail >](#)

Protocol

Deriving genotypes from RAD-seq short-read data using Stacks

Nicolas C Rochette & Julian M Catchen [✉](#)

Nature Protocols 12, 2640–2659 (2017)

Published online: 30 November 2017

Methods in Ecology and Evolution

[Explore this journal >](#)

Research Article

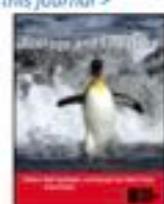
Lost in parameter space: a road map for STACKS

Josephine R. Paris, Jamie R. Stevens, Julian M. Catchen [✉](#)

First published: 18 April 2017 [Full publication history](#)

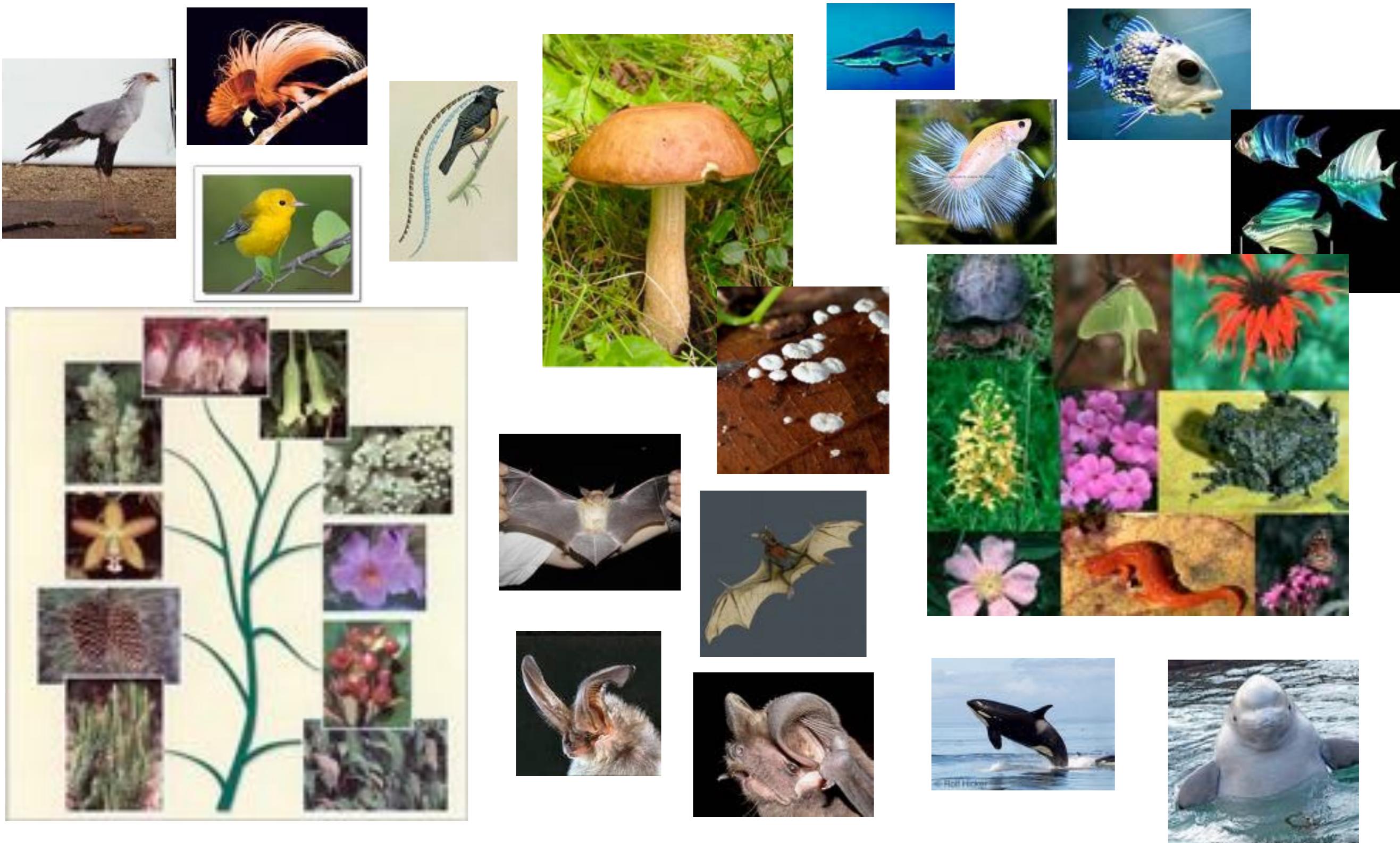
DOI: 10.1111/2041-210X.12775 [View/Save citation](#)

Cited by (CrossRef): 0 articles [Check for updates](#) [Citation tools](#)

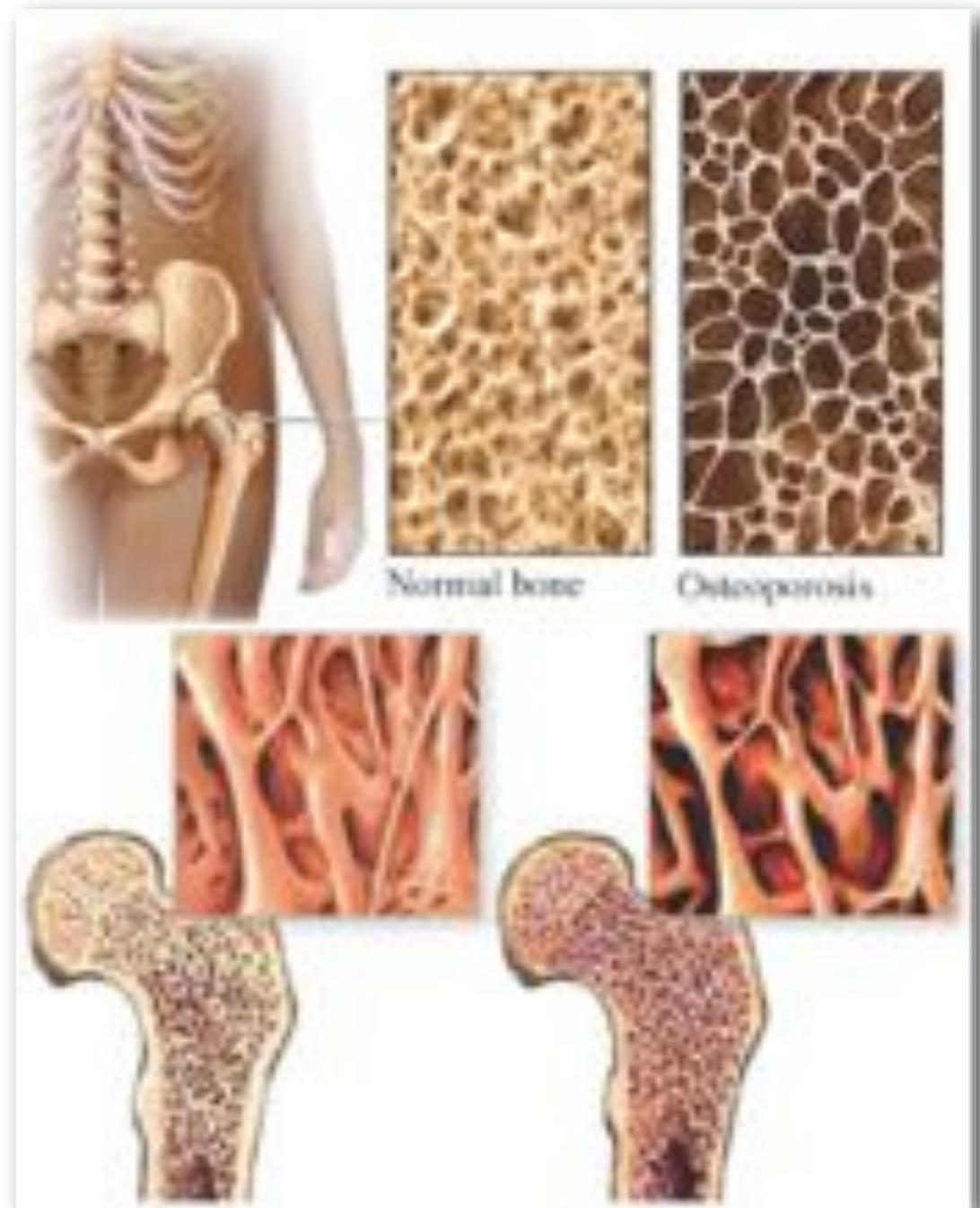


[View Issue TOC](#)
Volume 8, Issue 10
October 2017
Pages 1360–1373

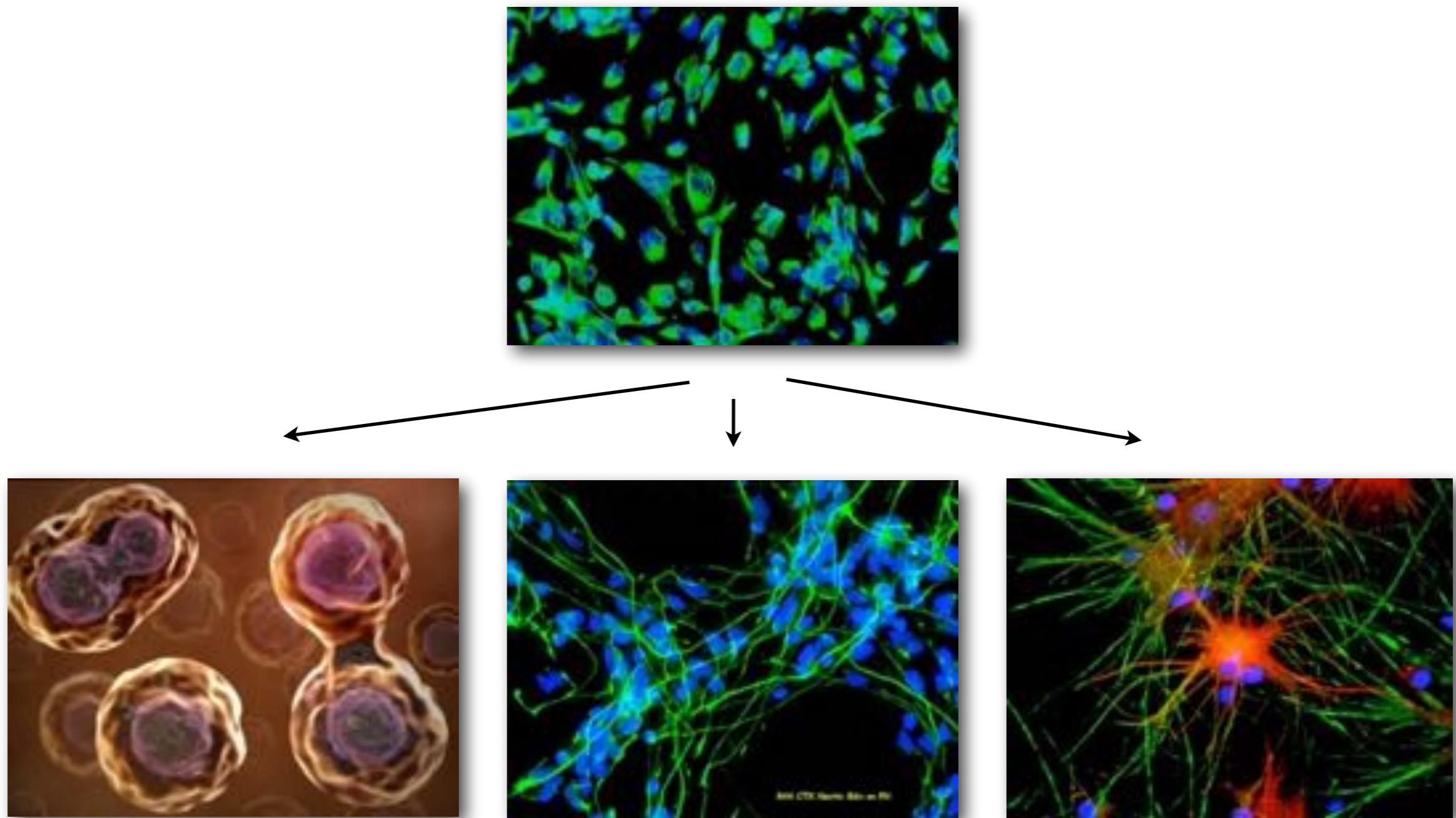
Why do species look the way that they do?



Why do organisms vary?

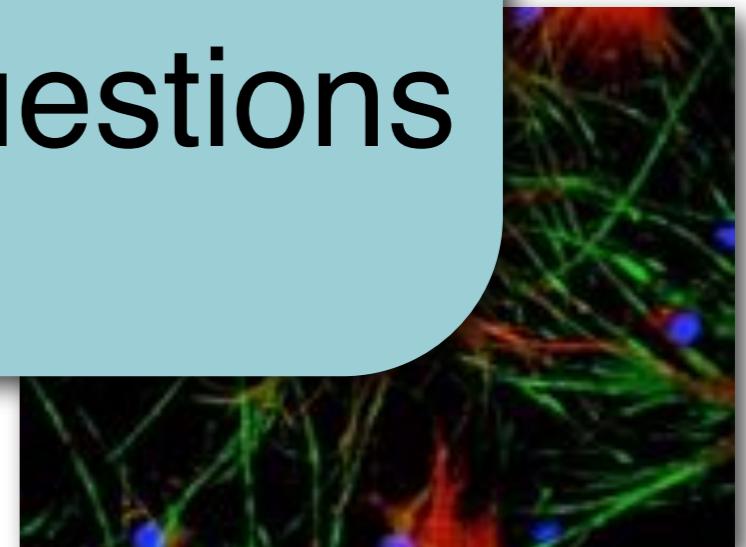


How is cellular functional diversity created?



How is cellular functional diversity created?

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





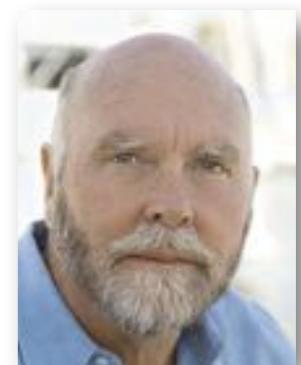
(1998)



(2000)



(2002)



(2006)



(2006)



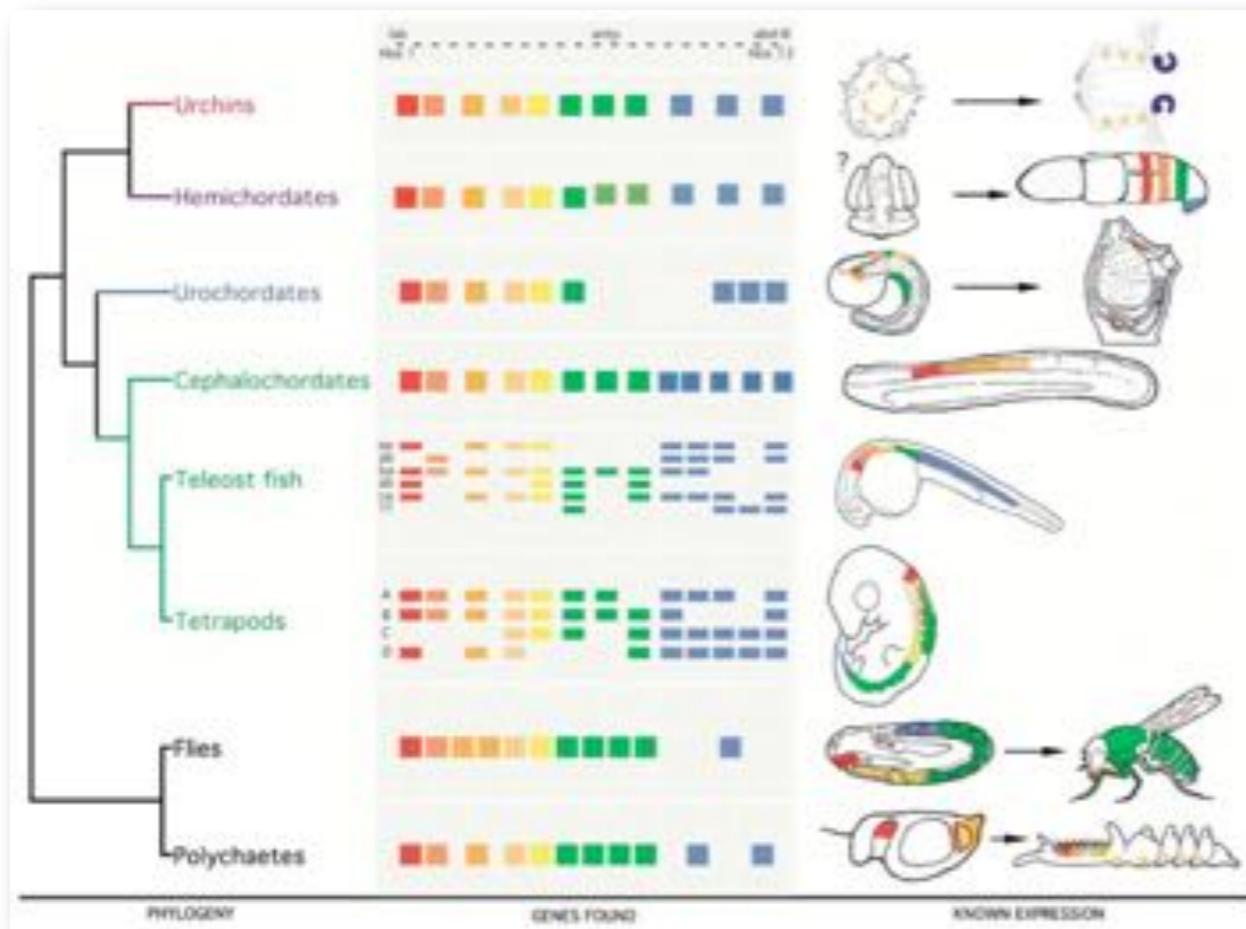
(2002)



(2007)



(2006)



Comparative Genomics

Vertebrate zygotes or embryos



28 day human



19h zebrafish



Vertebrate zygotes or embryos



28 day human



19h zebrafish





Population Genomics

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

Why not just sequence entire genomes??

- Still prohibitively expensive for many studies
- A full sequence for every sample may not be necessary
- Genetic maps are very useful in genomic studies

Alternative - reduced representation sequencing

- Focus sequencing of numerous samples on a subset of 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 RAD-seq?

(Restriction-site Associated DNA)



2007

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

Michael R. Miller,¹ Joseph P. Dunham,² Angel Amores,³ William A. Cresko,² and Eric A. Johnson^{1,4}

¹Institute for Molecular Biology, University of Oregon, Eugene, Oregon 97403, USA; ²Center for Ecology & Evolutionary Biology, University of Oregon, Eugene, Oregon 97403, USA; ³Institute 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^{1*}, Paul D. Etter^{1*}, Tressa S. Atwood², Mark C. Currey³, Anthony L. Shiver¹, Zachary A. Lewis¹, Eric U. Selker¹, William A. Cresko³, Eric A. Johnson^{1*}

¹Institute of Molecular Biology, University of Oregon, Eugene, Oregon, United States of America, ²Floragenex, Eugene, Oregon, United States of America, ³The Center for Ecology and Evolutionary Biology, University of Oregon, Eugene, Oregon, United States of America

Chr I

22,830 *SbfI* sites in threespine stickleback genome

~ 45,000 RAD-Tags

1

TGCAGG

TCCATACACGCCGTCTGTTCACTGAACAGAA

CCTGCAGG

TTGTGACTAACAGGCAATAAAGTAGTAAACAAAC



TGCAGG

TTCTGTTCACTGAAGCAGACGCGCGTGTATGGA

TCCATACACGCCGTCTGTTCACTGAACAGAA

CCTGCAGG

TTGTGACTAACAGGCAATGAAGTAGTAAACAAAC



TGCAGG

TTGTGACTAACAGGCAATGAAAGTAGTAAACAAAC

TGCAGG

TTCTGTTCACTGAAGCAGACGCGCGTGTATGGA

SbfI

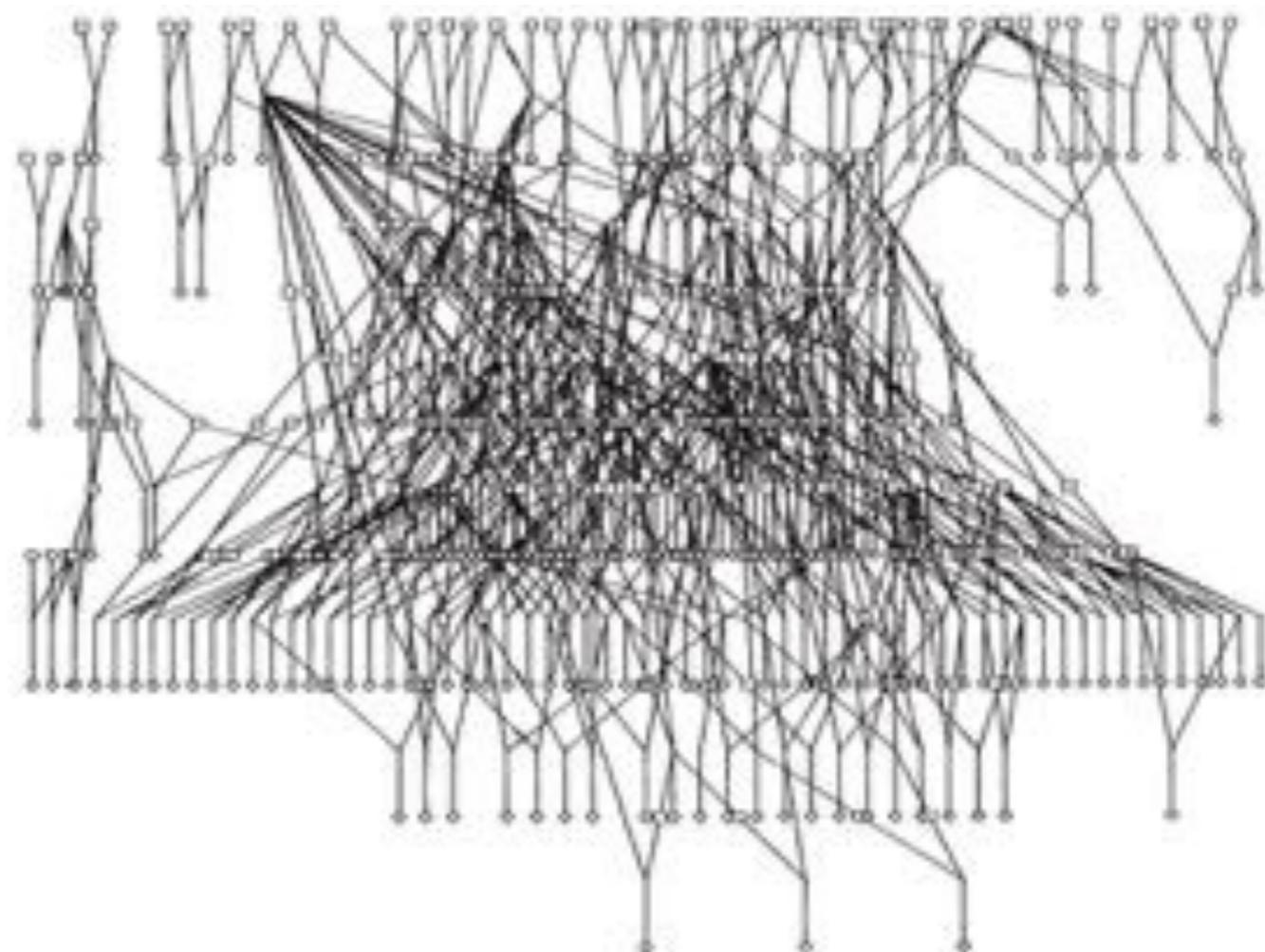
For what types of studies can
RAD-seq be useful?

Identifying genetically distinct individuals and estimating genetic diversity



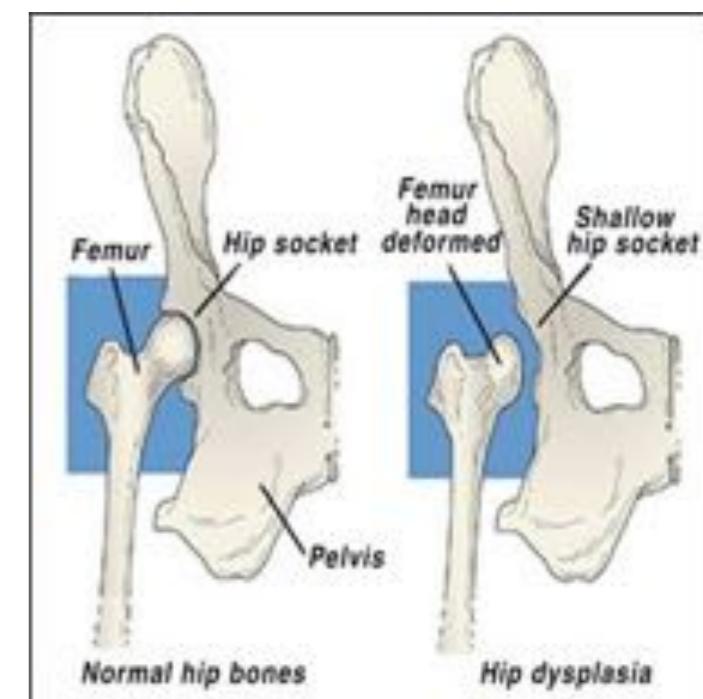
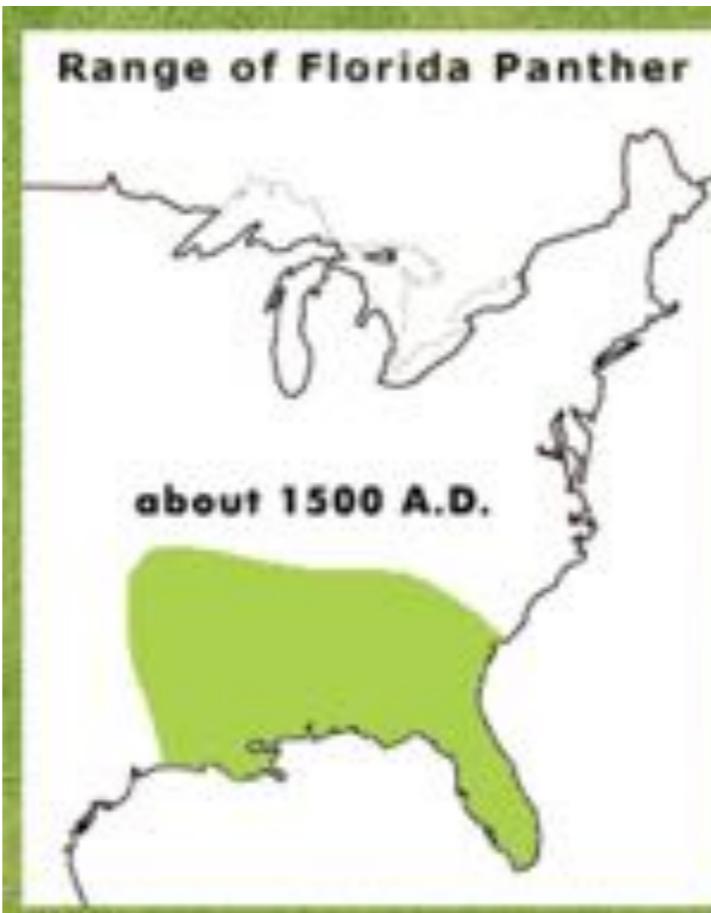
Quaking Aspens

Defining the relationships among individuals and populations

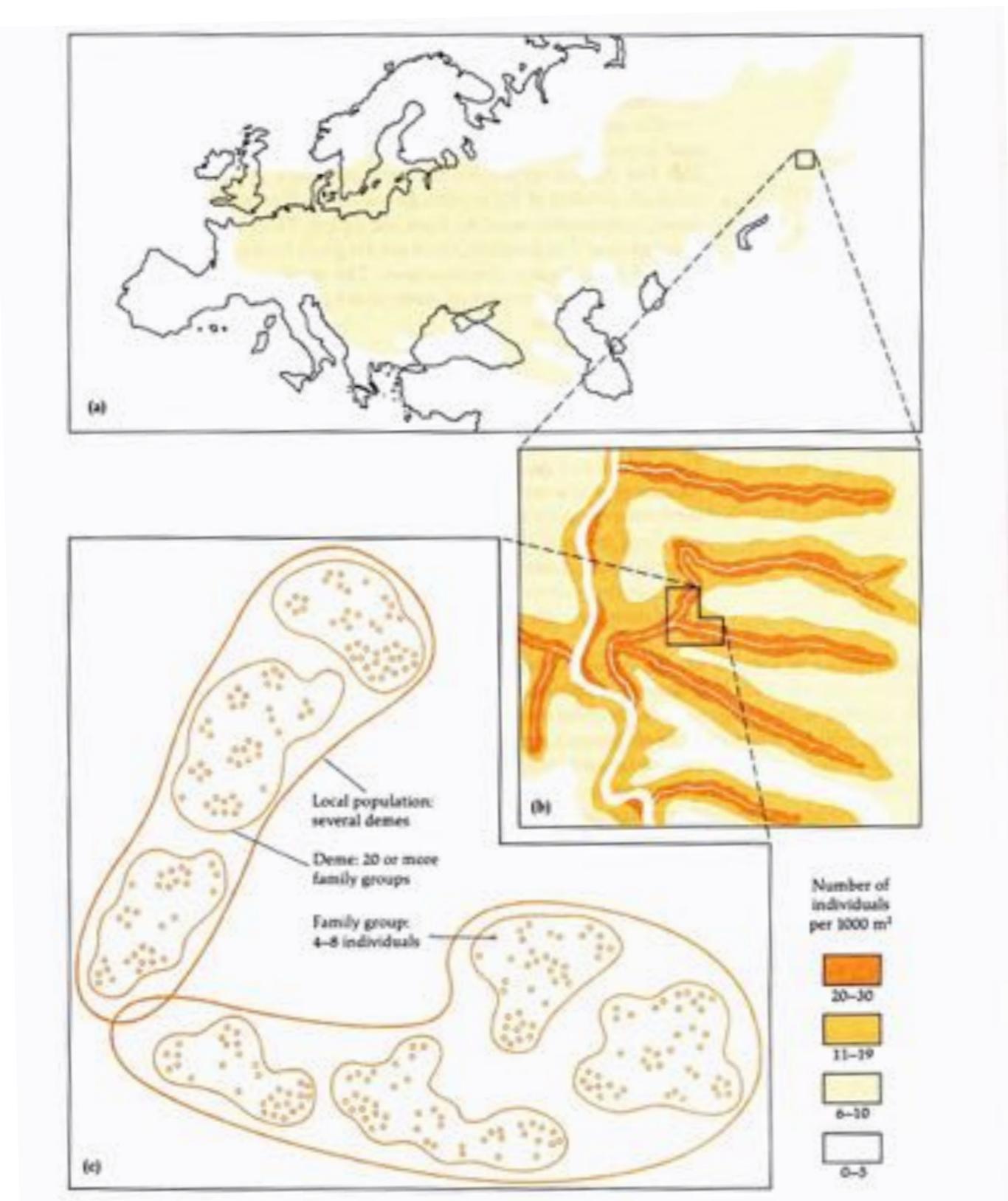


The golden lion tamarin (*Leontopithecus rosalia*).

Precisely quantifying the amount of inbreeding in wild and captive populations



Defining the relationships among individuals and populations



Phylogenetic relationships

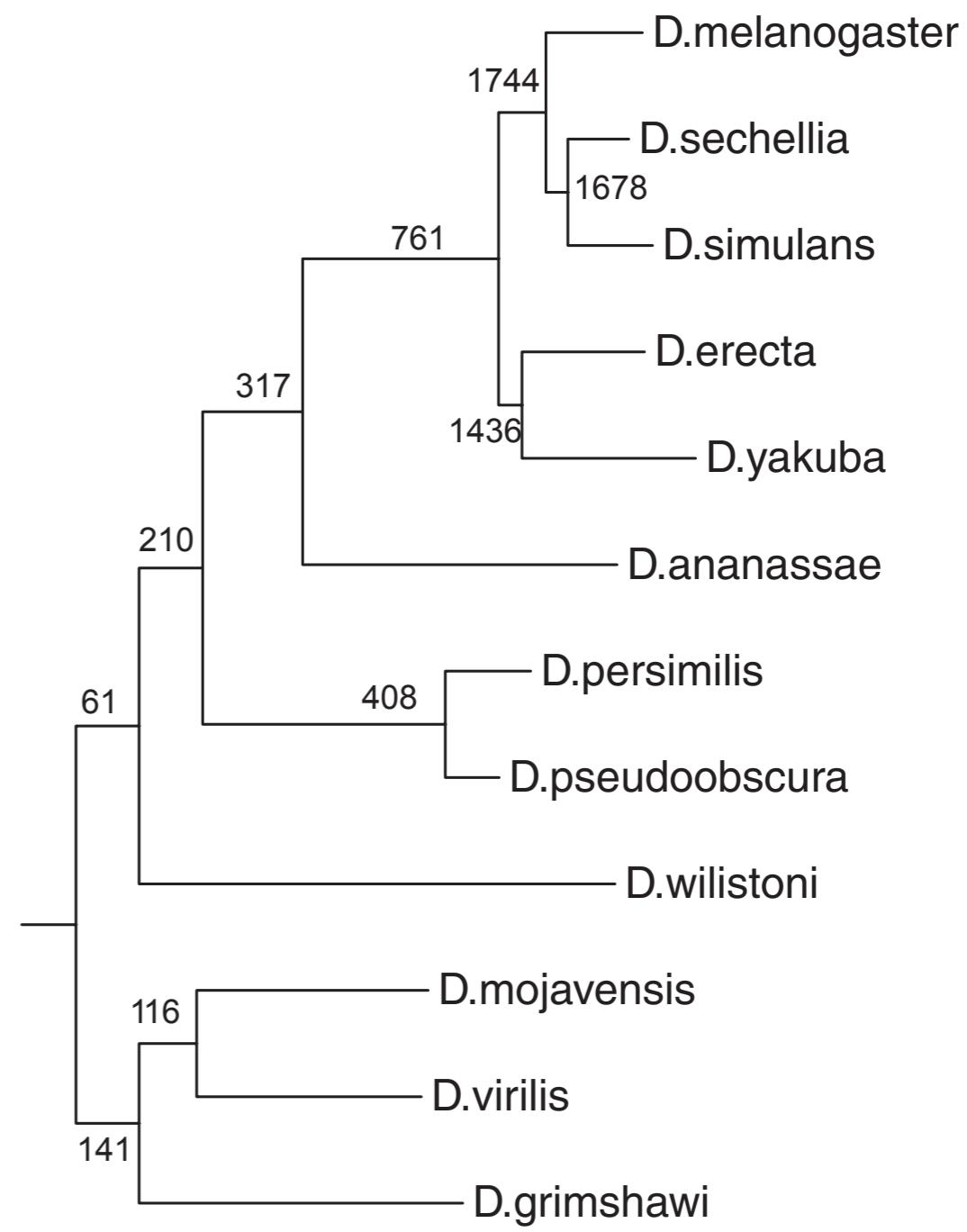
Ecology and Evolution

Open Access

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

Marie Cariou, Laurent Duret & Sylvain Charlat

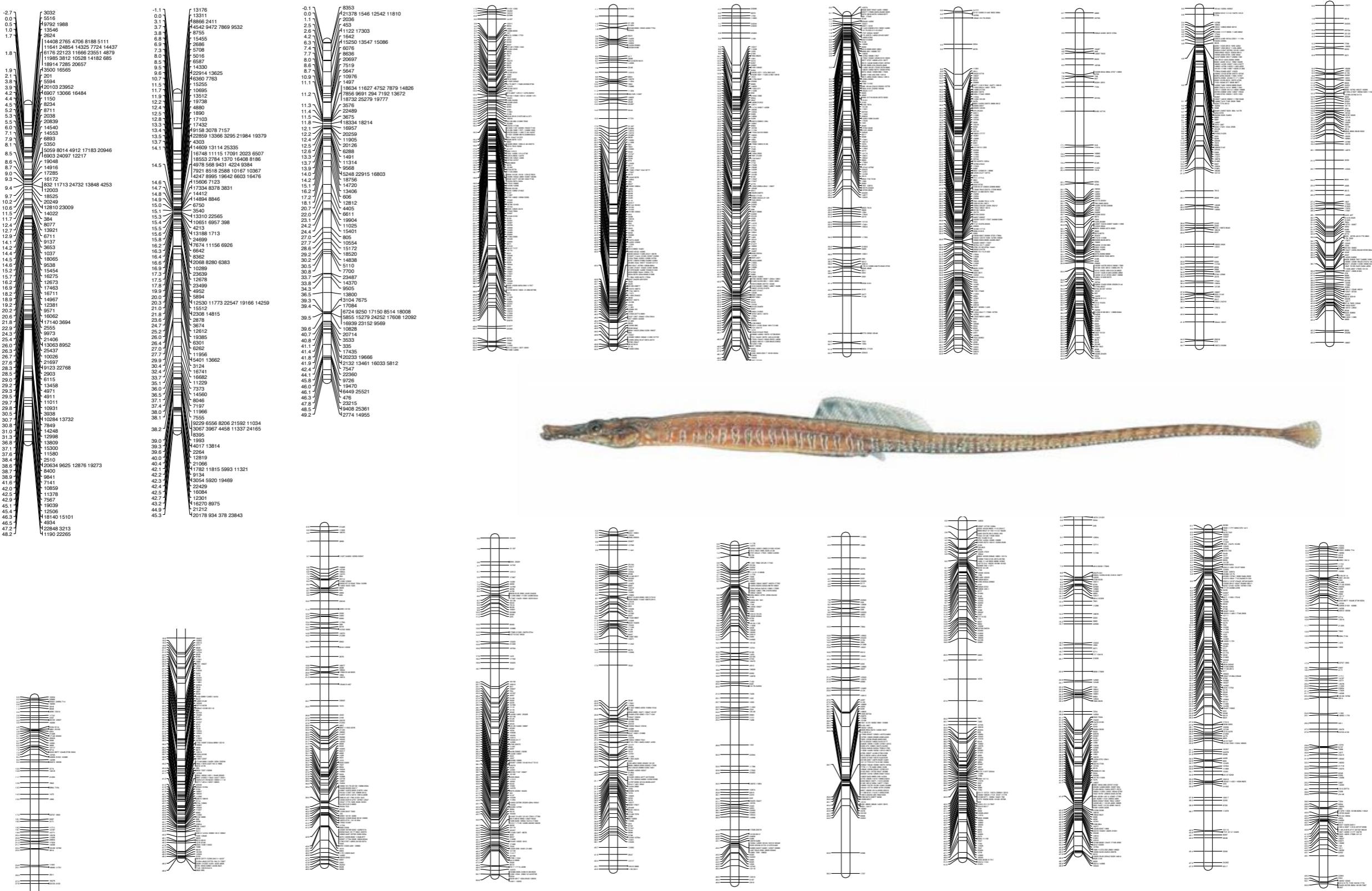
Université de Lyon, Université Lyon 1, CNRS, UMR 5558, Laboratoire de Biométrie et Biologie Evolutive, 43 bou
Villeurbanne F-69622, France



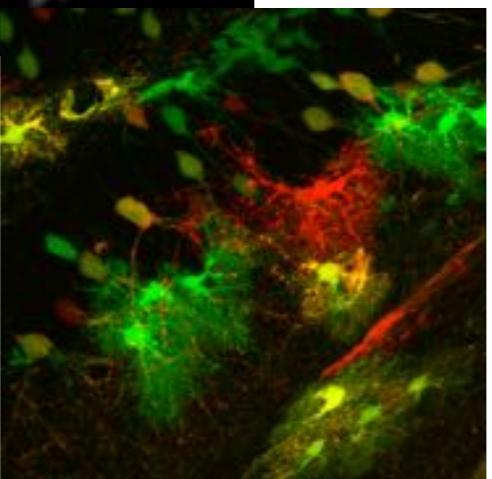
Phylogenetic relationships

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

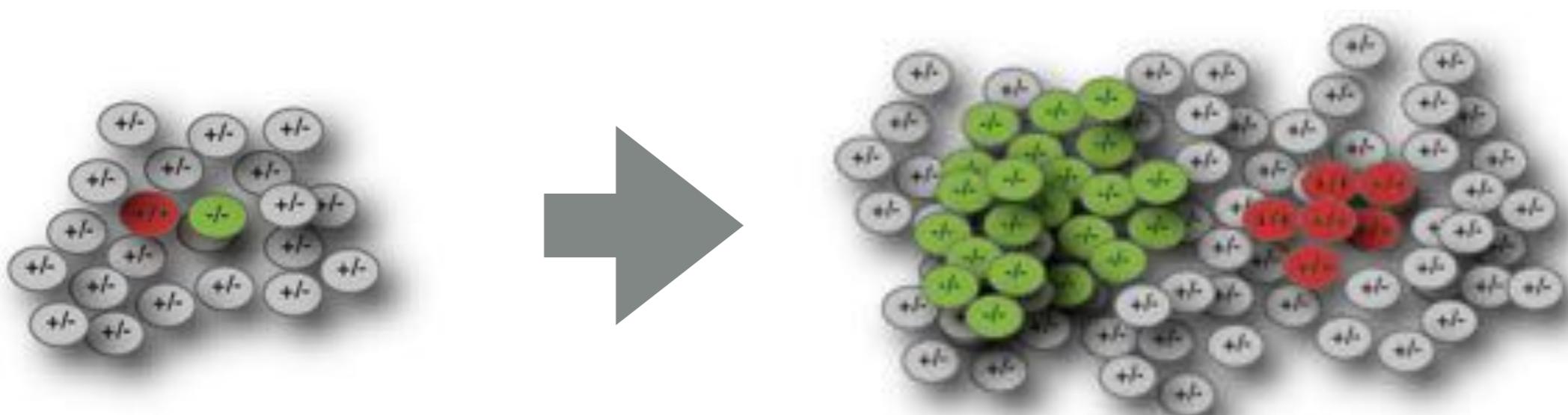
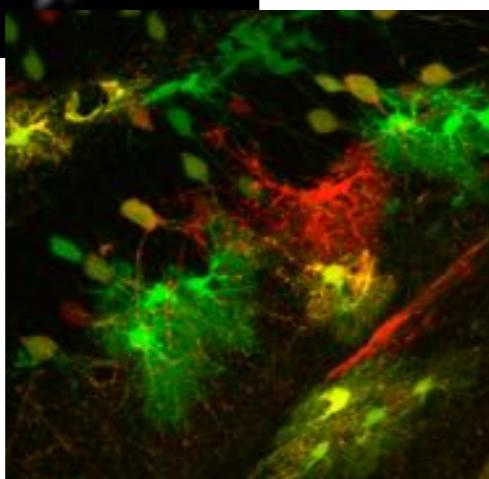
Improve genome assemblies with genetic maps



Cancer genomics



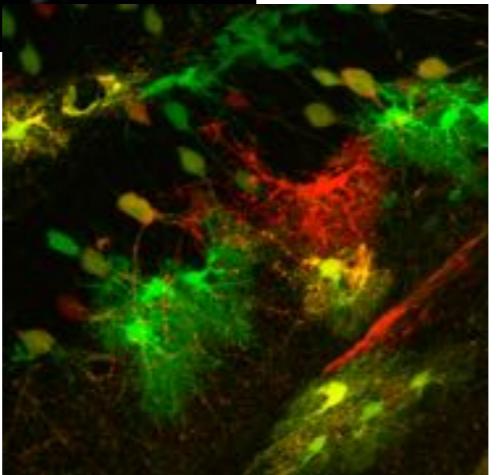
Cancer genomics



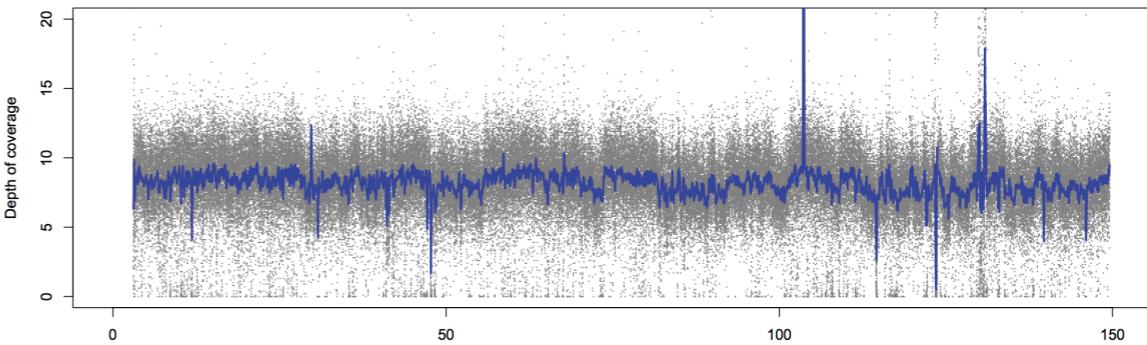
Cancer genomics



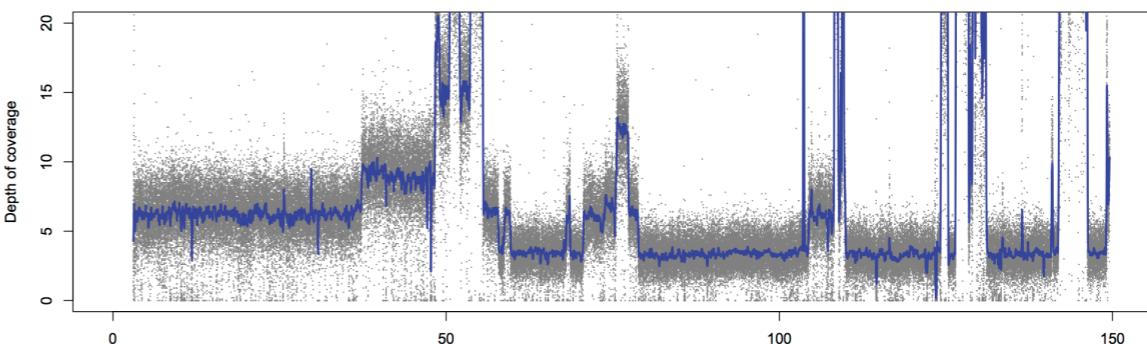
Wild
Type



Mutant

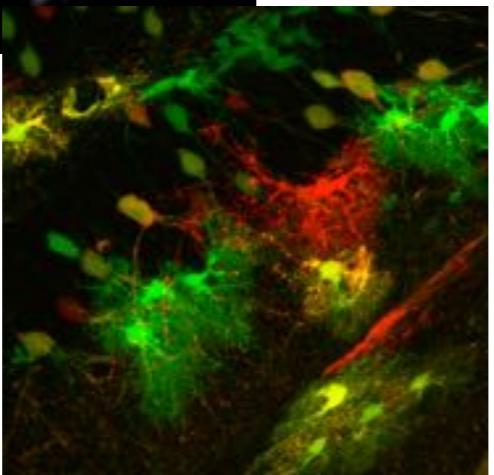


WGS

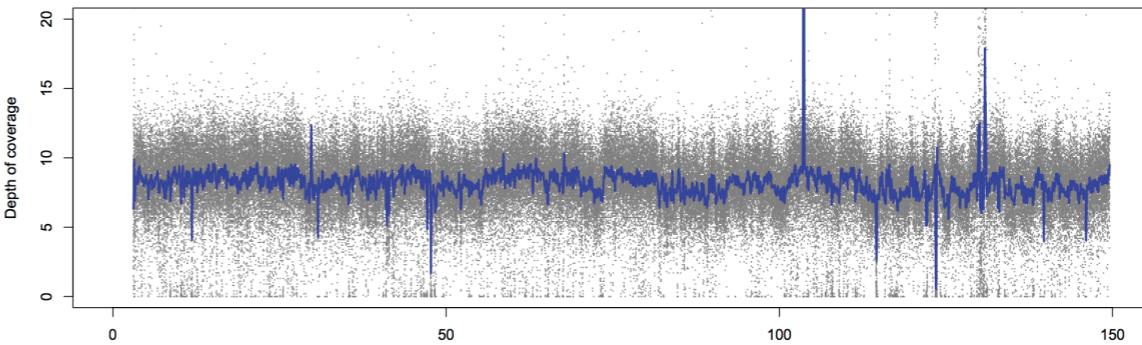


WGS

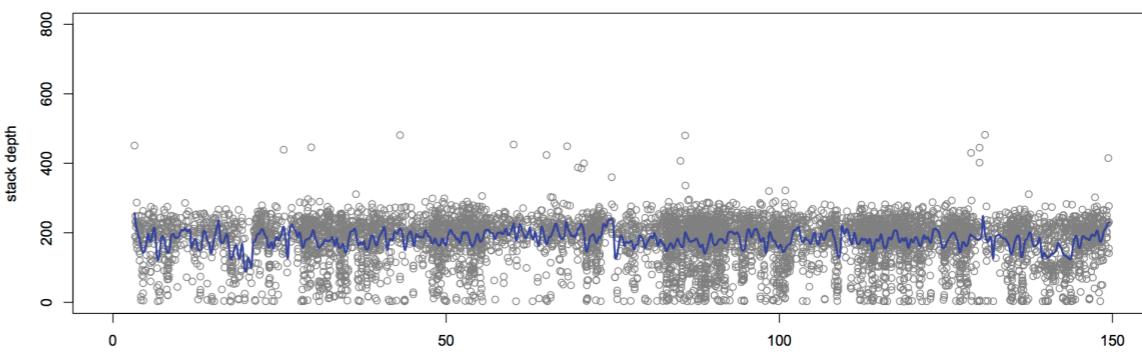
Cancer genomics



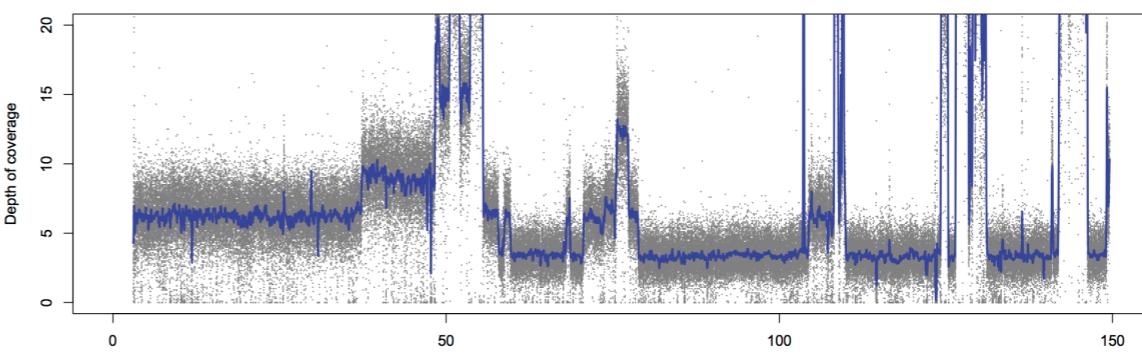
Wild
Type



WGS

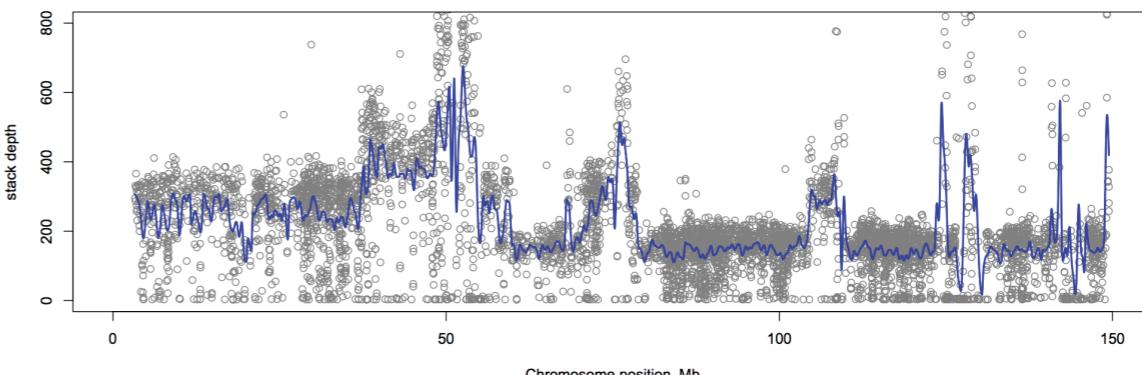


RAD



WGS

Mutant



RAD

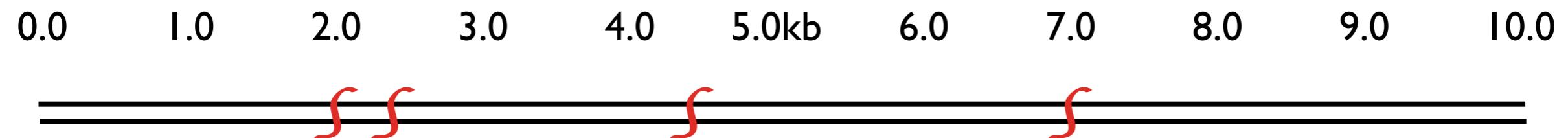
Let's go into a little more detail
on the methods

Restriction enzyme digestion and first adaptor ligation

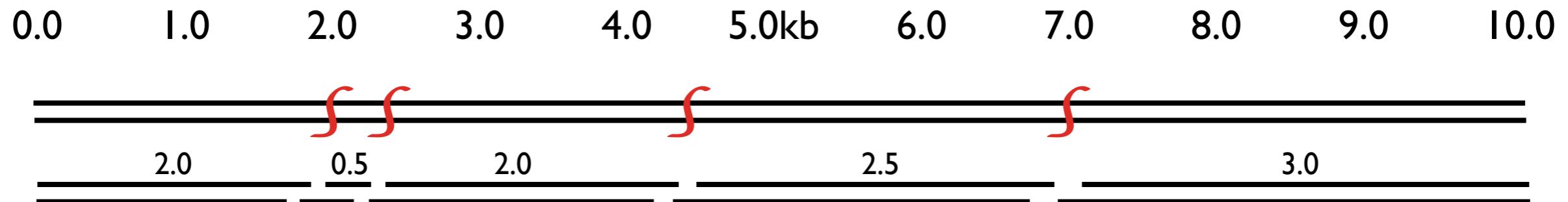
0.0 1.0 2.0 3.0 4.0 5.0kb 6.0 7.0 8.0 9.0 10.0



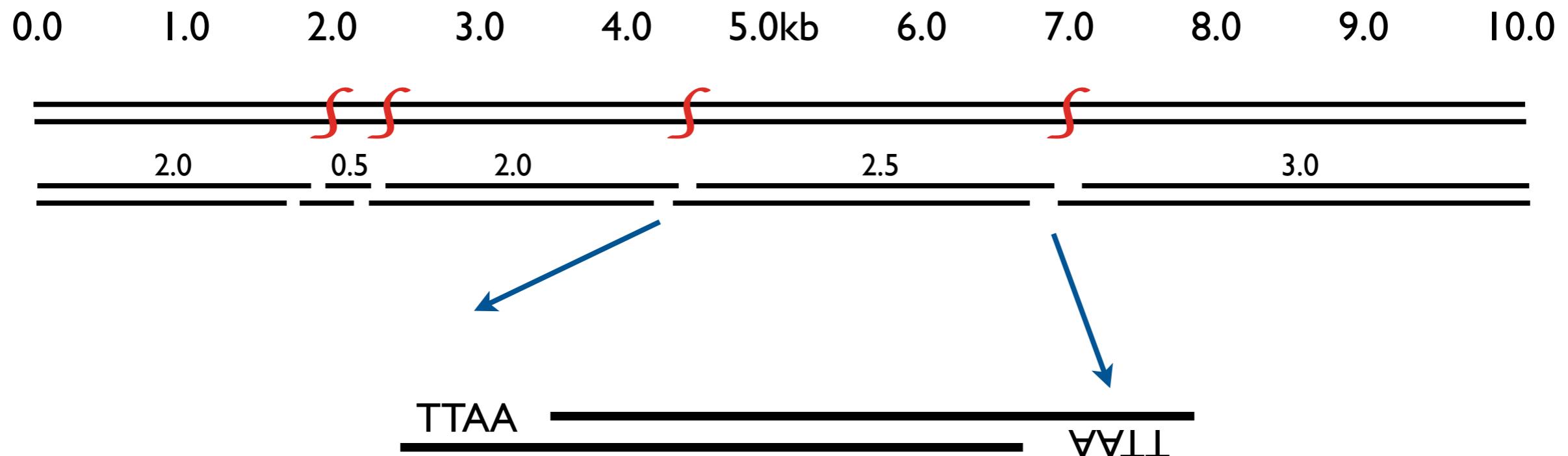
Restriction enzyme digestion and first adaptor ligation



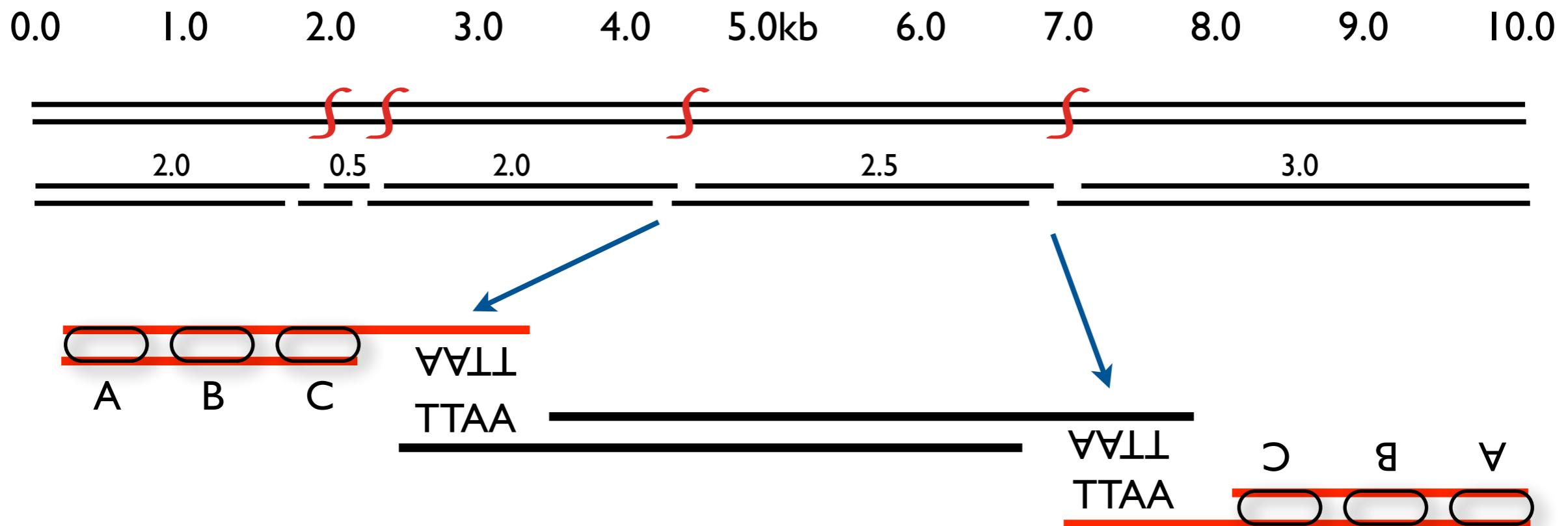
Restriction enzyme digestion and first adaptor ligation



Restriction enzyme digestion and first adaptor ligation

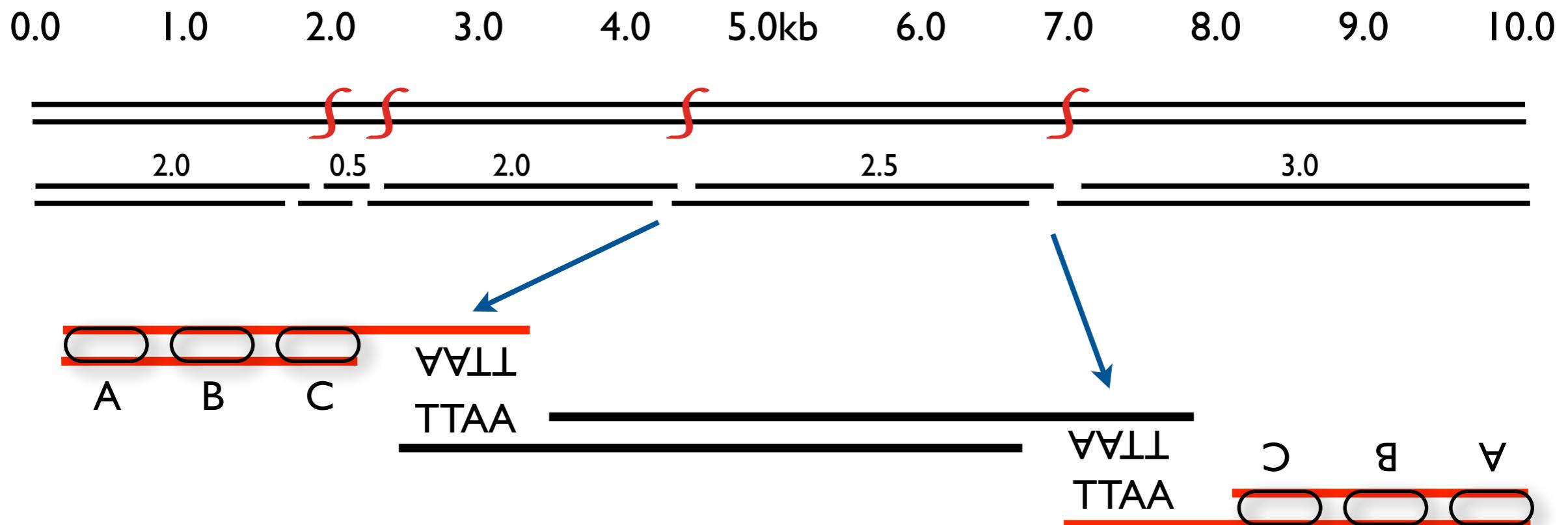


Restriction enzyme digestion and first adaptor ligation



A = Amplification primer
B = Sequencing primer
C = Barcode

Restriction enzyme digestion and first adaptor ligation

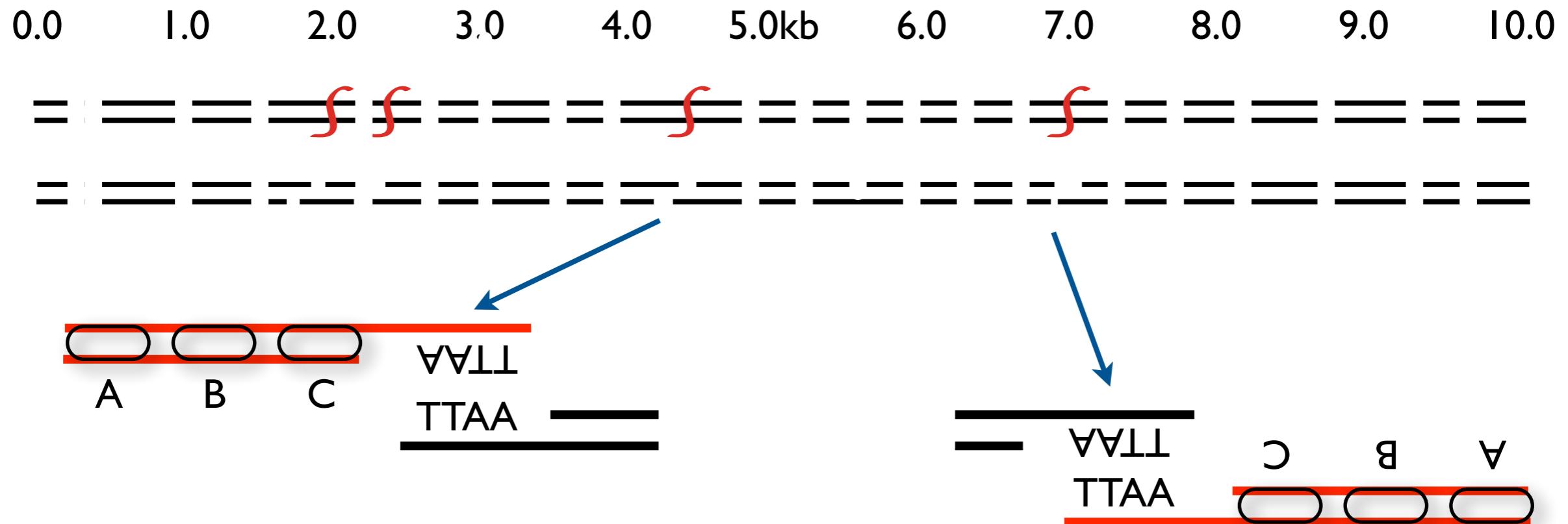


Note - there are now commonly two levels of barcodes used:

Sample Barcodes and
Molecular Identification Barcodes (MIPs)

A = Amplification primer
B = Sequencing primer
C = Barcode

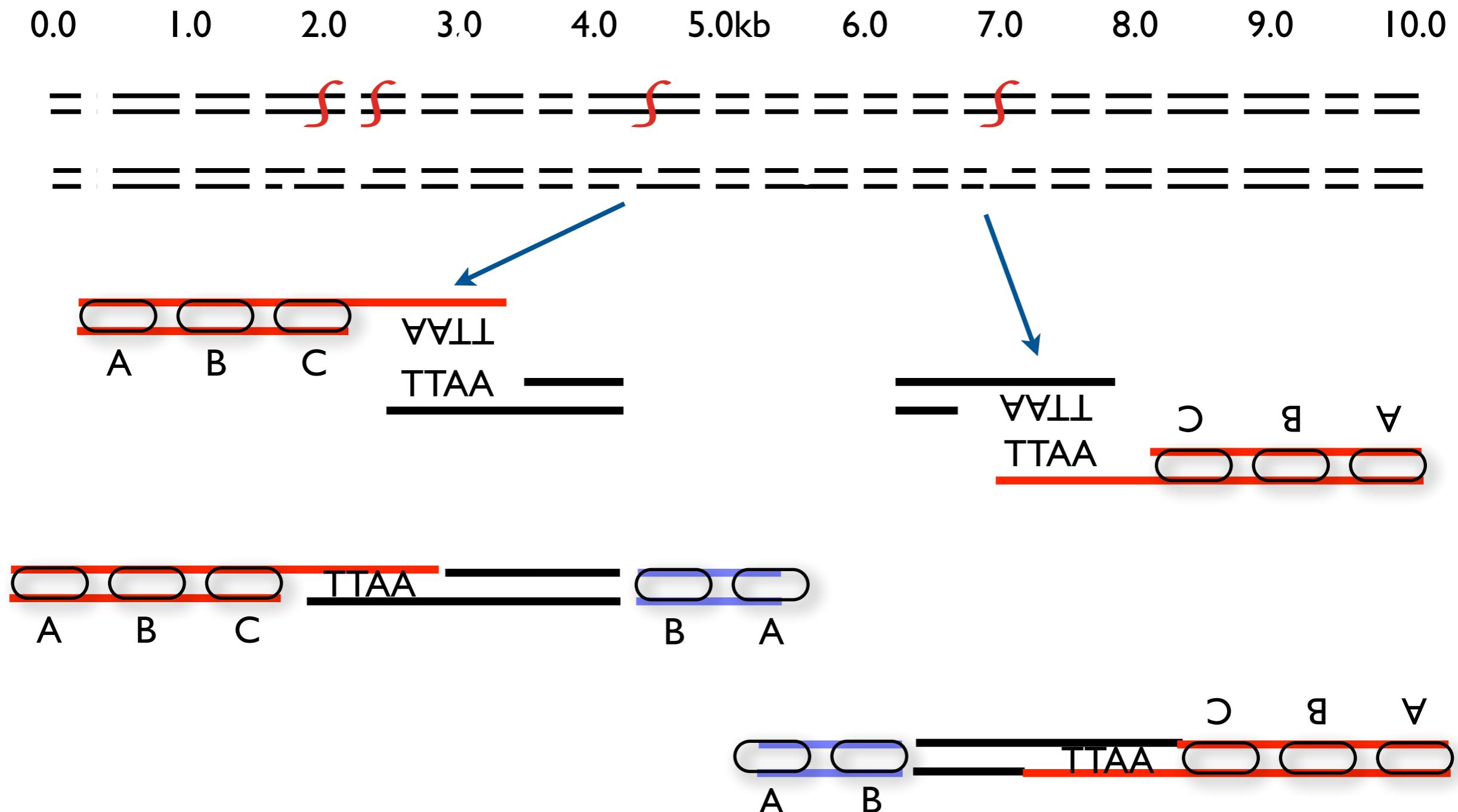
Shearing of DNA and second adaptor ligation



* Defining difference between original RAD and other approaches*

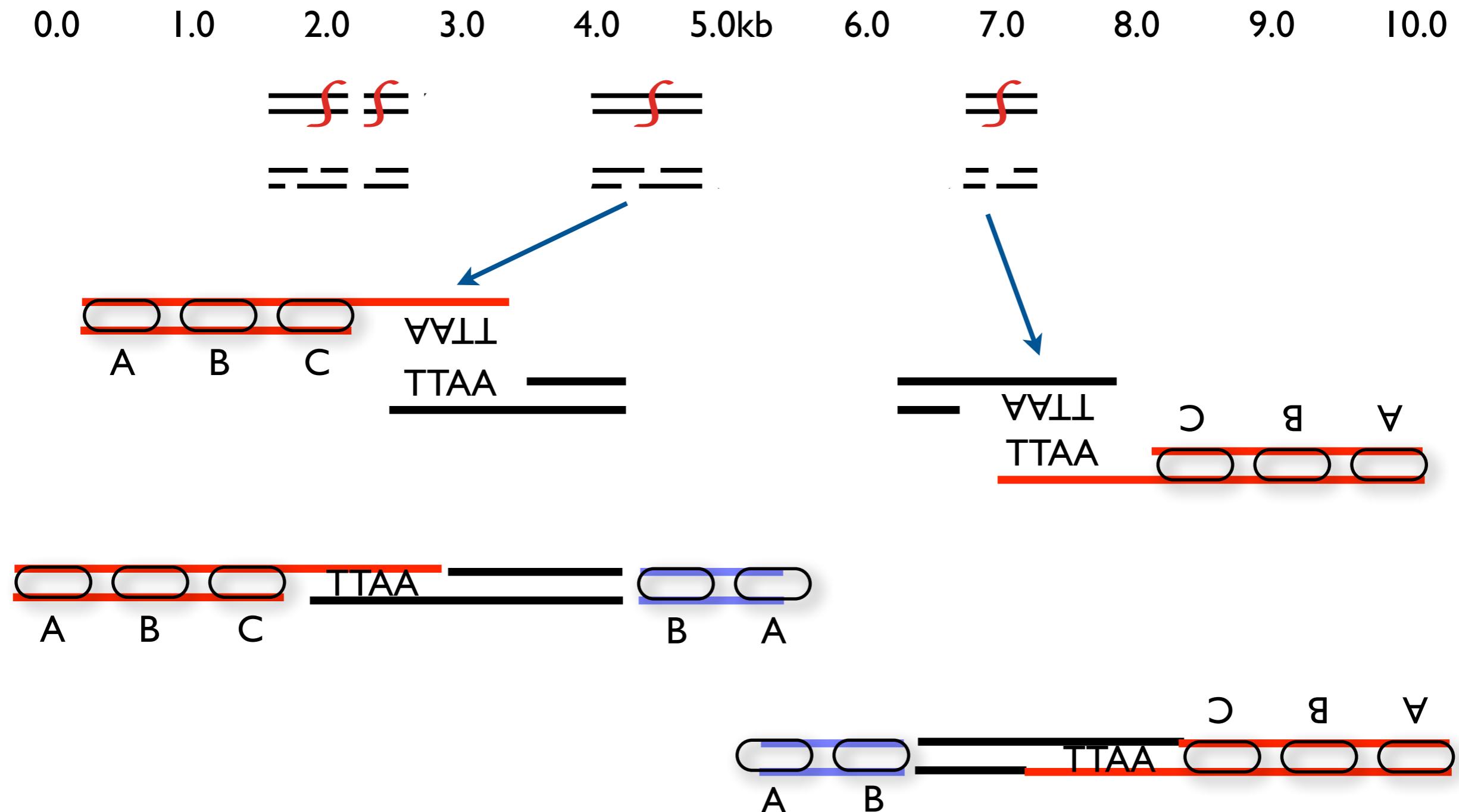
- A = Amplification primer
- B = Sequencing primer
- C = Barcode

Shearing of DNA and second adaptor ligation



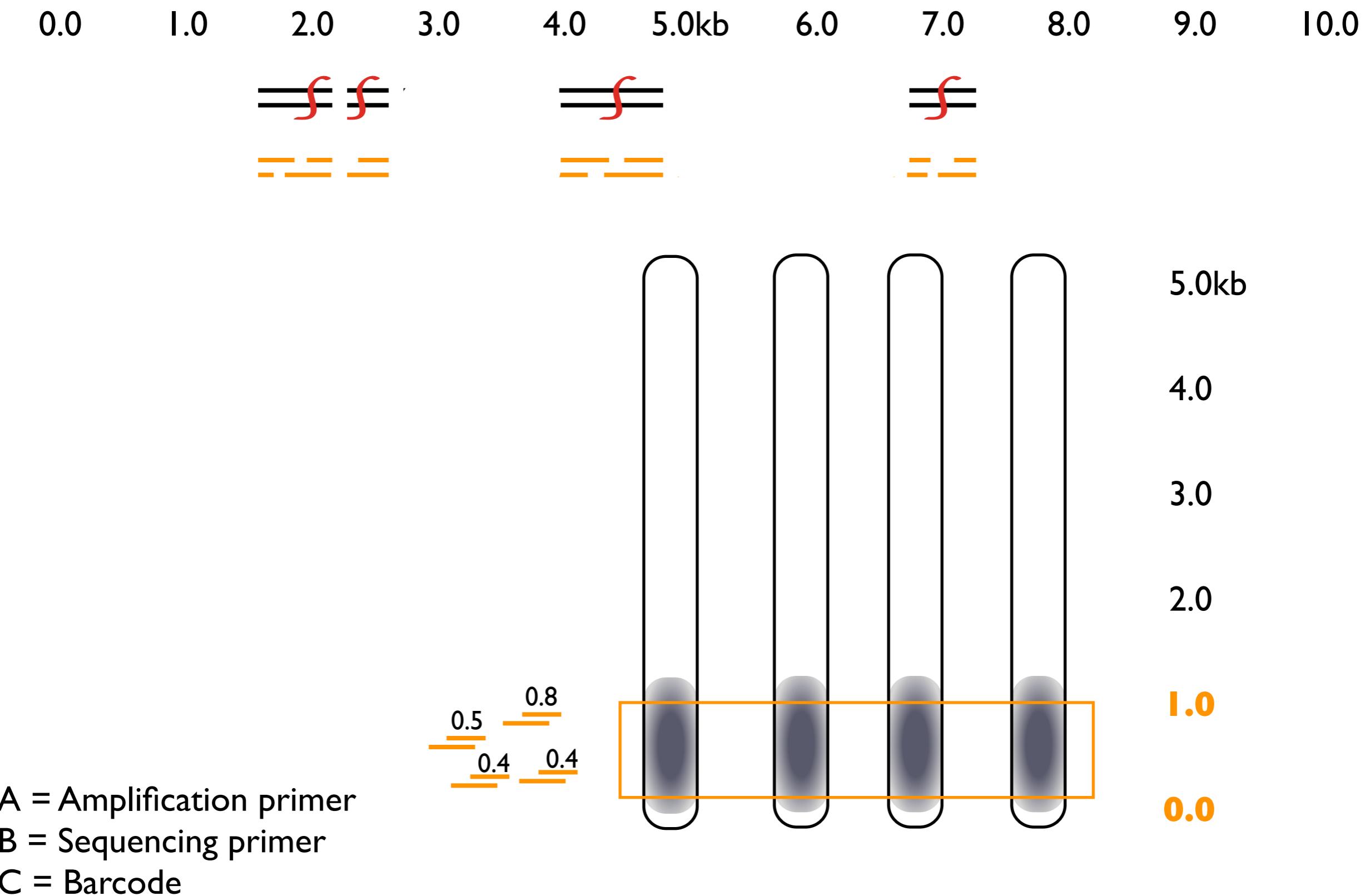
A = Amplification primer
B = Sequencing primer
C = Barcode

Shearing of DNA and second adaptor ligation

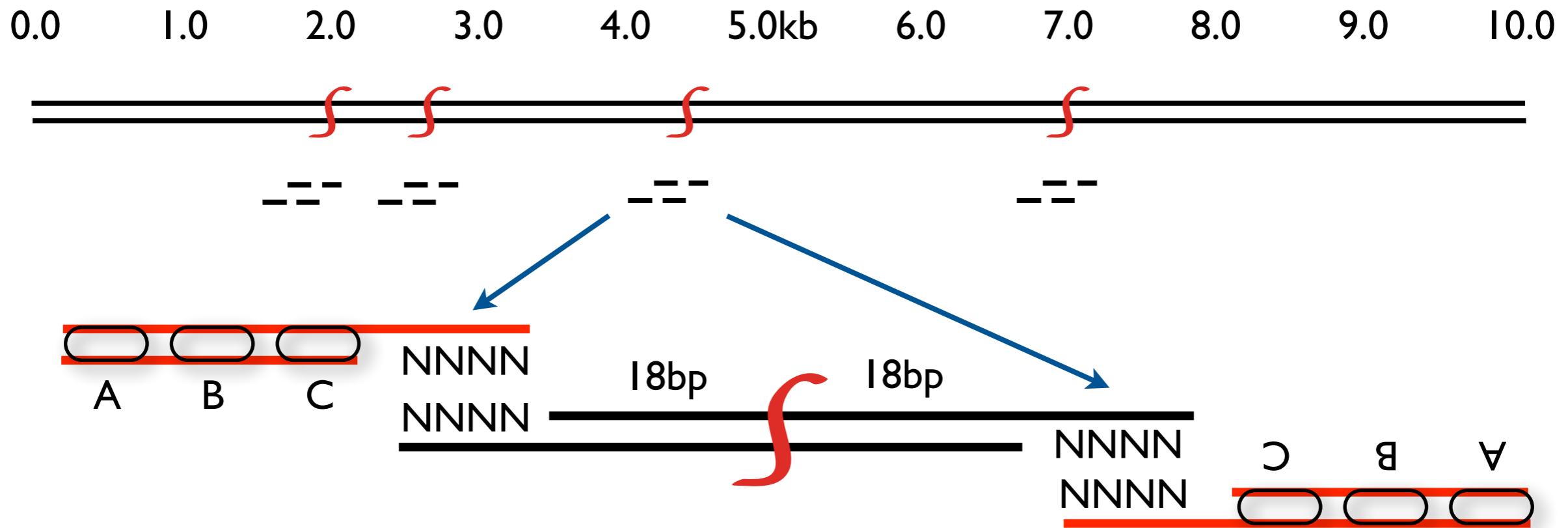


A = Amplification primer
B = Sequencing primer
C = Barcode

Shearing of DNA and second adaptor ligation

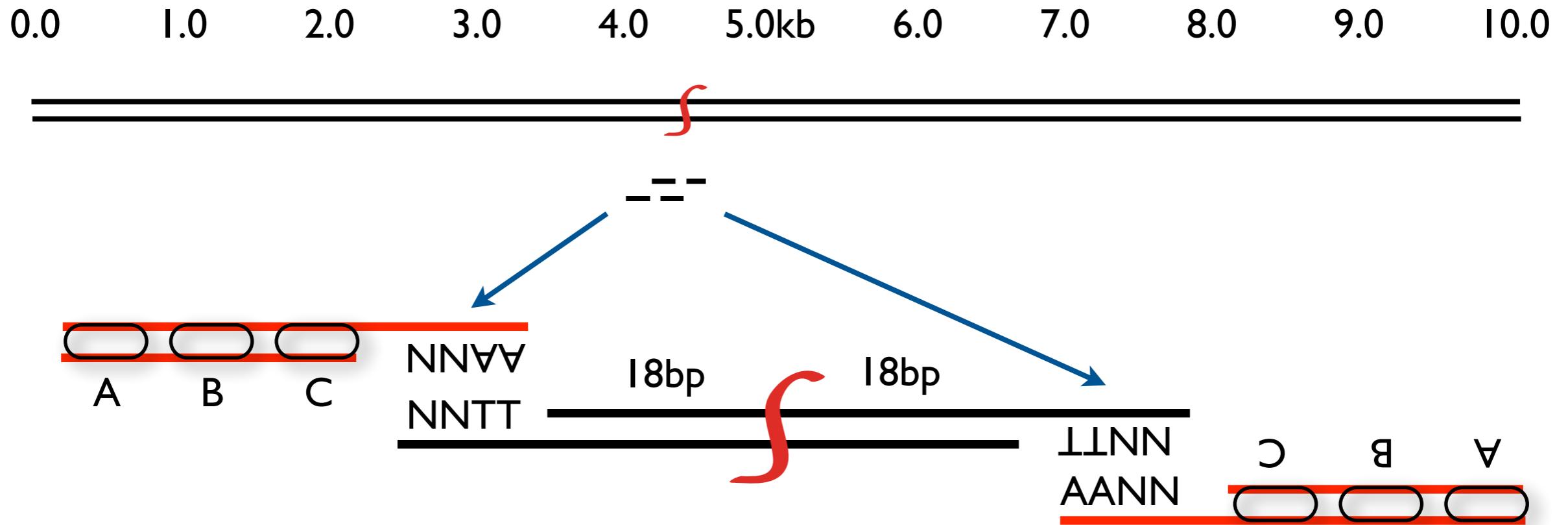


2bRAD - type 2b restriction enzyme



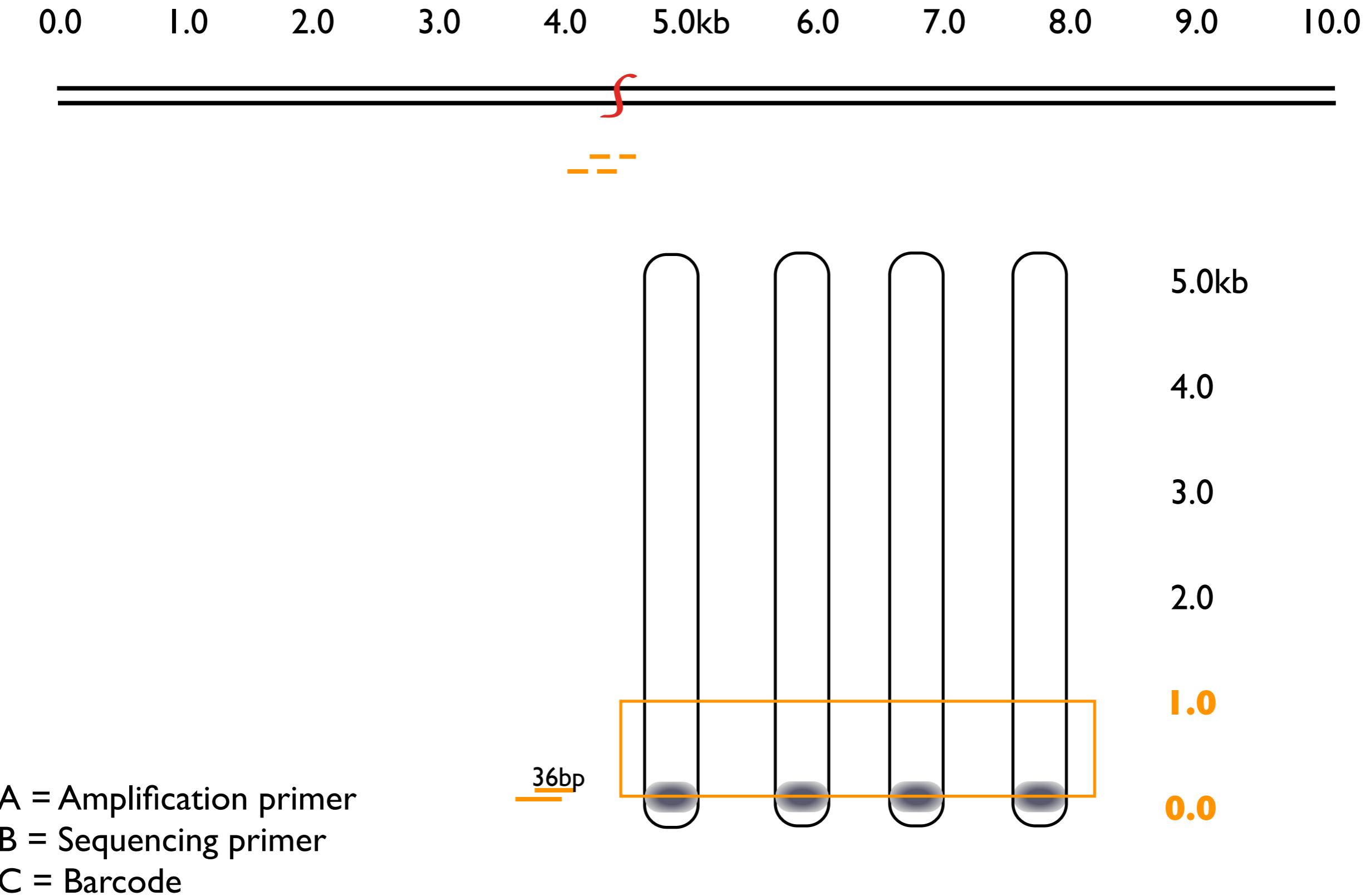
A = Amplification primer
B = Sequencing primer
C = Barcode

2bRAD - can scale number of markers easily

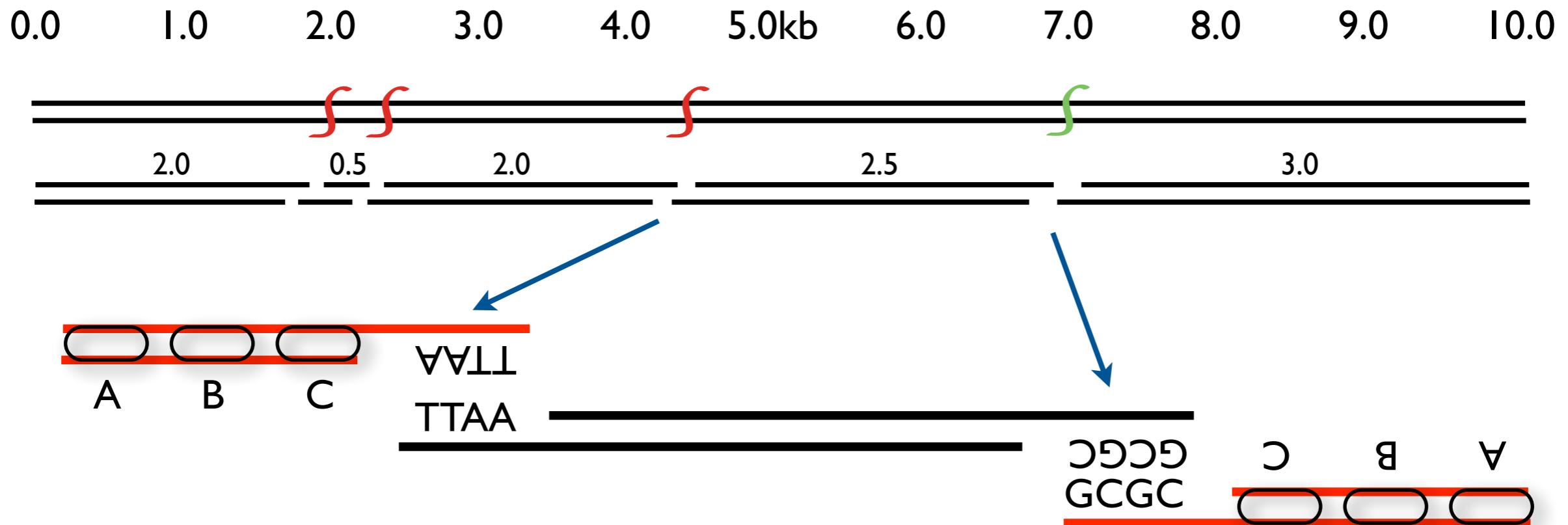


A = Amplification primer
B = Sequencing primer
C = Barcode

2bRAD - size selection is more difficult and reads are shorter

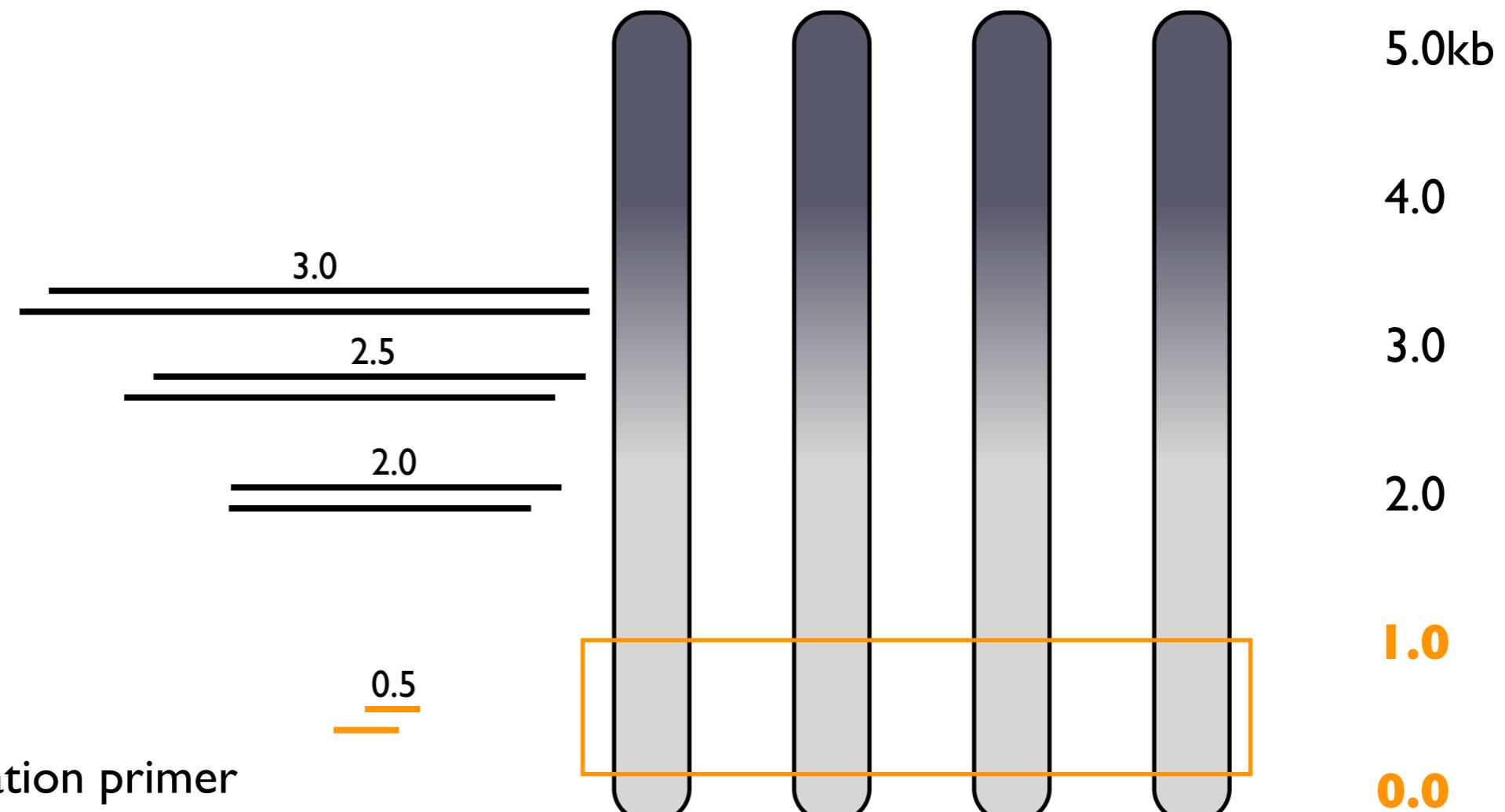
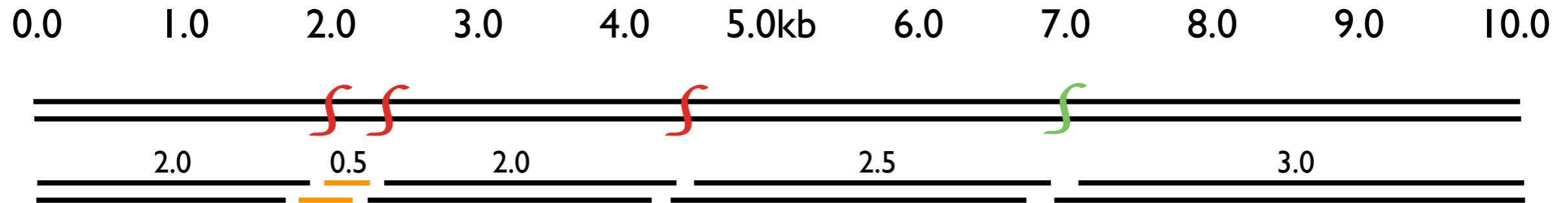


Single (GBS) or Double Digest RAD (ddRAD)



A = Amplification primer
B = Sequencing primer
C = Barcode

Size selection is more problematic without shearing



A = Amplification primer
B = Sequencing primer
C = Barcode

Summary of plusses and minuses of RAD family

Harnessing the power of RADseq for ecological and evolutionary genomics

Kimberly R. Andrews¹, Jeffrey M. Good², Michael R. Miller³, Gordon Luikart⁴
and Paul A. Hohenlohe⁵

	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	≤1 kb if building contigs; otherwise ≤300 bp [†]	33–36 bp	<300 bp [†]	≤300 bp [†]	≤300 bp [†]
Cost per barcoded or indexed sample	Low	Low	Low	Low	High
Effort per barcoded or indexed sample [‡]	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 [§]	Pippin Prep [§]
Suitability for large or complex genomes [¶]	Good	Poor	Moderate	Good	Good
Suitability for de novo locus identification (no reference genome) [¶]	Good	Poor	Moderate	Moderate	Moderate
Available from commercial companies	Yes	No	Yes	Yes	No

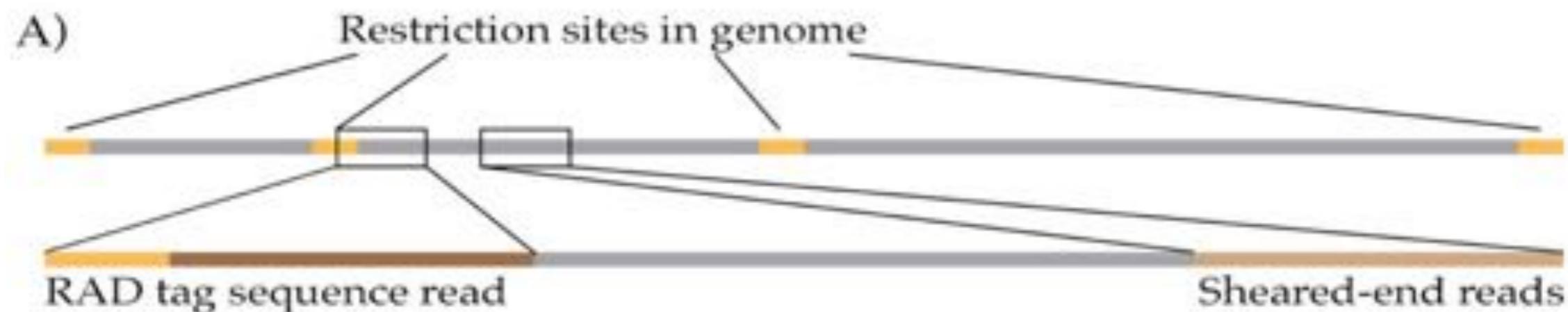
Summary of plusses and minuses of RAD family

Harnessing the power of RADseq for ecological and evolutionary genomics

Kimberly R. Andrews¹, Jeffrey M. Good², Michael R. Miller³, Gordon Luikart⁴
and Paul A. Hohenlohe⁵

	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 ^b	33–36 bp	<300 bp ^b	≤300 bp ^b	≤300 bp ^b
Cost per barcoded or indexed sample	Low	Low	Low	Low	High
Effort per barcoded or indexed sample ^b	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 ^d	Pippin Prep ^d
Suitability for large or complex genomes ^b	Good	Poor	Moderate	Good	Good
Suitability for de novo locus identification (no reference genome) ^b	Good	Poor	Moderate	Moderate	Moderate
Available from commercial companies	Yes	No	Yes	Yes	No

Local paired end assemblies to create haplotypes



OPEN ACCESS Freely available online

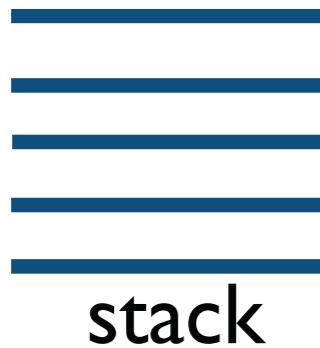
PLOS one

Local *De Novo* Assembly of RAD Paired-End Contigs Using Short Sequencing Reads

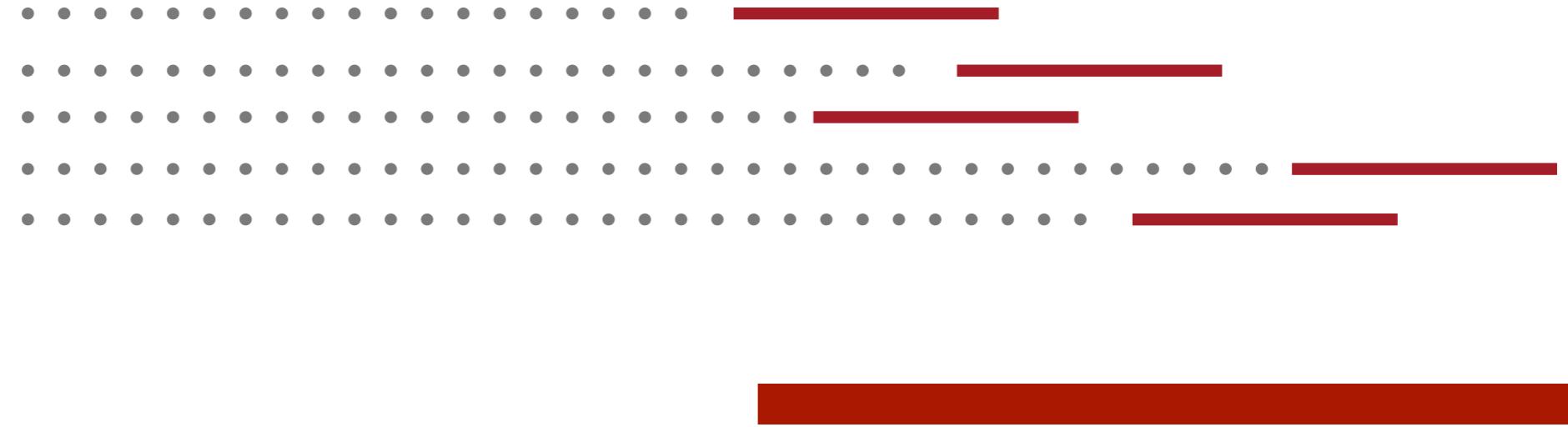
Paul D. Etter¹, Jessica L. Preston¹, Susan Bassham², William A. Cresko², Eric A. Johnson^{1*}

Local paired end assemblies to create haplotypes
One RAD site within an individual

Read 1
'Anchored' to
Restriction Site



Read 2
Randomly Sheared



Haplotype
200 - 800bp in length

‘Bias’ in RADseq

MOLECULAR ECOLOGY

Molecular Ecology (2013) 22, 3179–3190

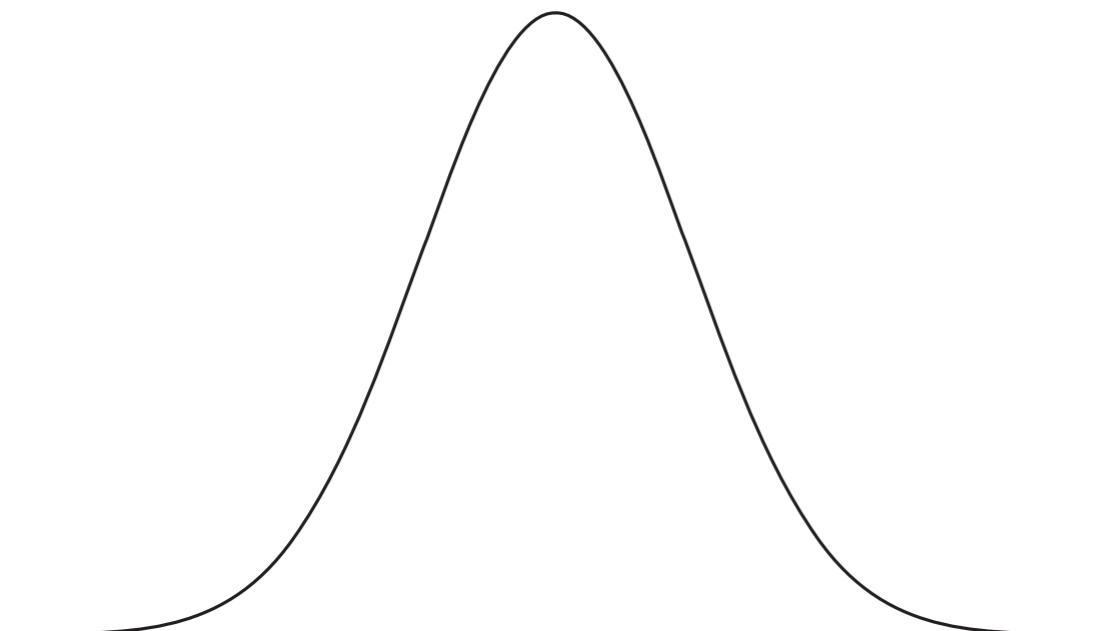
doi: 10.1111/mec.12276

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

B. ARNOLD,¹ R. B. CORBETT-DETIG,¹ D. HARTL and K. BOMBLIES

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

‘Bias’ in RADseq

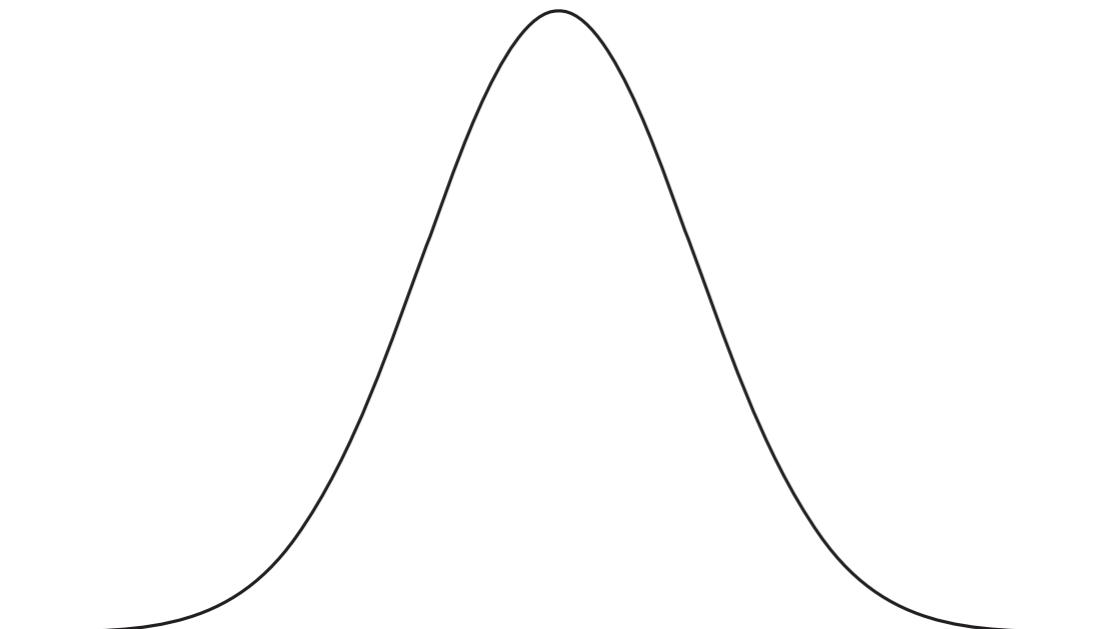


$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

$e = 2.7182\dots$

$\pi = 3.1415\dots$

‘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	θ per bp	Mean	
		θ_{we}/θ_{wa}	π_e/π_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

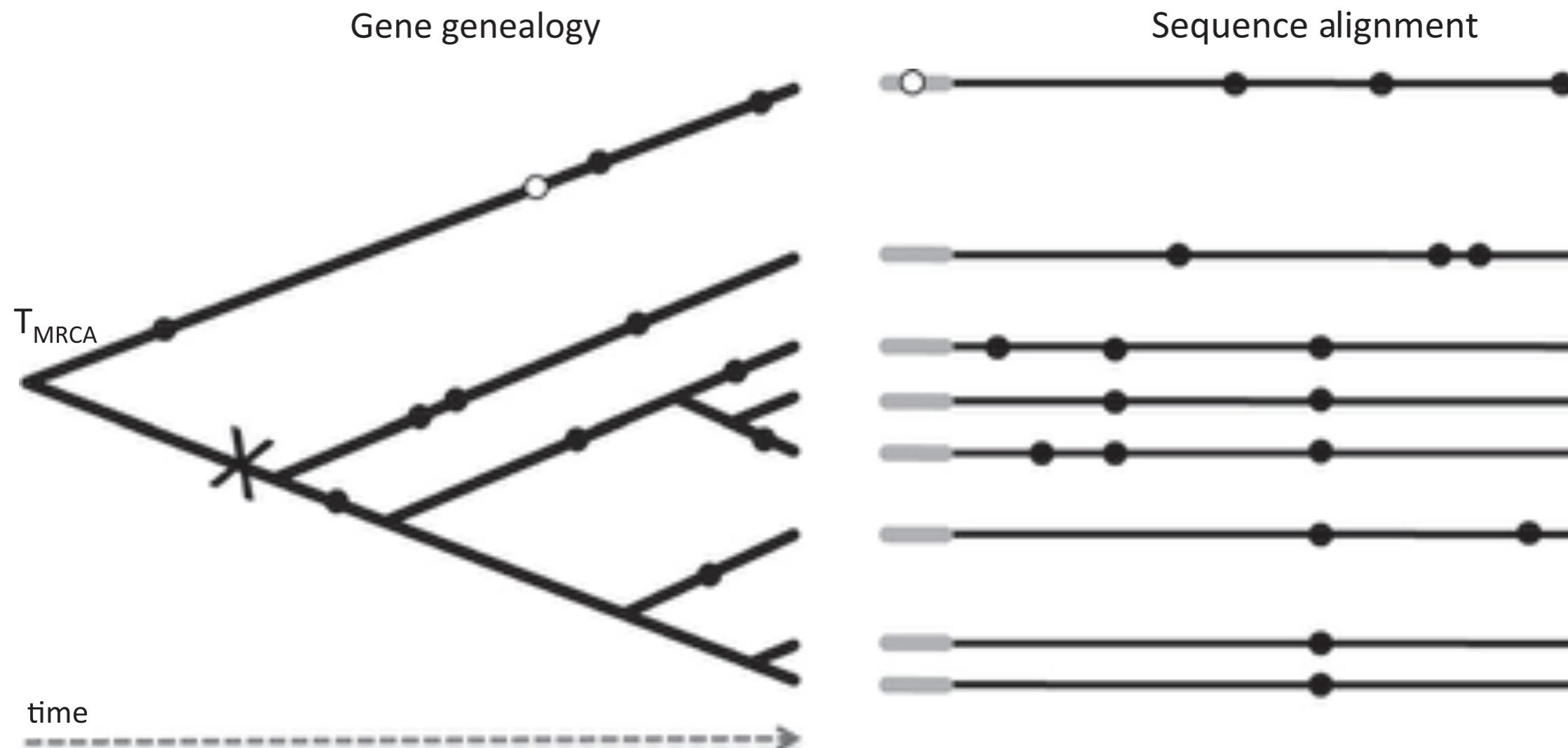
‘Bias’ in RADseq is increased in some RAD protocols

Protocol	θ per bp	Mean	
		θ_{we}/θ_{wa}	π_e/π_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	θ per bp	Mean	
		θ_{we}/θ_{wa}	π_e/π_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



Biological studies that benefit from 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 RAD-seq

- Defining individuality, parentage and pedigrees
 - Performing genome wide association studies
 - Fine scale evolution and adaptation
 - Identifying transmissible elements
 - Making maps of genetic variation
 - Genome Wide Population Studies
 - Building genomic resources for model organisms
 - Estimating species and higher level phylogenetic relationships
 - Population genomics - identifying the signatures of natural selection
- Any study where increasing biological sample size would be beneficial

Outline for today's lecture

RAD-seq for ecological and evolutionary genomics

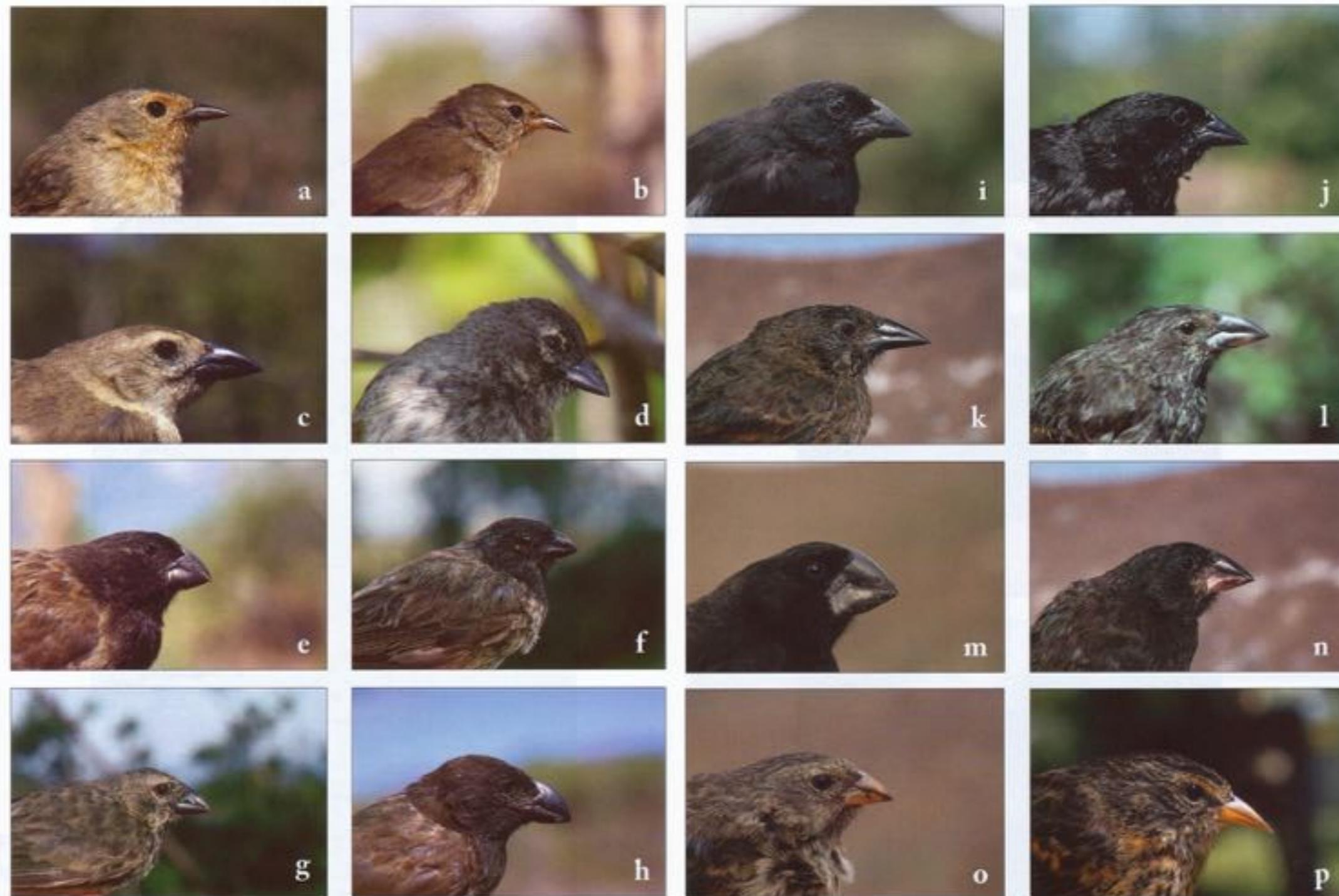
Primer on Population Genomics

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

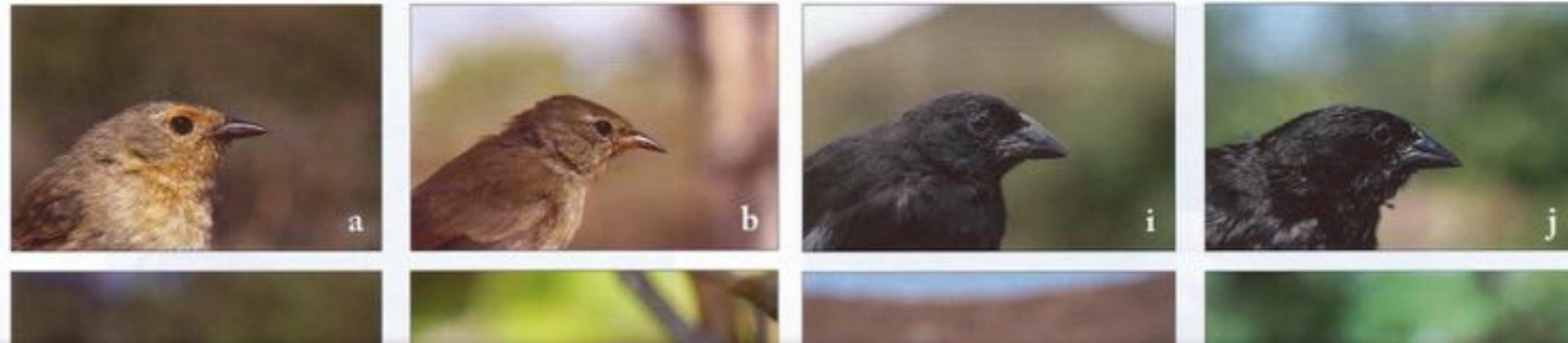
Genomically enabling the Gulf pipefish

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?

What genes are linked to evolving phenotypes?



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

Four fundamental processes in evolution

Origin of genetic variation

mutation

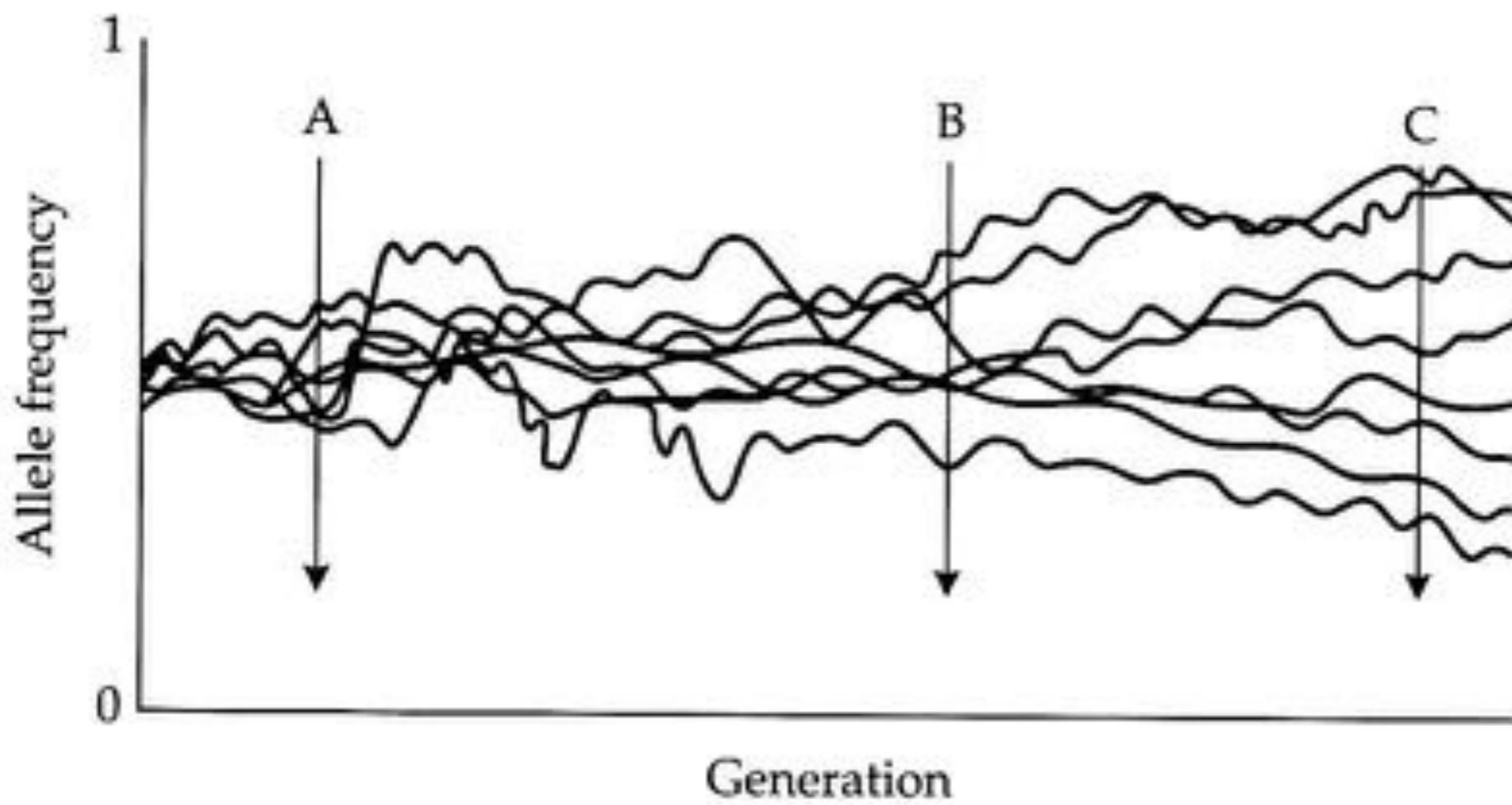
migration

Sorting of variation

genetic drift

natural selection

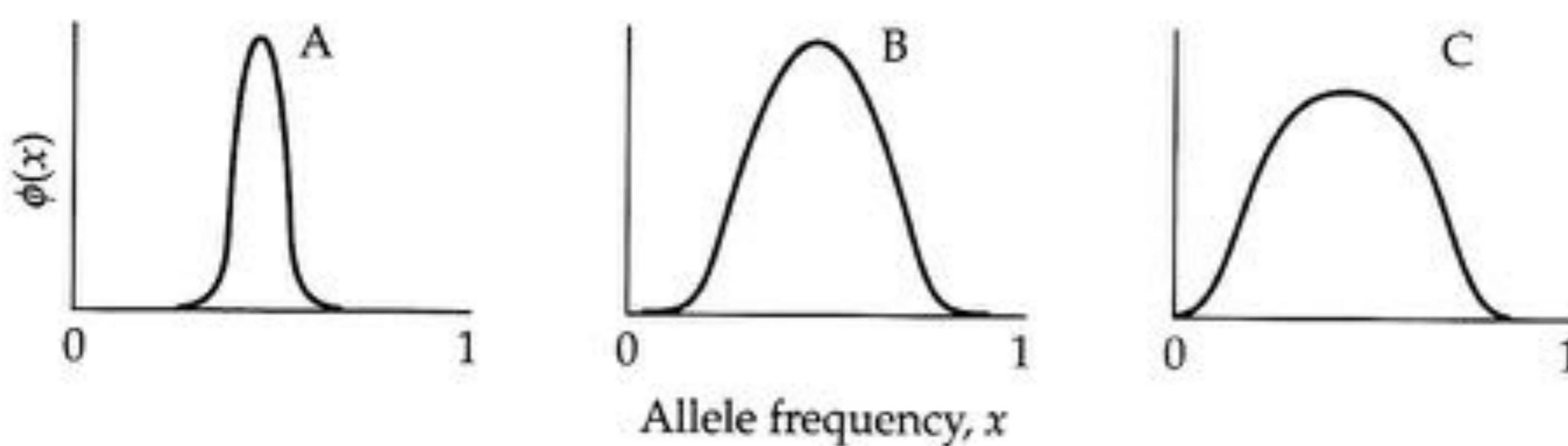
Genetic drift is a null model



R.A. Fisher

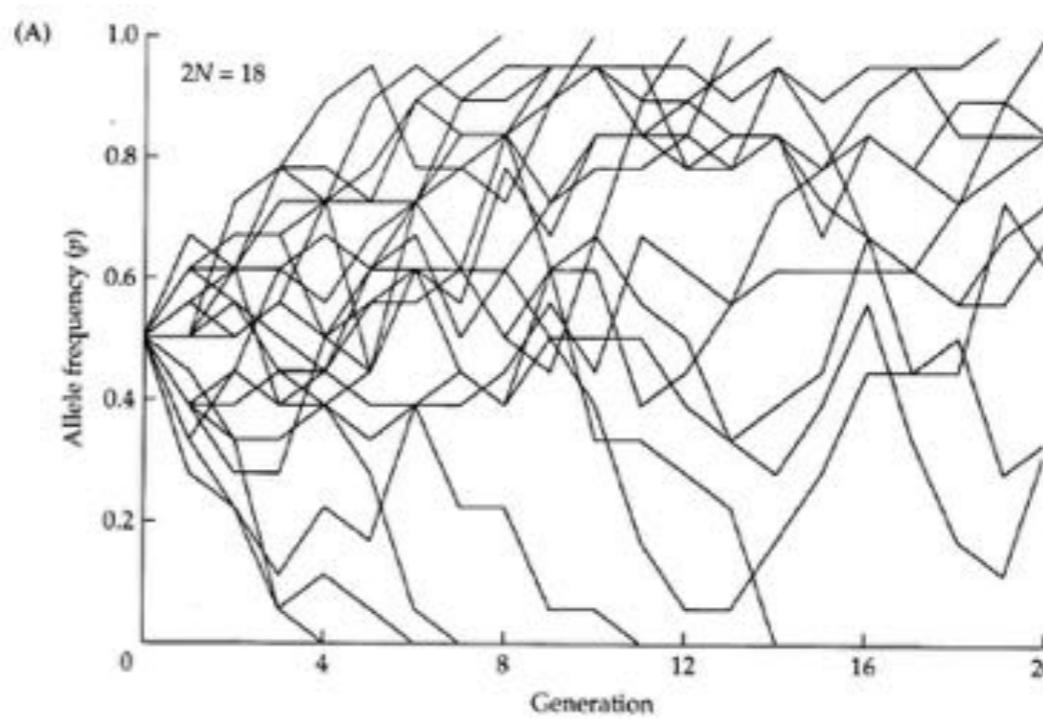


Sewall Wright

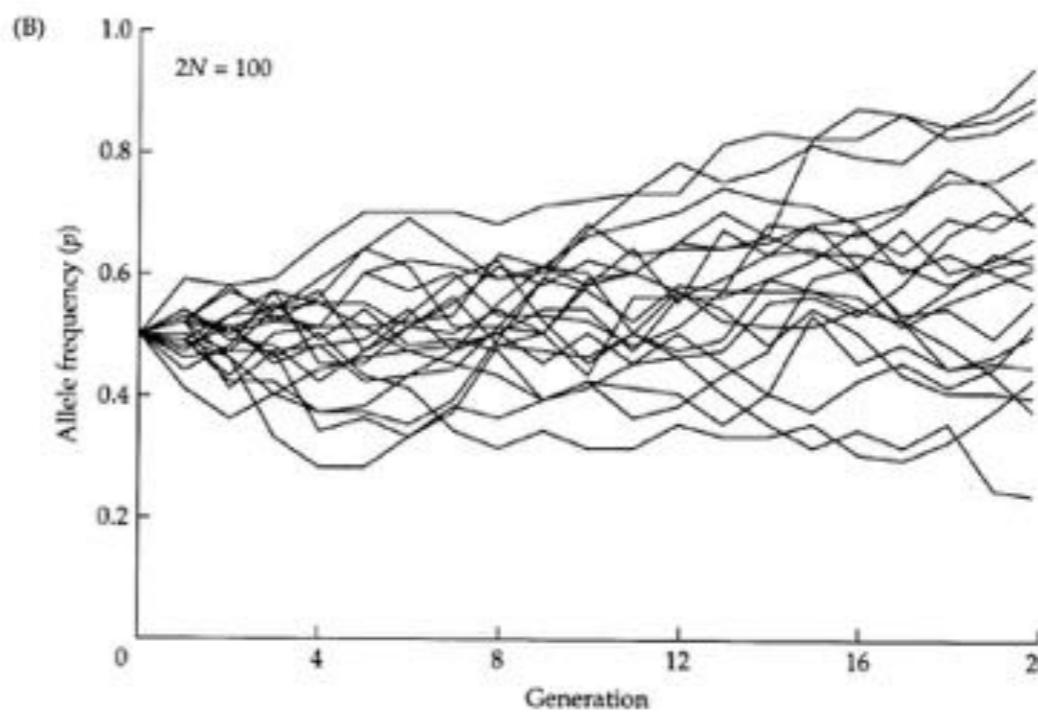


Population size affects rate of diffusion, and genetic drift affects the entire genome equally (on average)

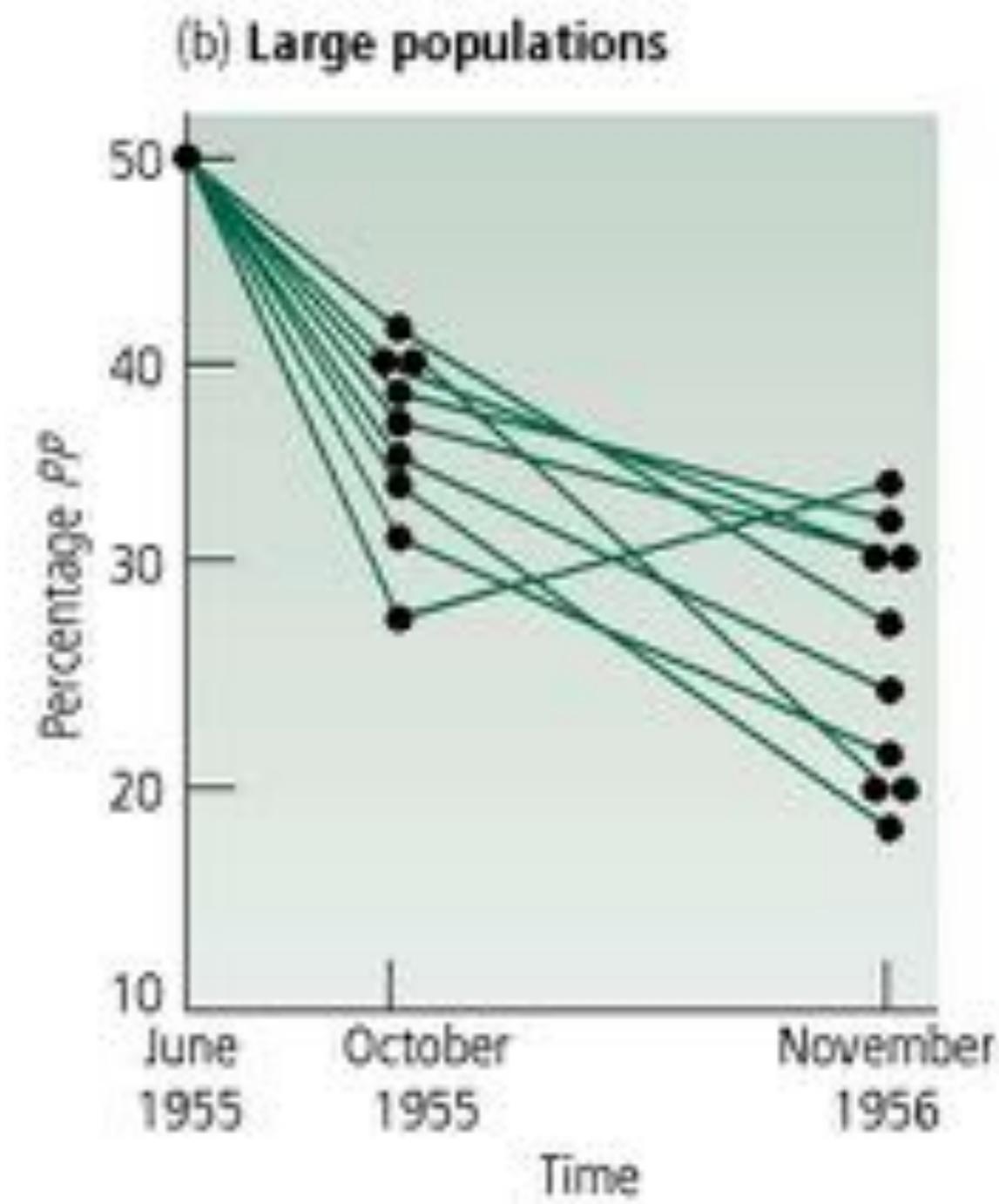
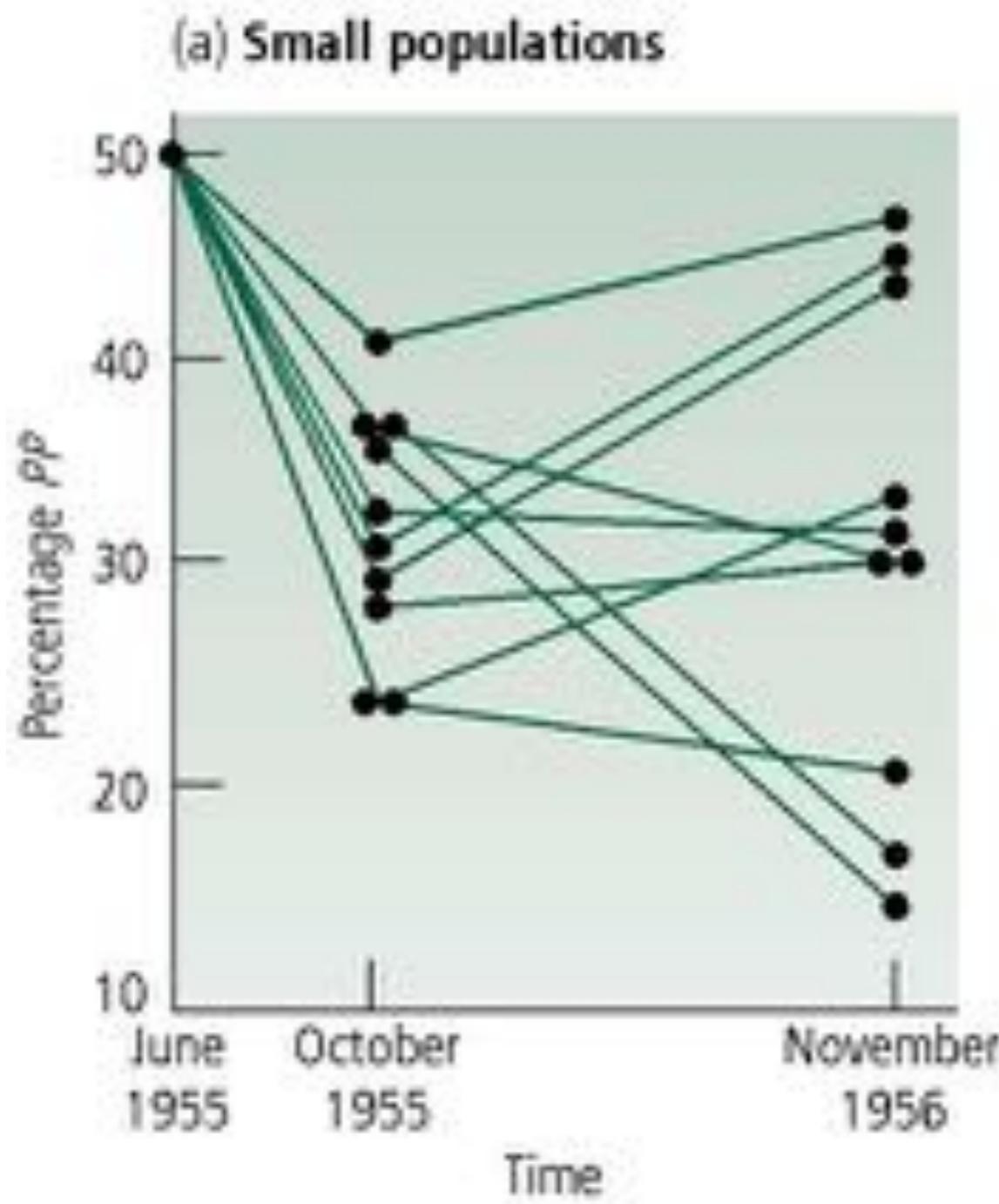
~10 individuals



~100 individuals



Natural selection biases the allele frequency changes and the effects can be genomically localized

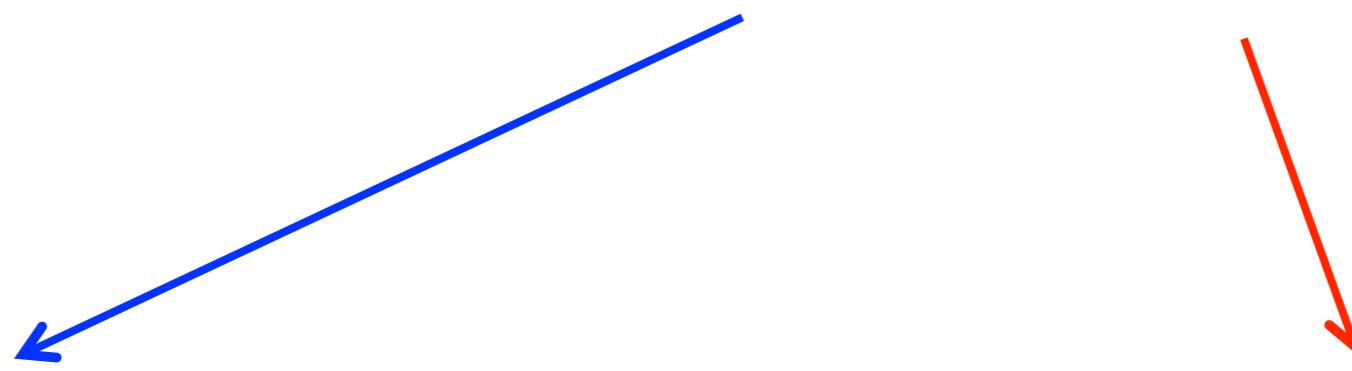


Population genomics

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

Genome-wide background provides more precise estimates:

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

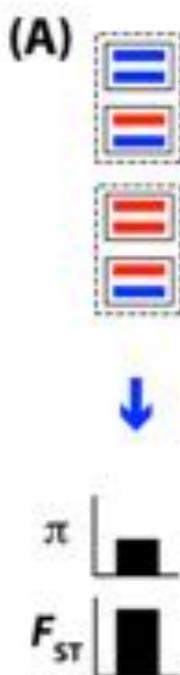


Outliers from background indicate:

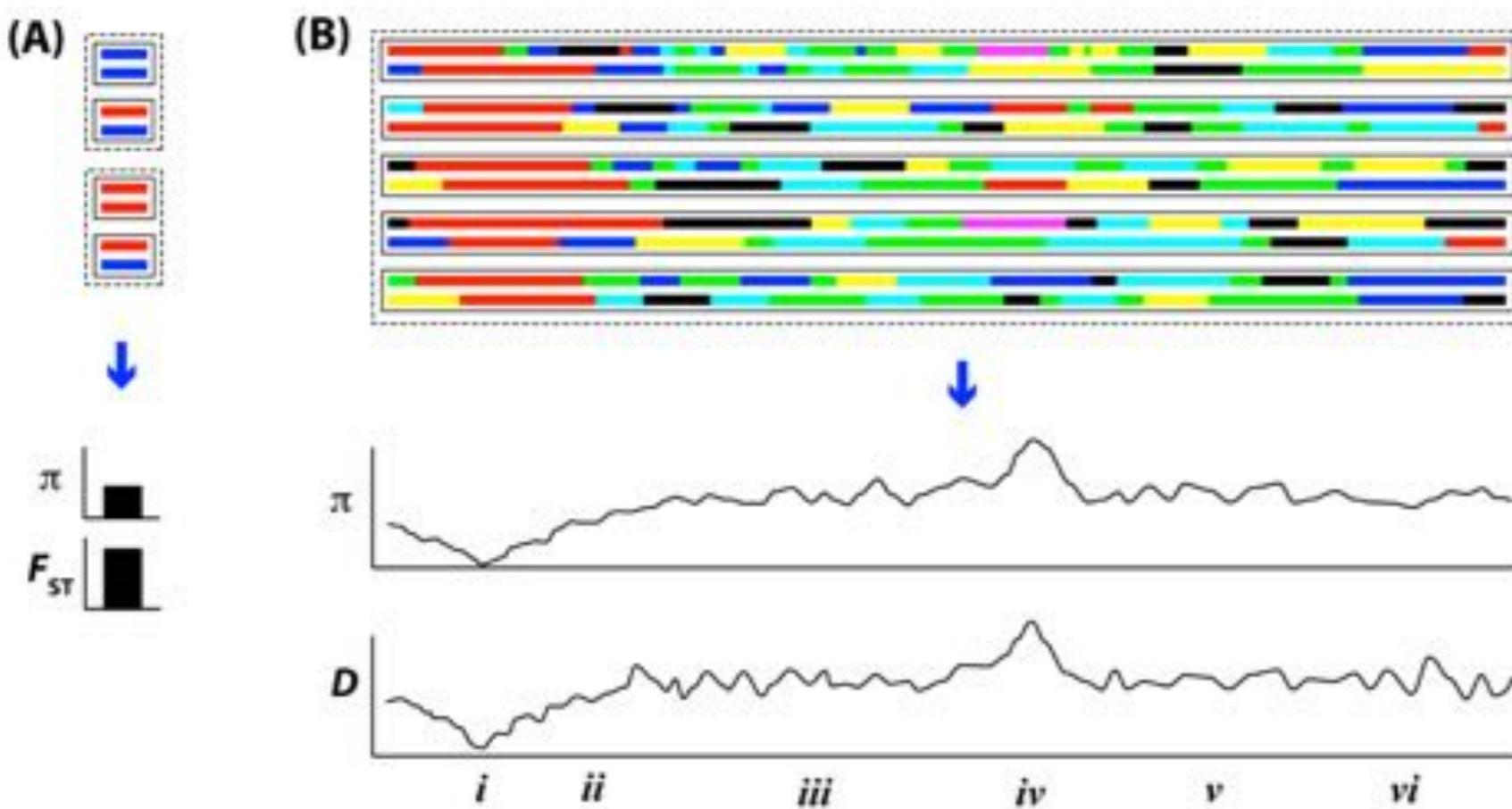
- Selective sweeps
- Local adaptation



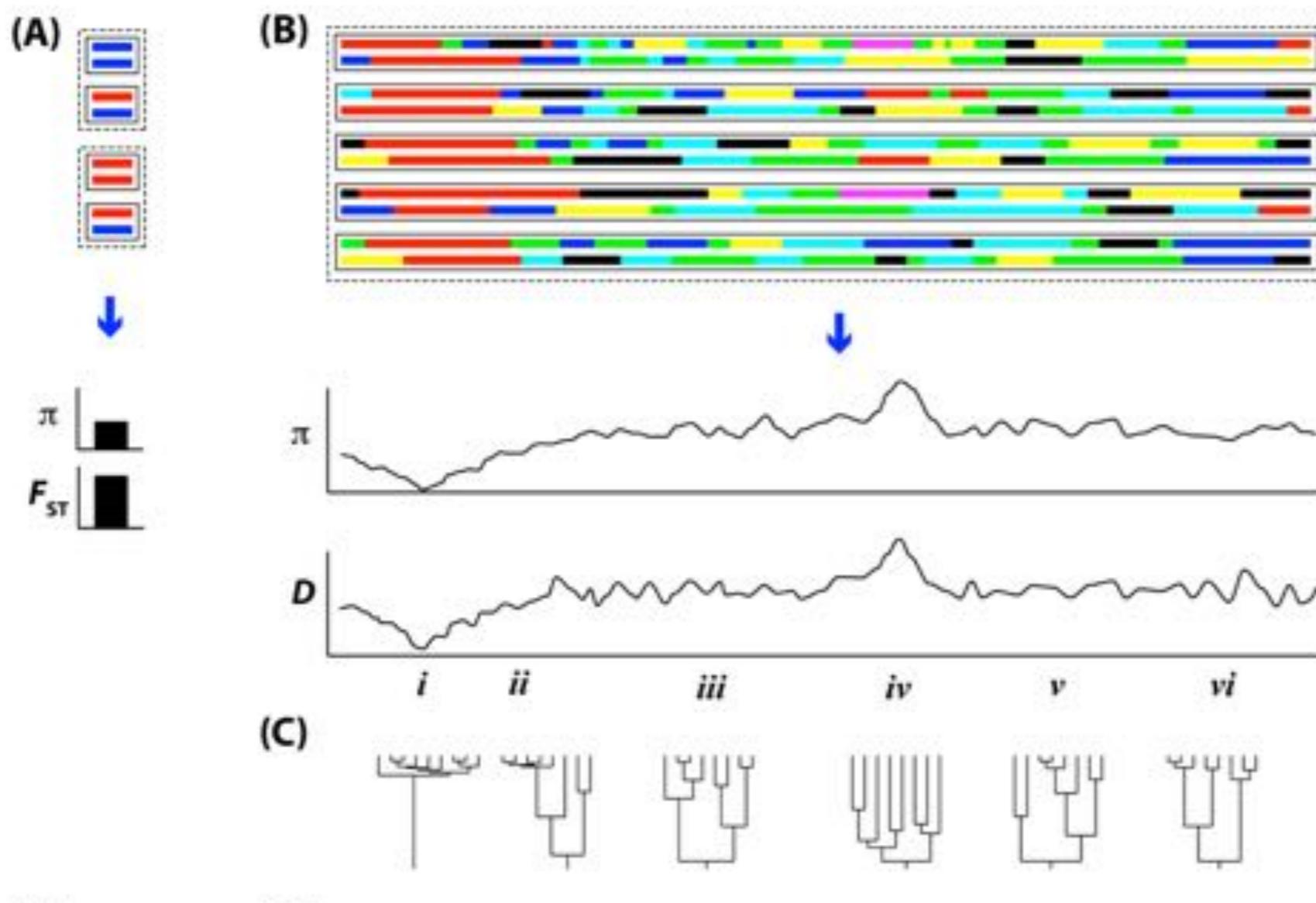
Population genomics of ordered markers



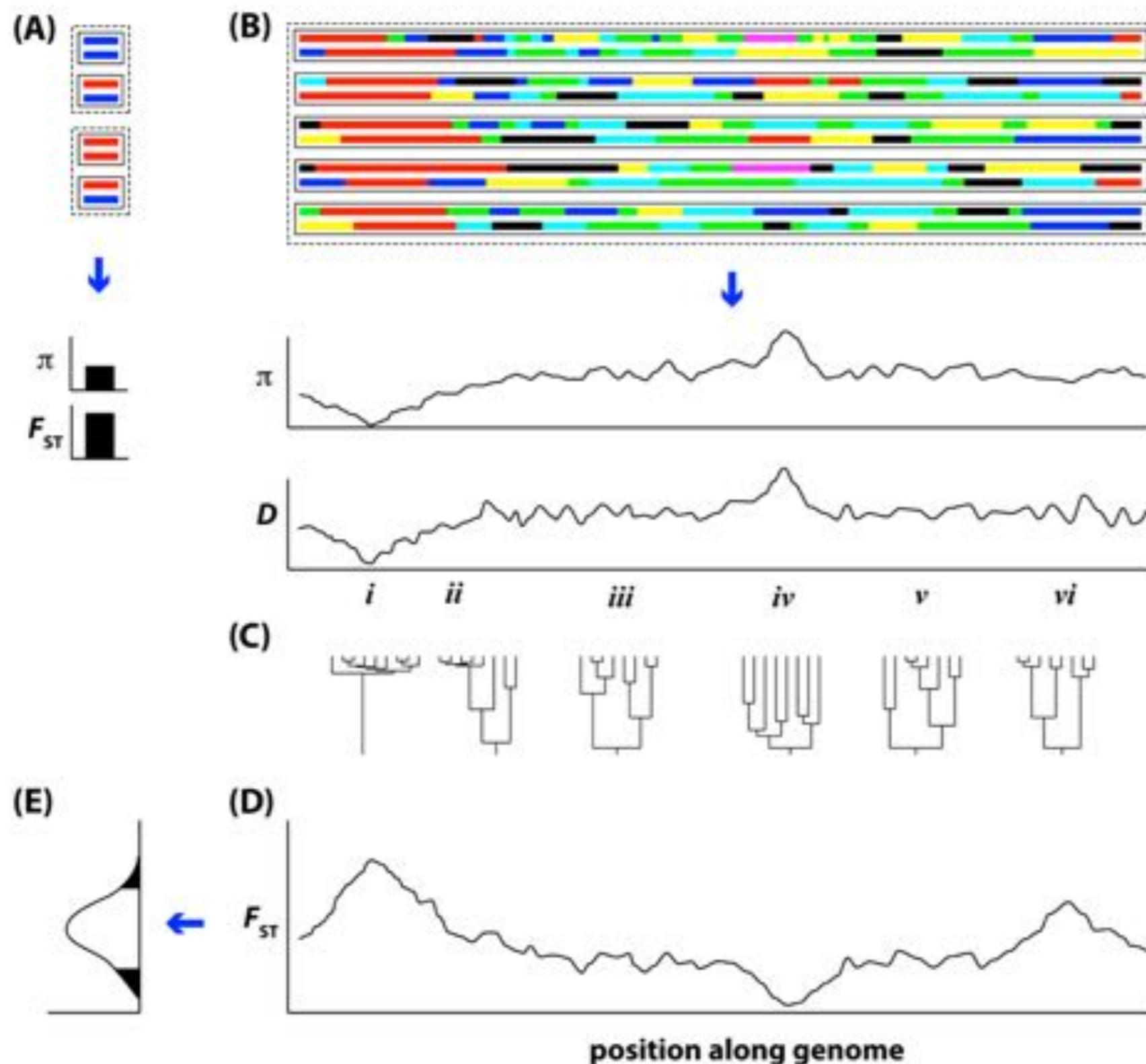
Population genomics of ordered markers



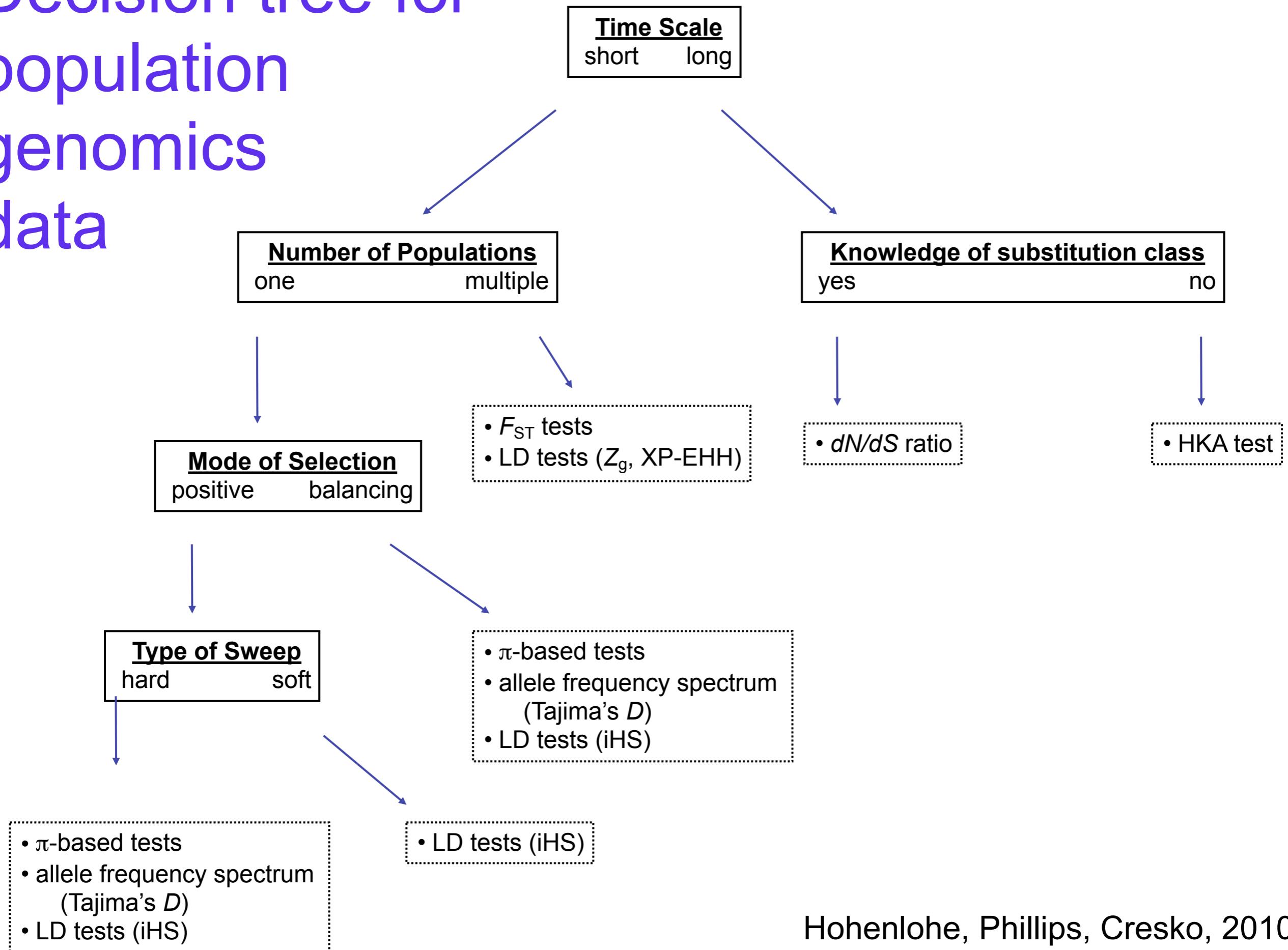
Population genomics of ordered markers



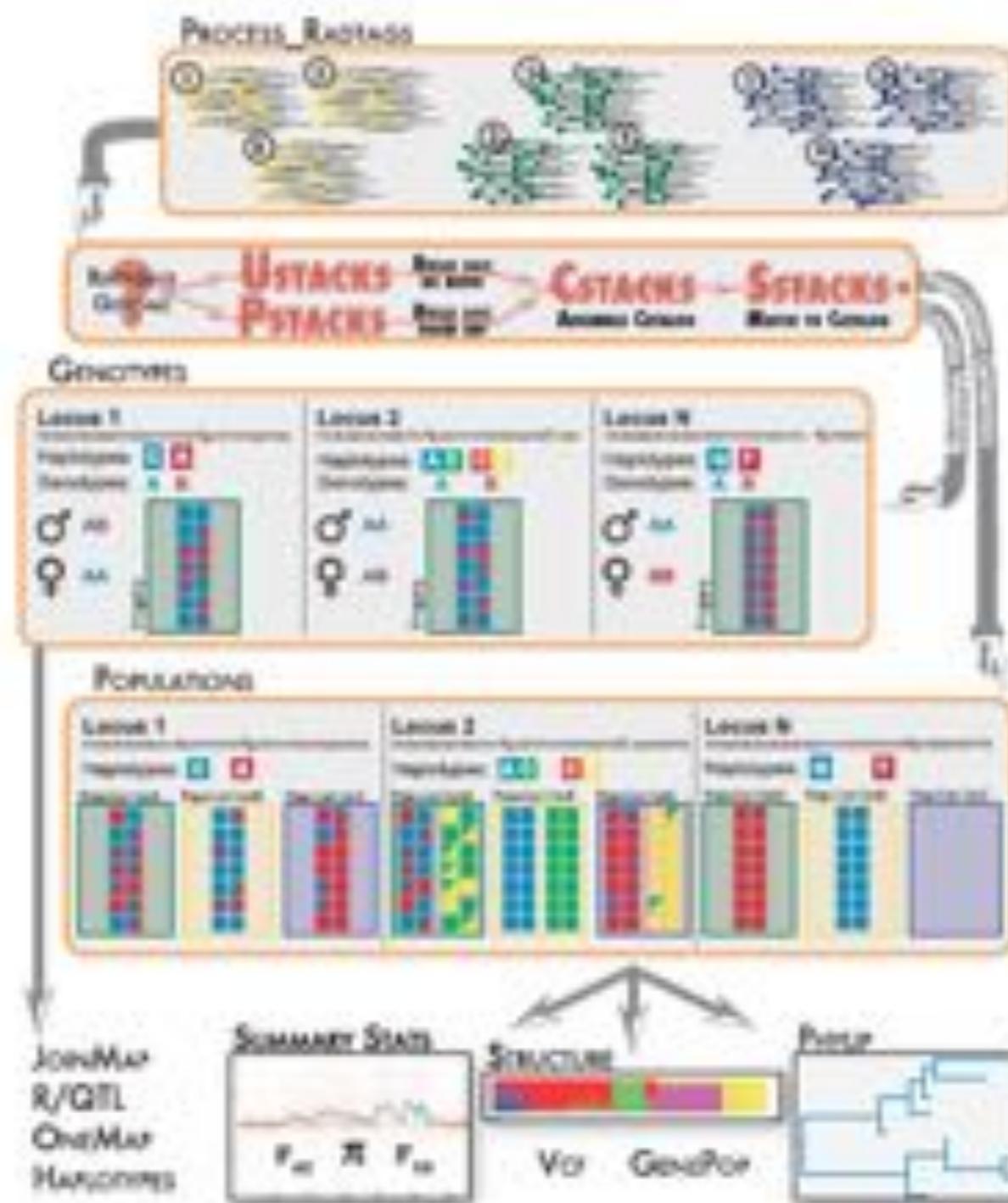
Population genomics of ordered markers



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,[†] Paul Hohenlohe,* William Cresko,* and John H. Postlethwait^{‡,1}

*Center for Ecology and Evolutionary Biology and [†]Institute of Neuroscience, University of Oregon, Eugene, Oregon 97403

Stacks: an analysis tool set for population genomics

JULIAN CATCHEN,* PAUL A. HOHENLOHE,*[†] SUSAN BASSHAM,* ANGEL AMORES[‡] and WILLIAM A. CRESKO*

^{*}Institute of Ecology and Evolution, University of Oregon, Eugene, OR 97403-5289, USA, [†]Biological Sciences, University of Idaho, Moscow, ID 83844-3051, USA, [‡]Institute of Neuroscience, University of Oregon, Eugene, OR 97403-1254, USA

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	Π	Mean value of π in this population.
22	Π Variance	
23	Π 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.

Stacks

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 π	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 3σ of sequence in both directions.
17	AMOVA F_{ST}	Analysis of Molecular Variance alternative F_{ST} calculation. Derived from Weir, <i>Genetic Data Analysis II</i> , 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 3σ 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 -e)

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	Φ_{ST}	
8	Smoothed Φ_{ST}	
9	Smoothed Φ_{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

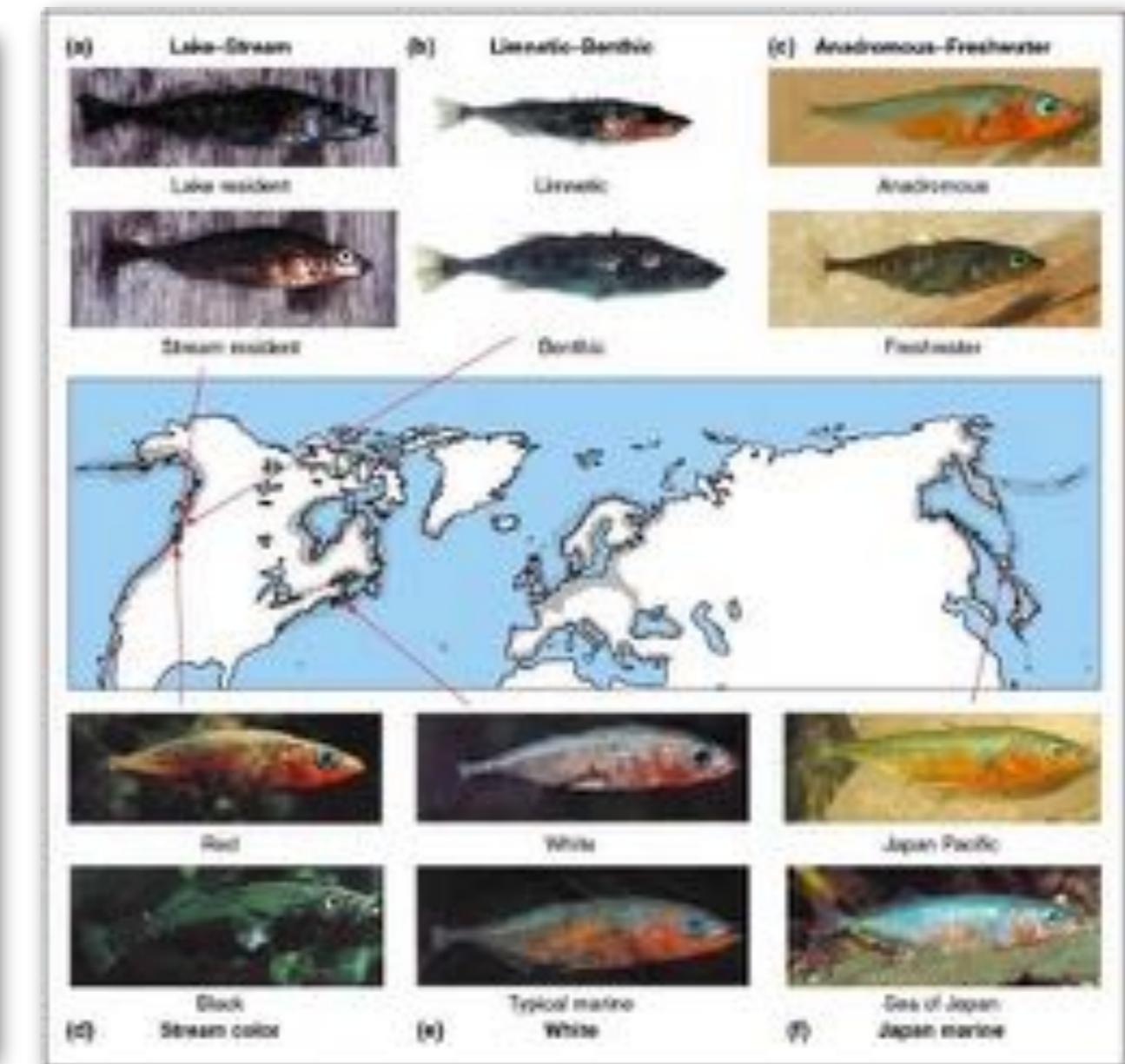
- Population genomics of rapid adaptation
- Using long read RAD-seq for coalescent analyses
- Genome Wide Association Studies using RAD-seq

Genomically enabling the Gulf pipefish

Threespine stickleback, *Gasterosteus aculeatus*



Threespine stickleback, *Gasterosteus aculeatus*



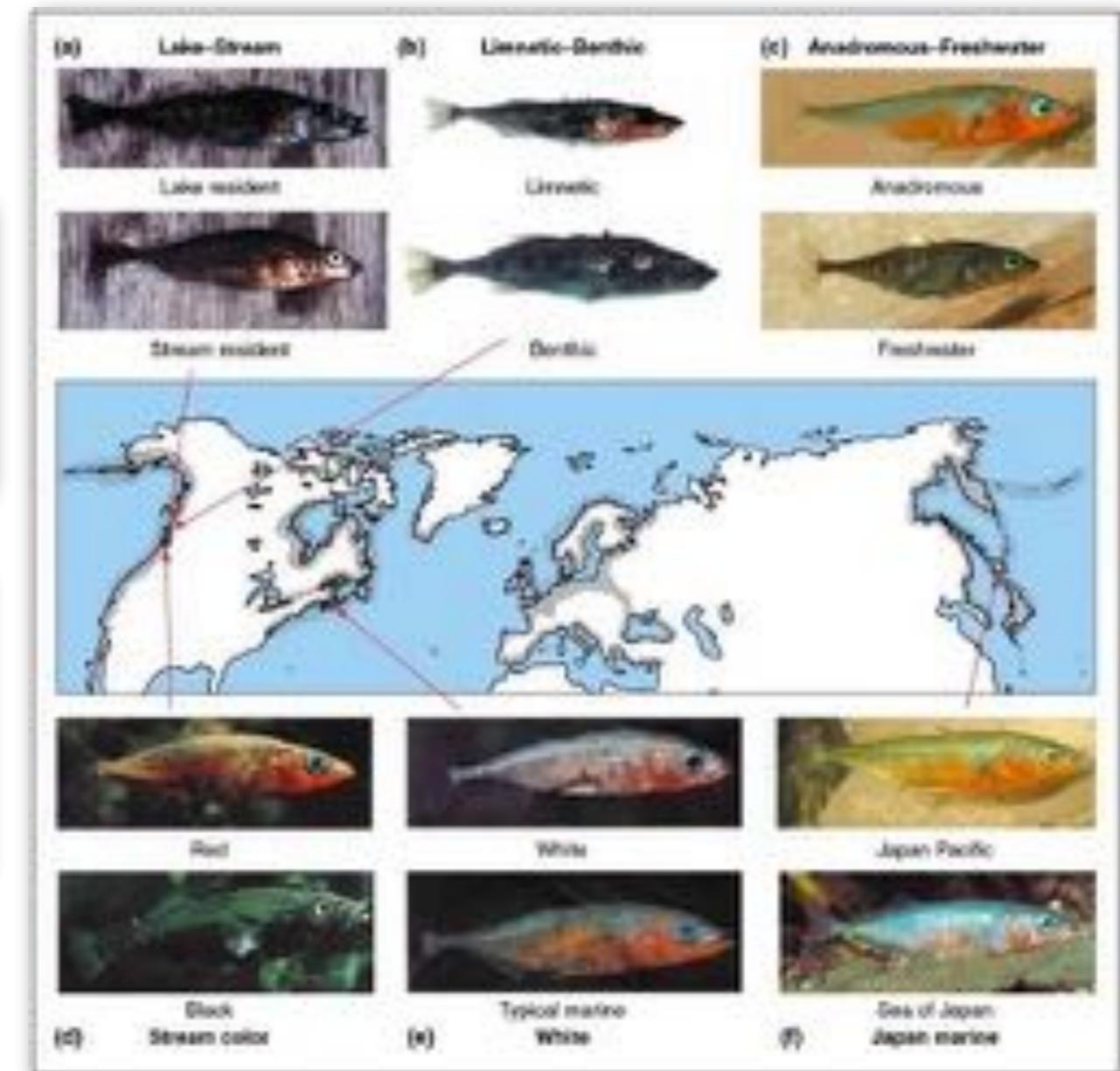
Rundle and McKinnon 2002

Threespine stickleback, *Gasterosteus aculeatus*

Pelvic
Structure



Lateral
Plates



Rundle and McKinnon 2002



Some stickleback phenotypes mapped in the lab far

Pelvic structure size and shape *** (Eda)

Lateral plate number *** (Pitx1)

Body coloration *** (KitL)

Opercle bone shape

Pelvic spine length

Body shape

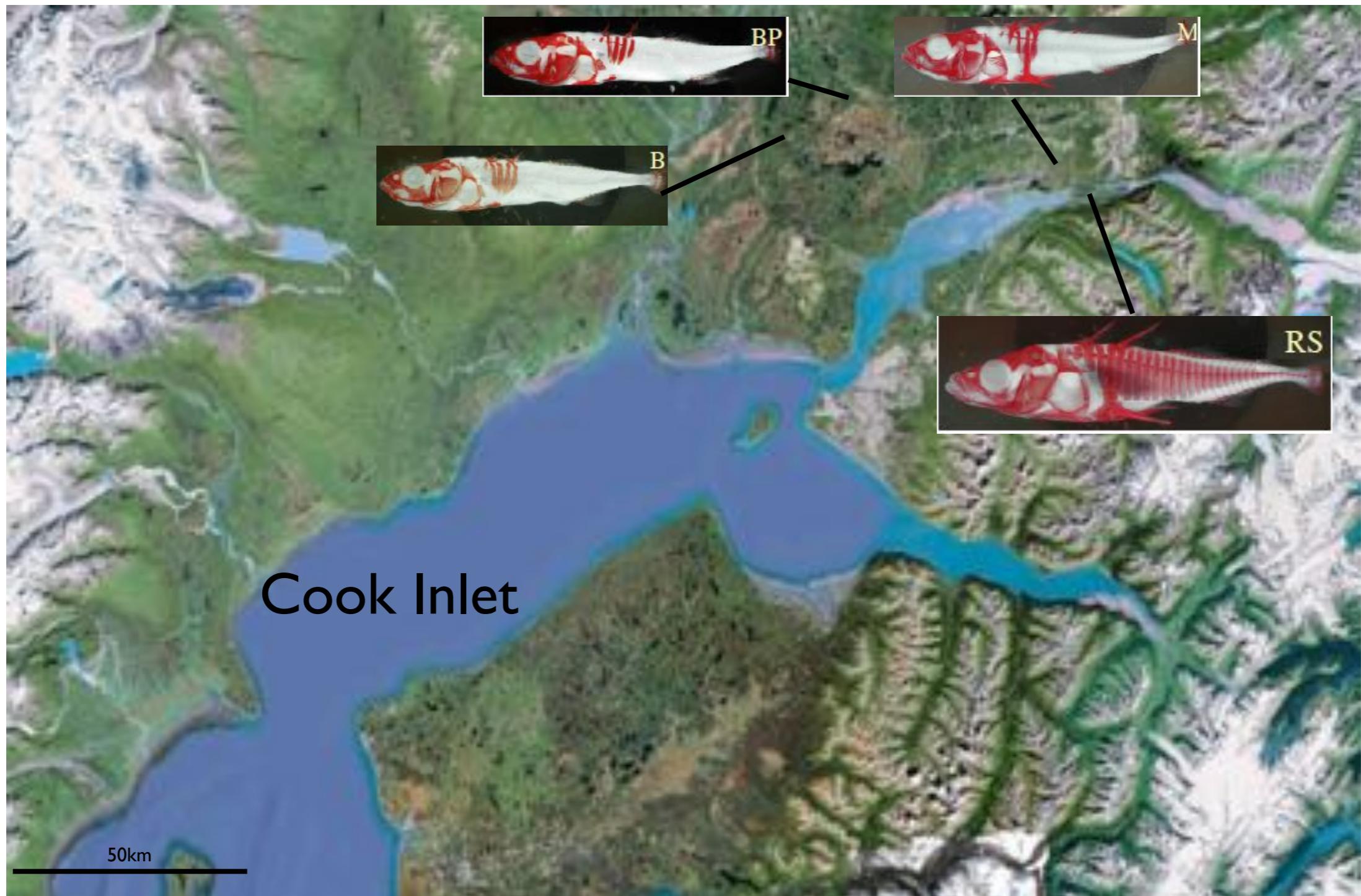
Courtship behavior

Gill raker size

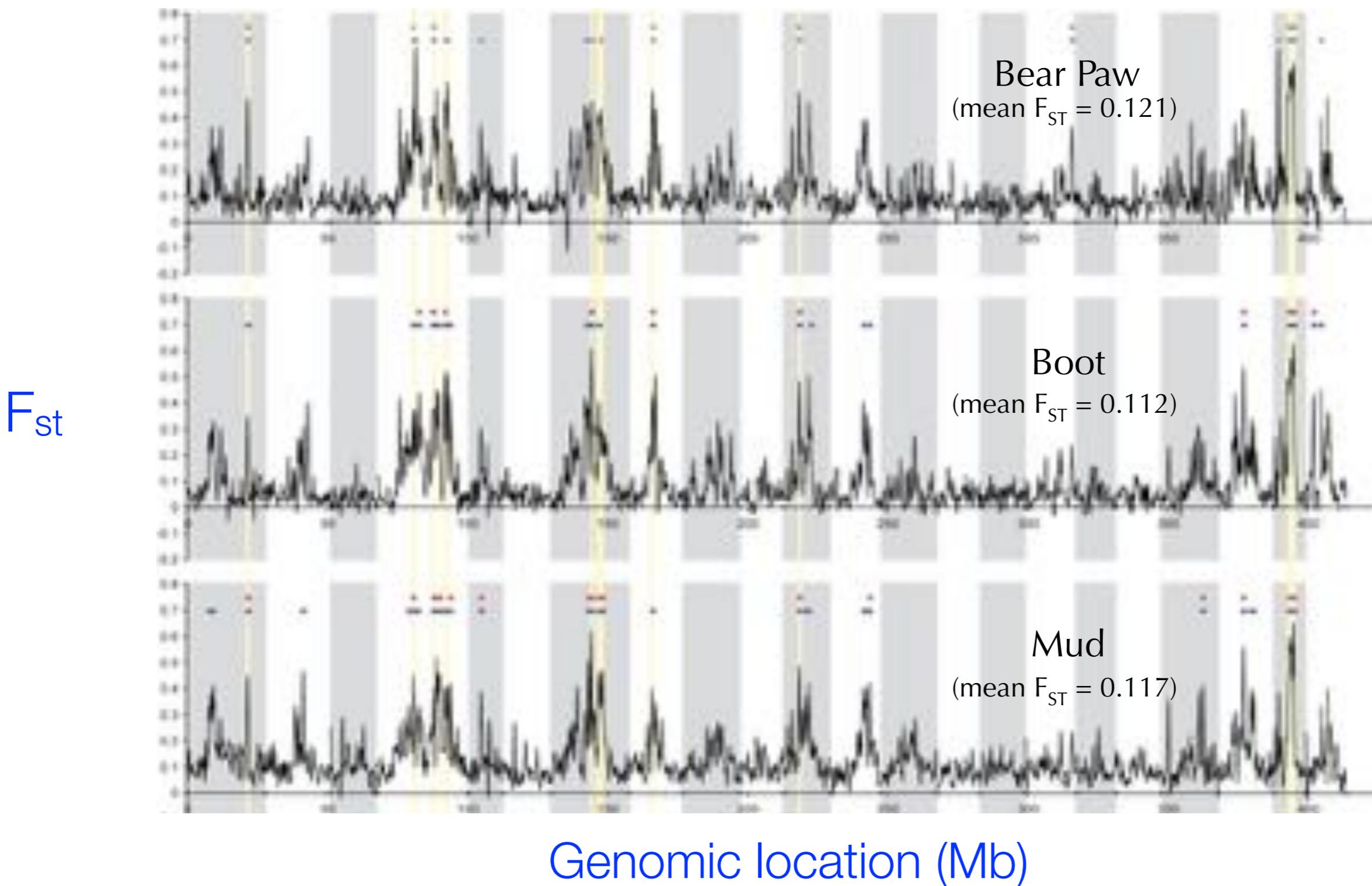
Dorsal spine length

-
- 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

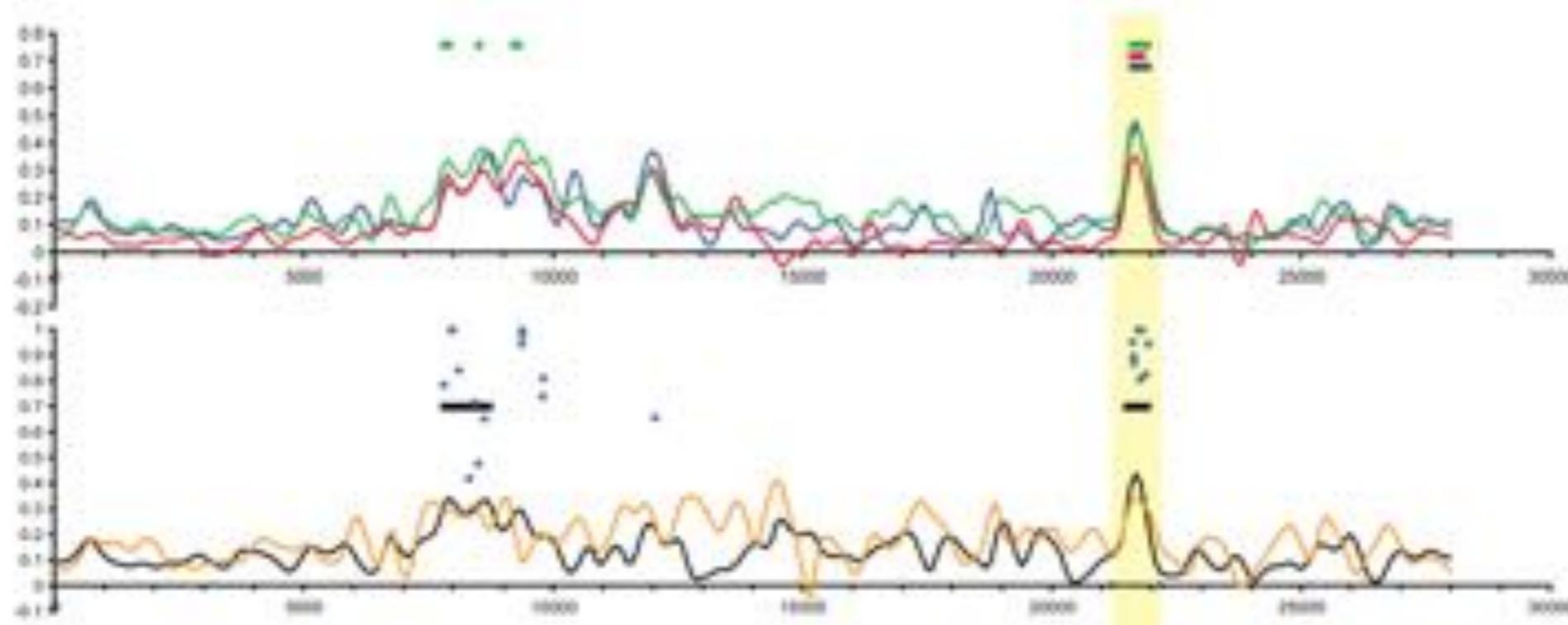


Signatures of natural selection in 13,000 years

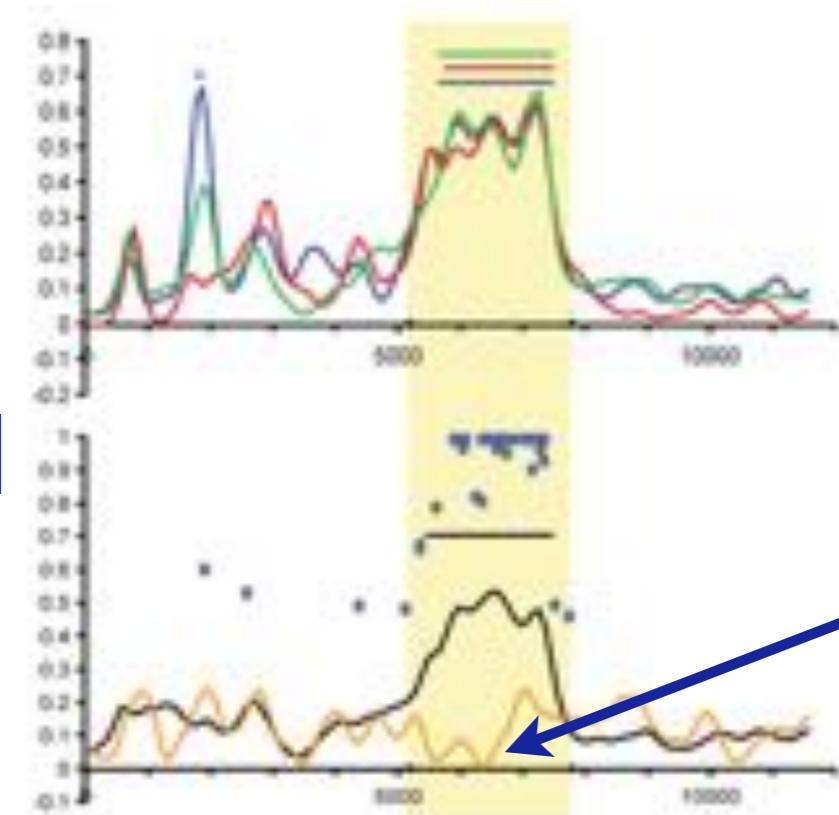


Numerous novel regions identified

LGI



LGXXI



Emily
Lescak



Julian
Catchen



Susan
Bassham



Mary
Sherbick



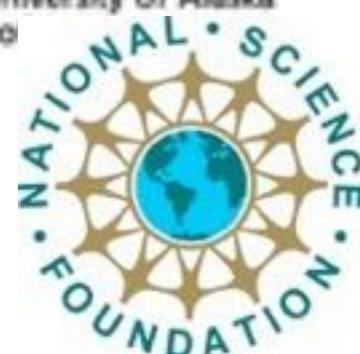
Frank
von Hippel

Evolution of stickleback in 50 years on earthquake-uplifted islands

Emily A. Lescak^{a,b}, Susan L. Bassham^c, Julian Catchen^{c,d}, Ofer Gelson^{b,1}, Mary L. Sherbick^b, Frank A. von Hippel^b, and William A. Cresko^{c,2}

^aSchool of Fisheries and Ocean Sciences, University of Alaska Fairbanks, Fairbanks, AK 99775; ^bDepartment of Biological Sciences, University of Alaska Anchorage, Anchorage, AK 99508; ^cInstitute of Ecology and Evolution, University of Oregon, Eugene, OR 97403; and ^dDepartment 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

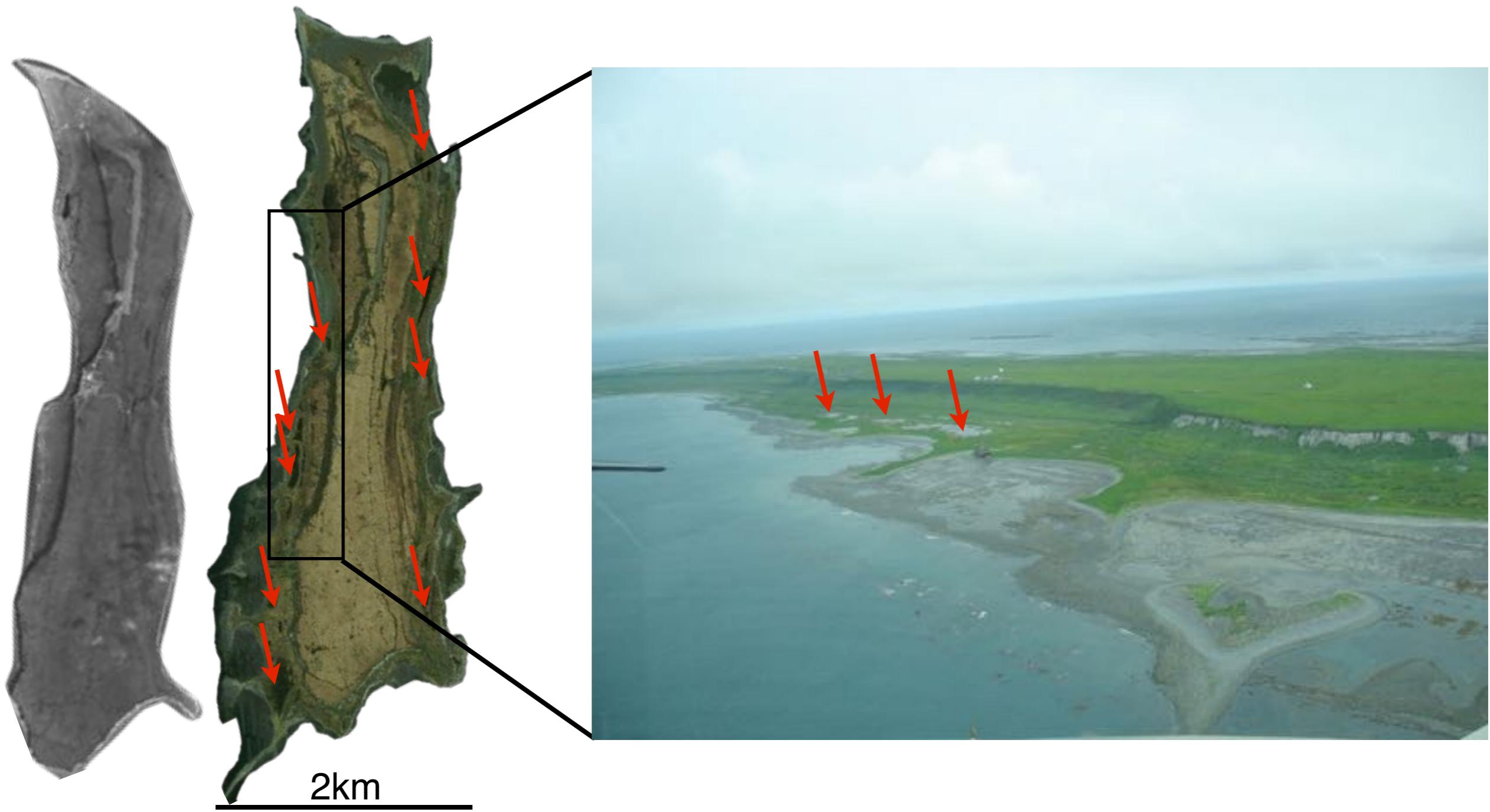
2008



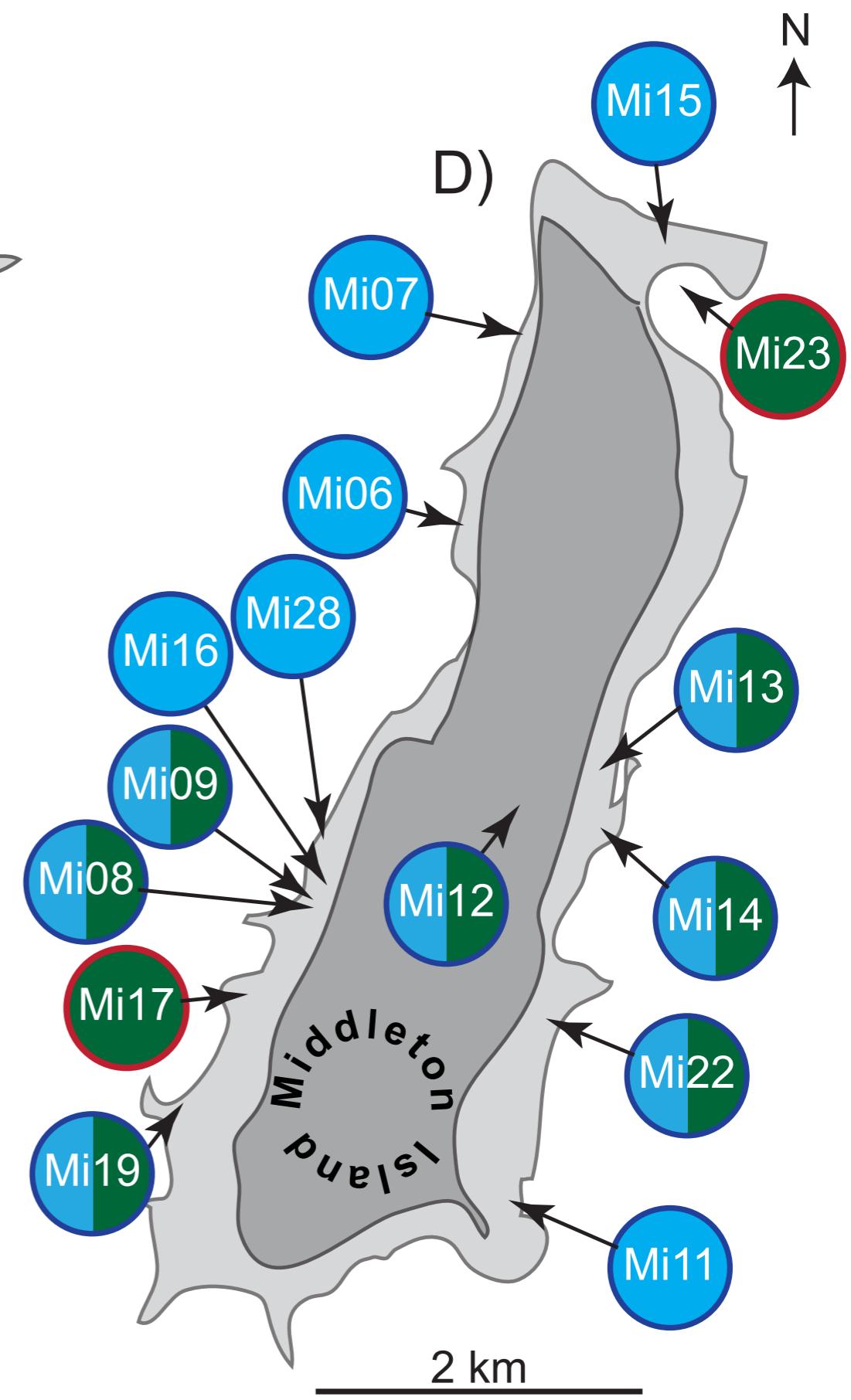
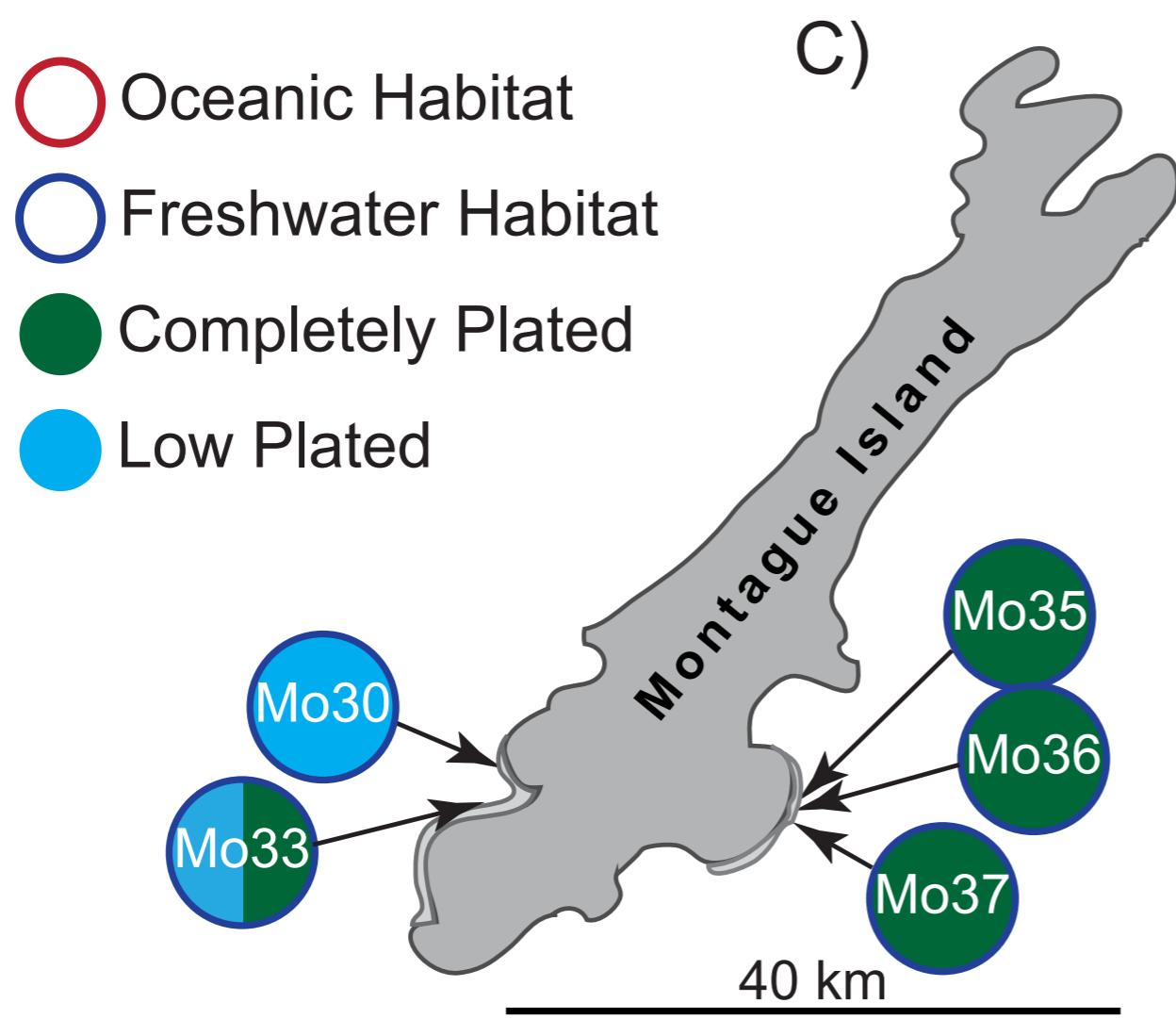
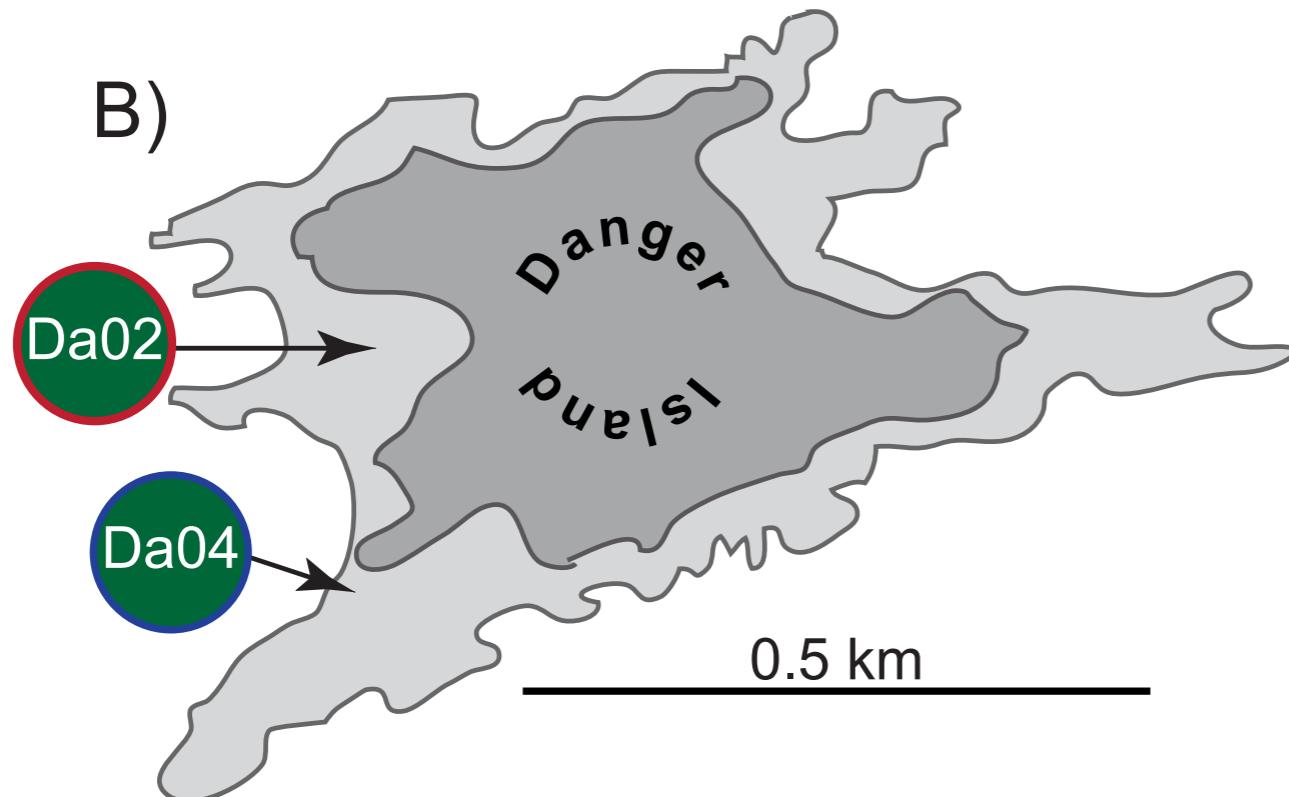
Middleton Island - 50 year old locations

1955

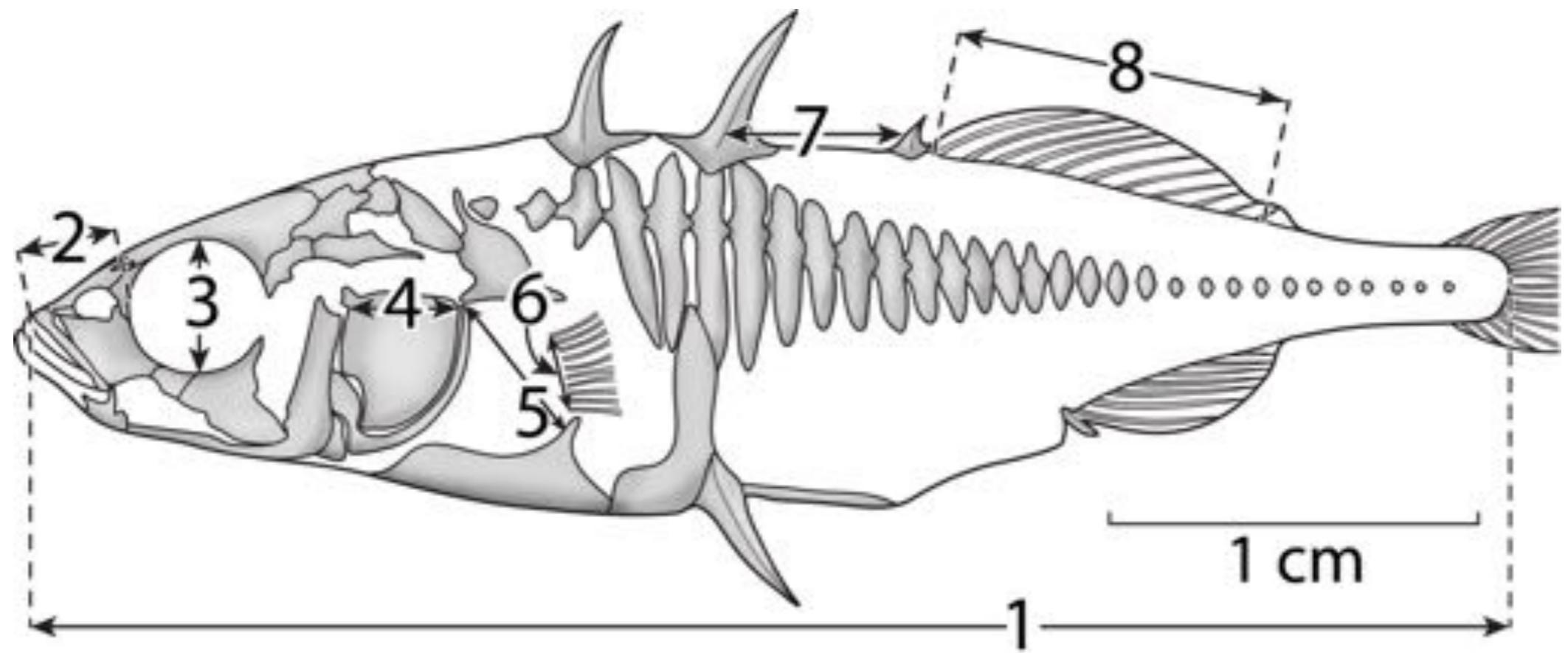
2008

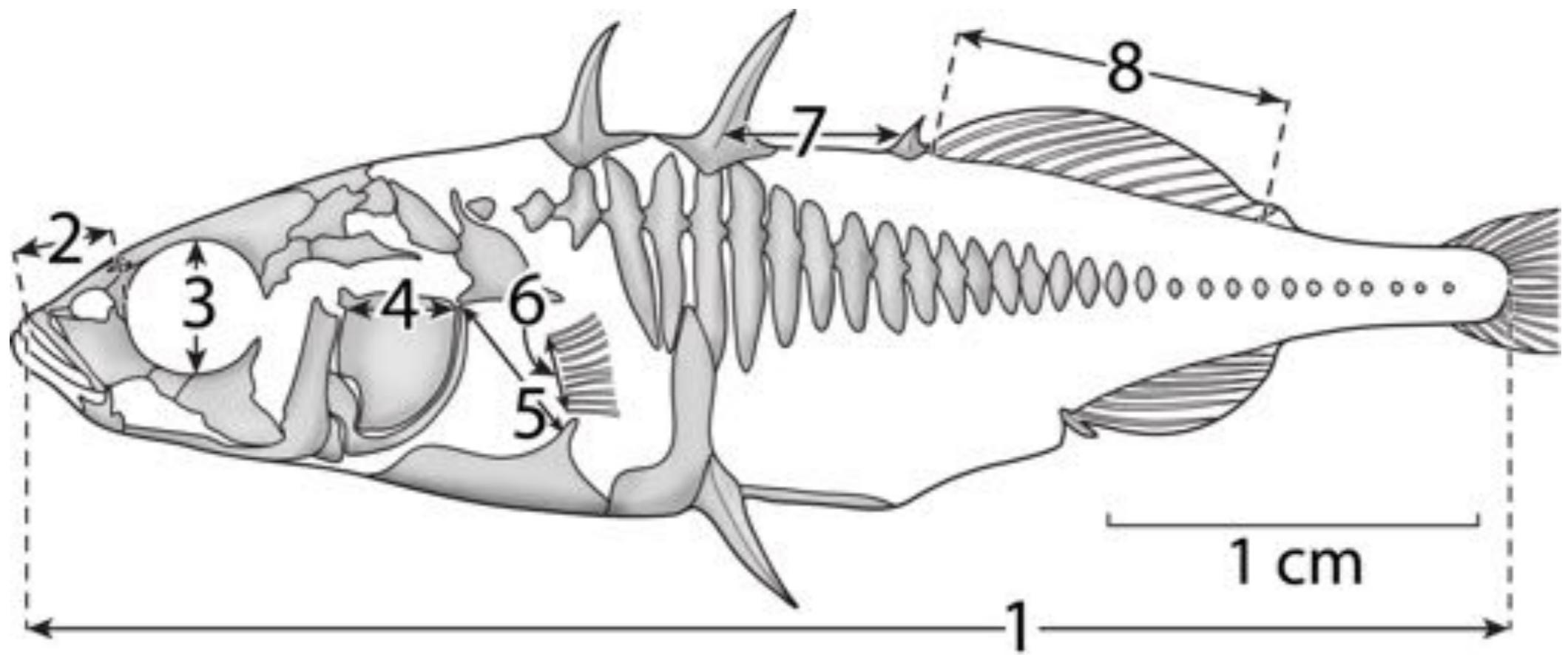






- Oceanic Habitat
- Freshwater Habitat
- Completely Plated
- Low Plated

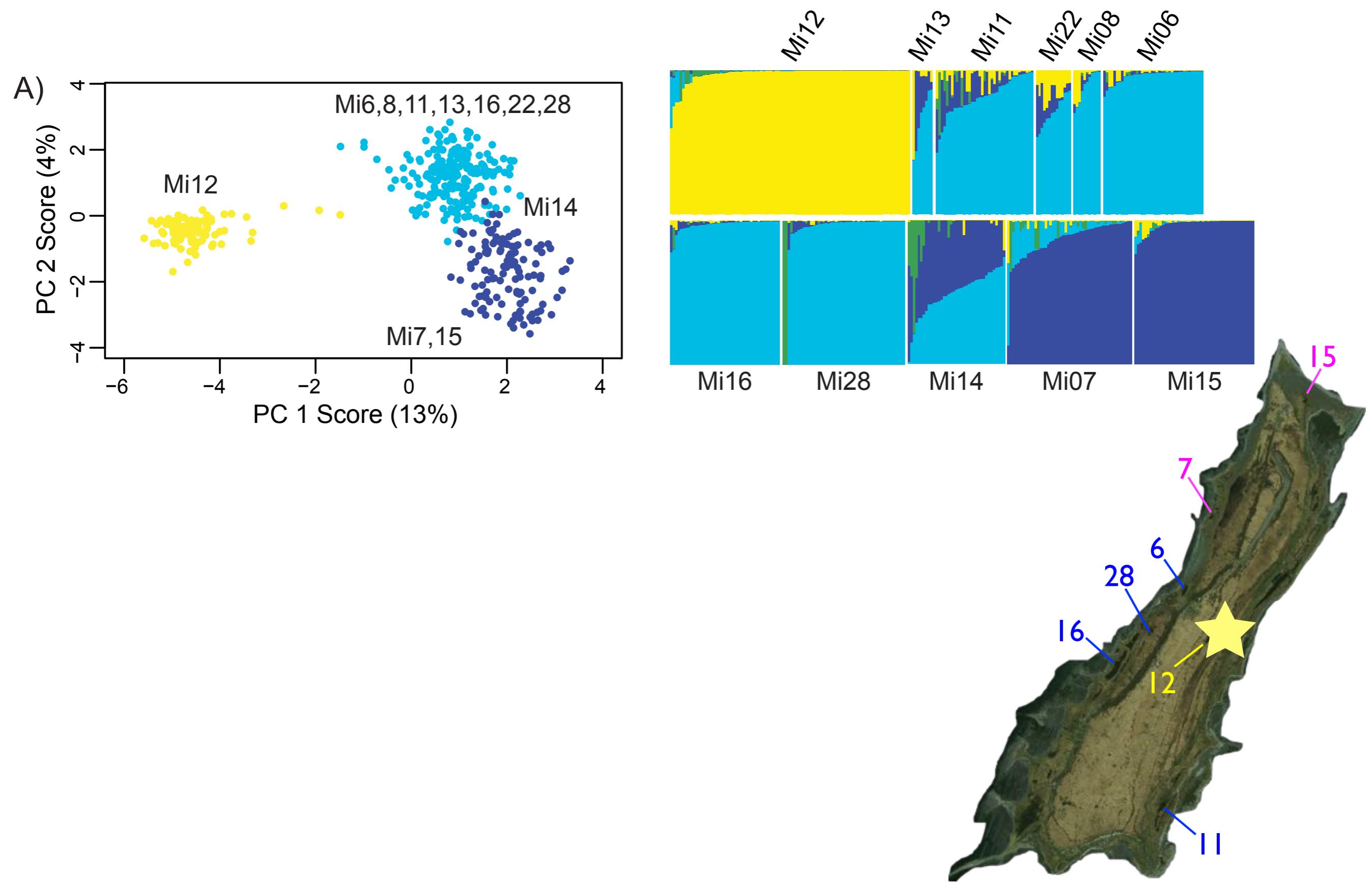




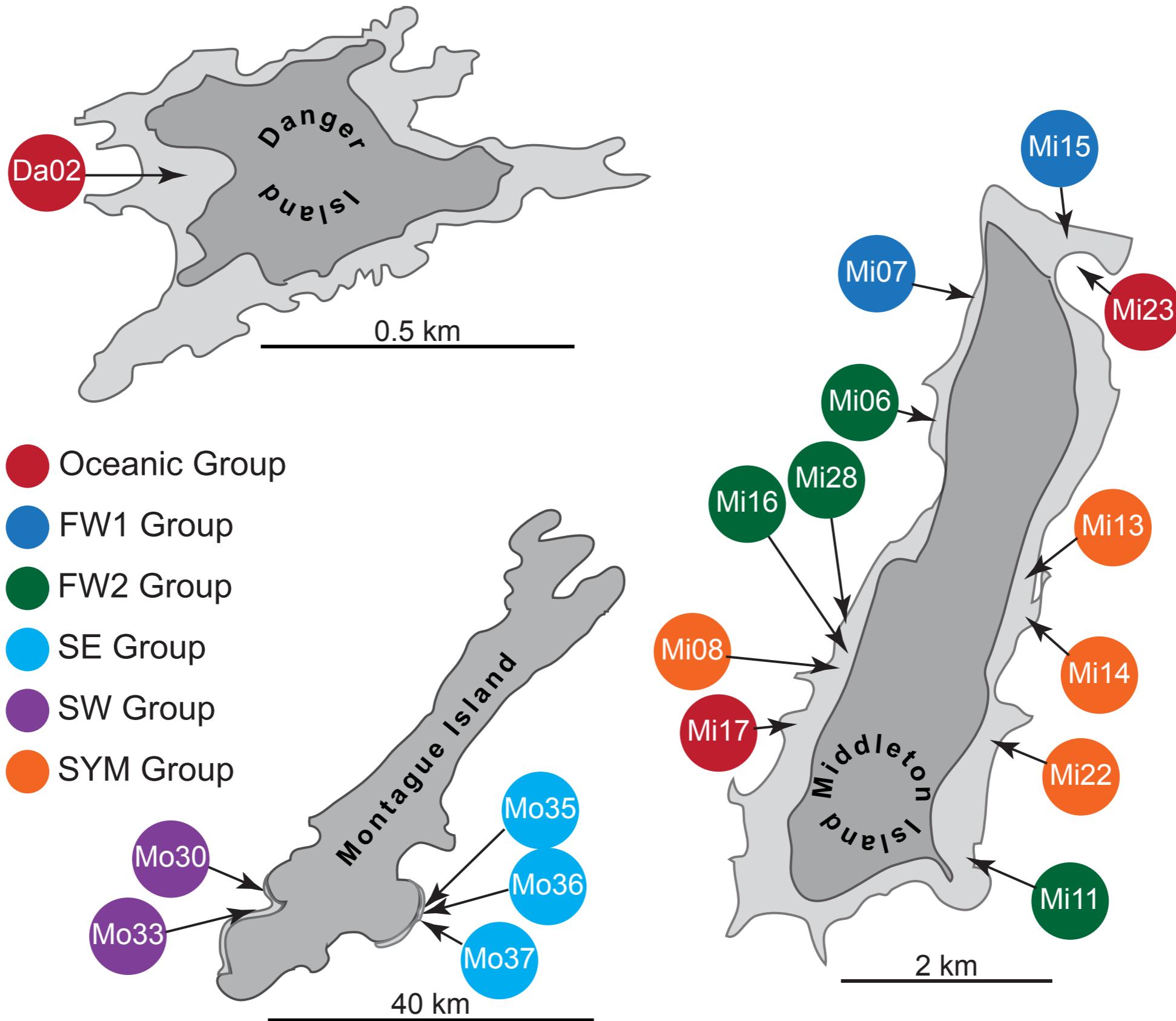
RAD-seq analysis

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

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



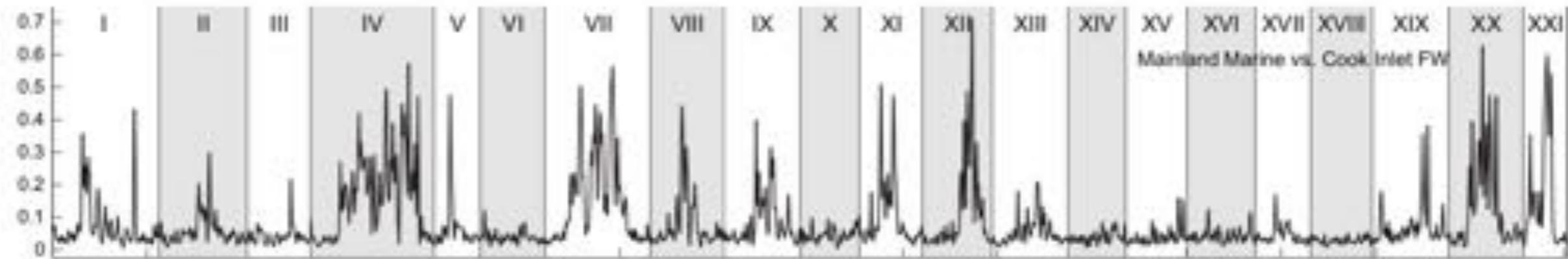
At least six independent evolutionary events in freshwater in the last 50 years



How much of the genome is differentiated?

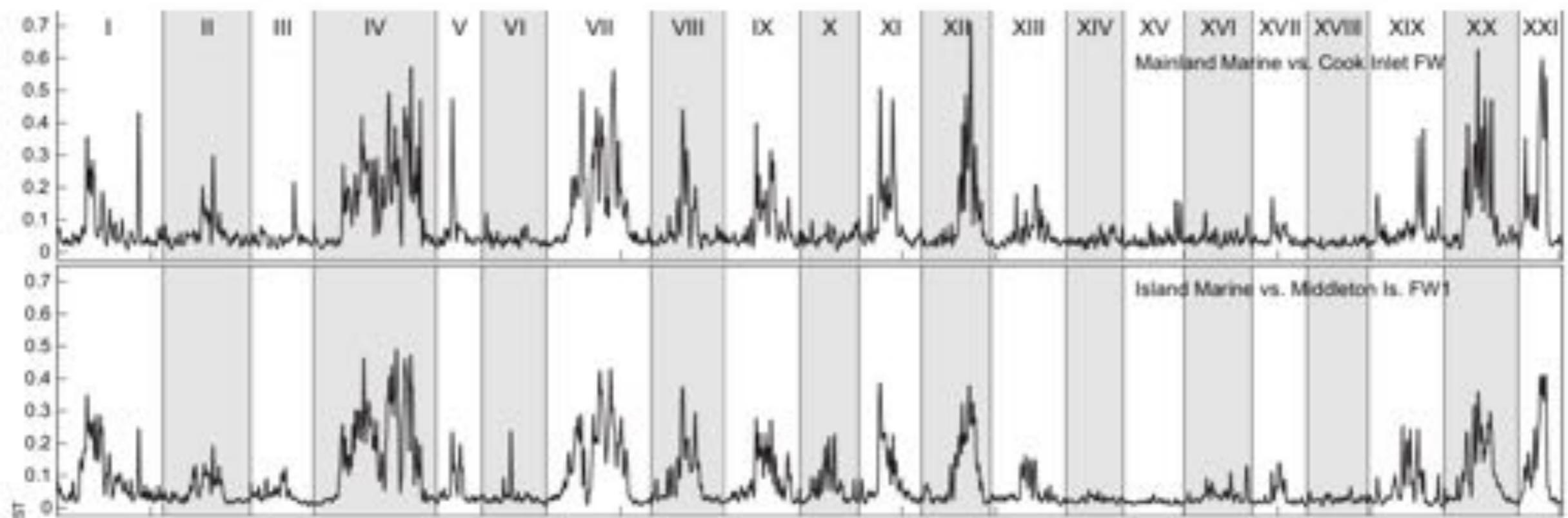
How similar are the genomic patterns of differentiation?

F_{ST}

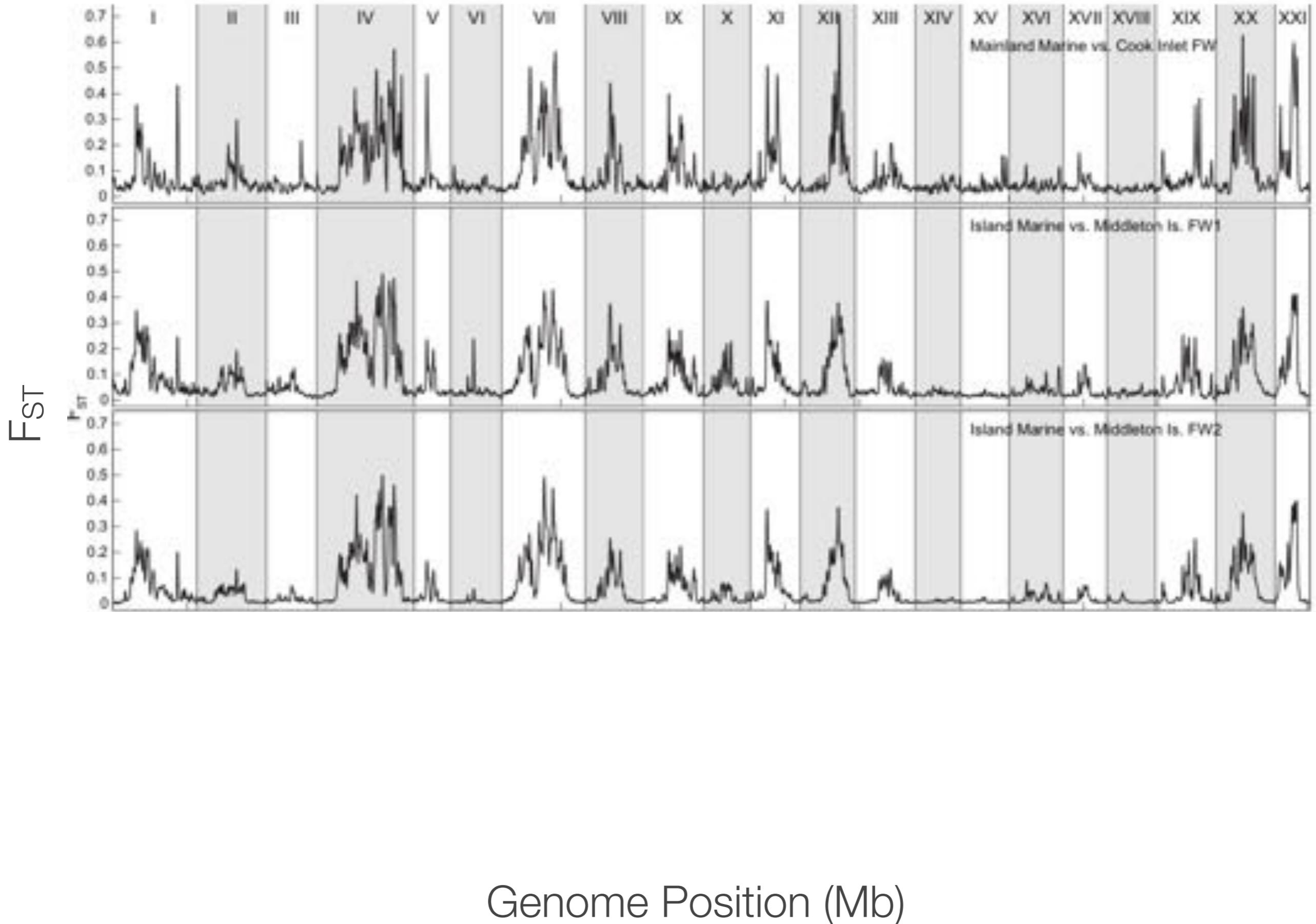


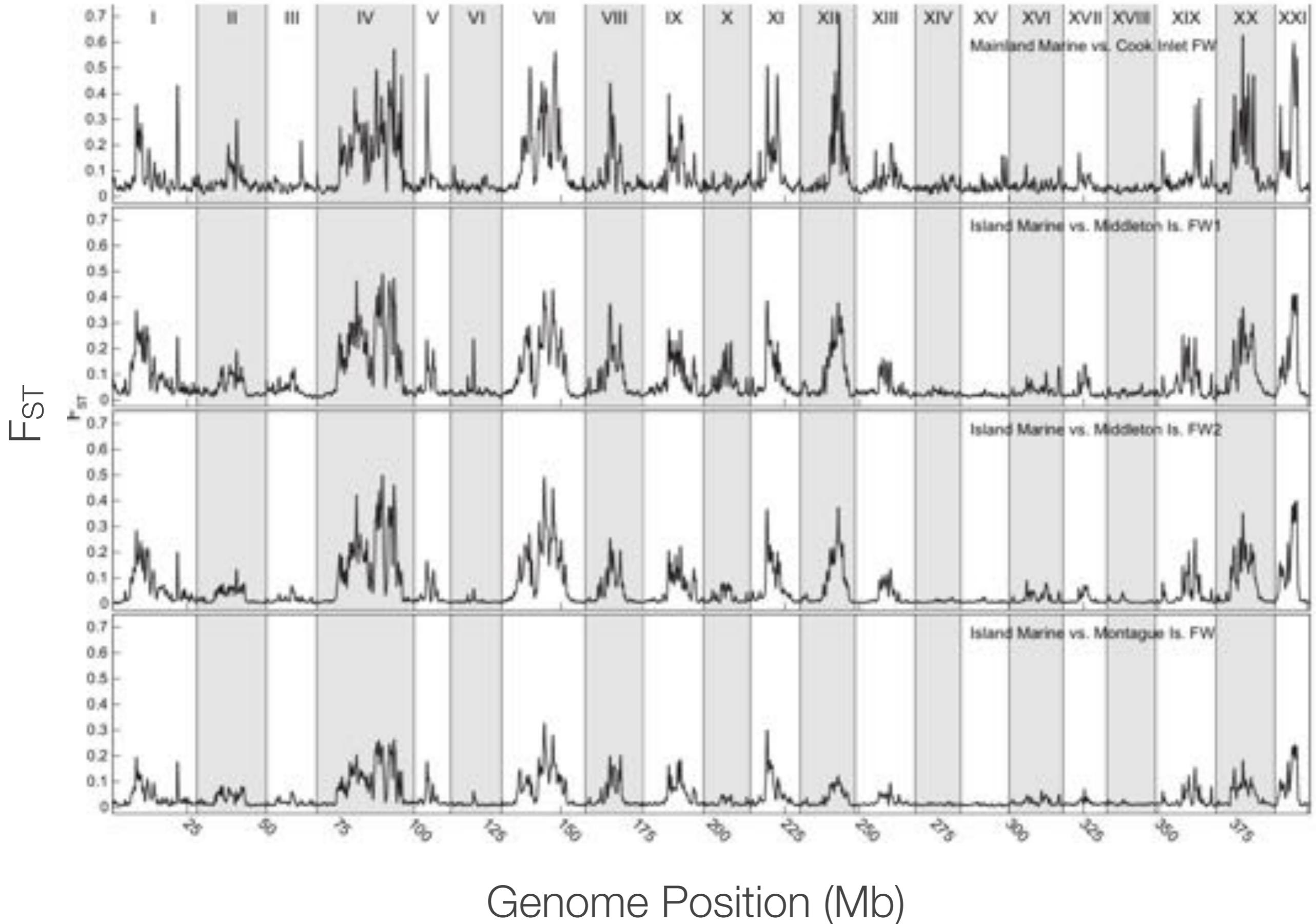
Genome Position (Mb)

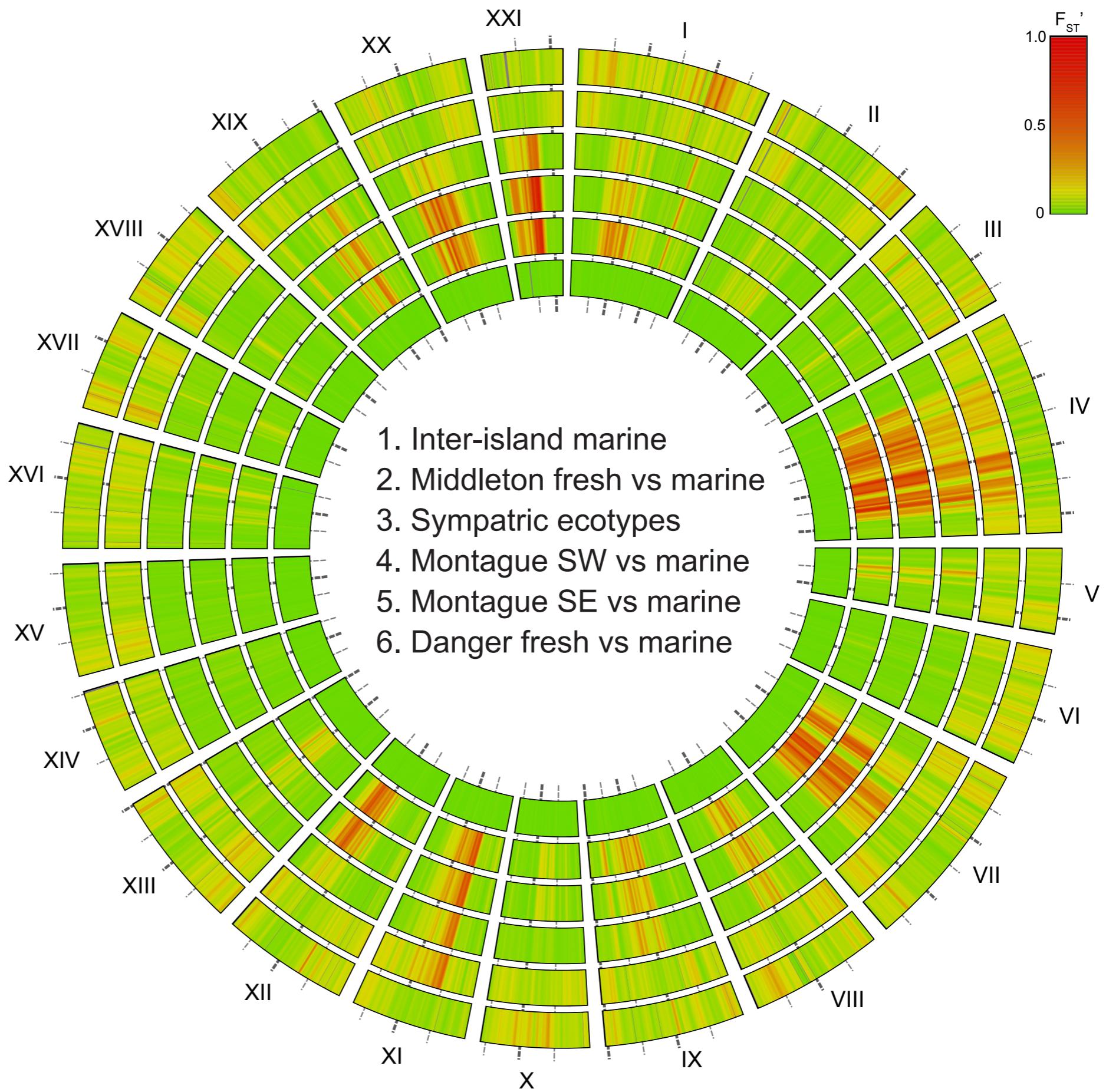
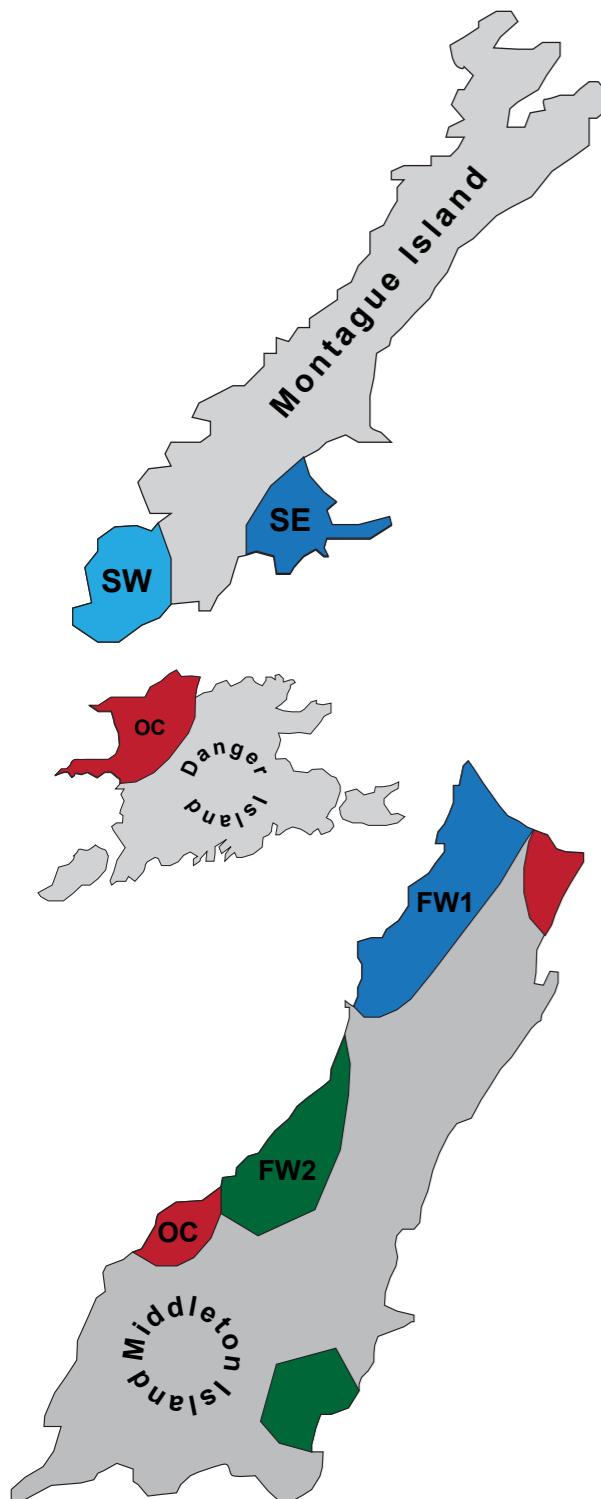
FST

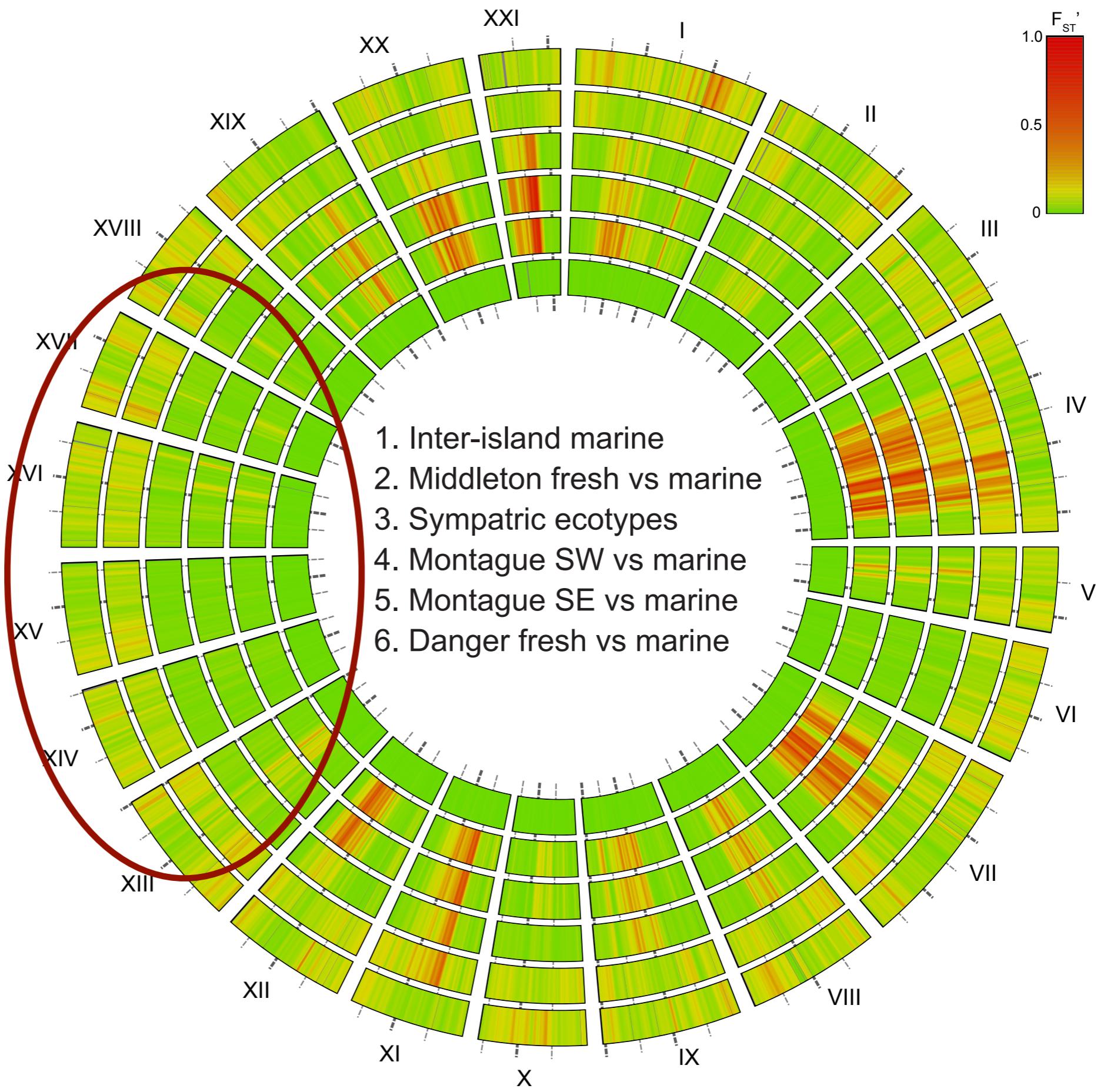
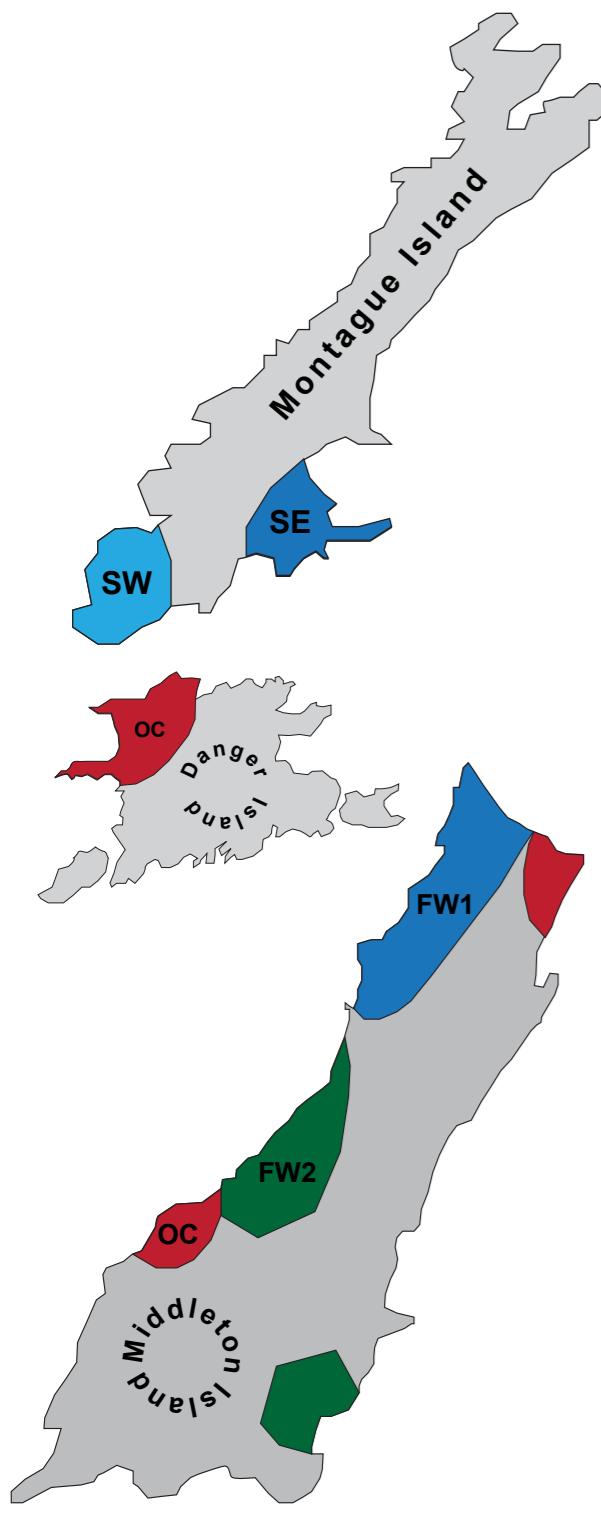


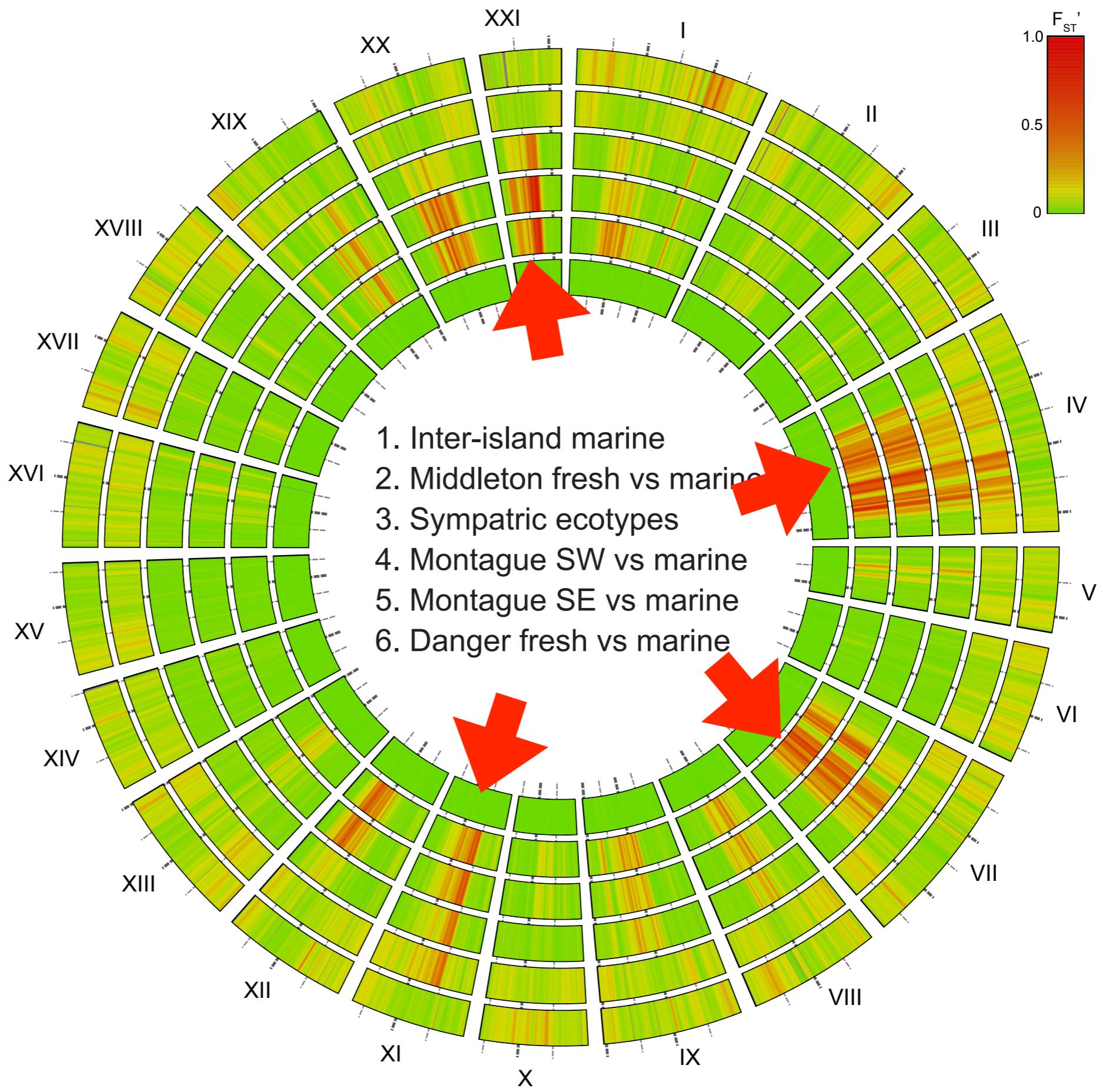
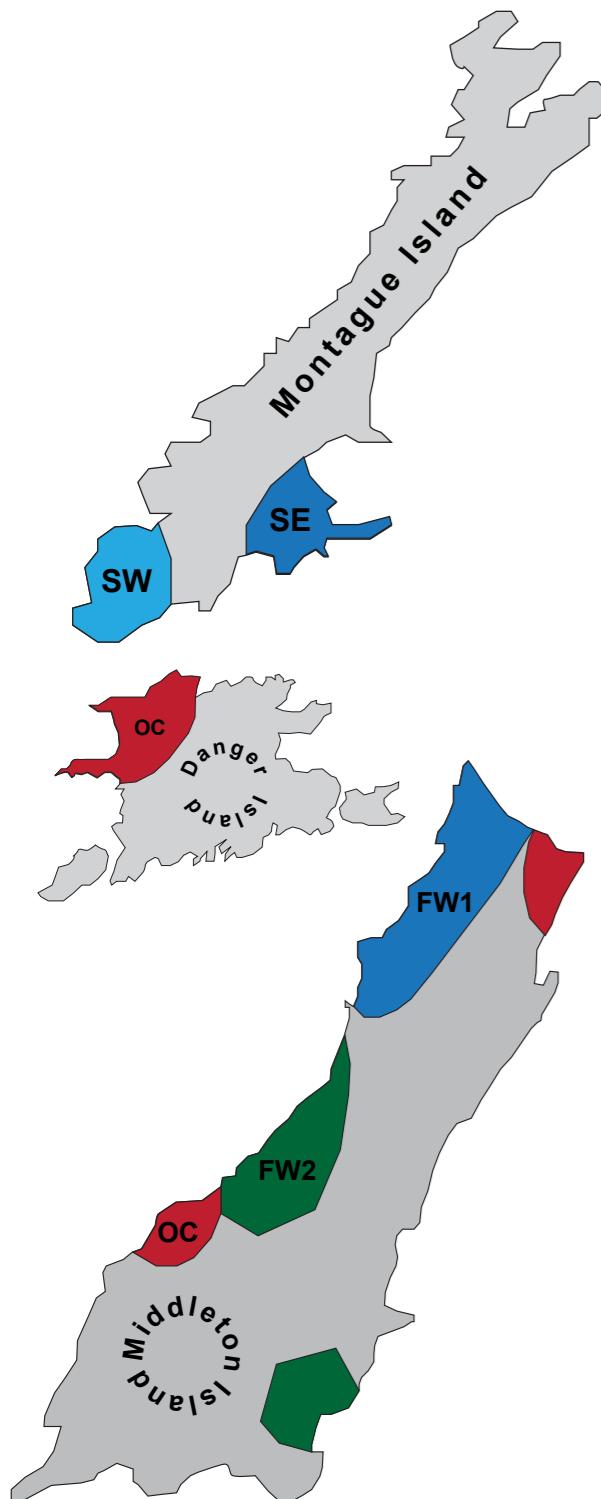
Genome Position (Mb)



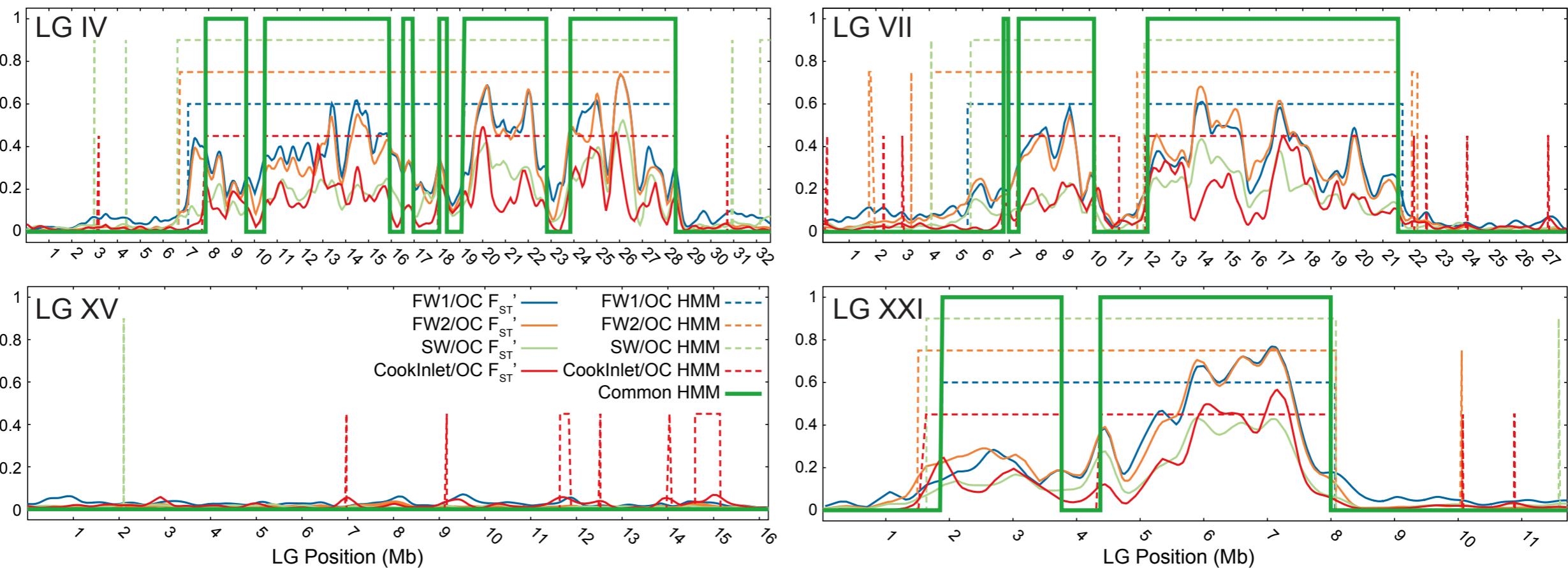








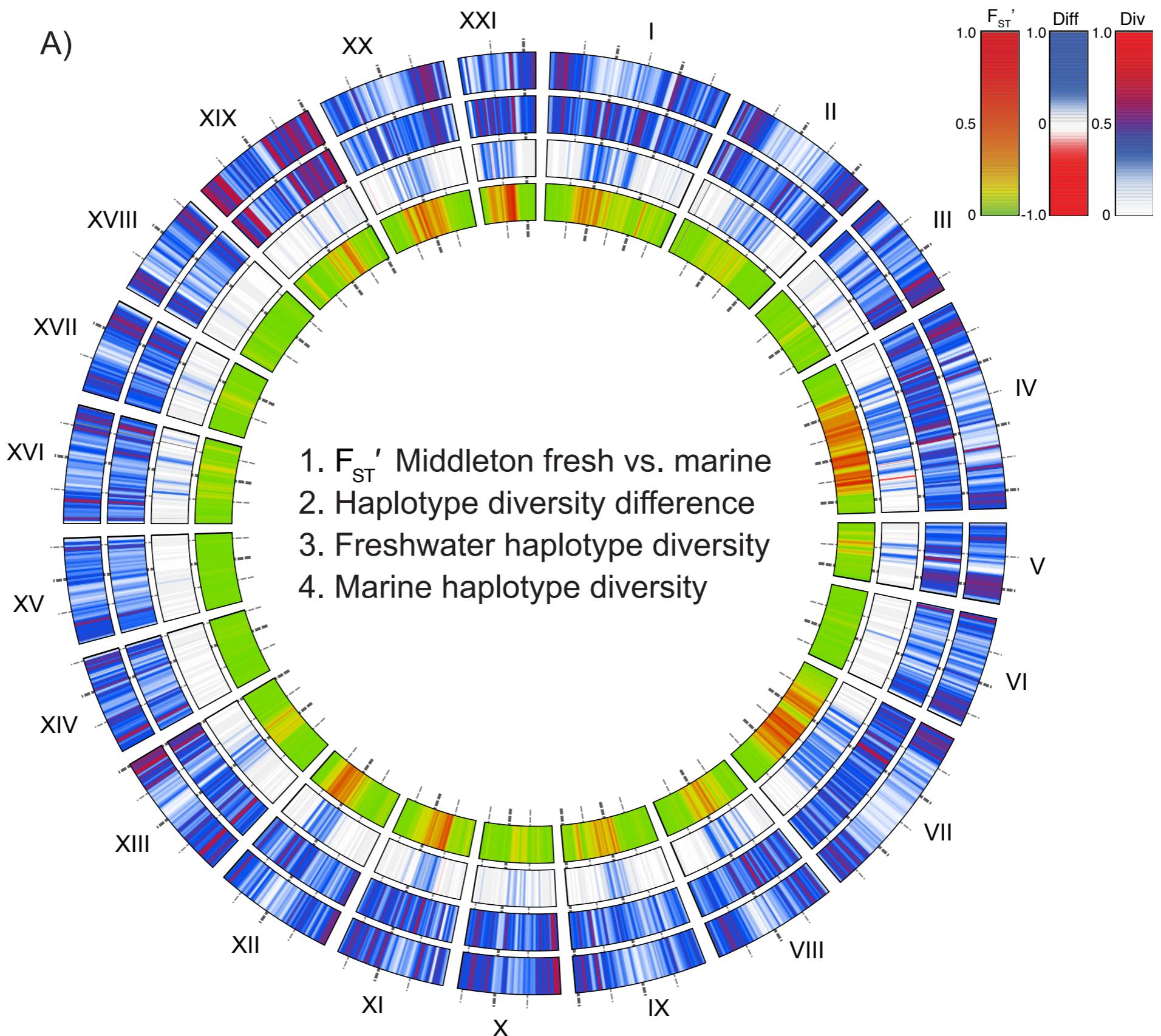
Significant overlap of divergence among populations



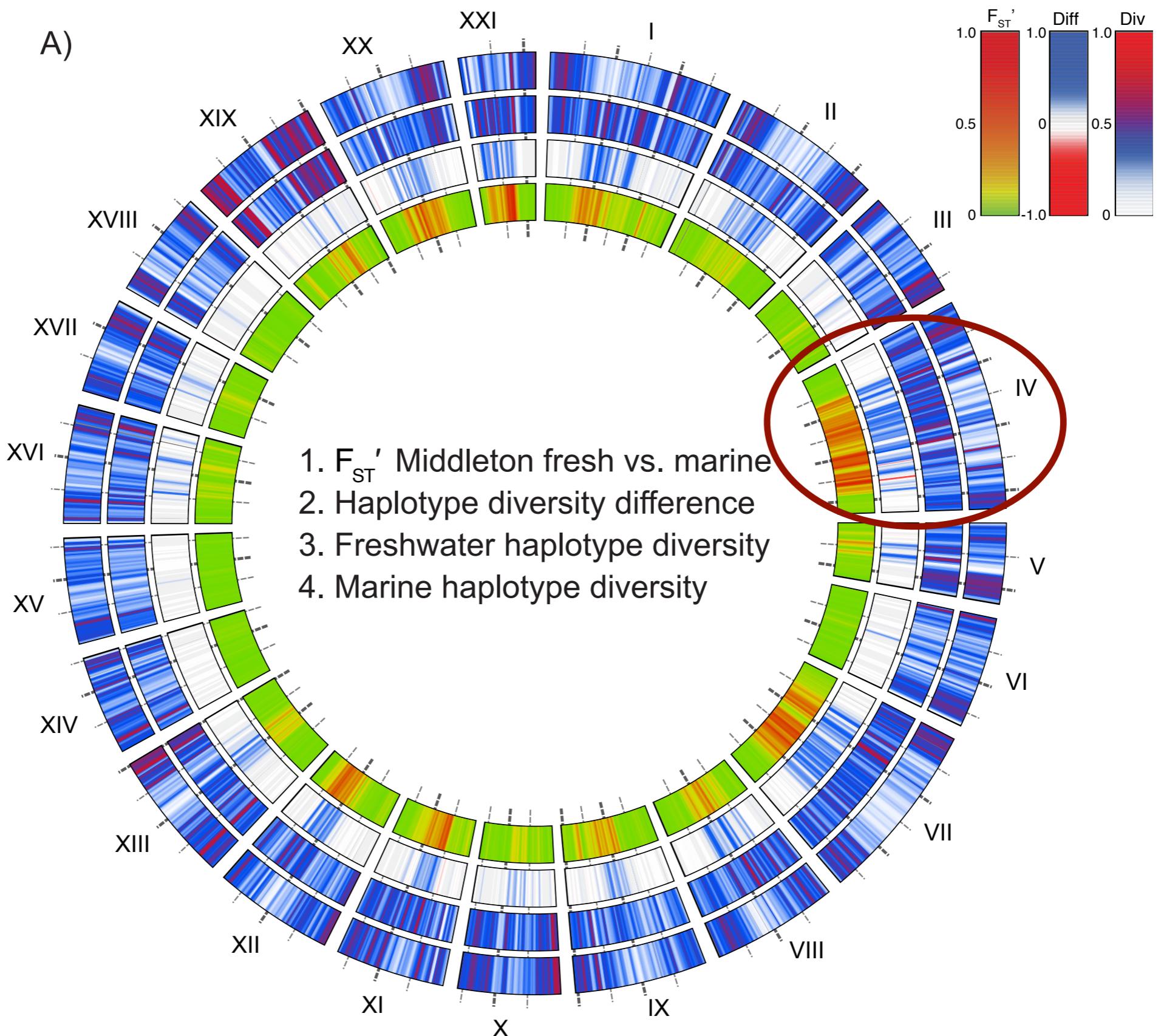
Significant overlap of divergence among populations

	MiFW1	MiFW2	MoSW	Cook Inlet
MiFW1	23.26%	22.12%	20.07%	15.15%
MiFW2	73.27%	29.06%	22.01%	16.27%
MoSW	72.94%	70.15%	24.33%	15.95%
Cook Inlet	56.65%	51.79%	59.02%	18.64%

Relative haplotype diversity is almost always elevated in marine populations in regions of divergence



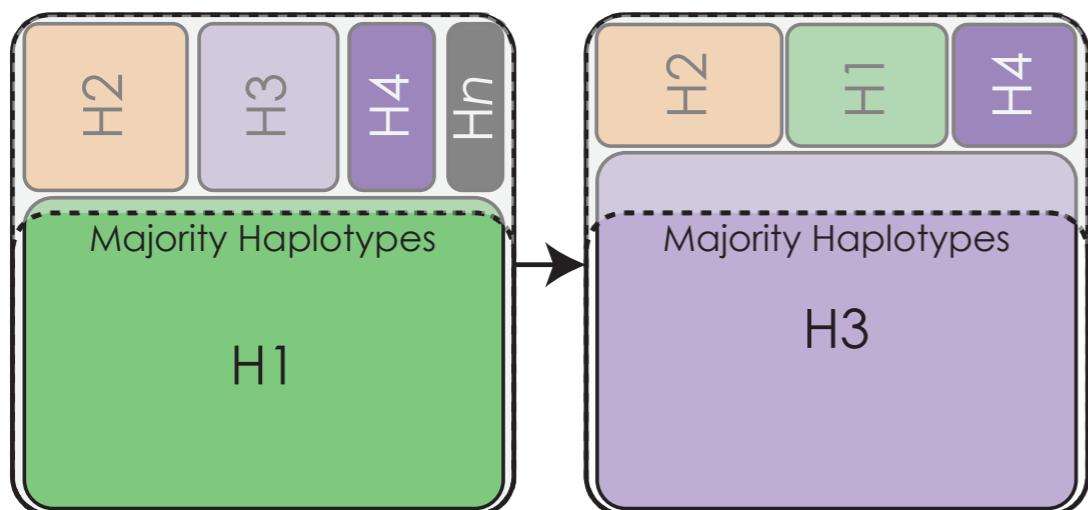
Relative haplotype diversity is almost always elevated in marine populations in regions of divergence



Regions of divergence contain alternative haplotypes

A)

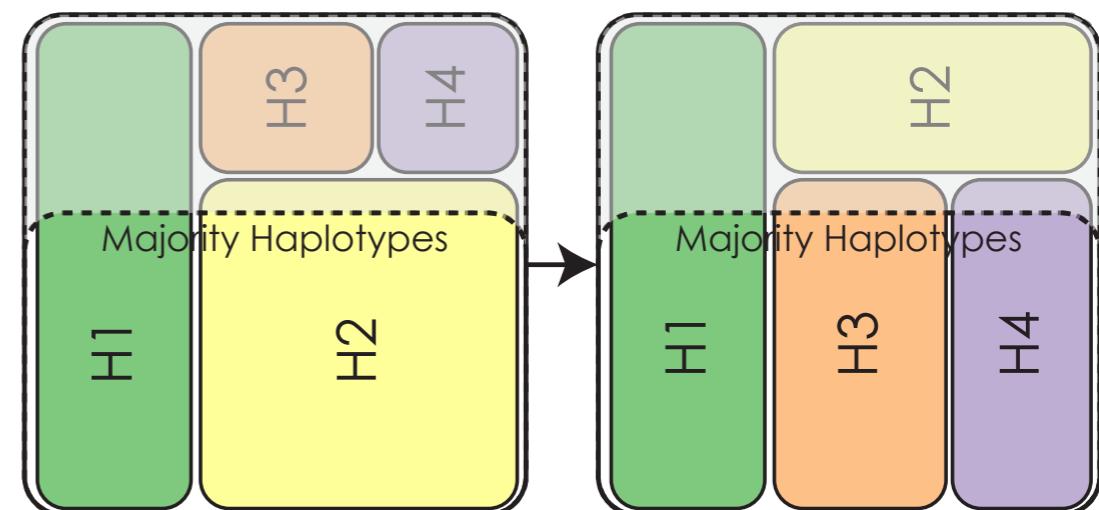
Case 1



Marine

Freshwater

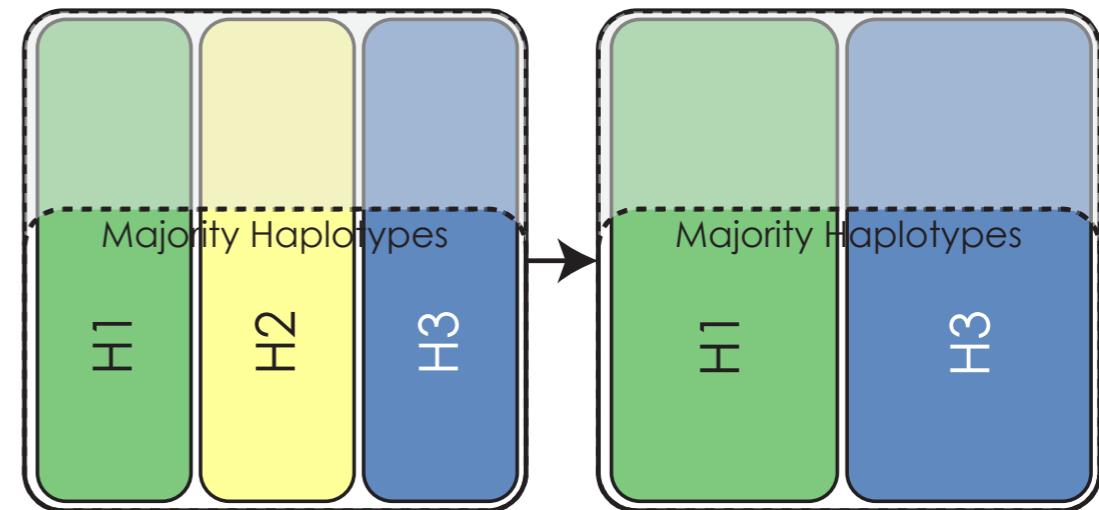
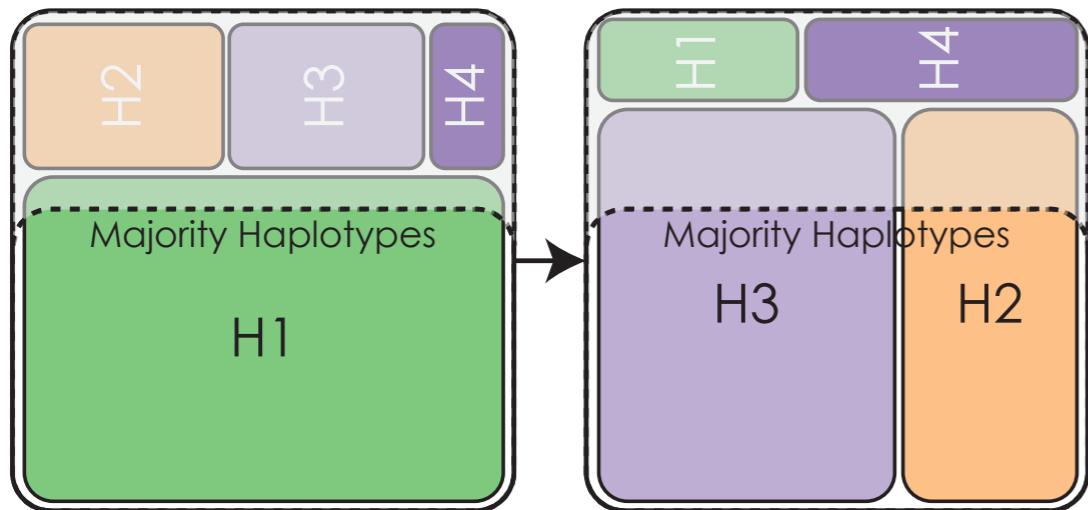
Case 2



Marine

Freshwater

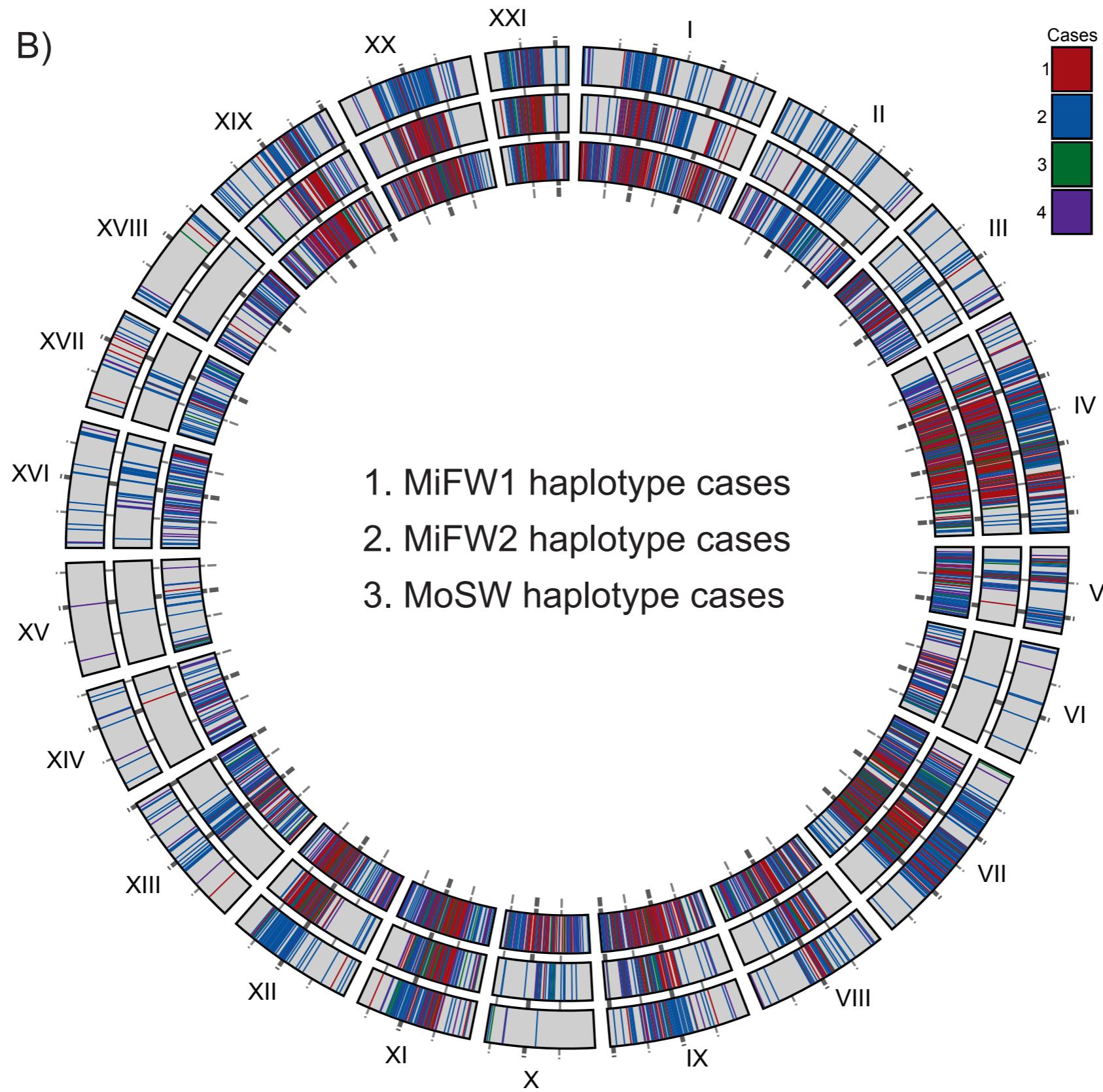
Case 3



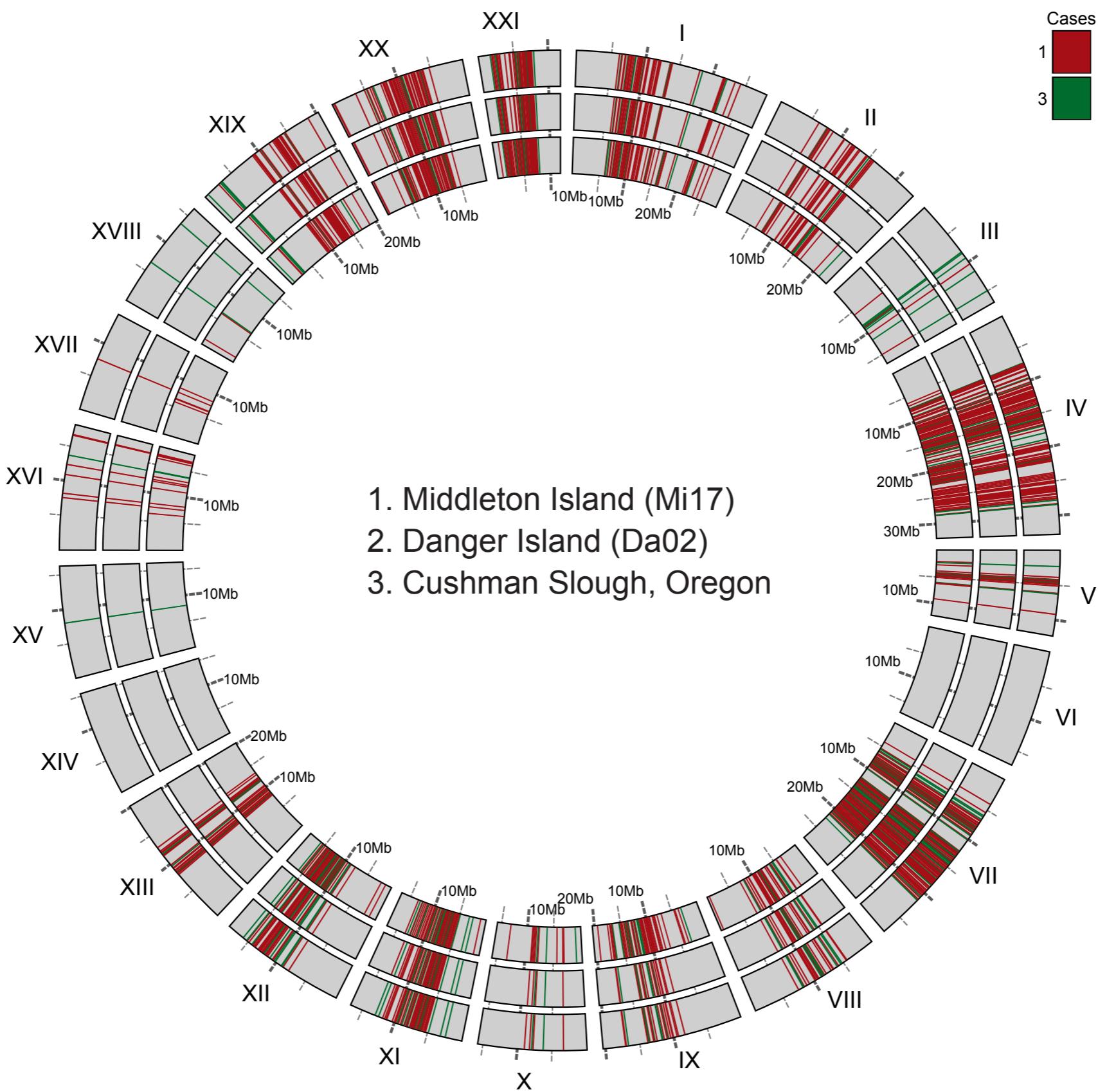
Marine

Freshwater

Regions of divergence contain alternative haplotypes

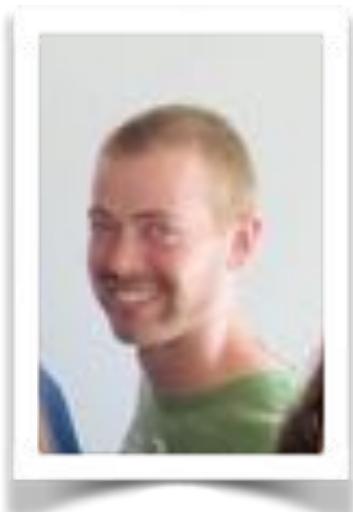


Marine stickleback carry haplotypes at low frequency
that are the majority alleles in freshwater populations



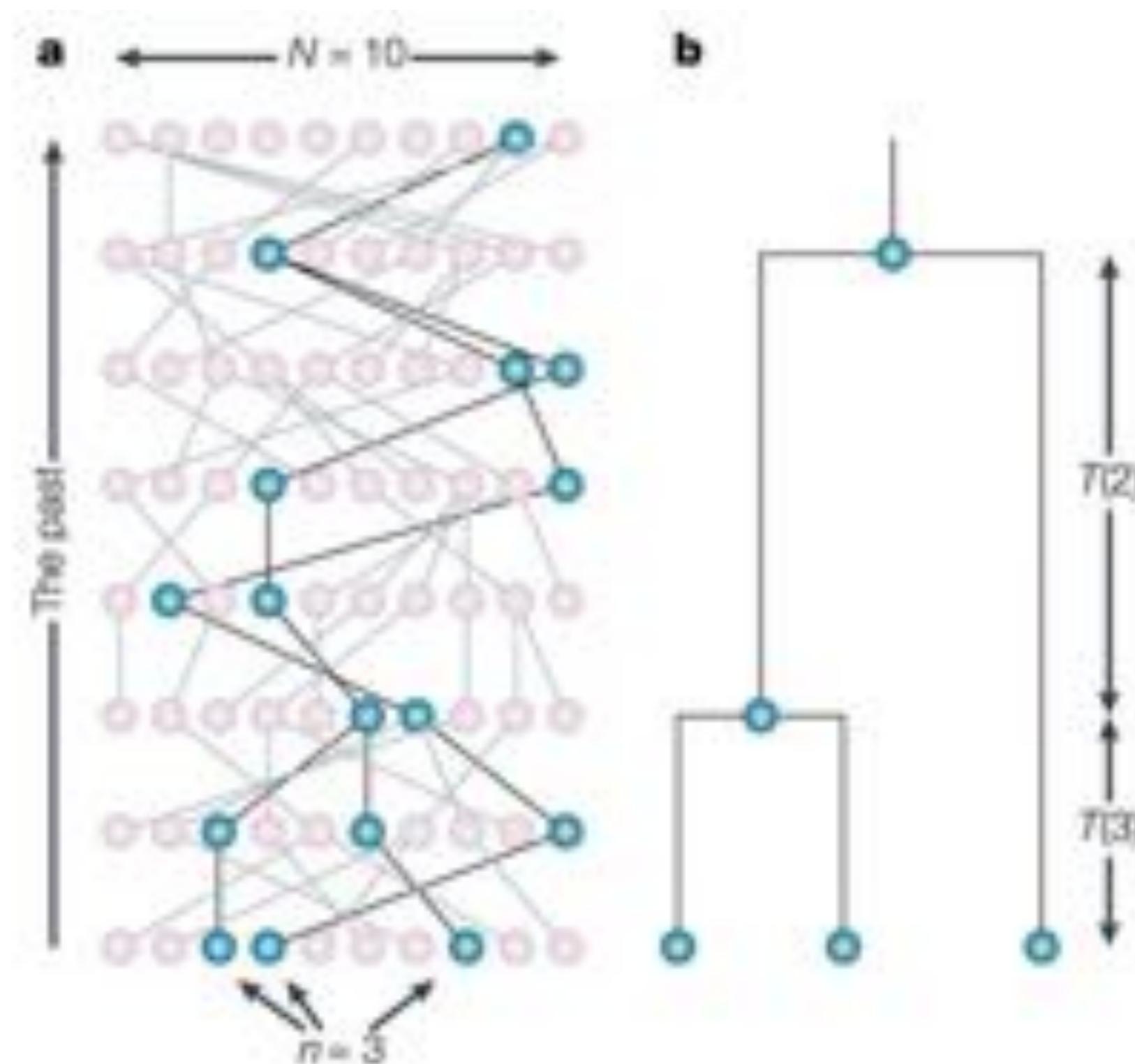
- Stickleback on uplifted islands evolve in decades
- ~25% of the genome in each population is affected by local adaptation, the majority of which is parallel across populations
- Many haplotypes are *habitat specific* in terms of majority allele frequency
- Freshwater-dominant haplotypes can be found at low frequency in marine populations across the genome
- The stickleback genome is crafted by migration-selection balance, and probably has been for millennia
- What is the molecular evolutionary history of the haplotypes involved in rapid local adaptation?

Coalescent analyses using RADseq



Nelson and Cresko. *Evolution Letters* (in press)

The coalescent in population genetics



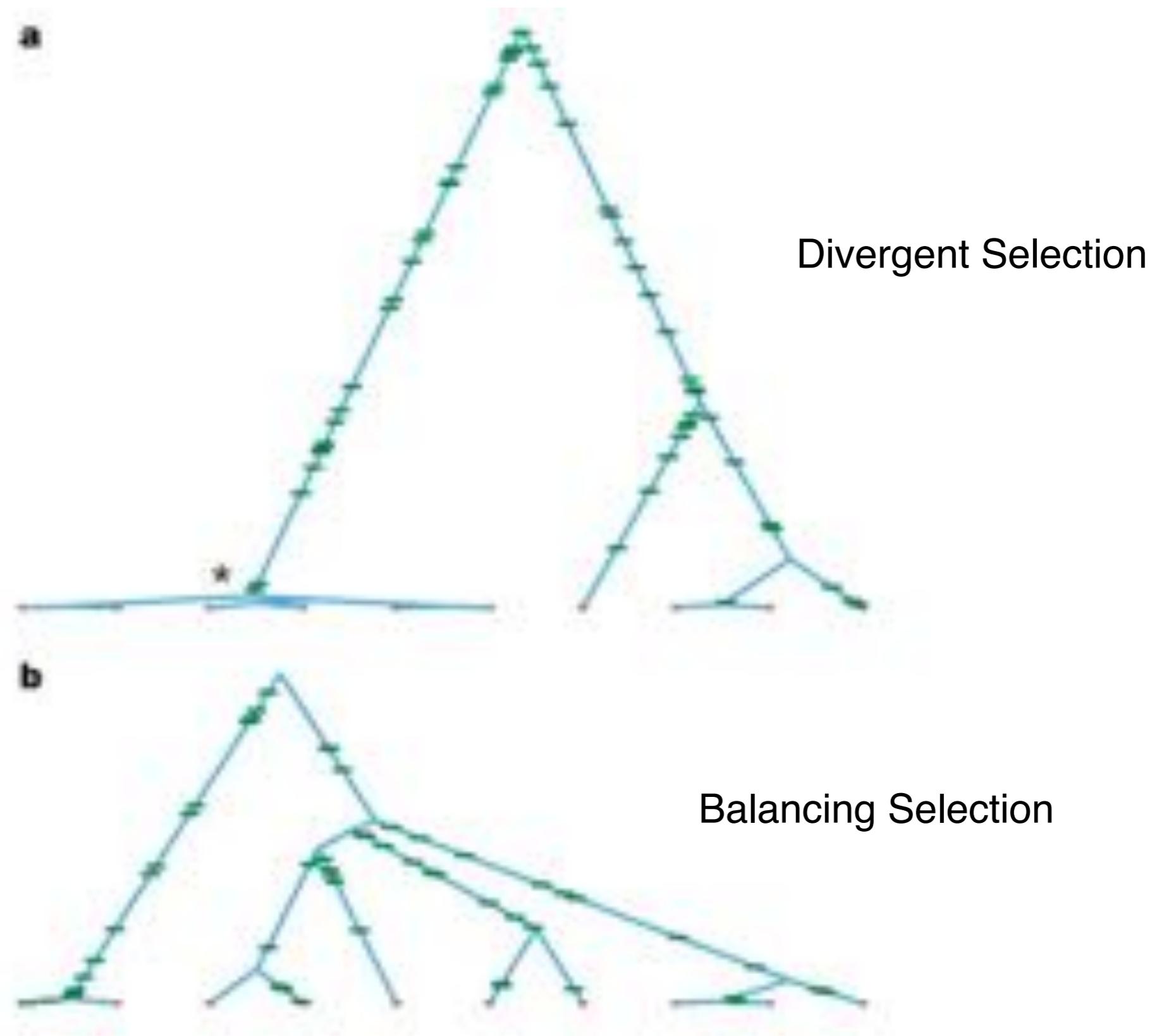
Noah A. Rosenberg & Magnus Nordborg
Nature Reviews | Genetics

Neutral coalescent expectations

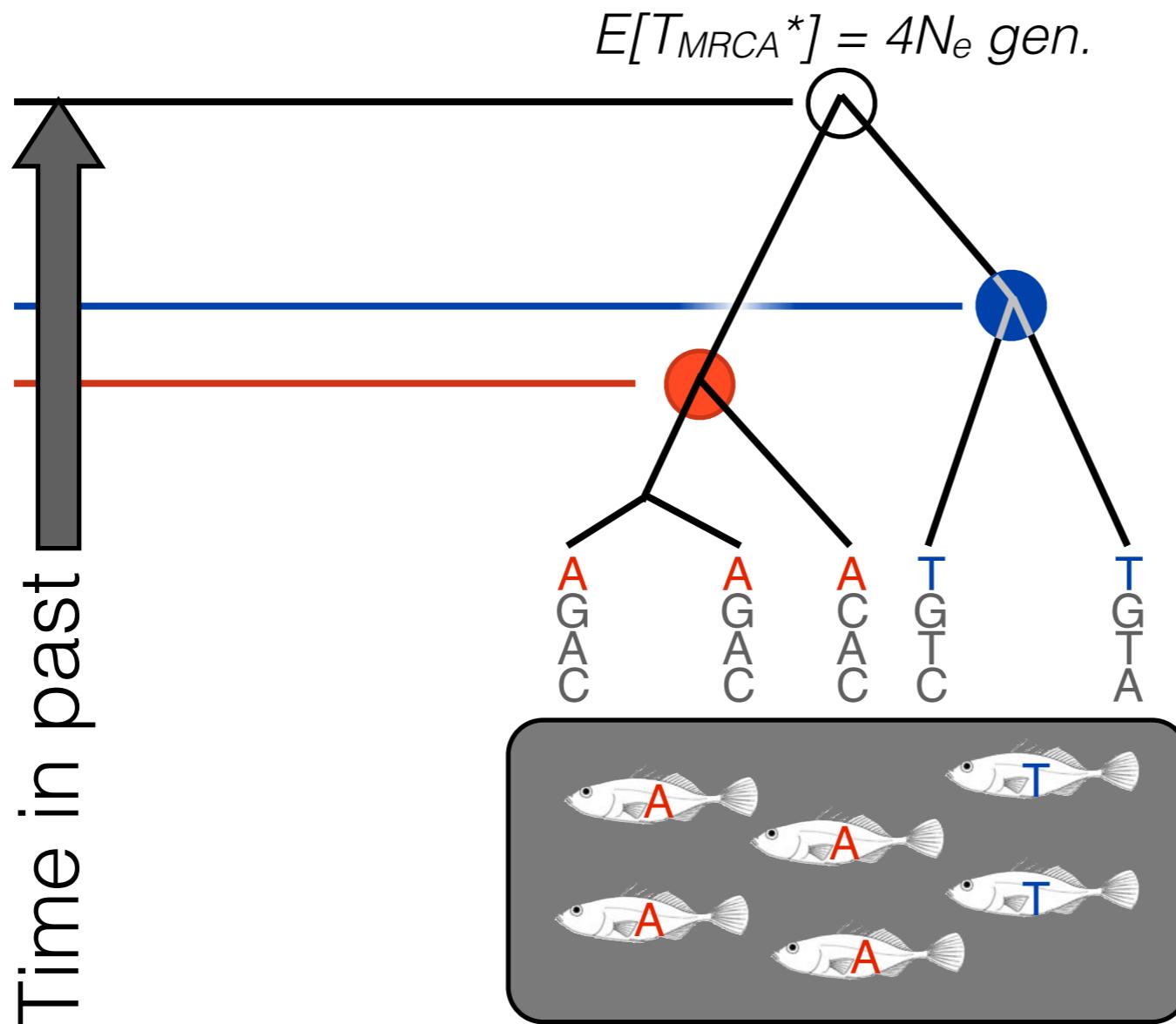


Nature Reviews | Genetics

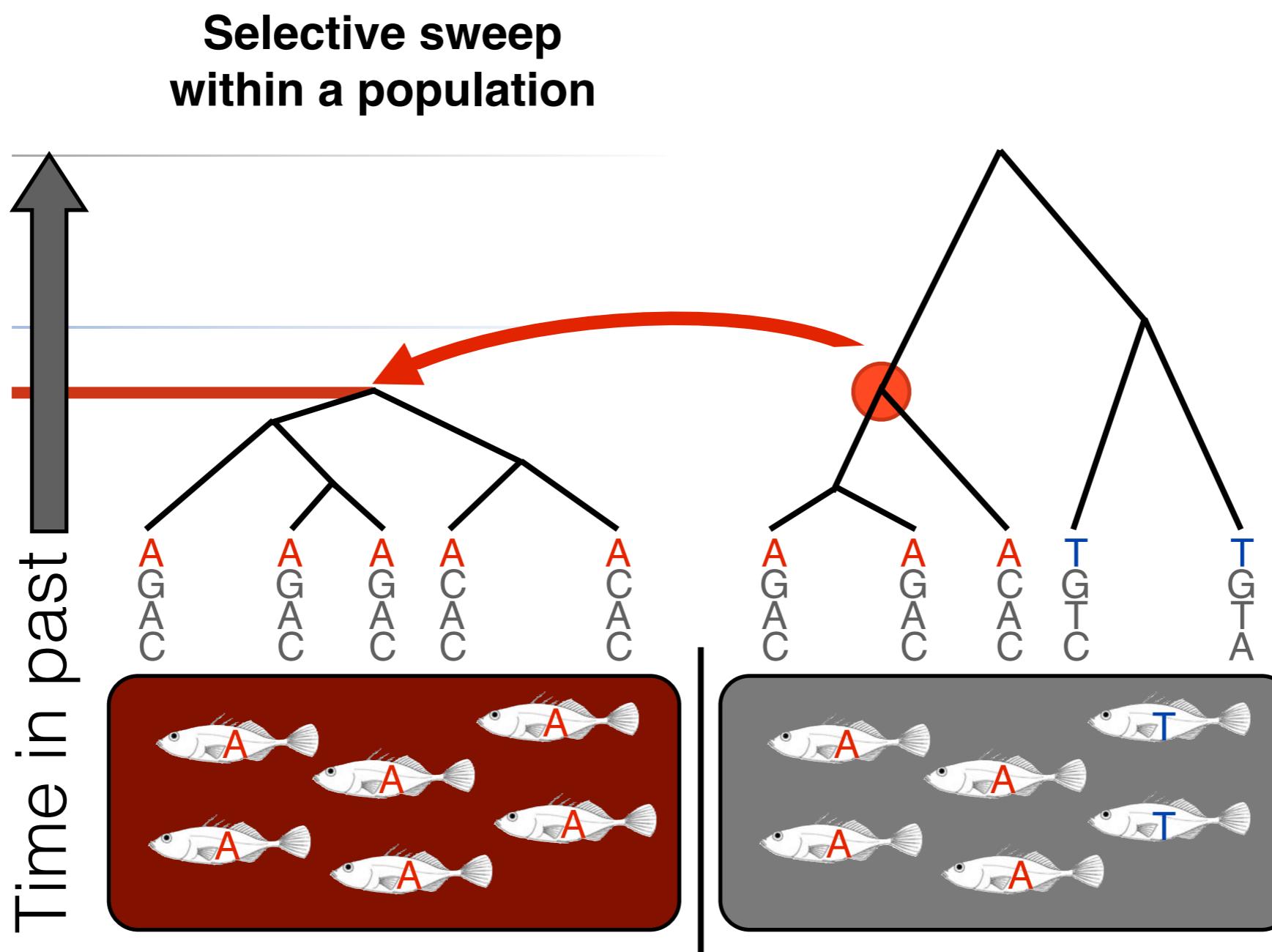
Natural selection and the coalescent



Coalescent analyses in stickleback



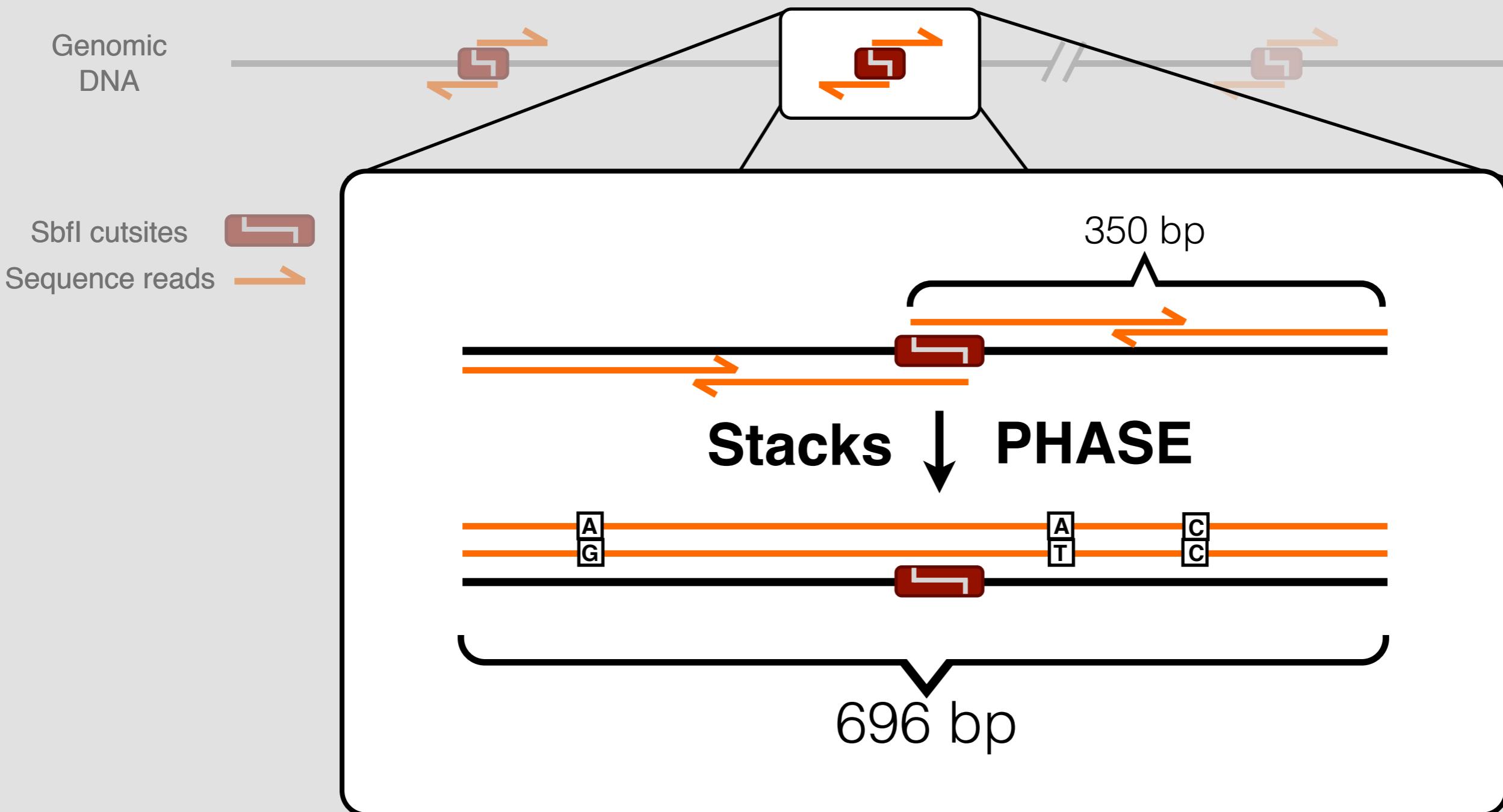
Coalescent analyses in stickleback



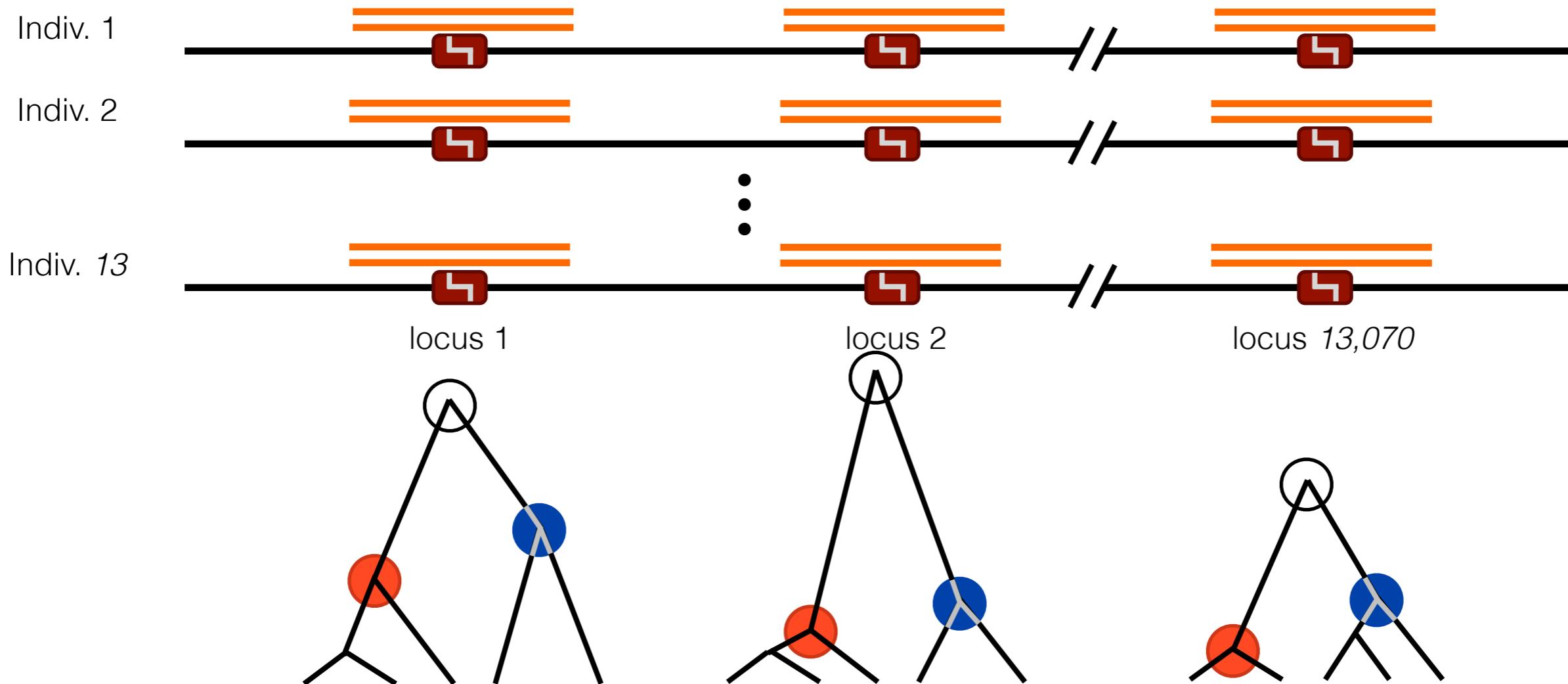
Coalescent analyses with RAD-seq



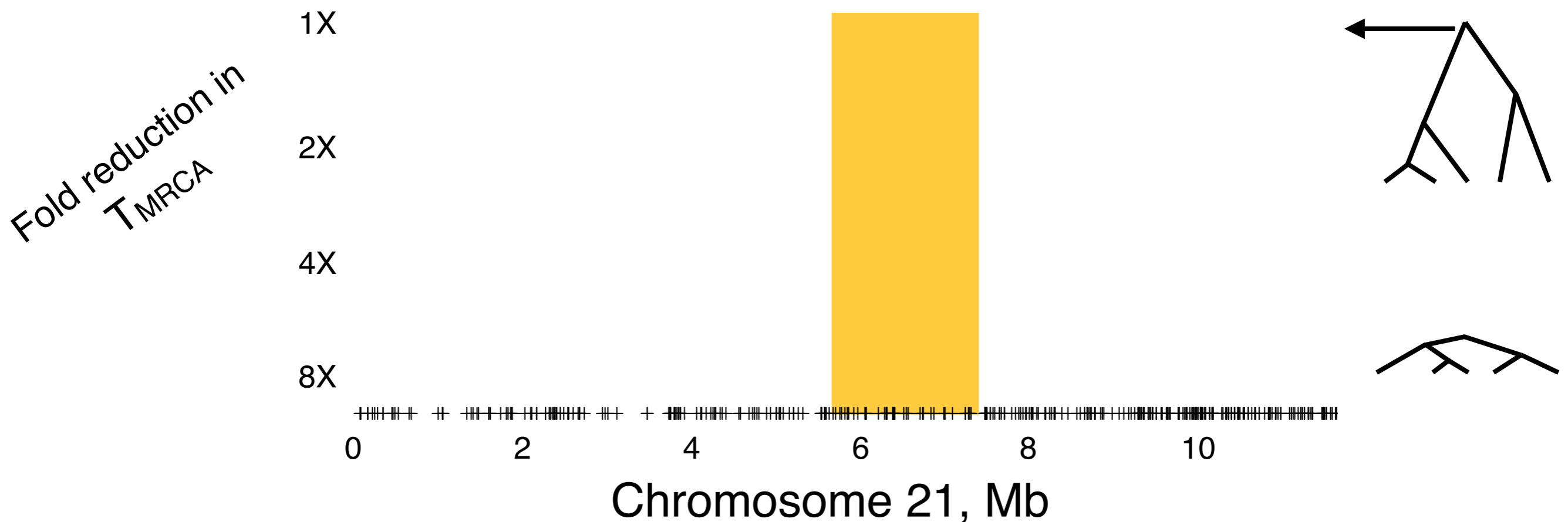
Coalescent analyses with RAD-seq



Coalescent analyses with RAD-seq

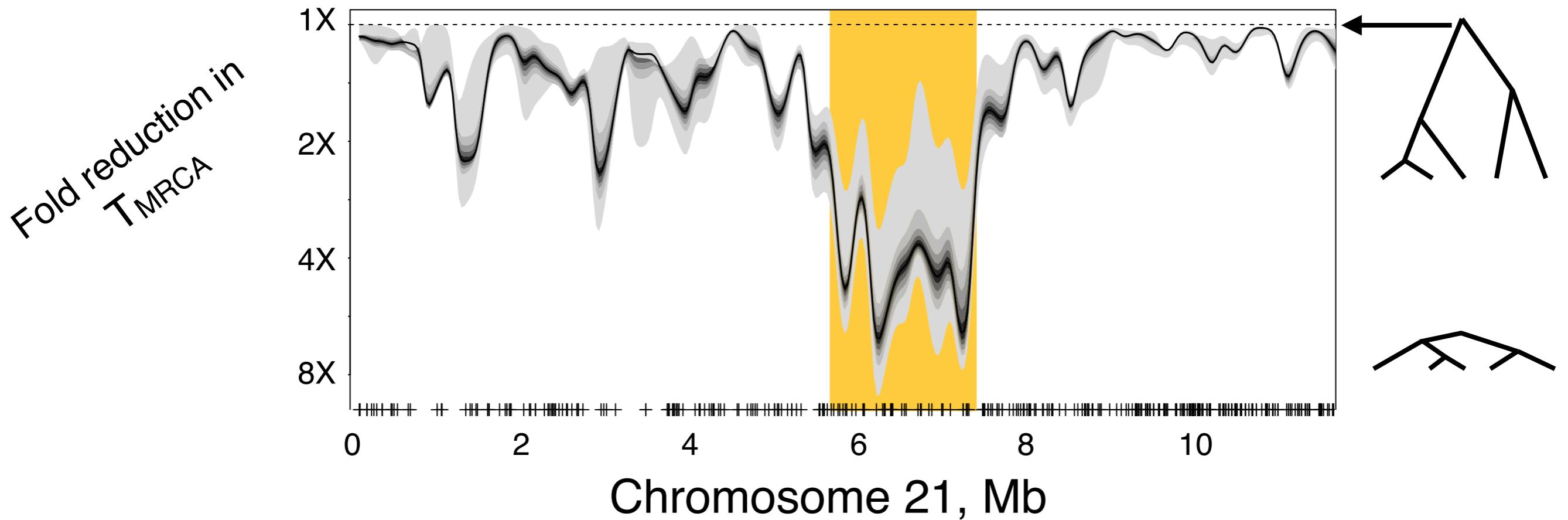


Selection in one population reduces coalescence time



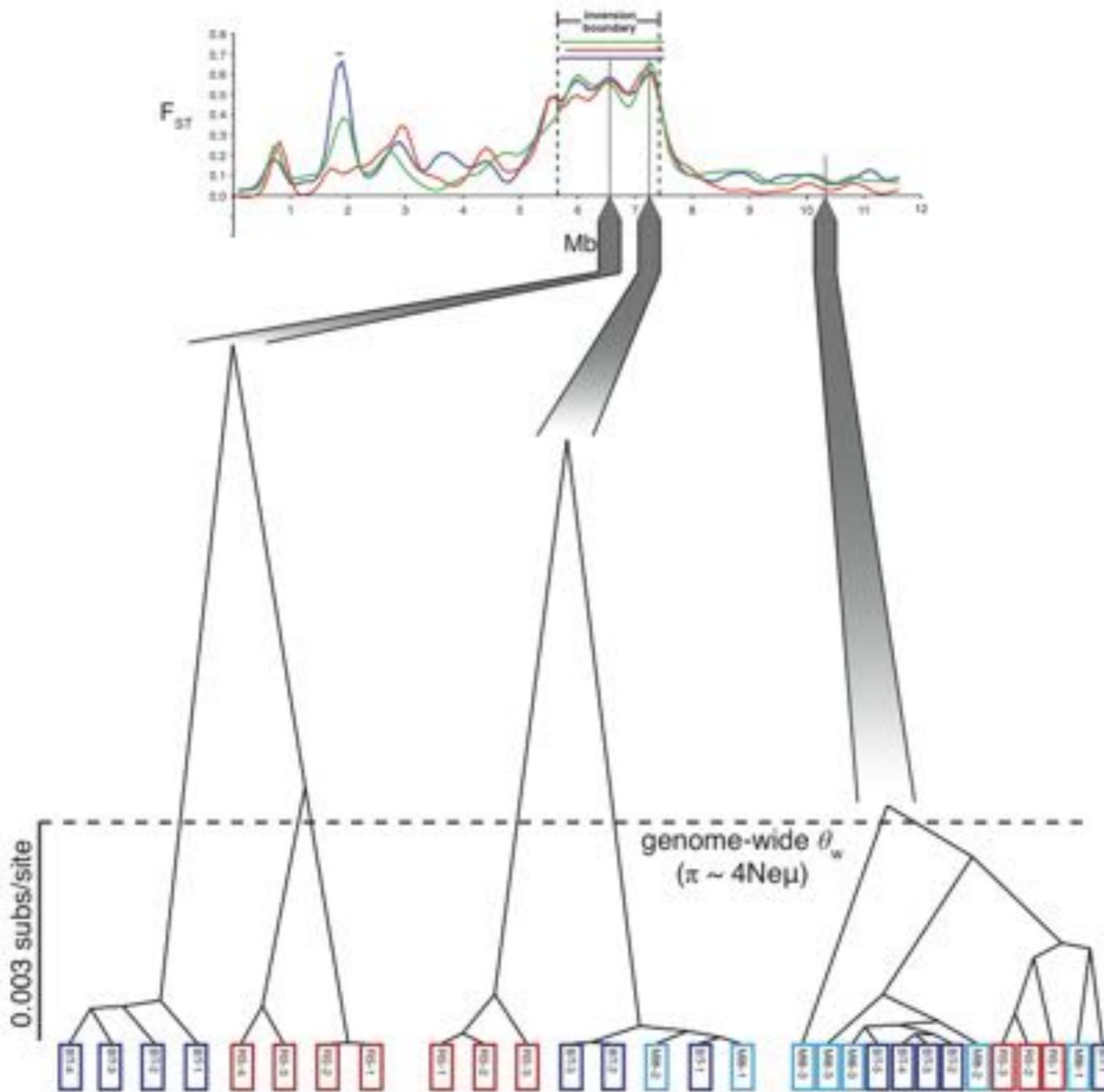
$$\log_2 \left(\frac{T_{\text{MRCA}} \text{FW}}{T_{\text{MRCA}} \text{ALL}} \right)$$

Selection in one population reduces coalescence time

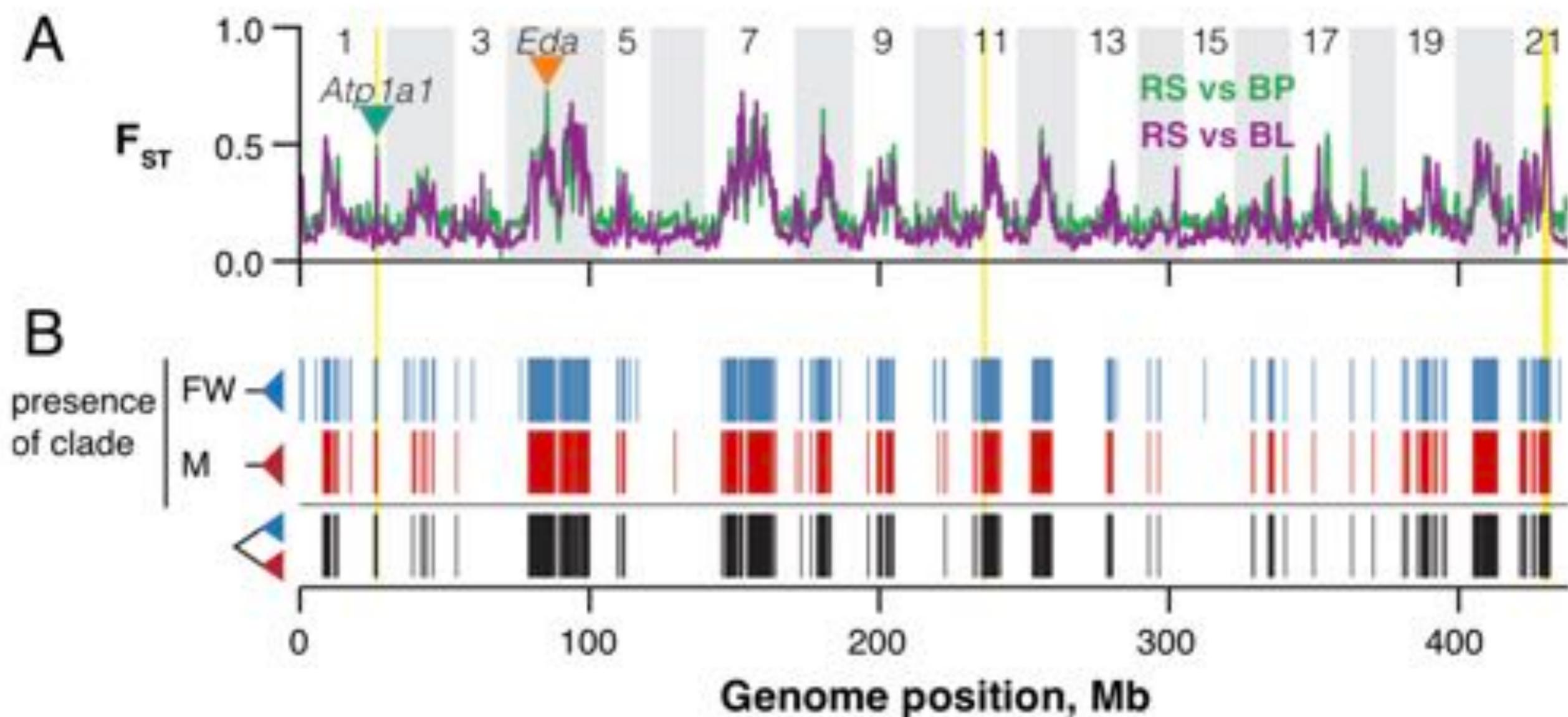


$$\log_2 \left(\frac{T_{\text{MRCA}} \text{FW}}{T_{\text{MRCA}} \text{ALL}} \right)$$

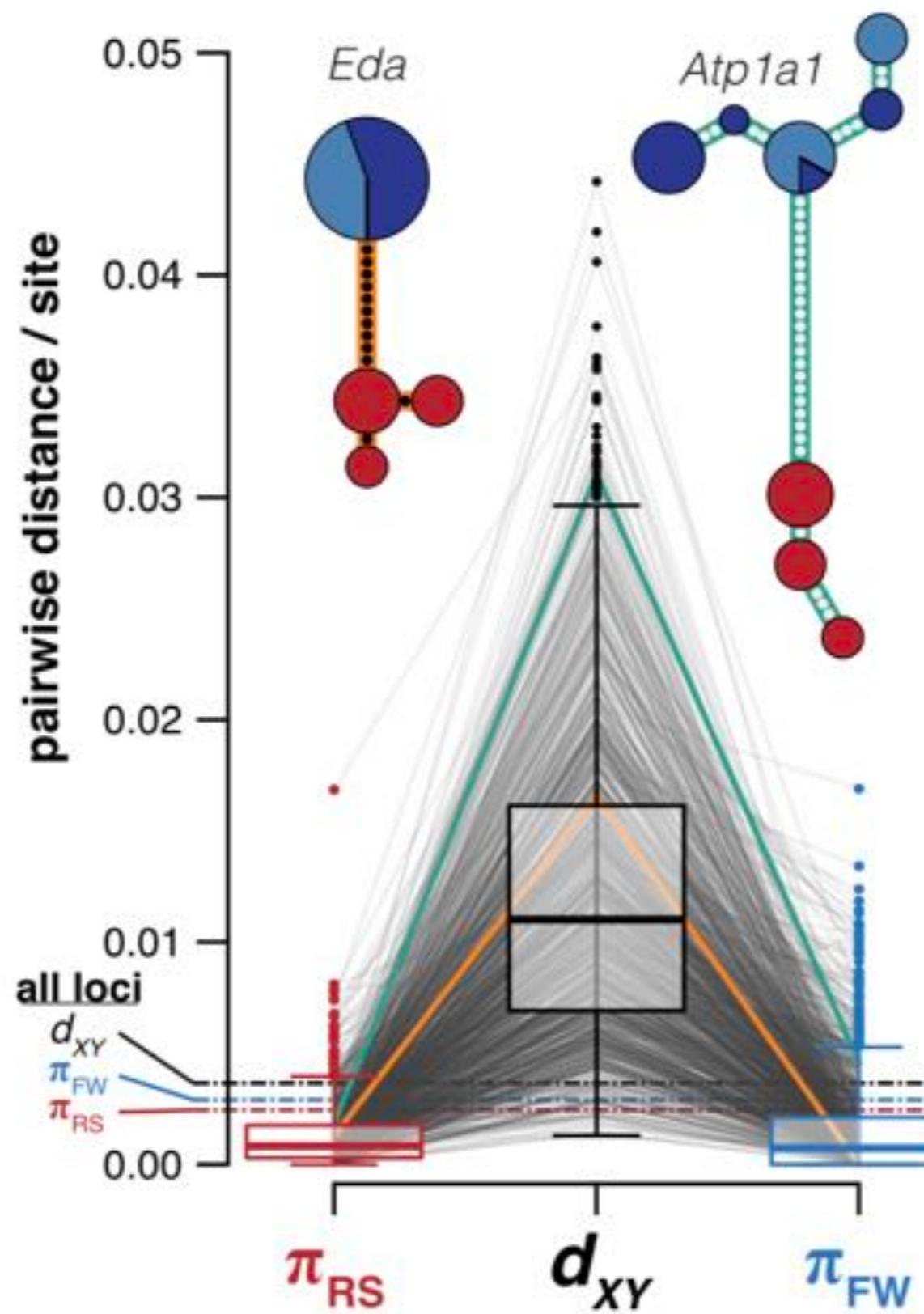
However, across populations the coalescence time can increases significantly in a genomic region



Reciprocal monophyly across different habitats

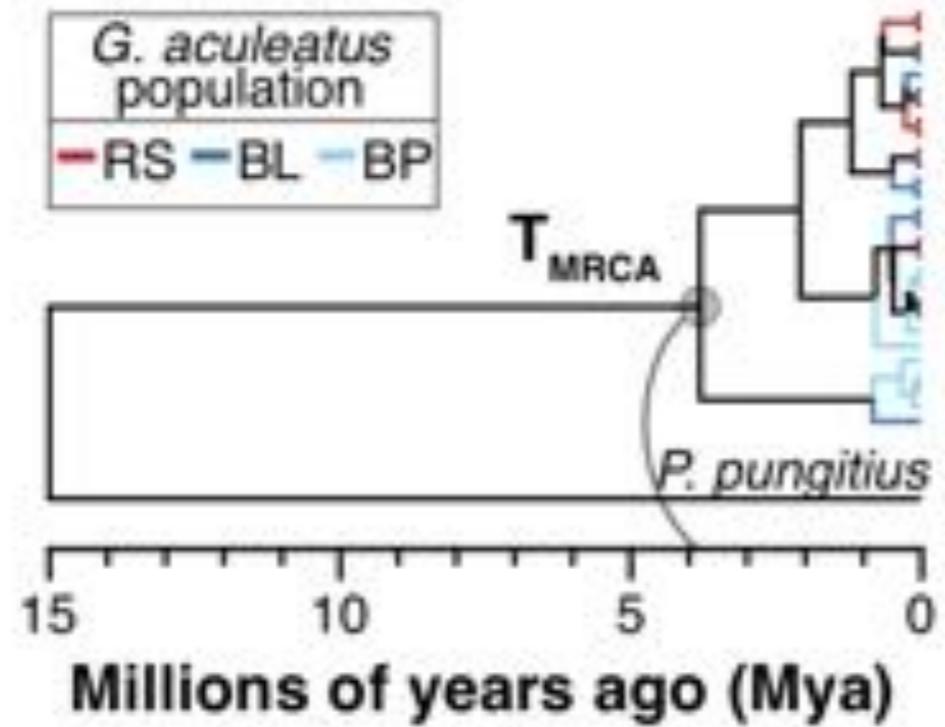


Increased absolute divergence between habitats

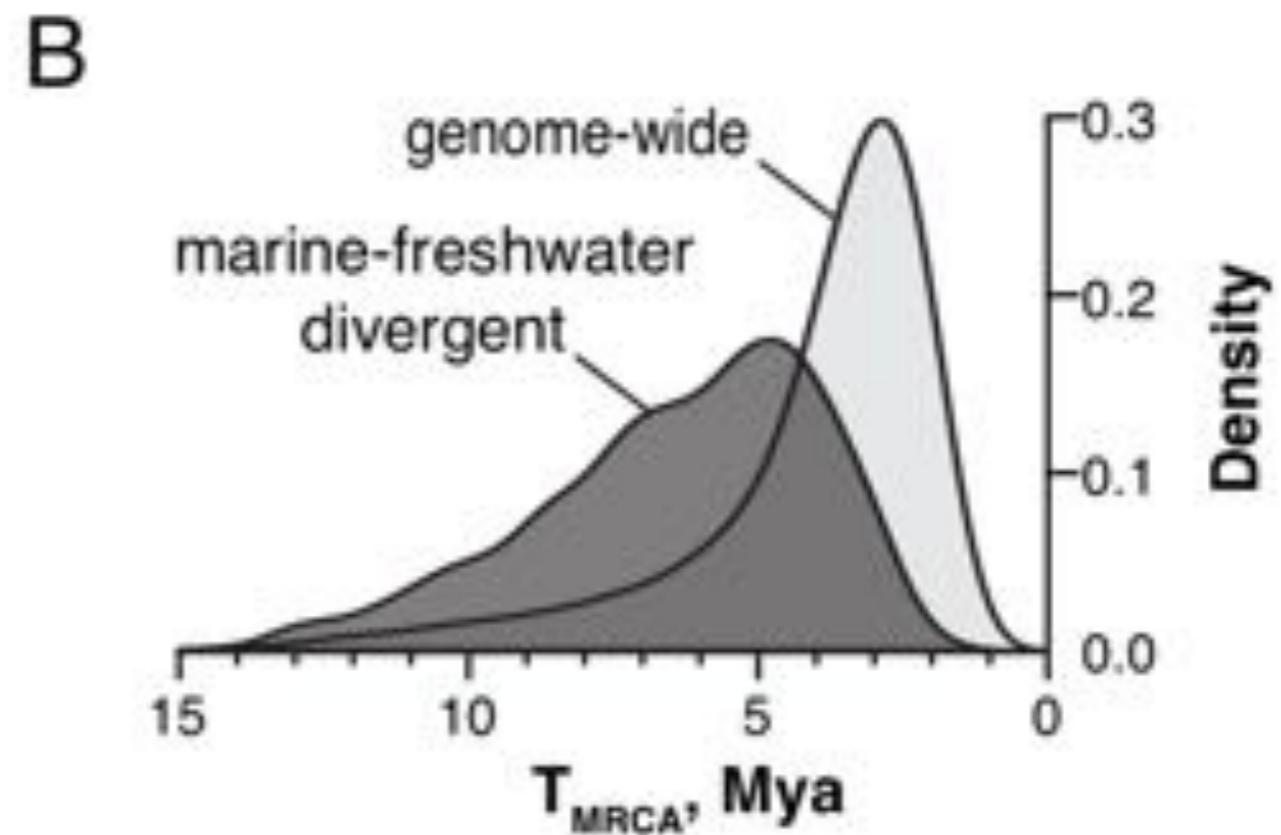
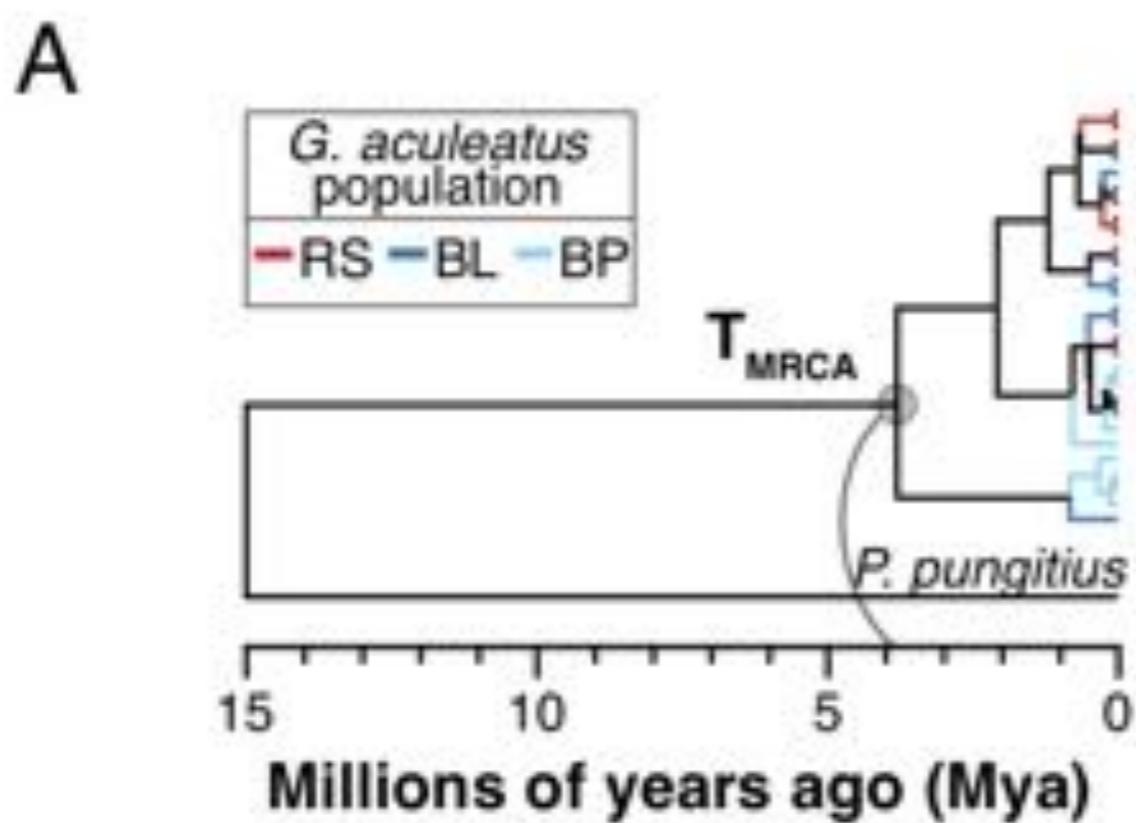


Habitat associated haplotypes evolved millions of years ago

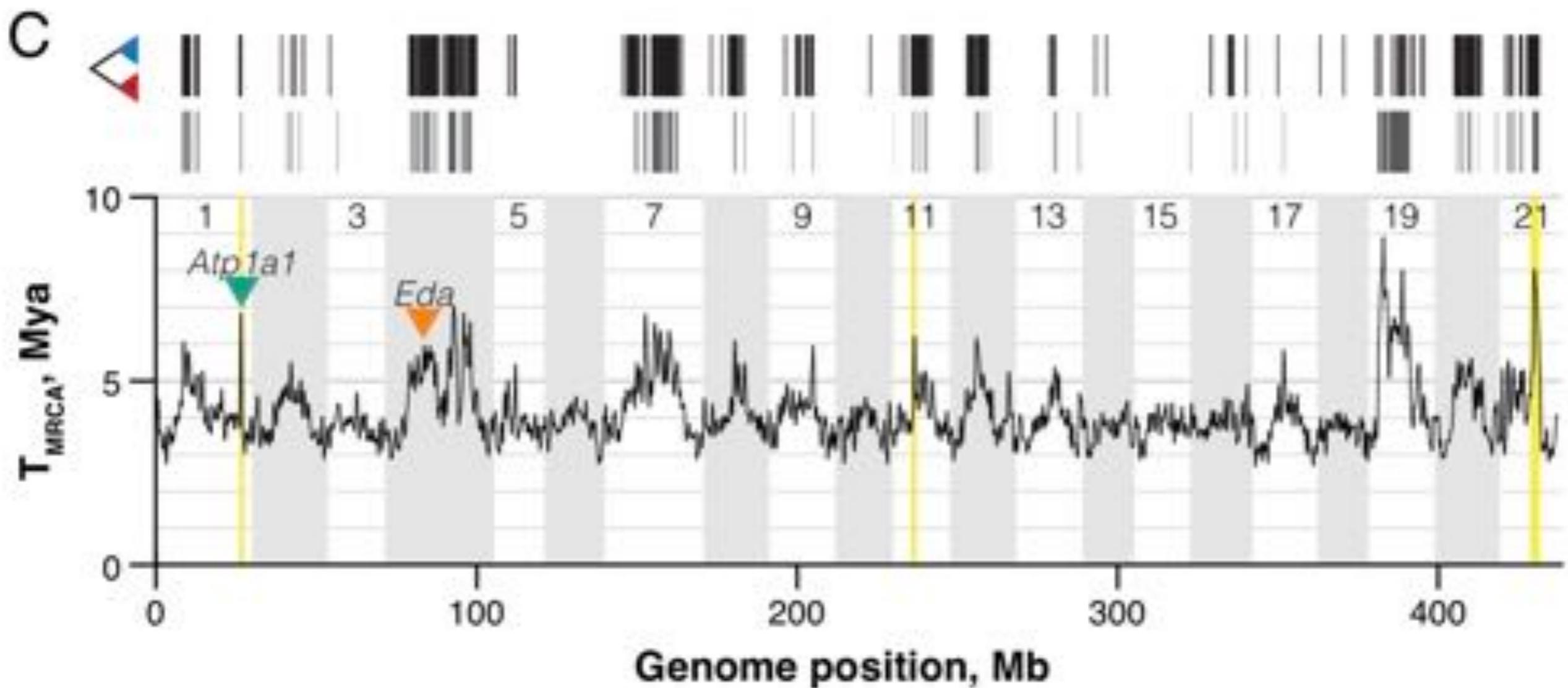
A



Habitat associated haplotypes evolved over millions of years ago



Habitat associated haplotypes evolved over millions of years

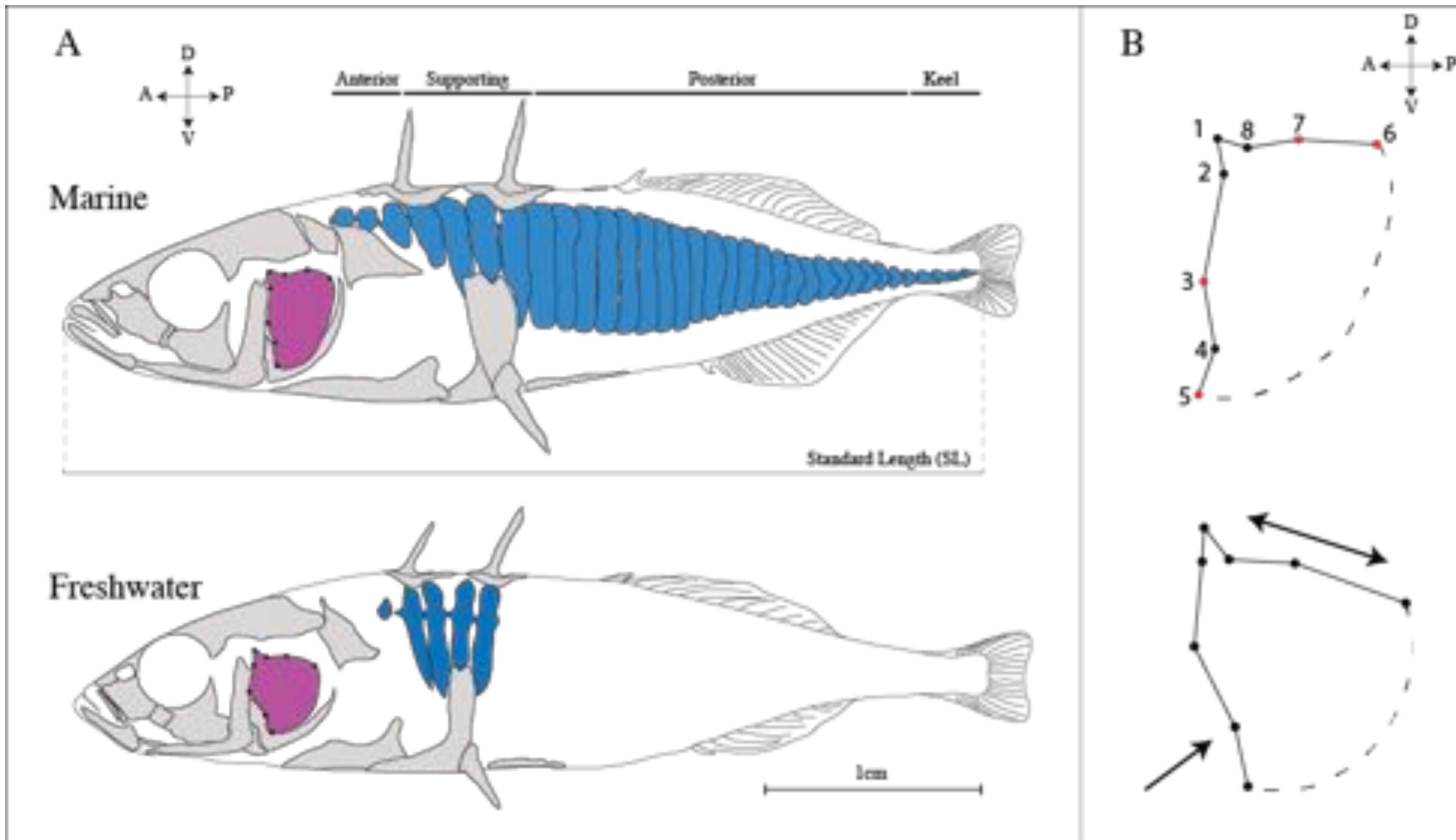


How do the genomic patterns of divergence link to phenotypic diversification?

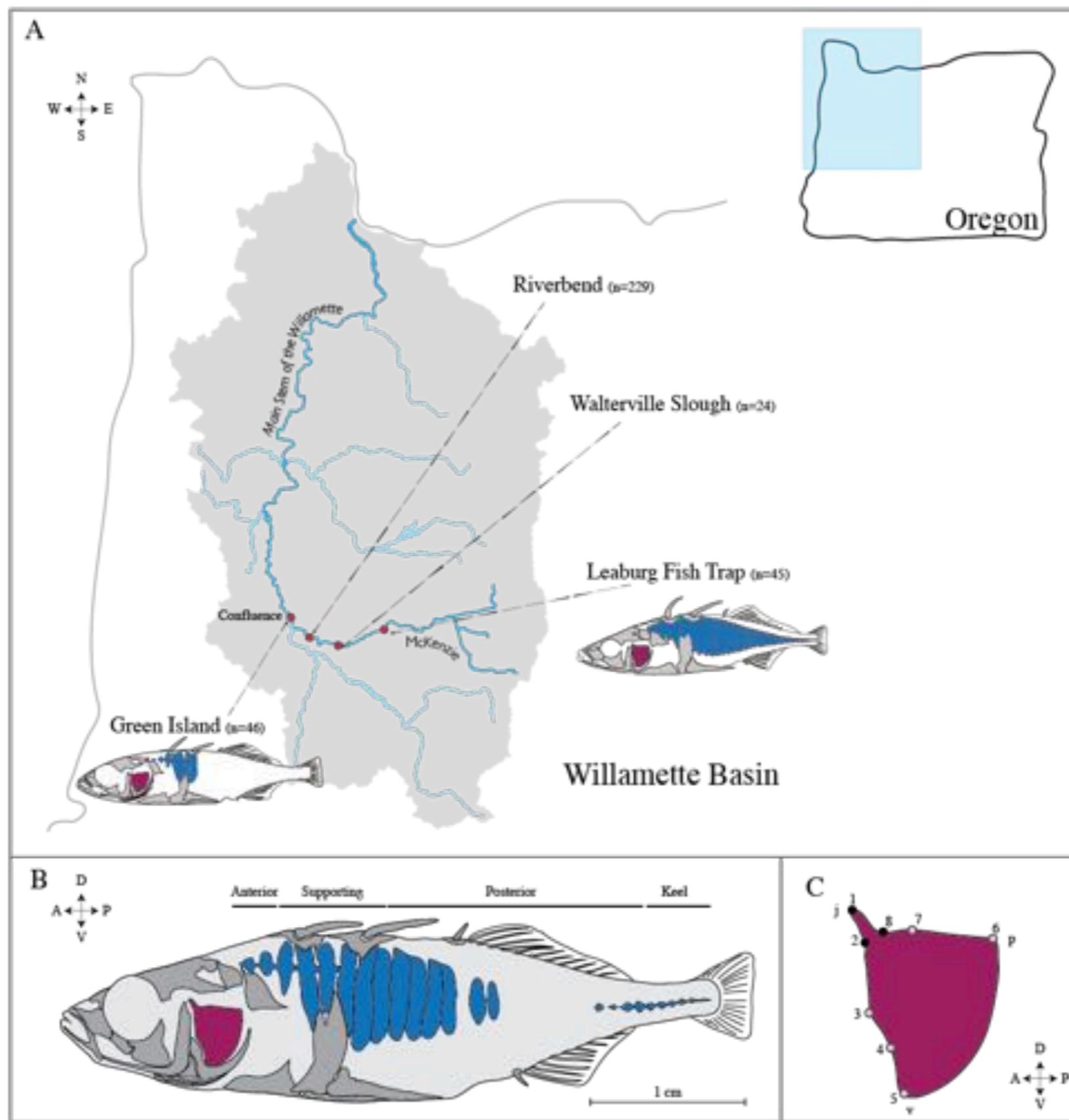


Kristin Alligood

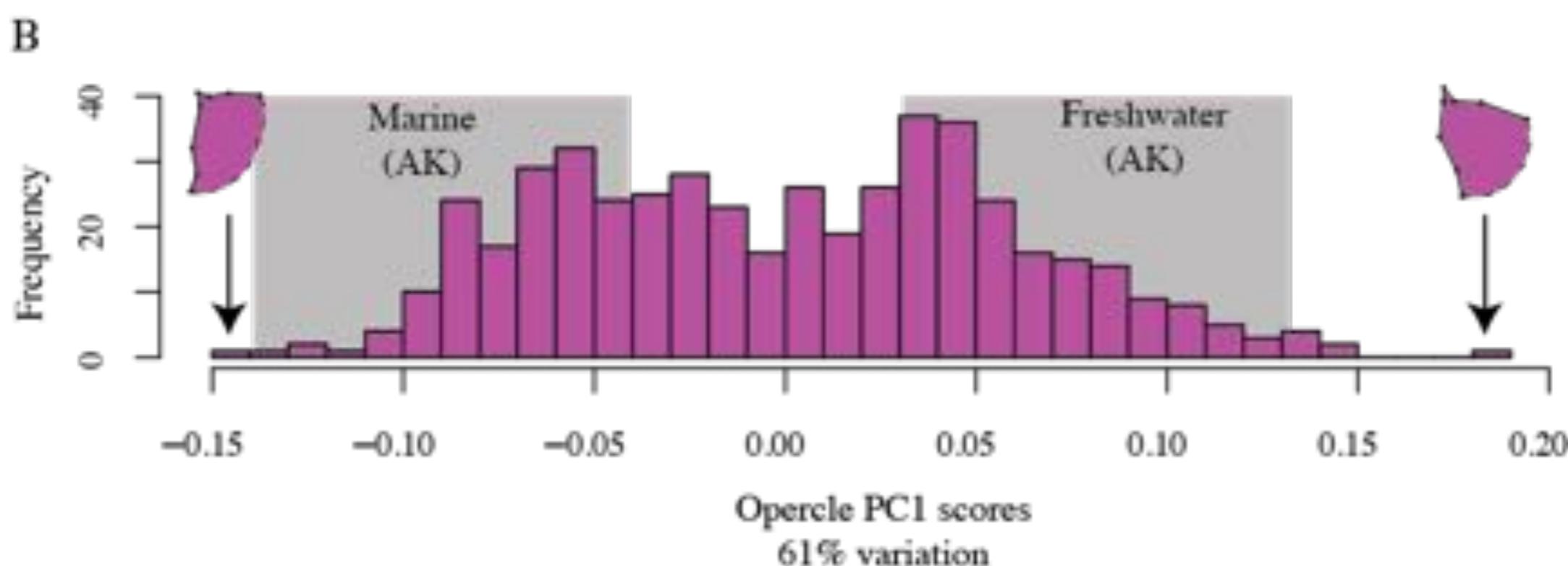
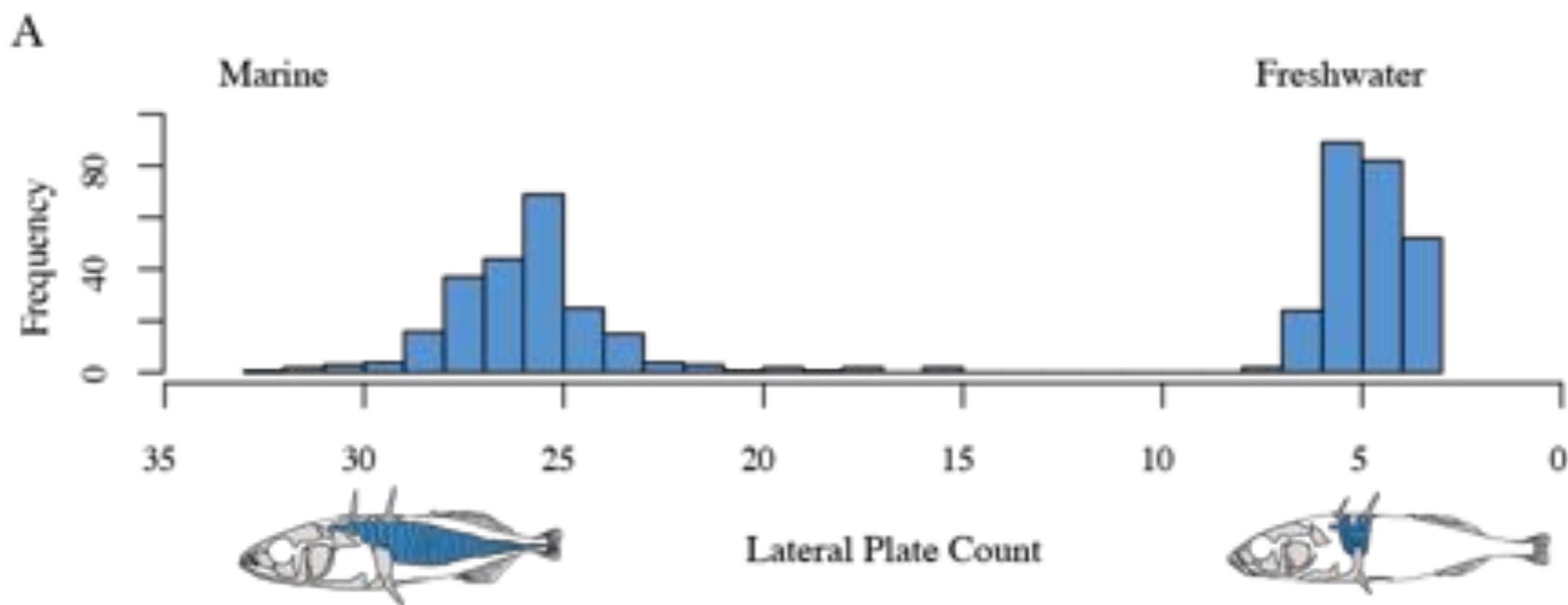
Lateral plate and opercle shape co-vary in the wild



An interesting stickleback hybrid population in Oregon

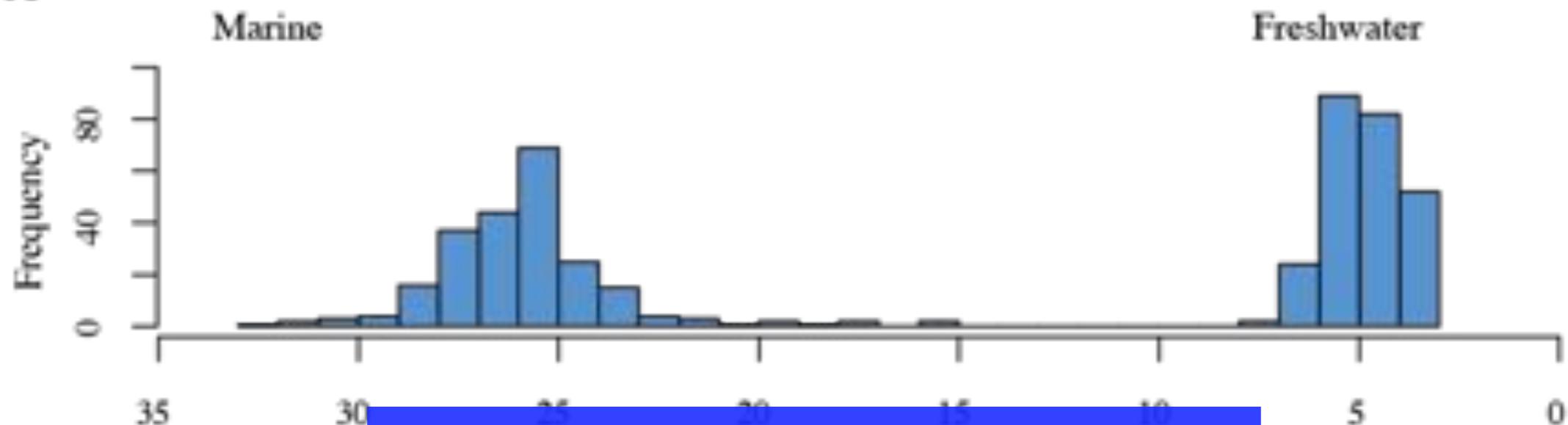


An interesting stickleback hybrid population in Oregon



An interesting stickleback hybrid population in Oregon

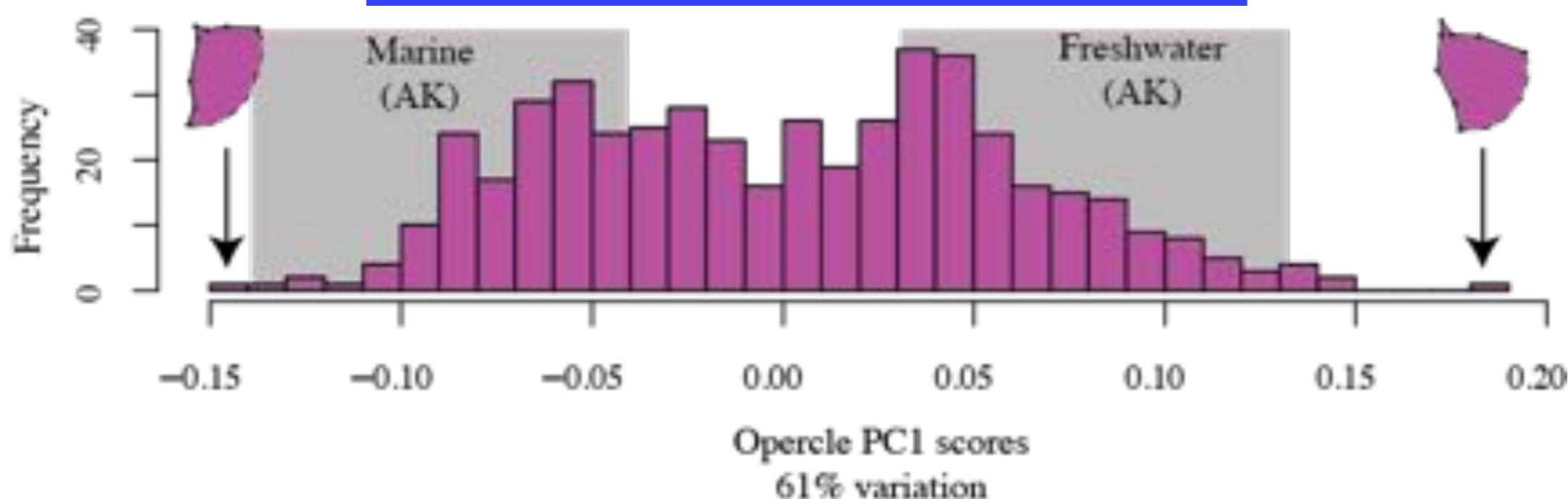
A



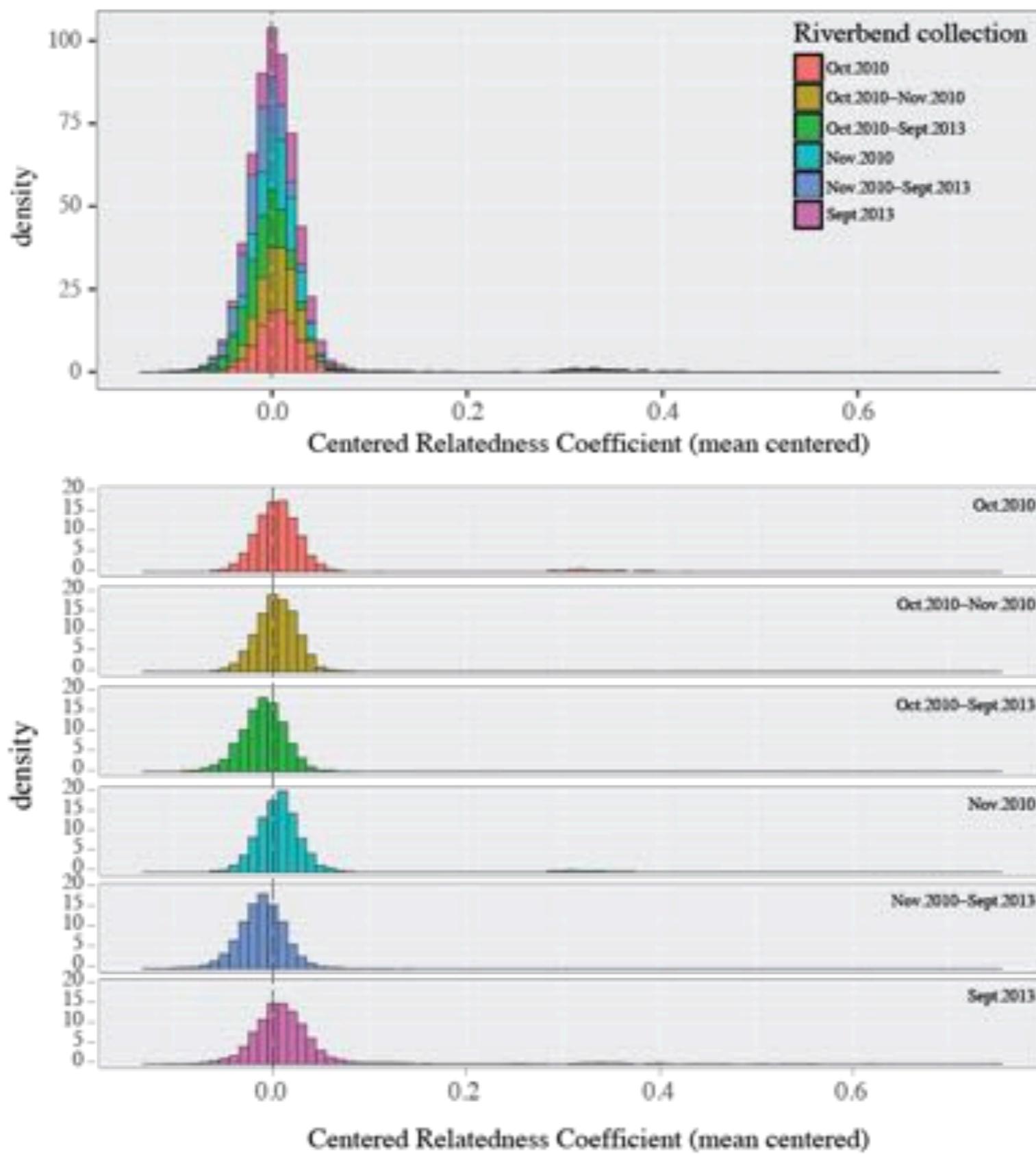
Genome Wide Association:
GEMMA

Zhou and Stephens 2012

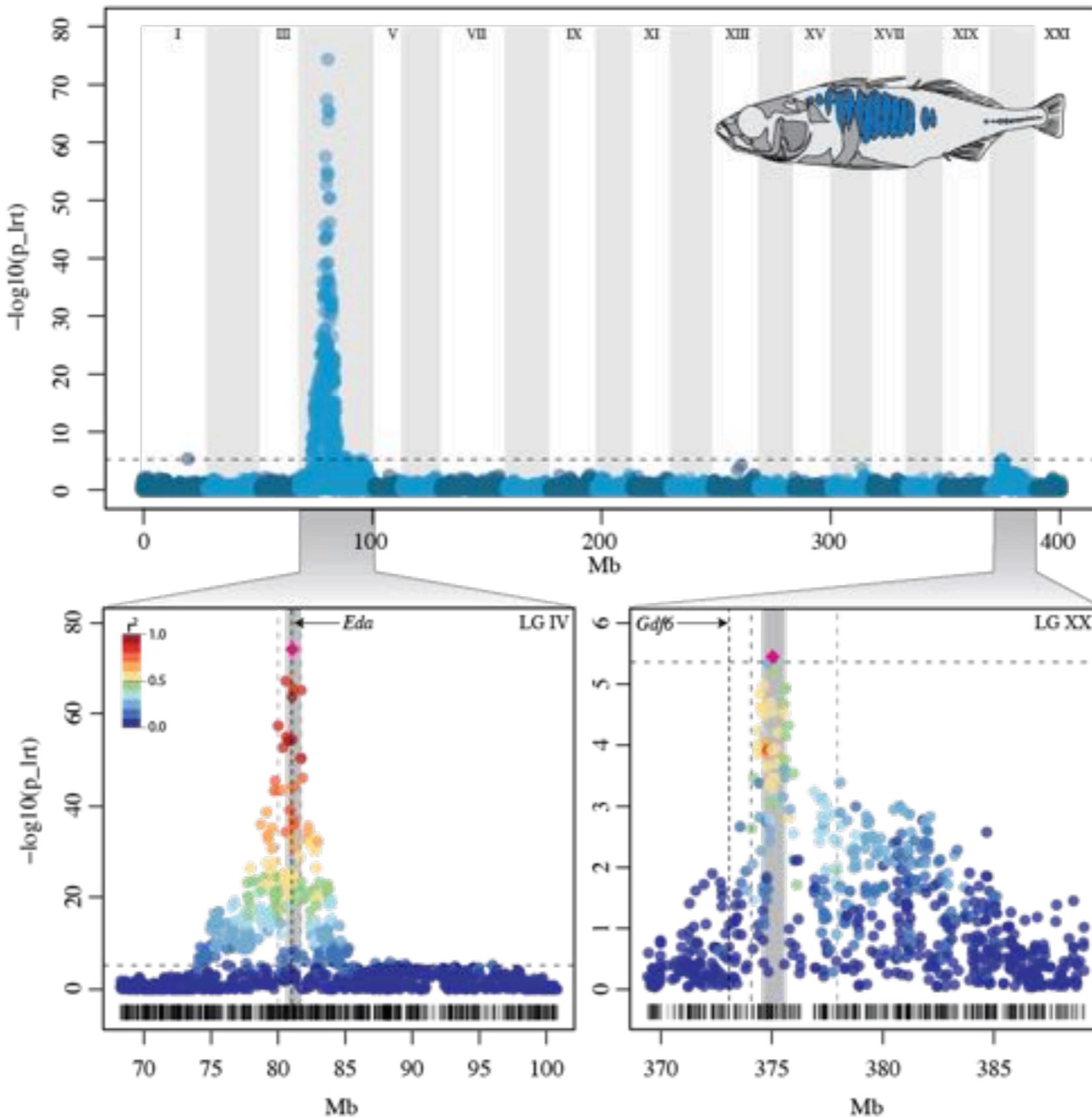
B



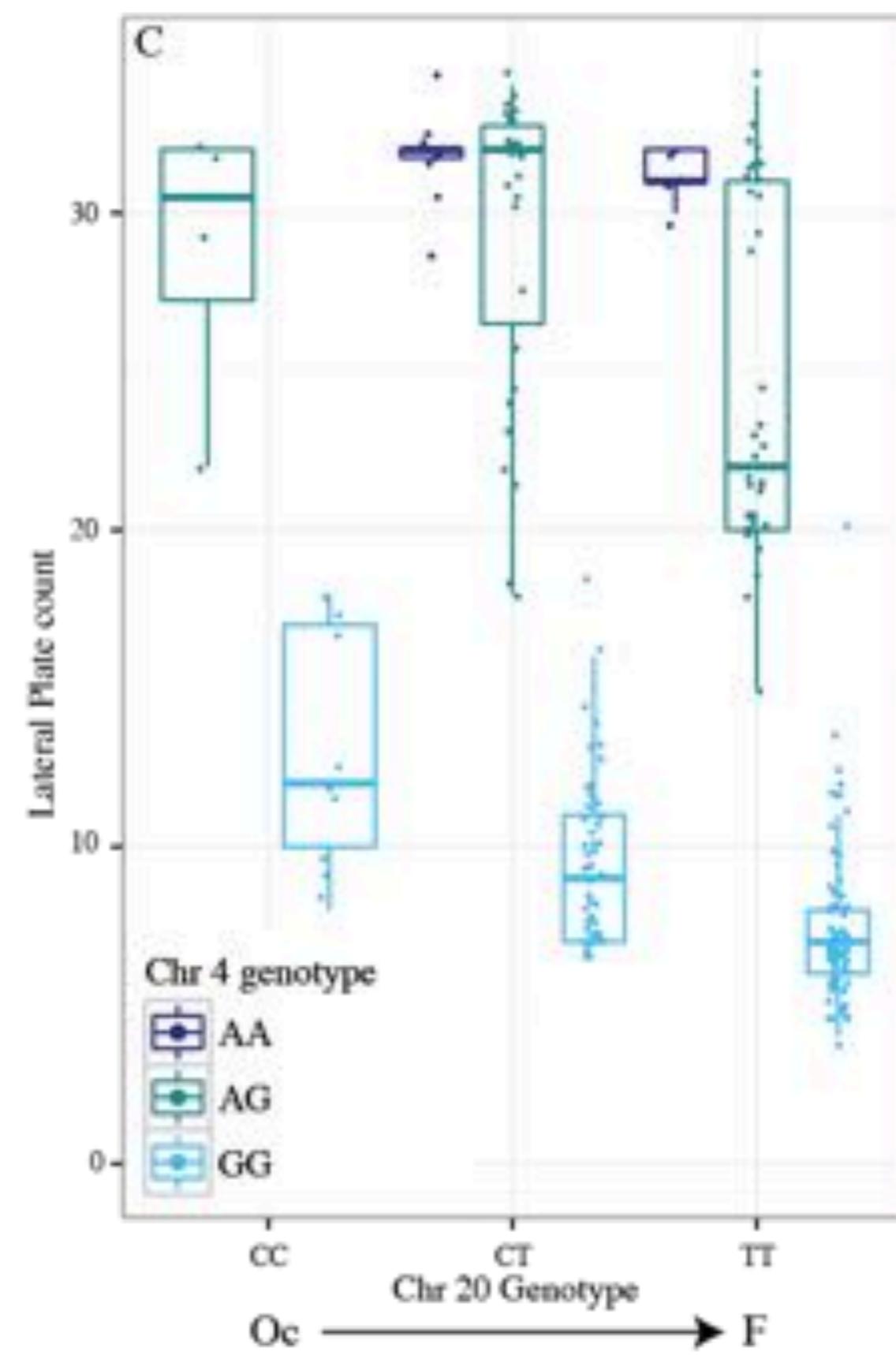
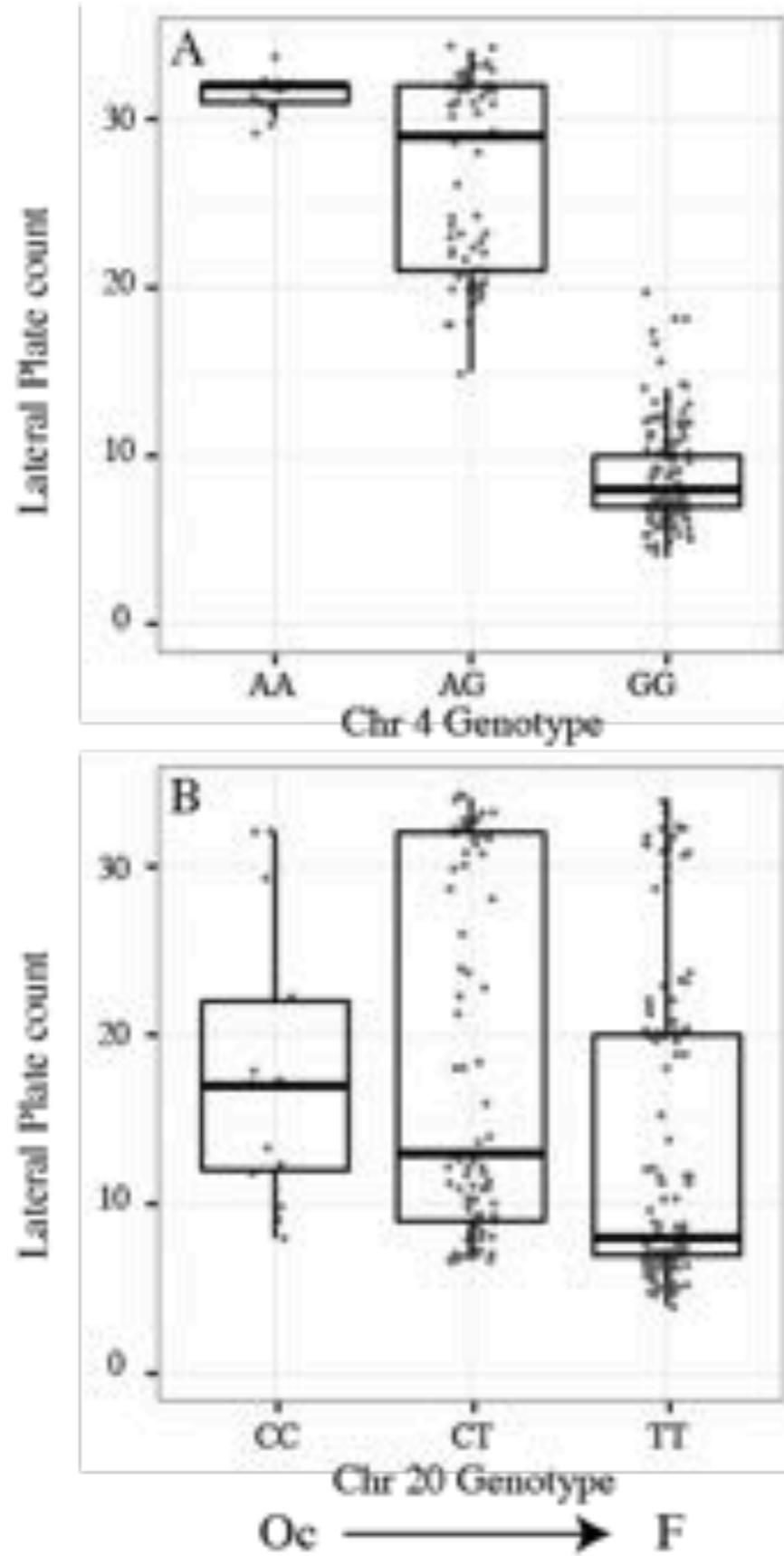
Pairwise relatedness among individuals



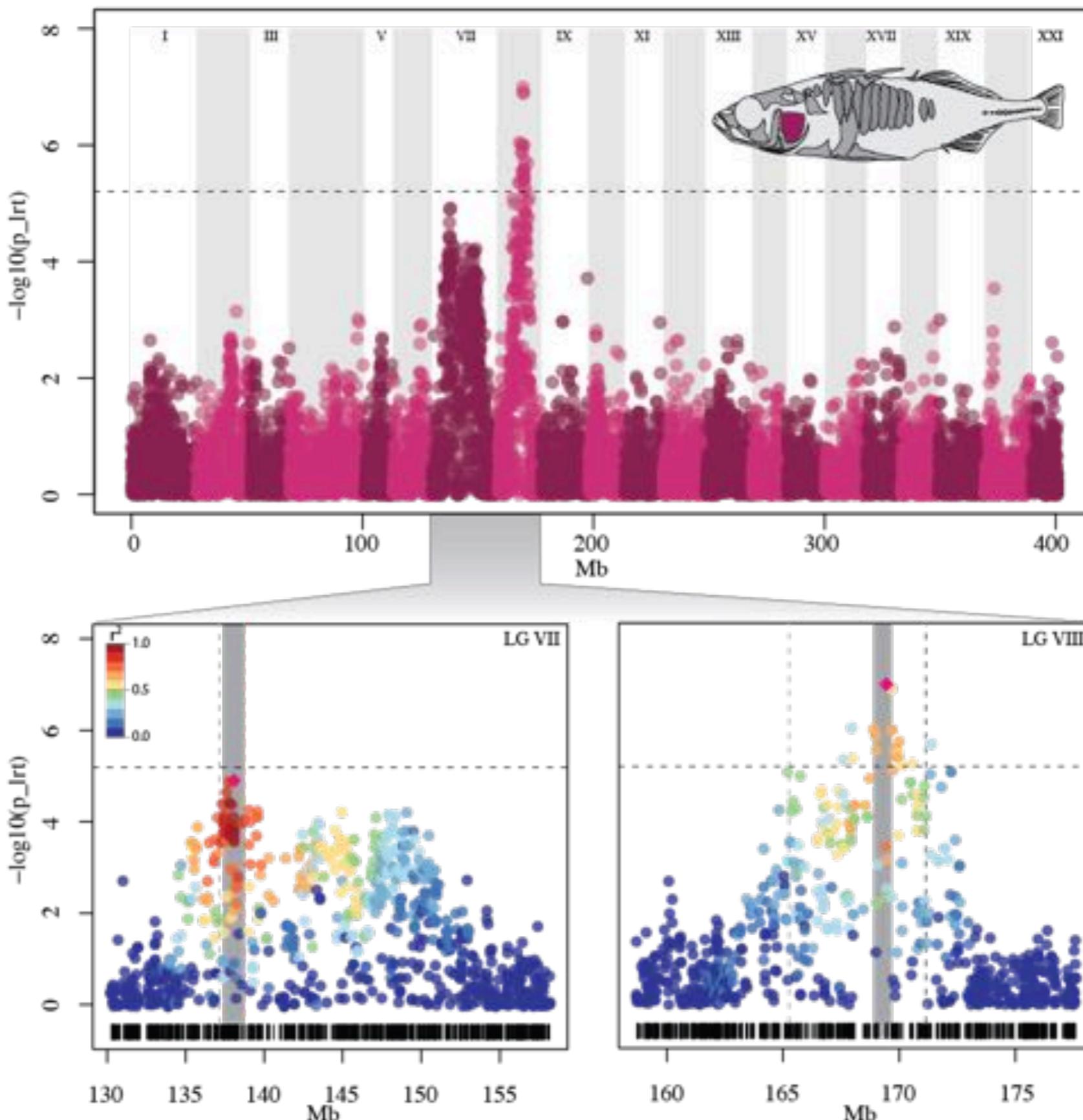
Lateral plate count: A good hit to a novel locus on LGXX



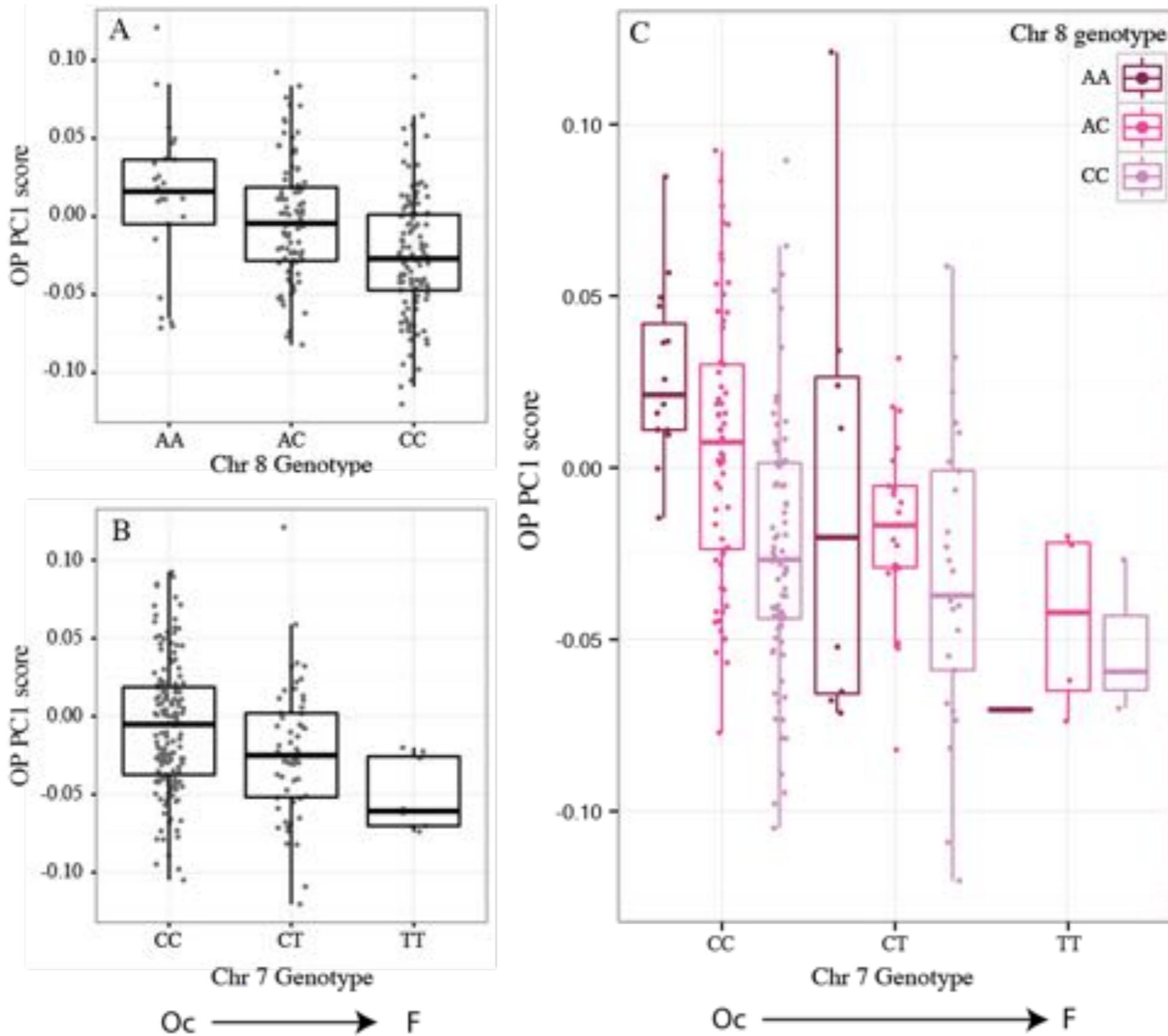
Lateral plate count: A good hit to a novel locus on LGXX



Opercle shape: Good hits on LGVII and LGVIII

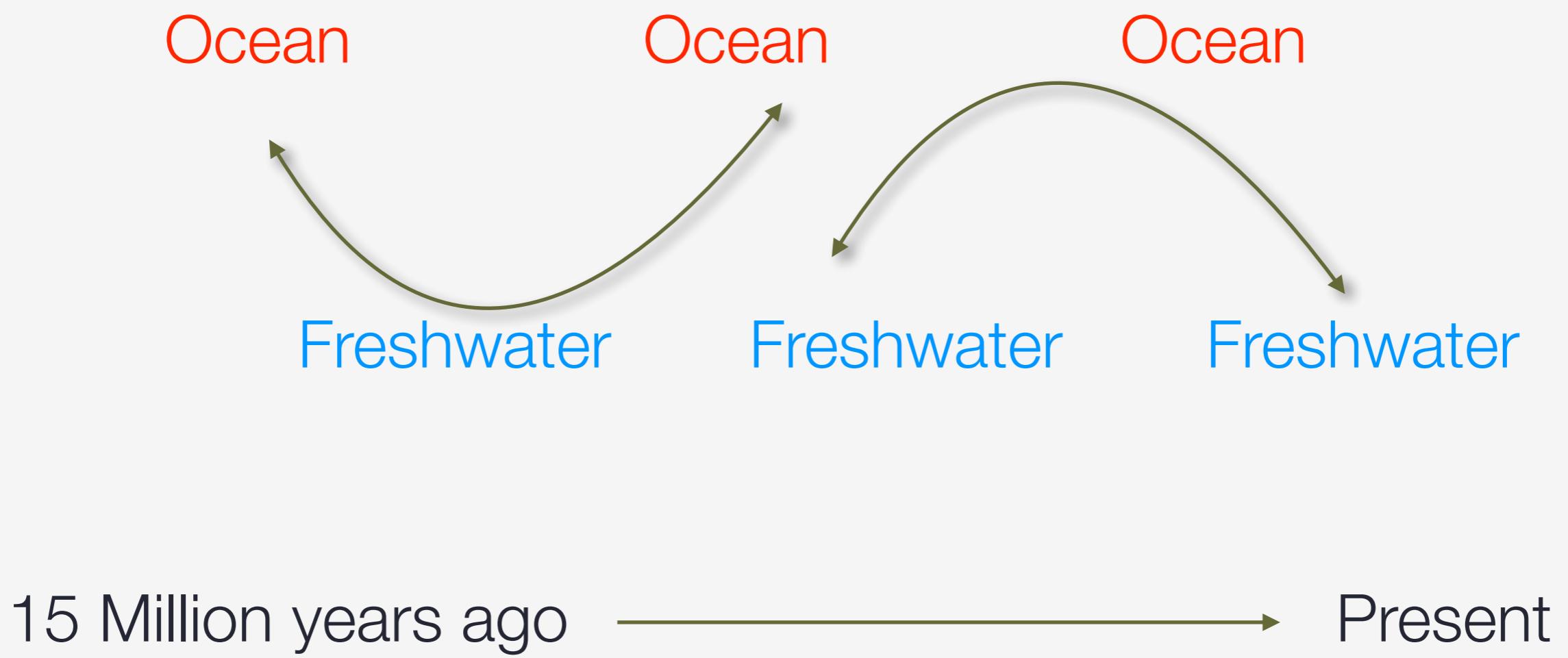


Opercle shape: Good hits on LGVII and LGVIII



- Previous work has shown that the freshwater genomes evolve in 13,000 years
- Uplift island data shows that the evolution can appear in as little as 50 years
- Much of the divergence involves soft sweeps of haplotypes
- 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

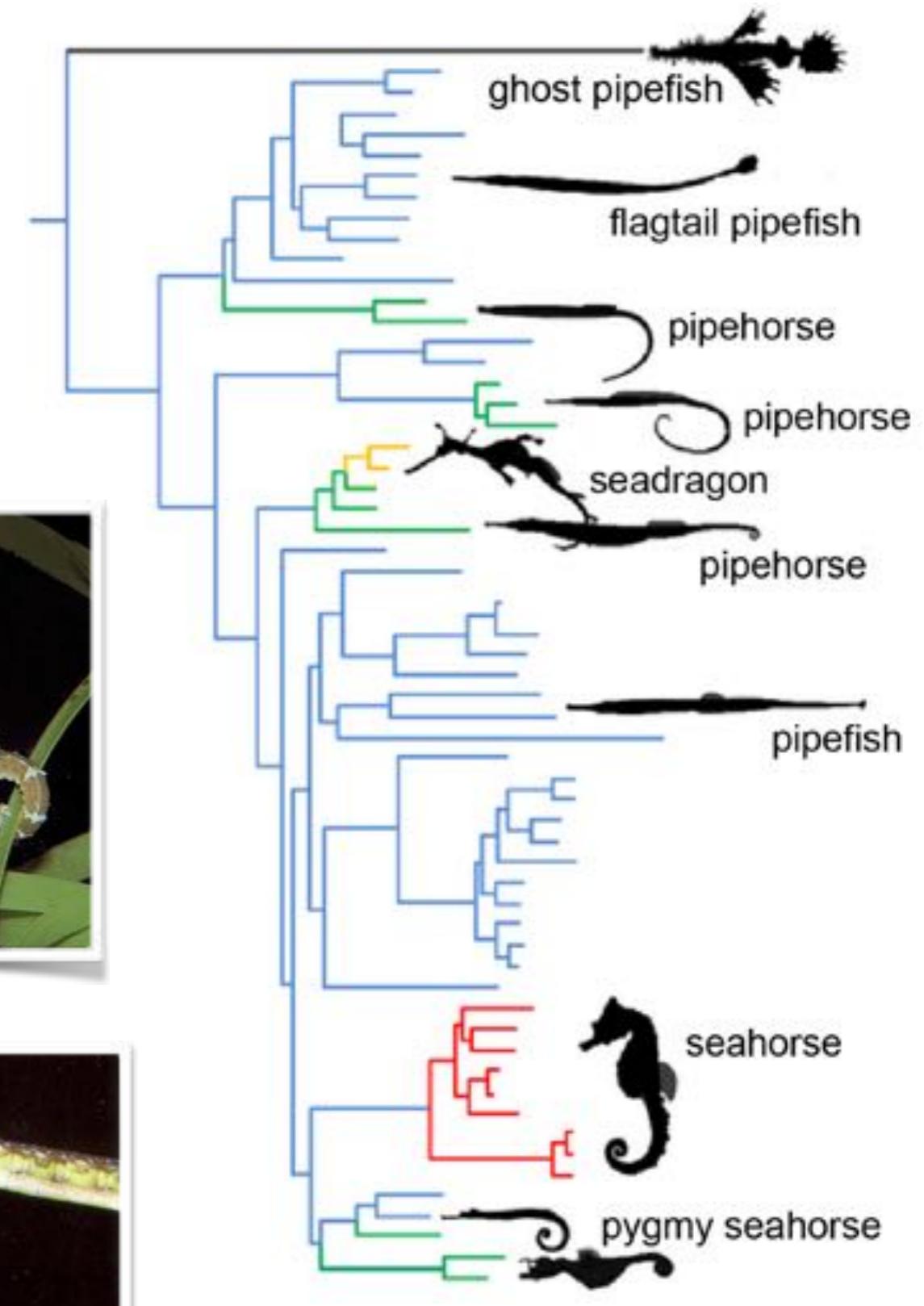
Genomically enabling the Gulf pipefish



What if you don't have a genome sequence?

A case study of RAD-seq and genome assembly

Seahorses, sea dragons and pipefishes



RESEARCH

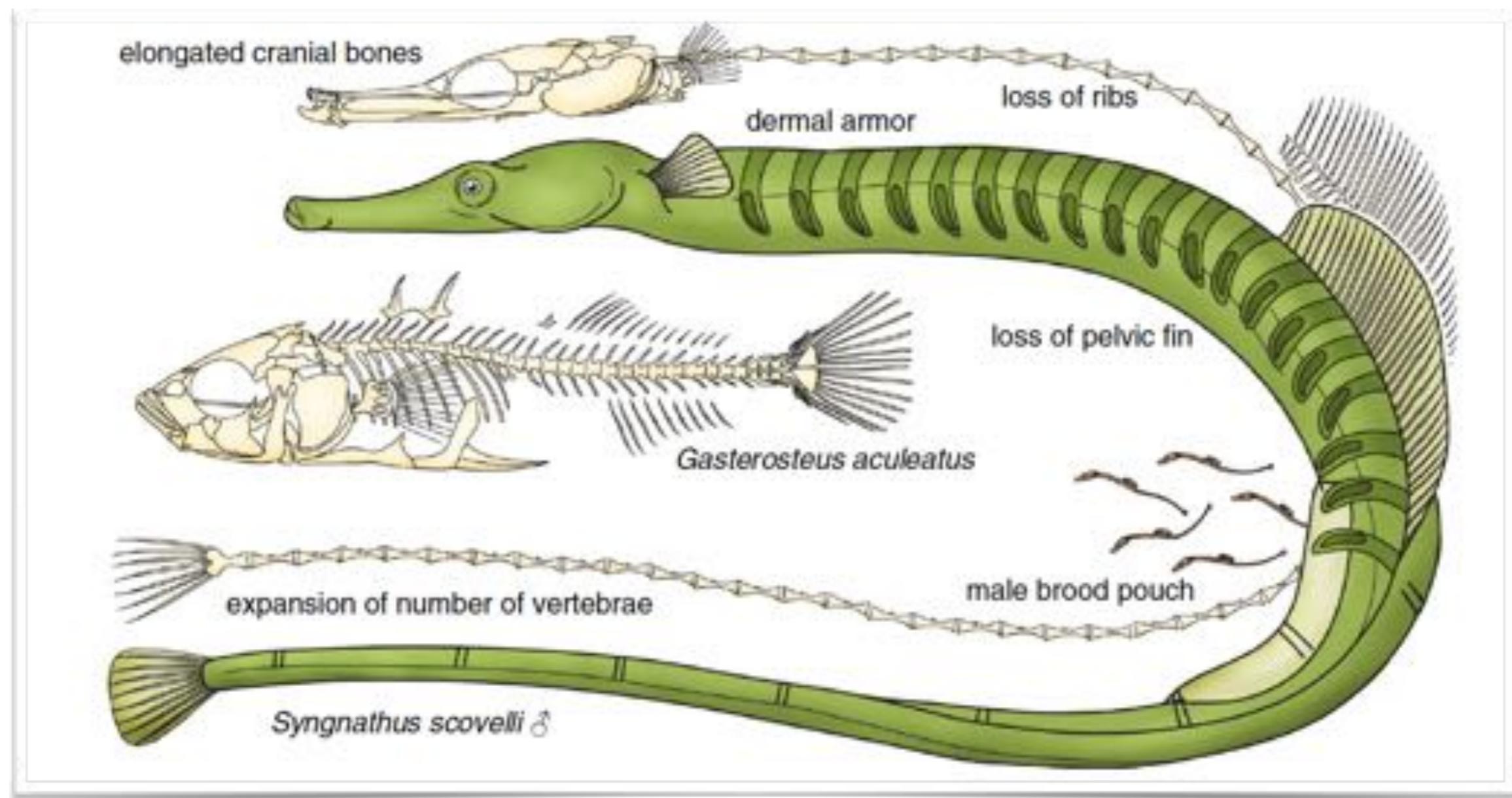
Open Access



CrossMark

The genome of the Gulf pipefish enables understanding of evolutionary innovations

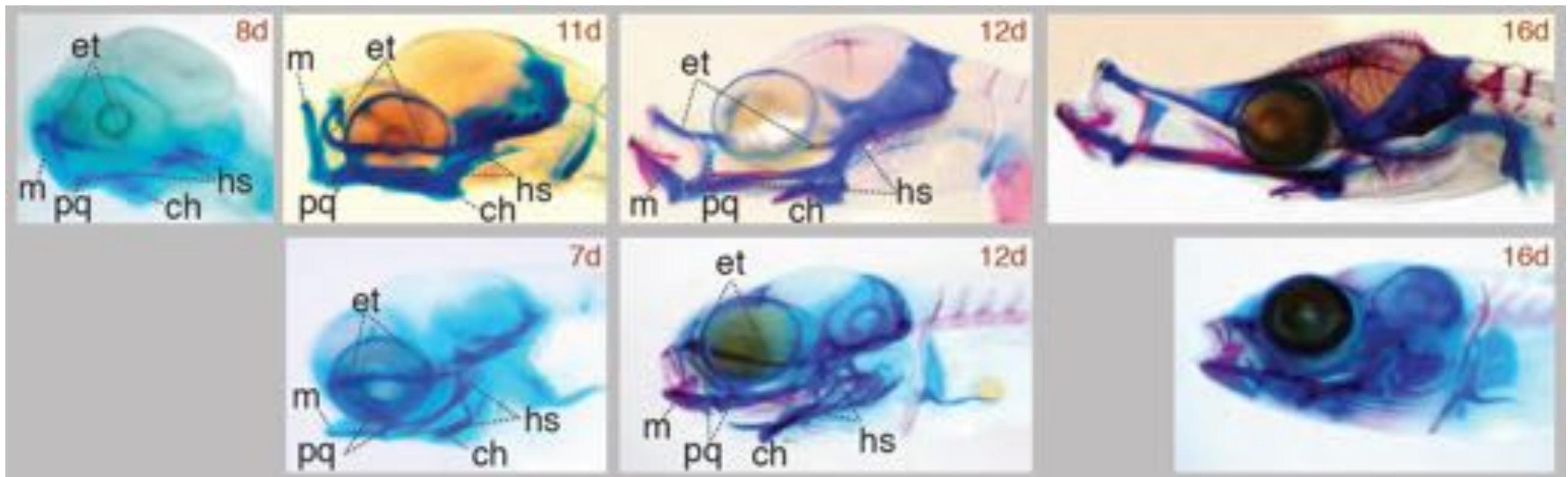
C. M. Small^{1†}, S. Bassham^{1†}, J. Catchen^{1,2†}, A. Amores³, A. M. Fuiten¹, R. S. Brown^{1,4}, A. G. Jones⁵ and W. A. Cresko^{1*}



We're really interested in head development



Pipefish



Stickleback

How did we genomically enable pipefish

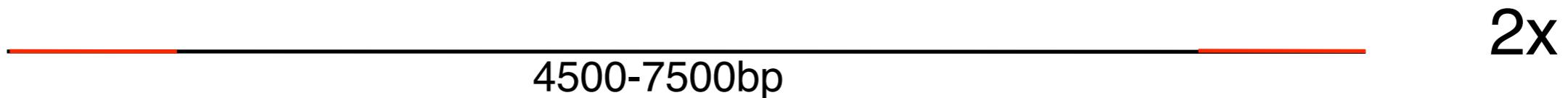
- 1) A high quality transcriptome
- 2) Very dense RAD genetic map
- 3) Deep shotgun sequencing of the genome
- 4) Order contigs against the RAD reference map

Illumina genomic libraries for pipefish genome

paired end 101bp



mate pair

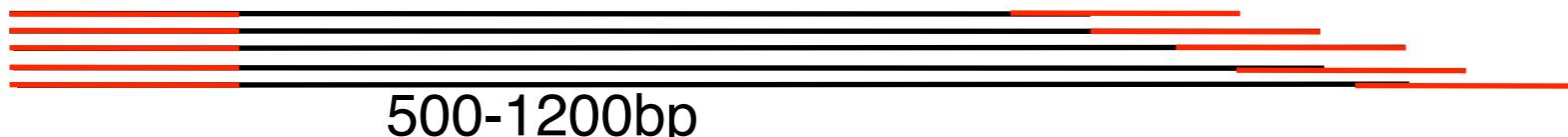


overlapping



paired end RAD

ACTCTC



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

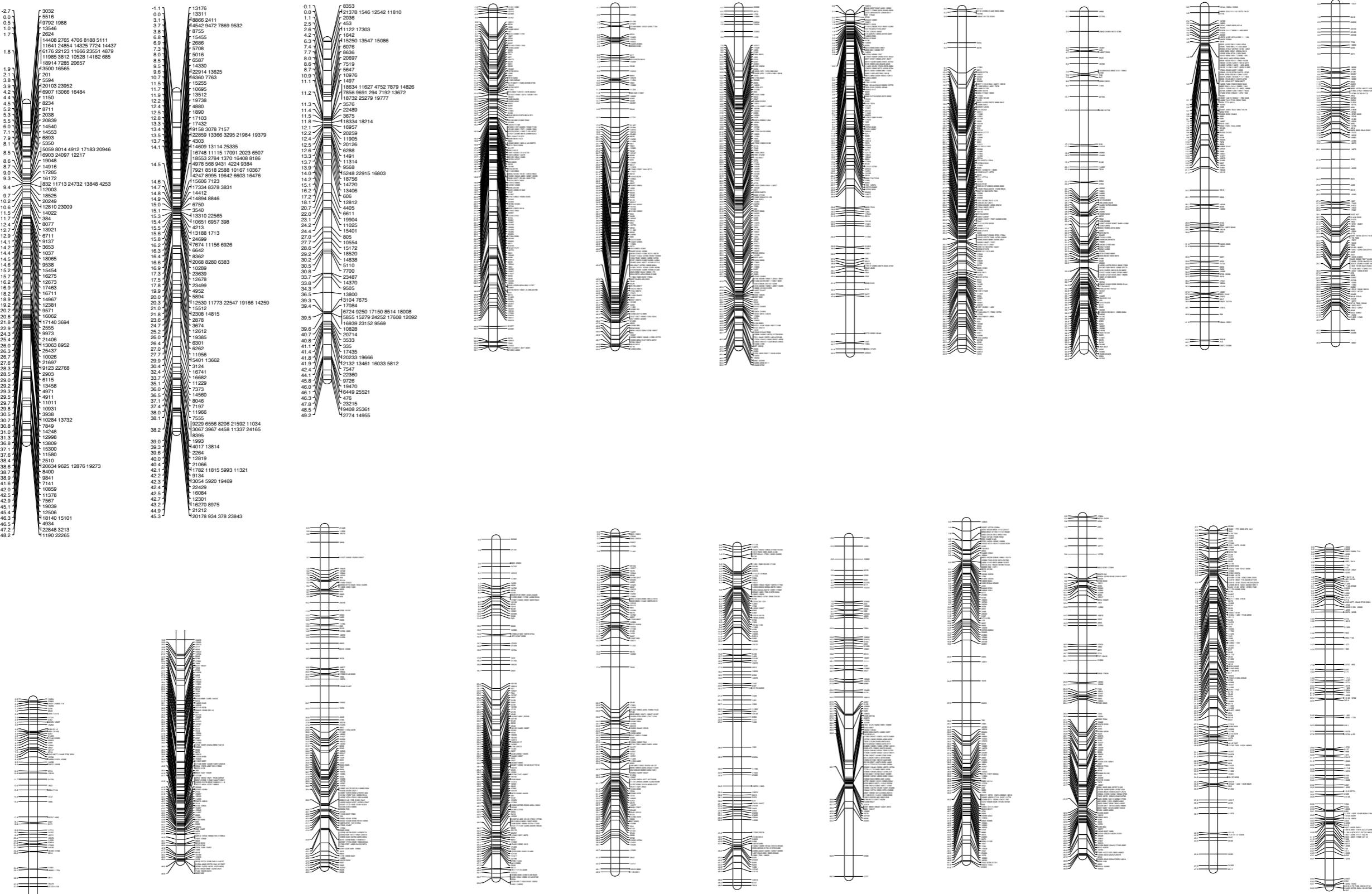
RAD sequenced the parents and offspring

Analyzed the data using *Stacks*

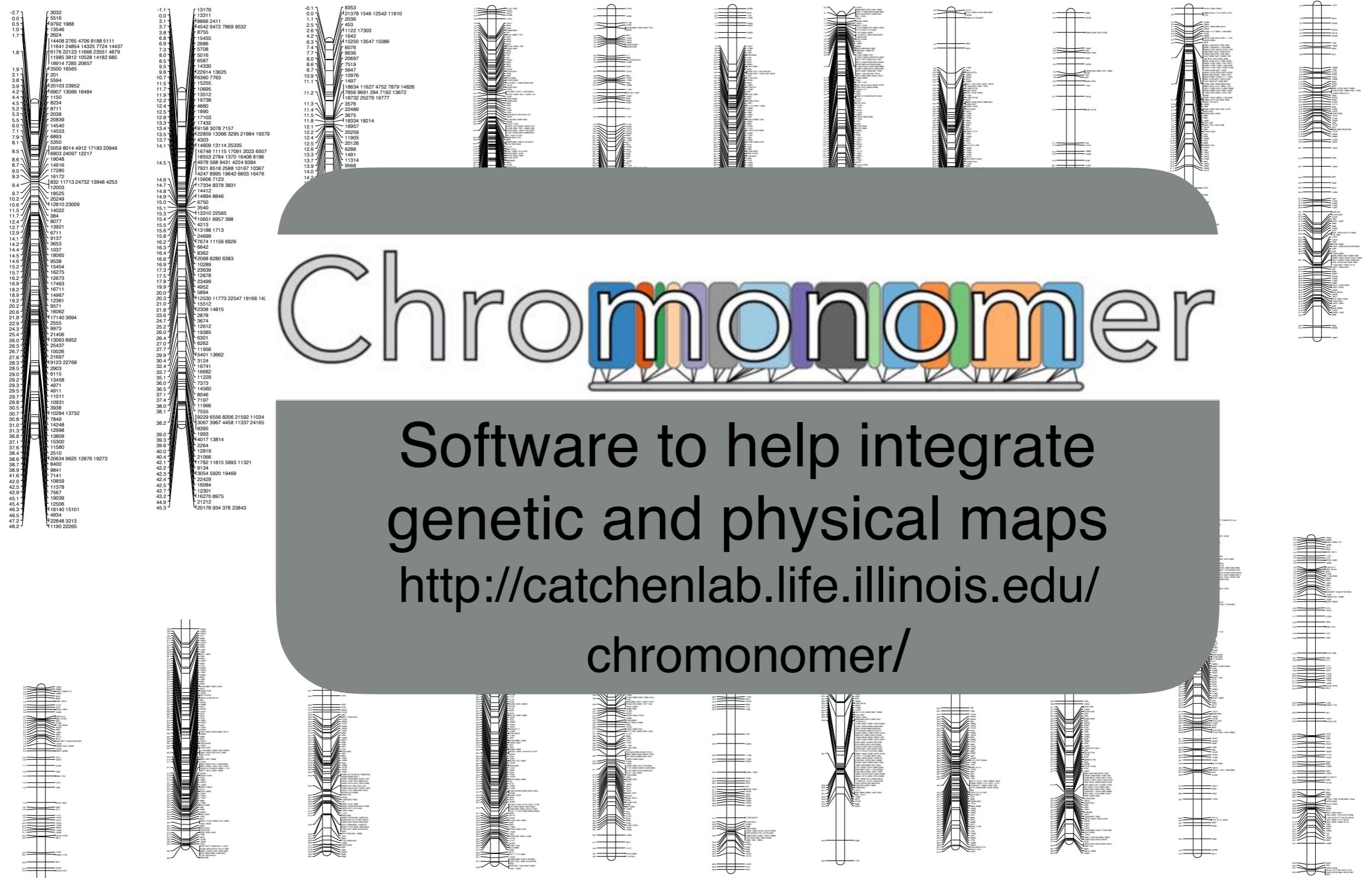
Output to JoinMap format

Created Linkage map

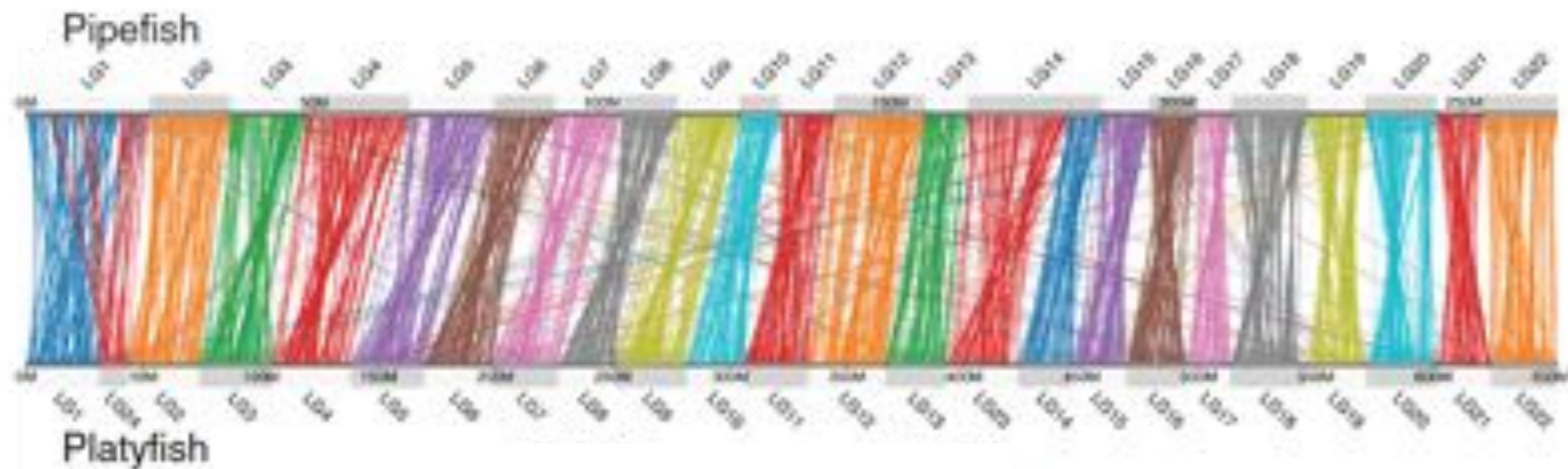
22 LGs; 6000 segregating SNPs; 30,000 RAD sites



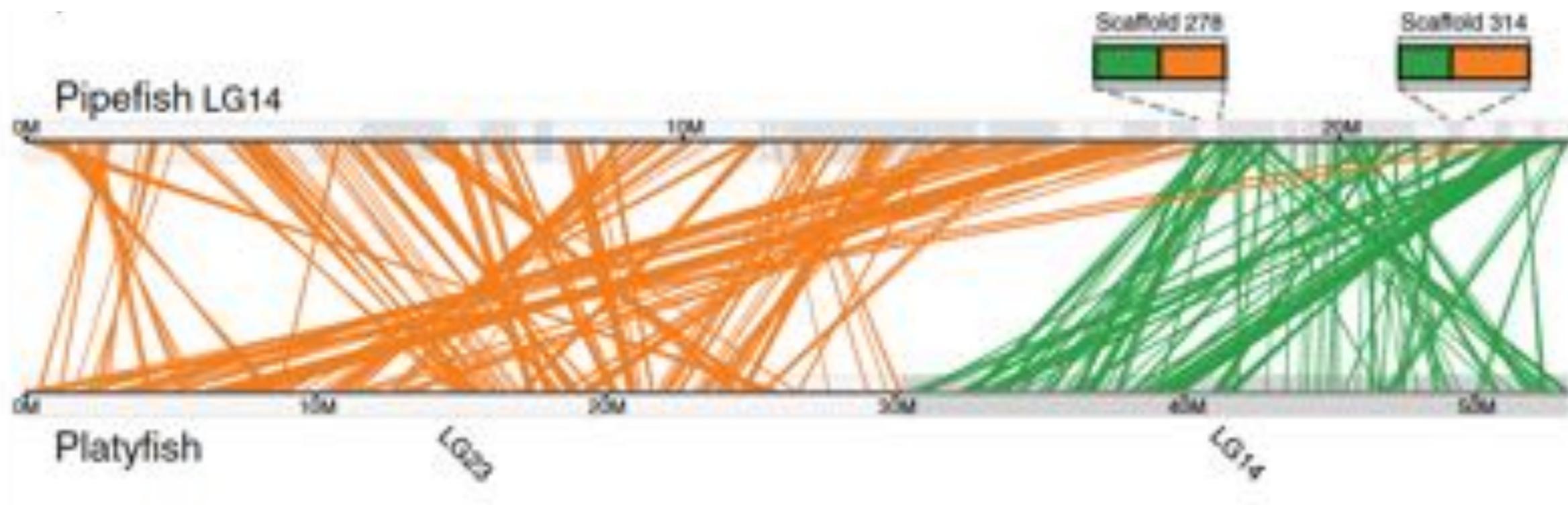
The pipefish genetic map and genome together



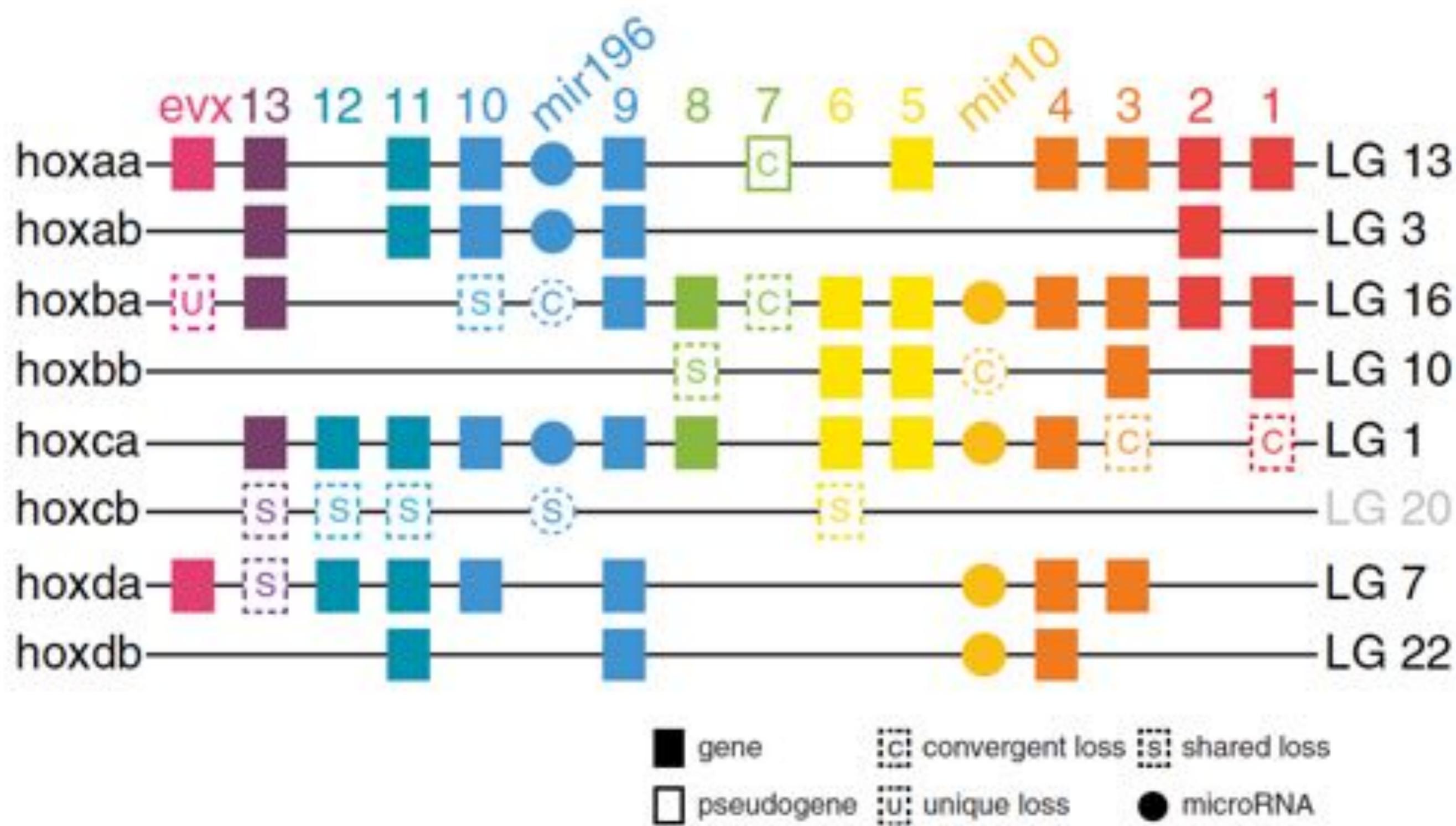
Two instances of chromosome fusion in pipefish



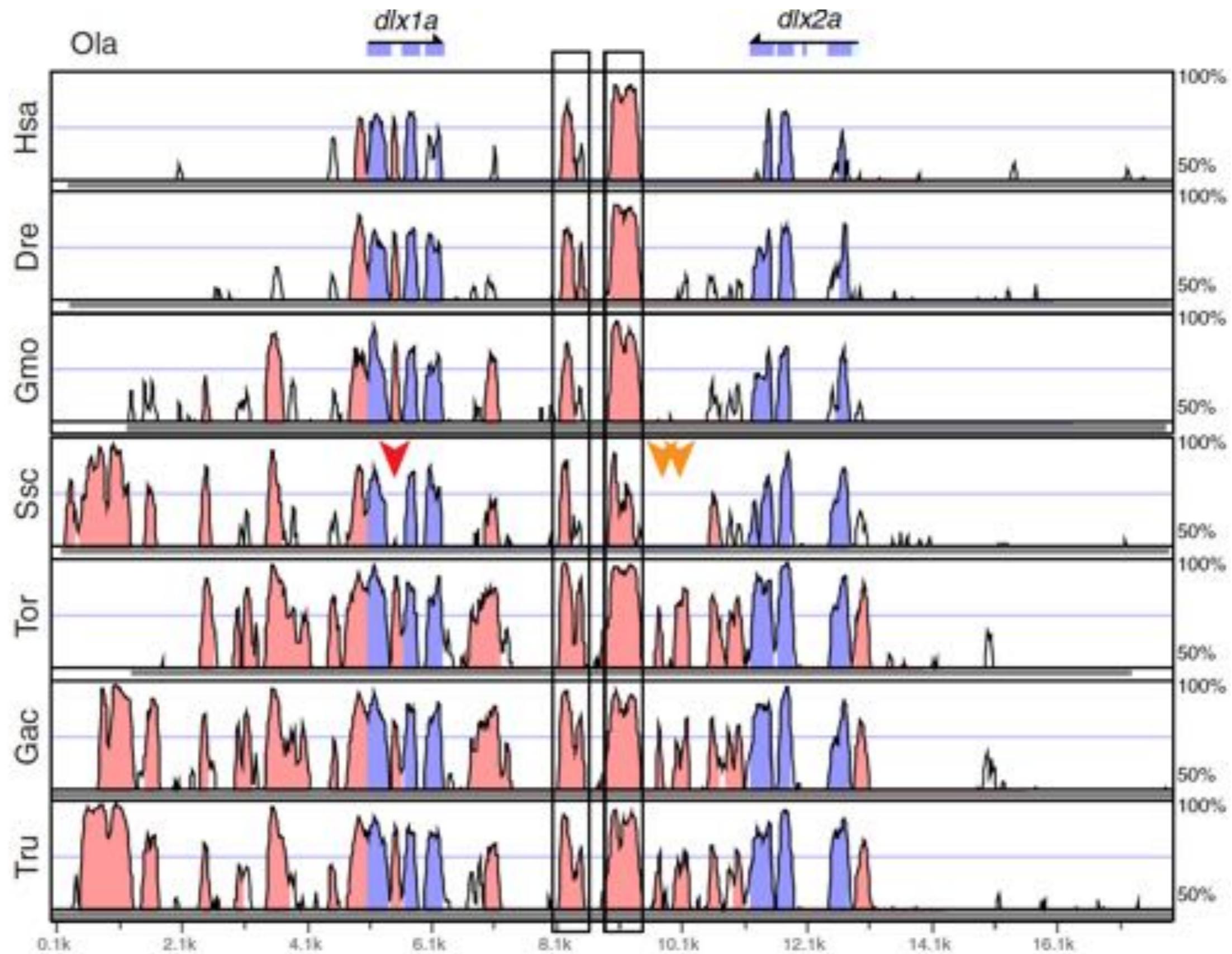
Two instances of chromosome fusion in pipefish



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

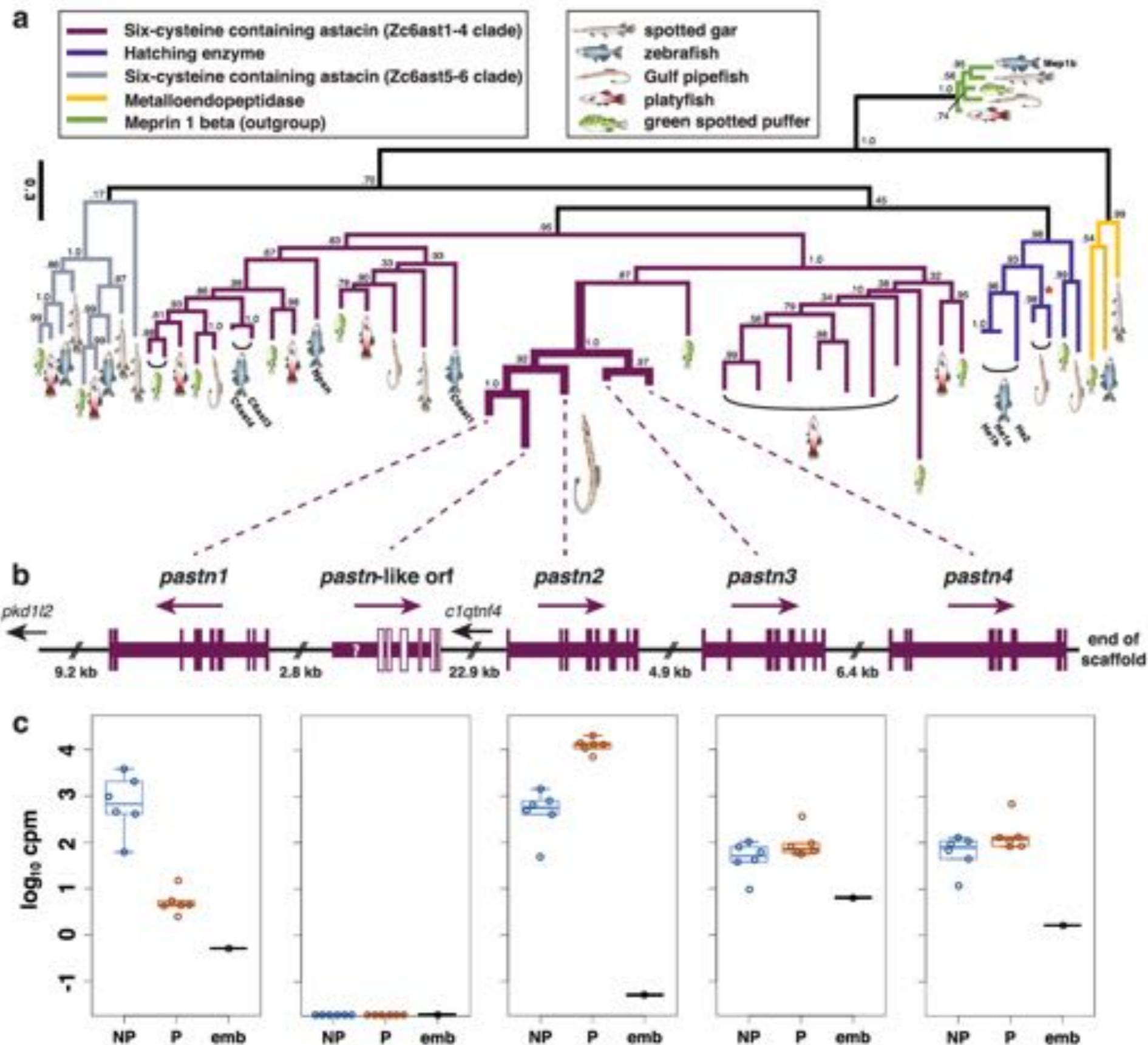
Disruption in coding sequence of *Pitx1*

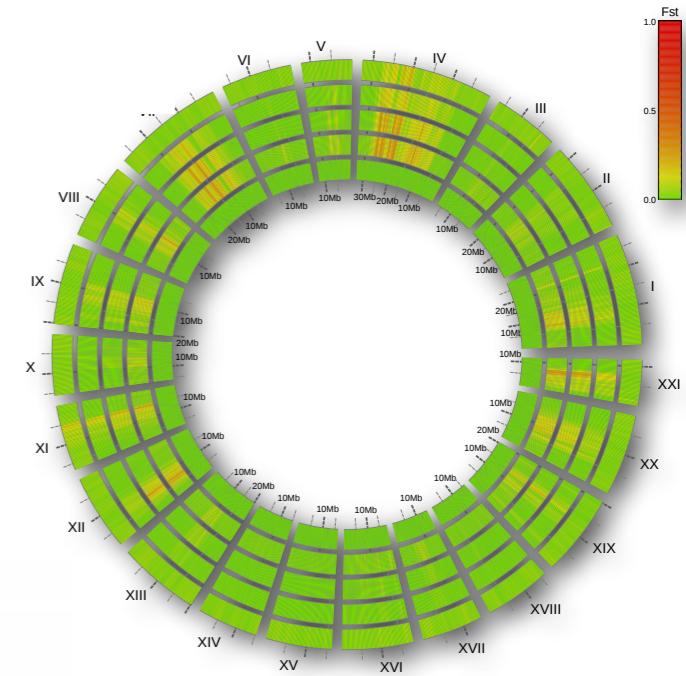


Pitx1

NRKREINQQMELCKNGYVPPSGLIVQPYED-MYA-----	AGYS-TNNWAAKSL-ADAPLSTKRPFT-PNSM--SPLSQ--SMF-SA-PSSISSMTH-----	PSIMGPAGV	human
NRKREINQQMELCKNGYVPPSGLIVQPYED-MYA-----	GYH-TNNWATKSL-TDAPLSTKRPFT-PNSM--SPLSQ--SMF-SA-PSSISSMH-----	SSTMGPSGV	coelacanth
NRKREINQQMELCKNSYLPPSGLIVQPYED-MYA-----	AYT-TNNWTKKL-ADAPLSTKRPFT-PNSM--SPLTSQ--SMF-SA-PSSISSMH-----	PPTMARTAV	gar
NRKREINQQMELCKNSYLPPSGLIVQPYED-VYP-----	TYT-TNNWTKKL-TDAPLSTKRPFT-PNSM--SPLTSQ--SMF-SA-PSSISSMH-----	AGOMRSAV	zebrafish
NRKREINQQMELCKNSYLPPSGLIVQPYED-MYA-----	TYT-TNNWTKKL-ADAPLSTKRPFT-PNSM--SPLTSQ--SVP-SA-PSSISSMH-----	ASGMRSAV	cavefish
NRKREINQQMELCKNSYLPPSGLIVQPYED-MYA-----	TYT-TNNWTKKL-ADAPLSTKRPFT-PNSM--SPLTSQ--SVP-SA-PSSISSMH-----	SSGMRSAV	ghost pipefish
NRKREINQQMELCKNSYLPPSGLIVQPYED-MYA-----	AYT-TNNWTKKL-ADAPLSTKRPFT-PNSM--SPLTSQ--SMF-SA-PSSISSMH-----	SSGMRSAV	Gulf pipefish
NRKREINQQMELCKNSYLPPSGLIVQPYED-MYA-----	AYT-TNNWTKKL-ADAPLSTKRPFT-PNSM--SPLTSQ--SMF-SA-PSSISSMH-----	ASGMRSAV	messmate pipefish
NRKREINQQMELCKNSYLPPSGLIVQPYED-MYA-----	AYT-TNNWTKKL-ADAPLSTKRPFT-PNSM--SPLTSQ--SMF-SA-PSSISSMH-----	ASGMRSAV	medaka
NRKREINQQMELCKNSYLPPSGLIVQPYED-MYA-----	AYT-TNNWTKKL-ADAPLSTKRPFT-PNSM--SPLTSQ--SMF-SA-PSSISSMH-----	ASGMRSAV	tilapia
NRKREINQQMELCKNSYLPPSGLIVQPYED-MYA-----	AYT-TNNWTKKL-ADAPLSTKRPFT-PNSM--SPLTSQ--SMF-SA-PSSISSMH-----	ASGMRSAV	pufferfish
NRKREINQQMELCKNSYLPPSGLIVQPYED-MYA-----	TYT-TNNWTKKL-ADAPLSTKRPFT-PNSM--SPLTSQ--SMF-SA-PSSISSMH-----	APOMCIPRAA	stickleback
NRKREINQQMELCKNGPDPGPMQPYED-MYA-----	GYS-TNNWAACGL-TSASLTKSFFFF-PNSMVNSPLSQ--SMF-SP-PNSISSEM-----	SSSMVPSAV	human
NRKREINQQMELCKNGPDPGPMQPYED-MYA-----	GYS-TNNWAACGL-TSASLTKSFFFF-PNSMVNSPLSQ--TMF-SP-PNSISSEM-----	SSSMVPS-V	coelacanth
NRKREINQQMELCKNGPDPGPMQPYED-MYA-----	STT-TNNWAACGL-TSASLTKSFFFF-PNSMVNSPLSQ--TMF-SP-PNSISSEM-----	SSSMVPSAV	gar
NRKREINQQMELCKNGPDPGPMQPYED-MYA-----	STT-TNNWAACGL-TSASLTKSFFFF-PNSMVNSPLSQ--TMF-SP-PNSISSEM-----	SSSMVPSAV	zebrafish
NRKREINQQMELCKNGPDPGPMQPYED-MYA-----	STT-TNNWAACGL-TSASLTKSFFFF-PNSMVNSPLSQ--AVF-SP-PTEISSEM-----	SSGMVPT--	cavefish
NRKREINQQMELCKNGPDPGPMQPYED-MYA-----	GYT-TNNWAACGL-TSASLTKSFFFF-PNSMVNSPLSQ--AMF-SP-PNSISSEM-----	TGGMVPSAV	Gulf pipefish
NRKREINQQMELCKNGPDPGPMQPYED-MYA-----	GYT-TNNWAACGL-TSASLTKSFFFF-PNSMVNSPLSQ--TMF-SP-PNSISSEM-----	TGGMVPAAV	medaka
NRKREINQQMELCKNGPDPGPMQPYED-MYA-----	GYT-TNNWAACGL-TSASLTKSFFFF-PNSMVNSPLSQ--TMF-SP-PNSISSEM-----	TGGMVPSAV	tilapia
NRKREINQQMELCKNGPDPGPMQPYED-MYA-----	GYT-TNNWAACGL-TSASLTKSFFFF-PNSMVNSPLSQ--TMF-SP-SPAPMSSEM-----	TGGMVPSAV	pufferfish
NRKREINQQMELCKNGPDPGPMQPYED-MYA-----	GYT-TNNWAACGL-TSASLTKSFFFF-PNSMVNSPLSQ--TMF-SP-SNSISEM-----	TGGMVPSAV	stickleback
NRKREINQQMELCKNGPAAFPGLIVQPYED-VYP-----	GYS-TNNWPKAL-A-PPLAAKTPPPAPNSVNSVPLSQ--SVP-SP-PSSIAAEM-----	PSLAAAPCT	human
NRKREINQQMELCKNGPAAFPGLIVQPYED-MYA-----	GYS-TNNWATKSL-ATPPLSAKSFFFF-PNSMVNSPLSQ--PMF-SP-PSSISSEM-----	PSSMVPSAV	coelacanth
NRKREINQQMELCKNGPAAFPGLIVQPYED-MYA-----	GYS-TNNWATKSL-ATPPLSAKSFFFF-PNSMVNSPLSQ--PMF-SP-PSSISSEM-----	ASGMVPSAV	gar
NRKREINQQMELCKNGPAAFPGLIVQPYED-MYA-----	GYS-TNNWATKSL-ASEPLSAKSFFFF-PNSMVNSPLSQ--PMF-SP-PSSIPSEM-----	ASGMVPSAV	zebrafish
NRKREINQQMELCKNGPAAFPGLIVQPYED-MYA-----	GYS-TNNWATKSL-ASEPLSAKSFFFF-PNSMVNSPLSQ--PMF-SP-PSSIPSEM-----	ASGMVPSAV	cavefish
NRKREINQQMELCKNGPAAFPGLIVQPYED-MYA-----	GYS-TNNWATKSL-ADQQLSAKSFFFF-PNSMVNSPLSQ--PMF-SP-PSSMPSEM-----	ASGMVPSAV	Gulf pipefish
NRKREINQQMELCKNGPAAFPGLIVQPYED-MYA-----	GYS-TNNWATKSL-ADQQLSAKSFFFF-PNSMVNSPLSQ--PMF-SP-PSSIPSEM-----	ASGMVPSAV	medaka
NRKREINQQMELCKNGPAAFPGLIVQPYED-VYP-----	GYS-TNNWATKSL-ASEPLSTKSFFFF-PNSMVNSPLSQ--PMF-SP-PSSIPSEM-----	ASGMVPTAV	tilapia
NRKREINQQMELCKNGPAAFPGLIVQPYED-VYS-----	GYS-TNNWAACL-ASEPLSAKSFFFF-PNSMVNSPLSQ--SMF-SP-PSSLPSEM-----	ASGMVPSAV	pufferfish
NRKREINQQMELCKNGPAAFPGLIVQPYED-VYT-----	GYS-TNNWATKSL-ASEPLSAKSFFFF-PNSMVNSPLSQ--PMF-SP-PSSIPSEM-----	ASGMVPSAV	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
- assisting in the assembly of genomes

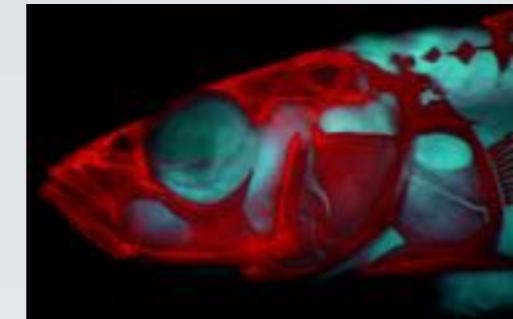
Open Source Genomics provides a suite of breakthrough technologies

- 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**
 - **Important to be comfortable with classical and modern statistics**

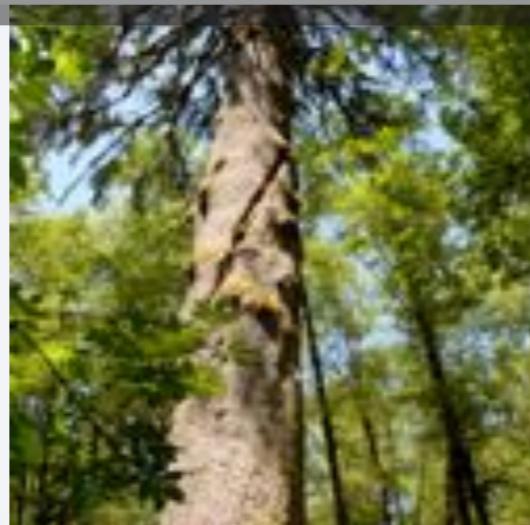
Acknowledgments



- *Past and present lab members*
- *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



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

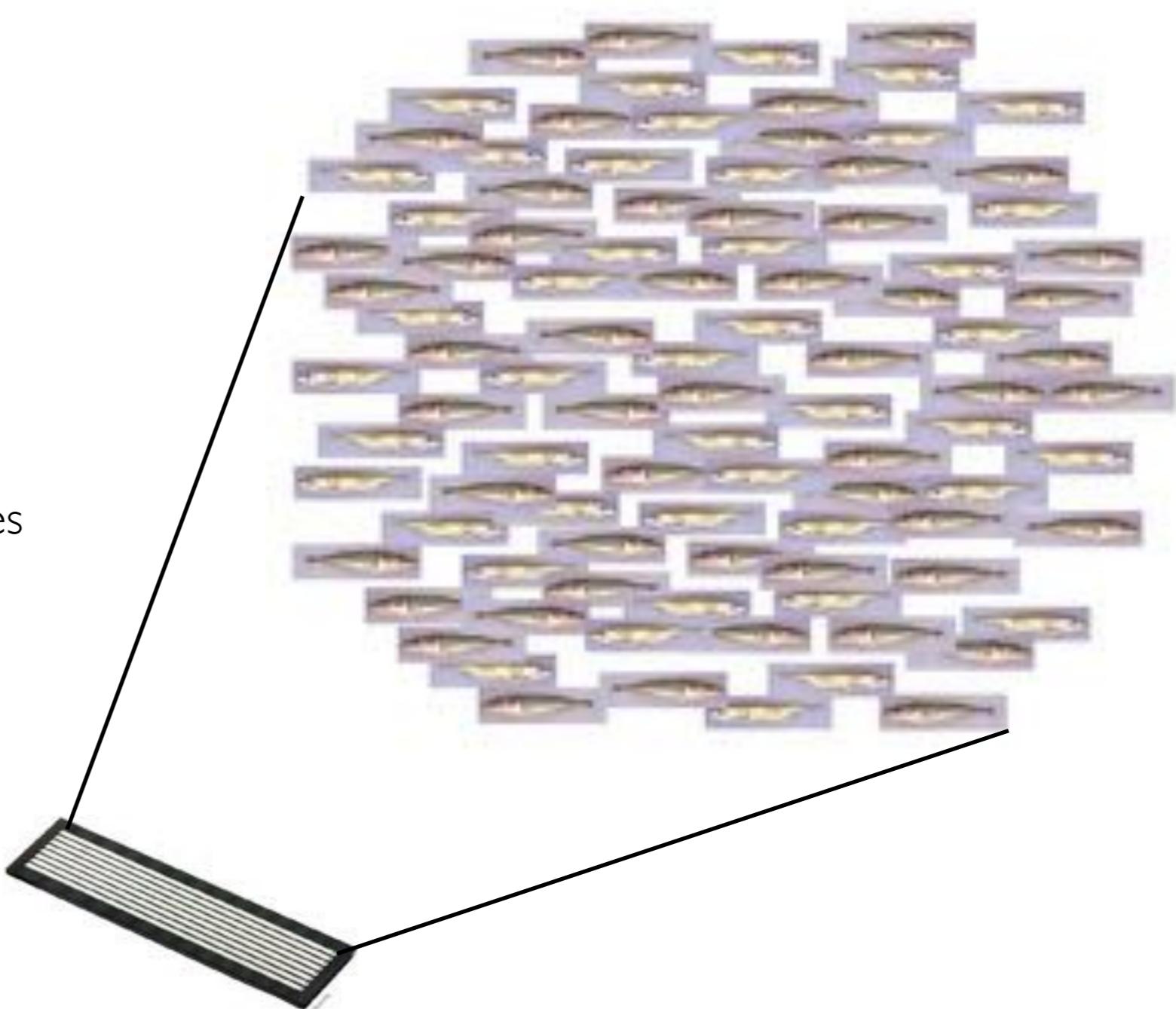
Barcoded adapters

5 to 8nt barcodes

Variable length barcodes

Combinatorial barcodes (PE)

Barcode distance - two mismatches



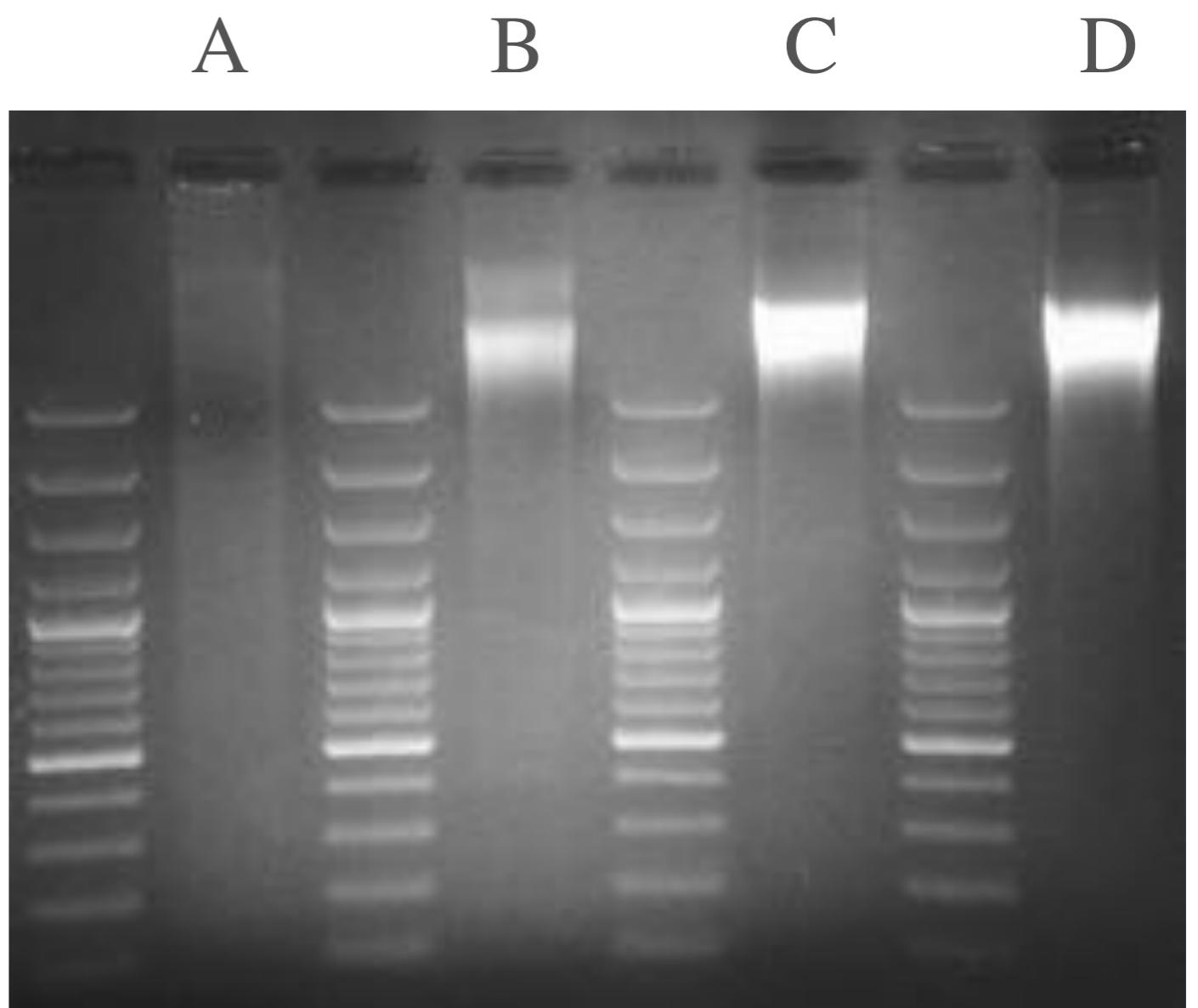
Molecular considerations in library building

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?

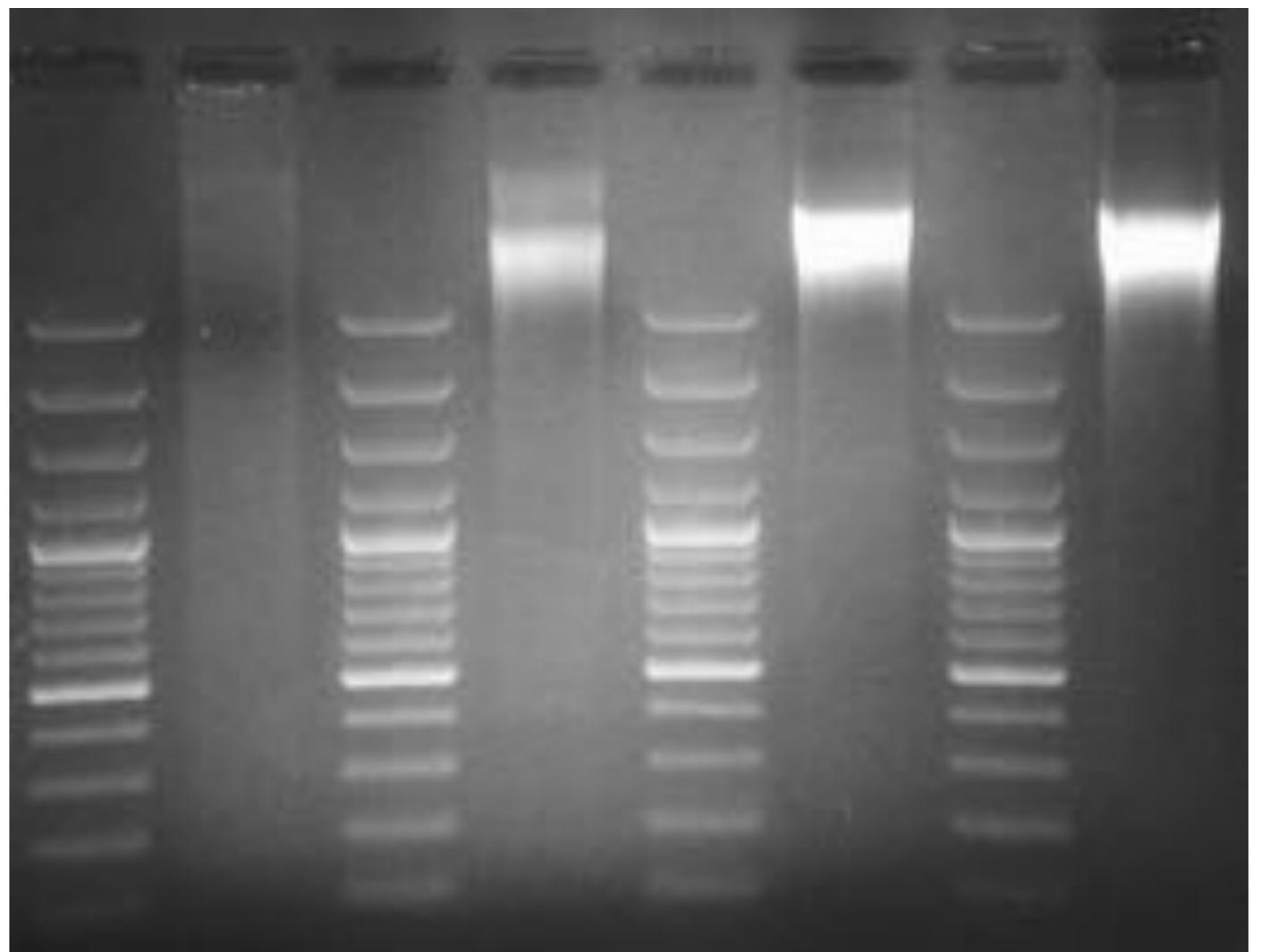
Things to consider

DNA Quality

[Diversify barcodes](#)

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

CGATA GTACA TAGCC ACTGC



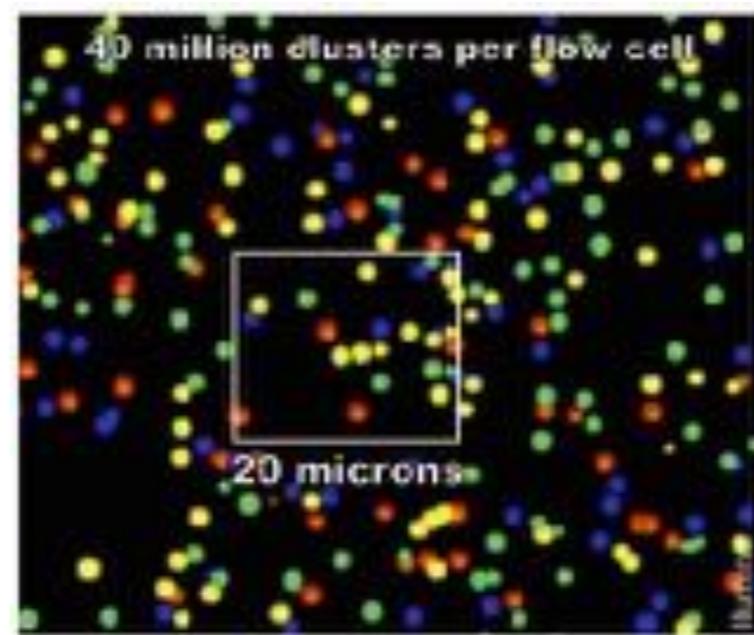
Molecular considerations in library building

How can I get the best depth of coverage?

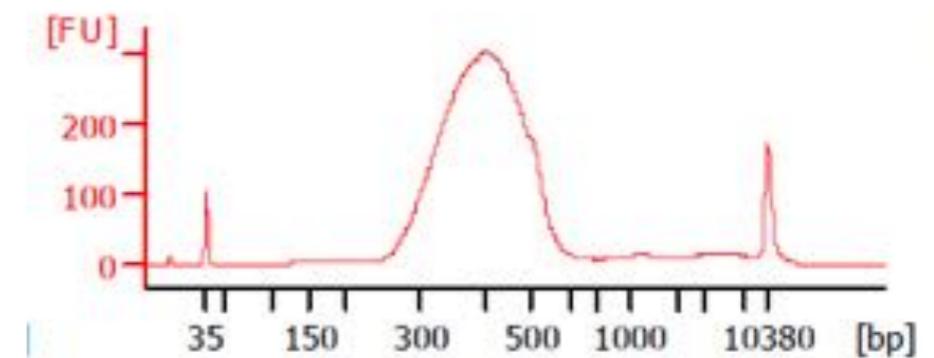
Things to consider

Fragment size

Smaller/tighter is better



Agilent Bioanalyzer



Molecular considerations in library building

How can I get the best depth of coverage?

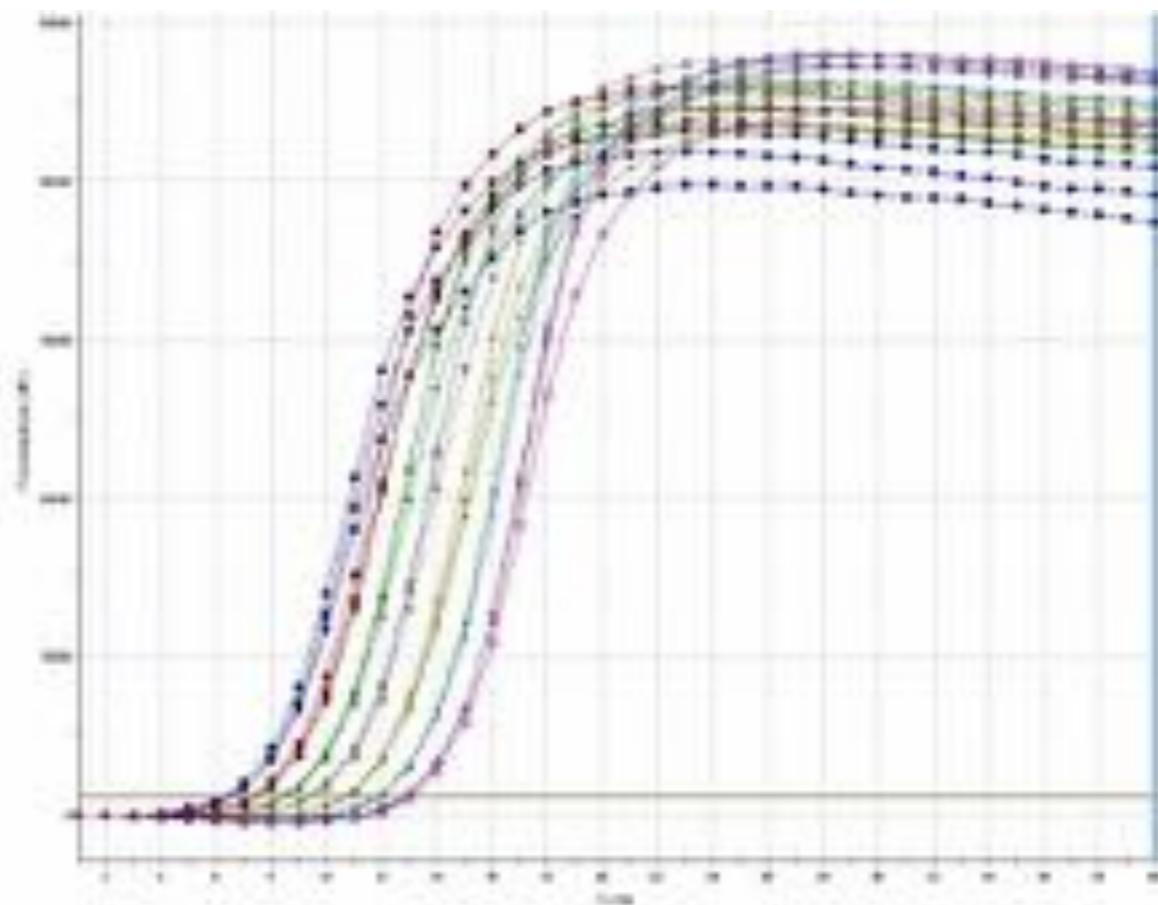
Things to consider

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

