

Short lectures

Short Lecture:

Alignment, reference genomes, denovo assemblies, stacks

Blast

BLAST®

Home Recent Results Saved Strategies Help

Basic Local Alignment Search Tool

BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance. [Learn more](#)

NEWS

BLAST+ 2.14.0 is here!
BLASTP, BLASTX, and TBLASTN are faster than before.
Fri, 28 Apr 2023 [More BLAST news...](#)

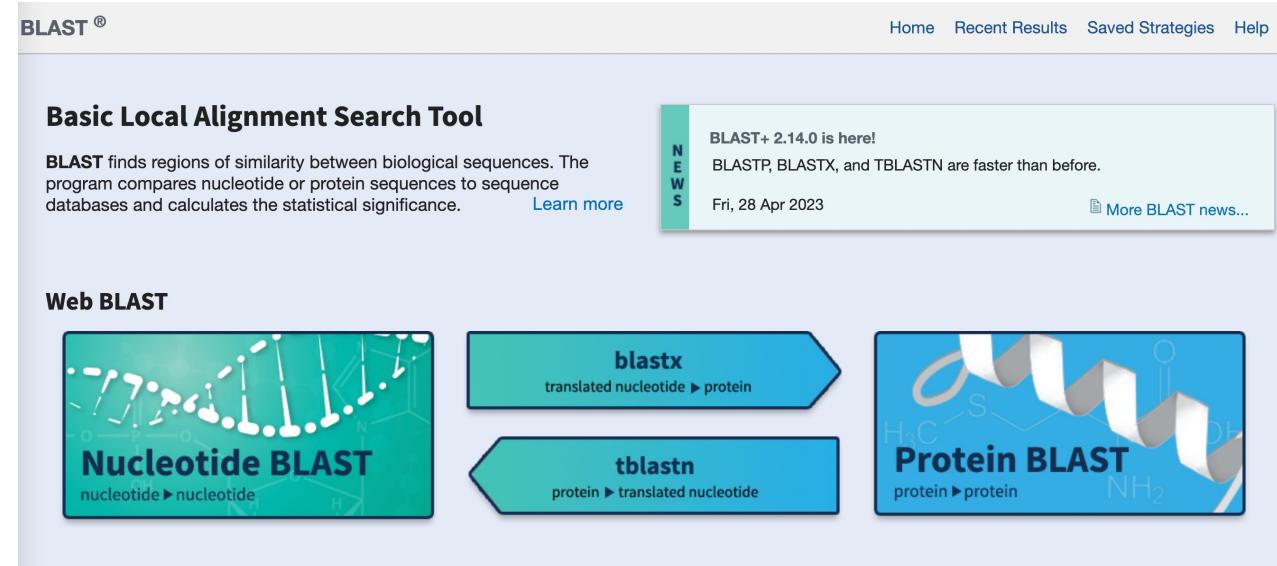
Web BLAST

Nucleotide BLAST
nucleotide ▶ nucleotide

blastx
translated nucleotide ▶ protein

tblastn
protein ▶ translated nucleotide

Protein BLAST
protein ▶ protein



- Finds conserved sequence (seed), then extends.
 - Mostly looking for homologous sequences, not exact matches
- Good for low #'s
- Slow!
 - Need to align hundreds of millions of reads

Short read aligners

- We want to find near exact matches within a genome for millions of reads
 - Usually working with one species
- Needs to be
 - Fast
 - Accurate
- Reference genomes are big
- We have MANY reads to align
- We don't want exact matches → why?

Short read aligners

Reference: AGTCGGATCACGGTTAGACTTCGACTAAGCGGACCGAATTGCG

- Basic approach: seed-and-extend

- Take a subset of your sequence (k ; usually ~ 19)
- Find the match in the reference genome (seeding)
- Try to extend these matches (extending)
- How might we alter k ?
 - shorter = more thorough but slower
 - longer = less thorough but faster.

Short read: CACGGTTACTTCTACT

Break into subsets: CACGGT TCGACT TCTACT

Seed:

CACGGT

AGTCGGATCACGGTTAGACTTCGACTAAGCGGACCGAATTGCG

Extend:

← CACGGT →

AGTCGGATCACGGTTAGACTTCGACTAAGCGGACCGAATTGCG

TCGACTTCTACT

CACGGT

AGTCGGATCACGGTTAGACTTCGACTAAGCGGACCGAATTGCG

CACGGTTCGACTTCTACT

Alignment: AGTCGGATCACGGTTAGACTTCGACTAAGCGGACCGAATTGCG

Short read aligners

- BWA
- Bowtie2
- STAR
- Pseudo aligners (for RNAseq)
 - Kallisto
 - Salmon
- And many others...
 - Depends on what you need
 - Think read length, speed, accuracy, RNA vs DNA, etc.
- BWA: Burrow-Wheeler Aligner for short-read alignment
 - Burrows-Wheeler Alignment → a compression method
 - Maximal Exact Match → longest match that can't be extended further
 - Most short read aligners use similar approaches

PREPRINT

Vol. 00 no. 00 2013
Pages 1–3

Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM

Heng Li

Broad Institute of Harvard and MIT, 7 Cambridge Center, Cambridge, MA 02142, USA

What if I don't have a genome?

- Option 1: Denovo assembly

- Best option
- Challenging:
 - Genomes are big
 - Repetitive
 - Inversions

- Long reads address these issues
 - PacBio, Nanopore

- Can also link scaffolds with other methods:
 - Hi-C, optical mapping, HiRise, etc.

TTCCGGAGAGGGAGCCTGAGAAATGGCTACCACATCCACGGAGAGG

GCCTGAGAAATGGCTACCACATC

CCACATCCACGGAGAGG

TTCCGGAGAGGGAGCCTGAG

Repeat regions:

AACCAACCAACCAACCAACCAACCAACCAACCC

What if I don't have a genome?

- Option 2: Reduced representation
 - Most commonly RAD-seq/GBS or similar
 - Still must generate a pseudo reference genome
 - Stacks, ipyrad, ddocent, others

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

Sequence analysis ipyrad: Interactive assembly and analysis of RADseq datasets

Deren A.R. Eaton^{1,*} and Isaac Overcast²

¹Department of Ecology, Evolution, and Environmental Biology, Columbia University, New York, NY 10027, USA and ²Department of Biology, Graduate School, University Center of the City University of New York, New York, NY 10016, USA

***dDocent*: a RADseq, variant-calling pipeline designed for population genomics of non-model organisms**

Jonathan B. Puritz, Christopher M. Hollenbeck and John R. Gold

Marine Genomics Laboratory, Harte Research Institute, Texas A&M University-Corpus Christi, Corpus Christi, TX, USA

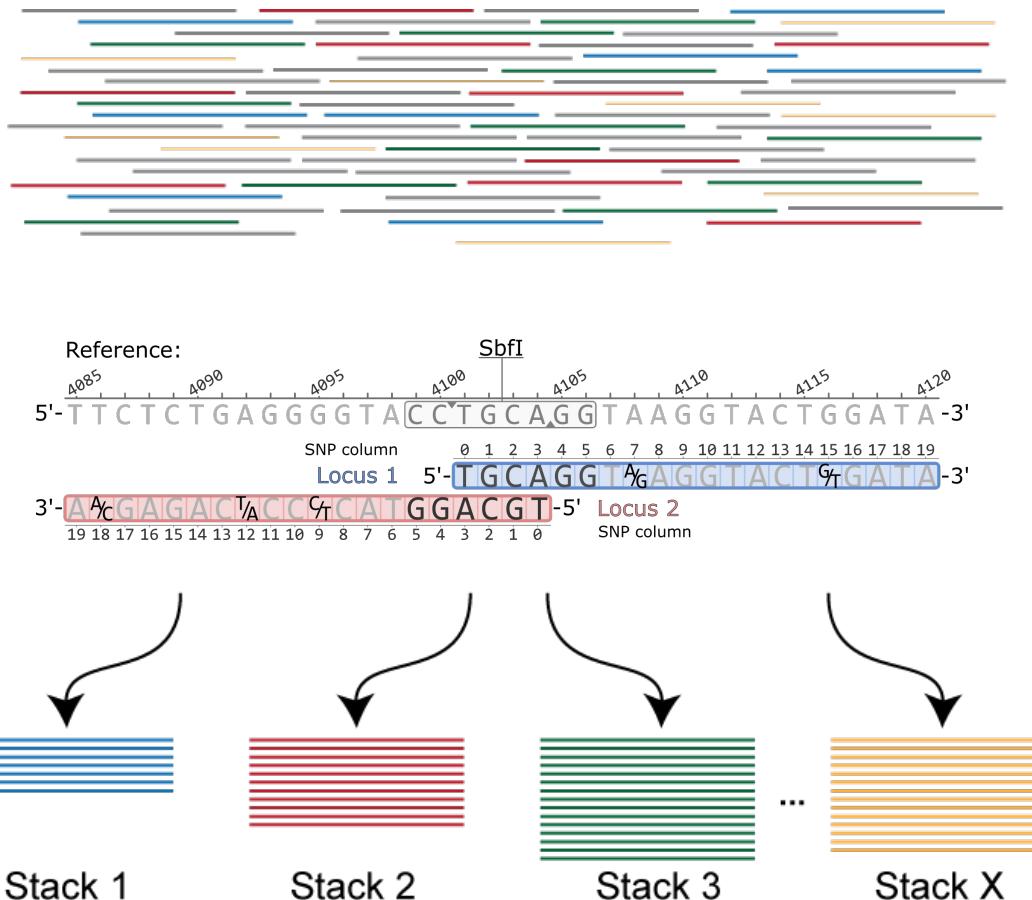
Stacks

Stacks Manual

Julian Catchen¹, Nicolas Rochette¹, Angel Rivera-Colón¹, William A. Cresko², Paul A. Hohenlohe³, Angel Amores⁴, Susan Bassham², John Postlethwait⁴

- Most popular method, good for within species comparisons
- Clustering by similarity:
 - Need to decide:
 - Minimum depth for a stack
 - # of mismatches to allow when making a stack
 - Allow for genetic variation
 - But not merge different genomic regions
- Harder than it seems...

pipeline for de novo rad assembly



Short Lecture:
Mapping results, PCR duplicates
Properly paired
multi-mapping

Alignment types

- Primary
- Secondary
- Supplementary
- Duplicates
- Properly paired

Alignment types

- Primary → When a read maps to multiple locations, the best alignment
- Secondary → When a read maps to multiple locations, the alternative, alignments (bwa mem doesn't use this flag, gives quality score of 0 instead)
- Supplementary
- Duplicates
- Properly paired

Alignment types

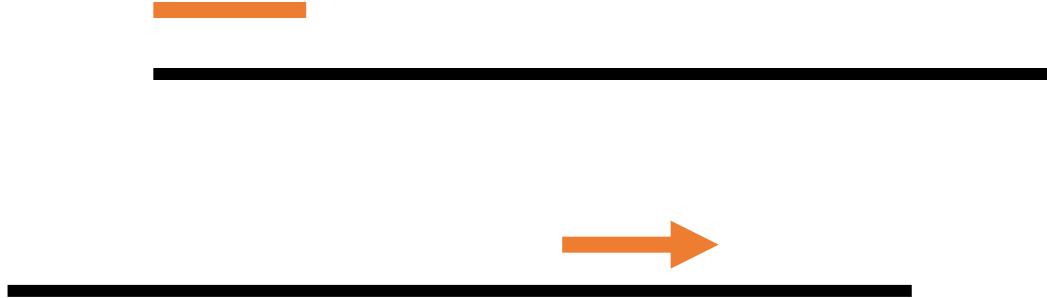
- Primary
- Secondary
- Supplementary
- Duplicates
- Properly paired

With reasonable distance (a few hundred bp)



Alignment types

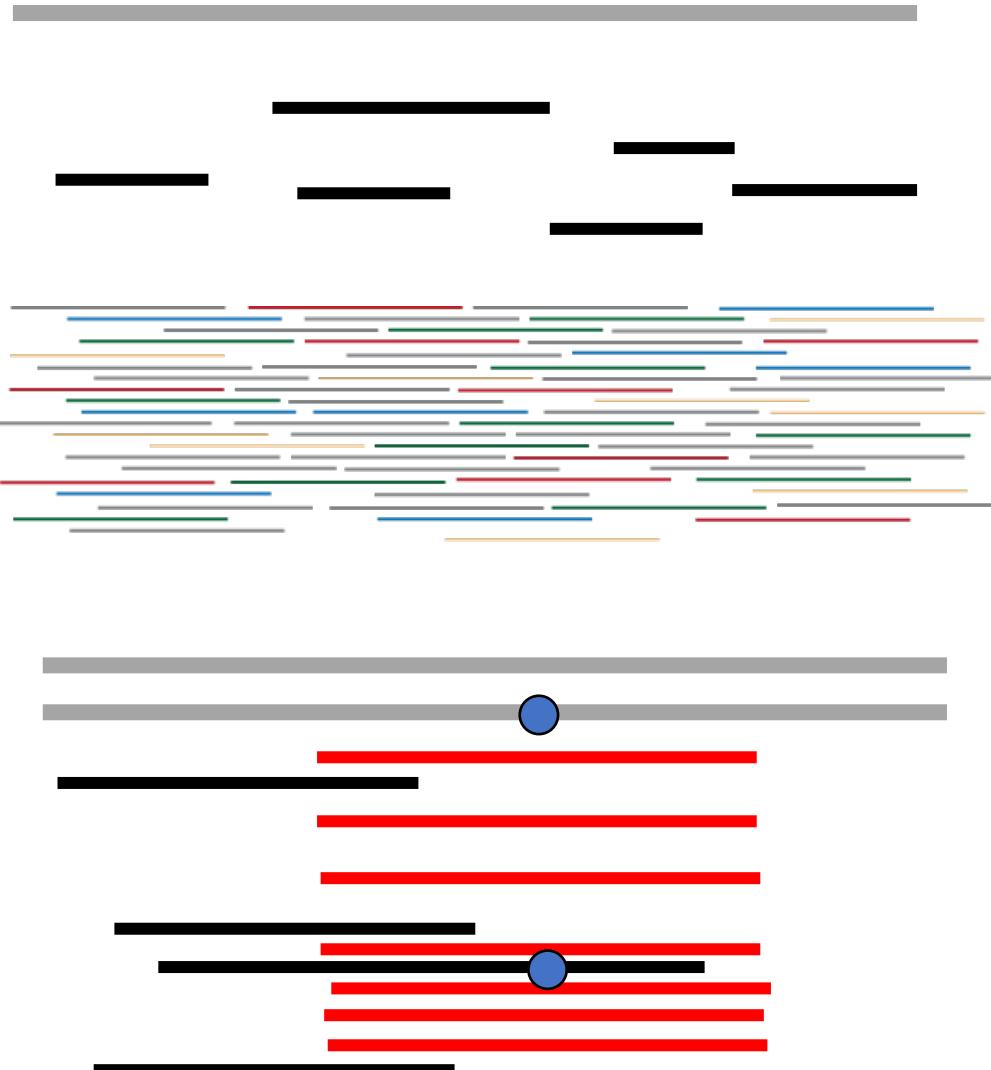
- Primary
- Secondary
- Supplementary → Split reads
- Duplicates
- Properly paired



Alignment types

- Primary
- Secondary
- Supplementary
- Duplicates
- Properly paired

Full genome sequencing



Randomly fragment DNA

PCR amplify

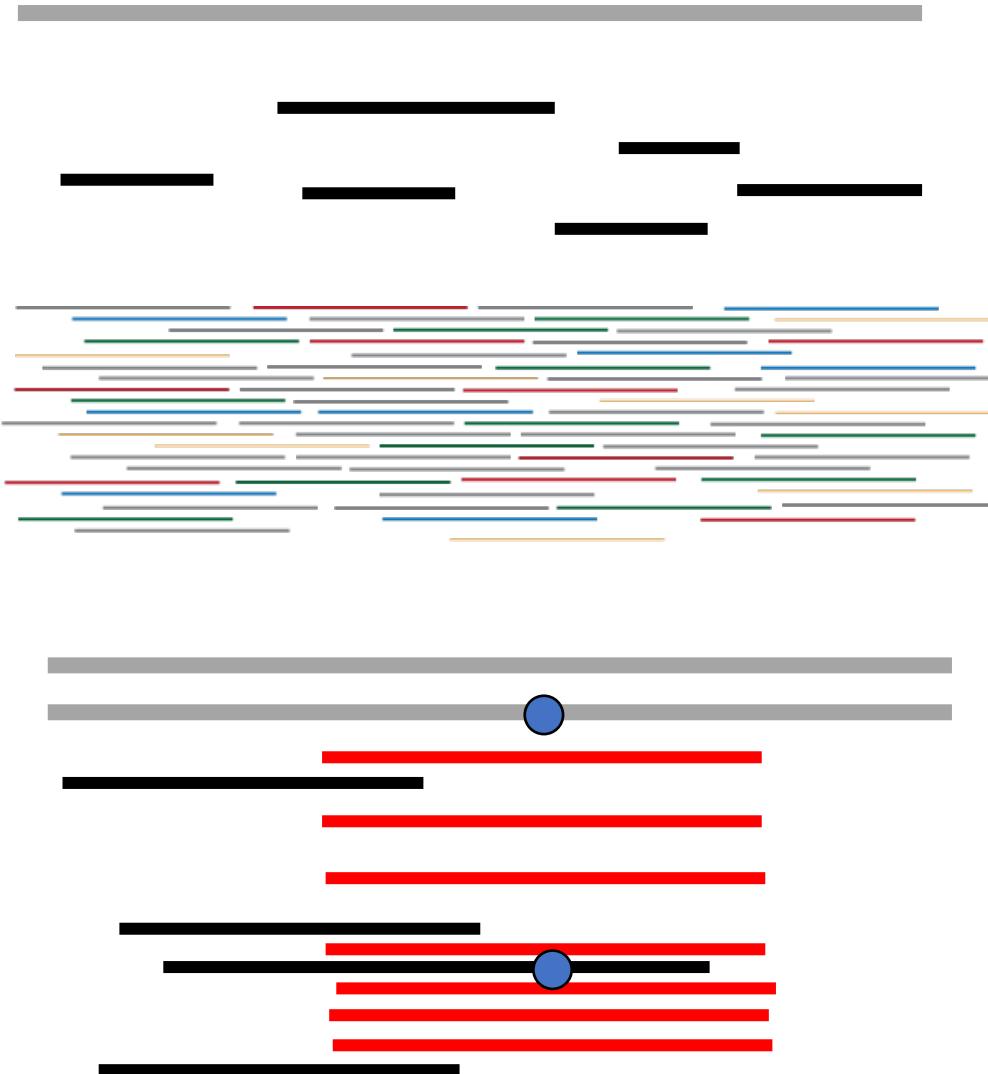
Map

False homozygote

Each DNA fragment should be:

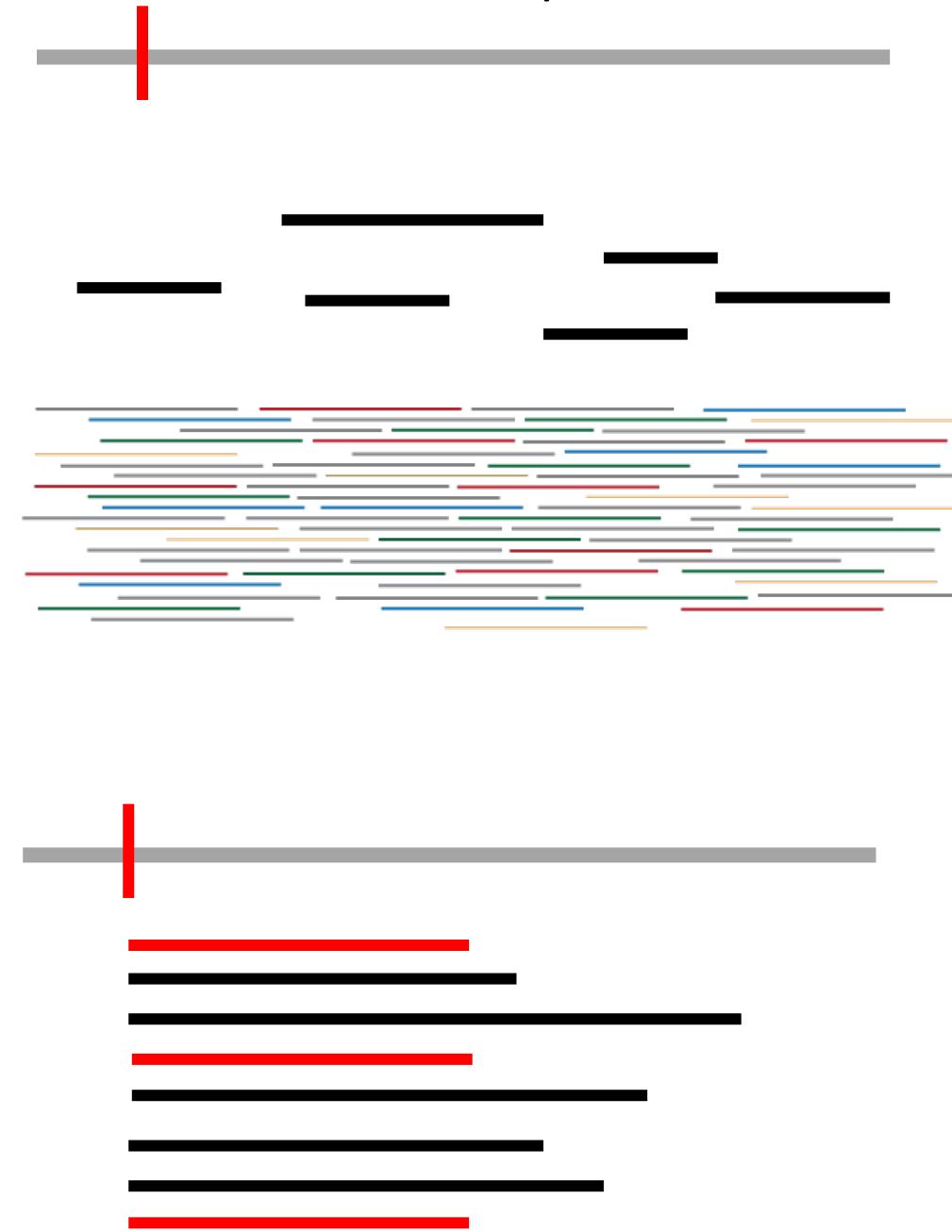
- Different length
- Different position

Full genome sequencing

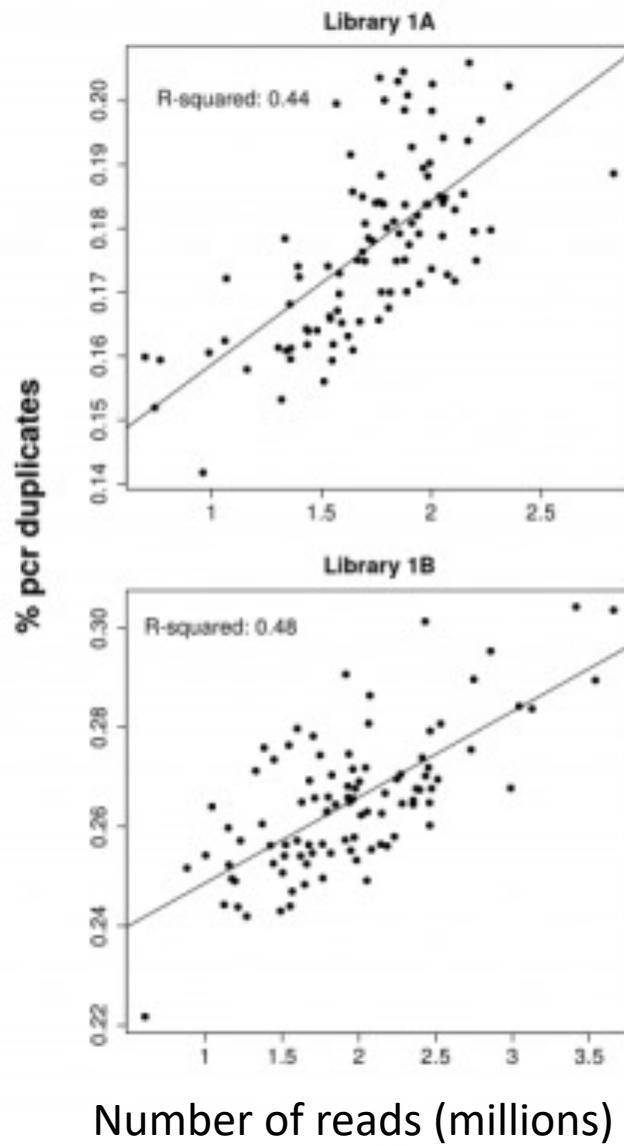


False homozygote

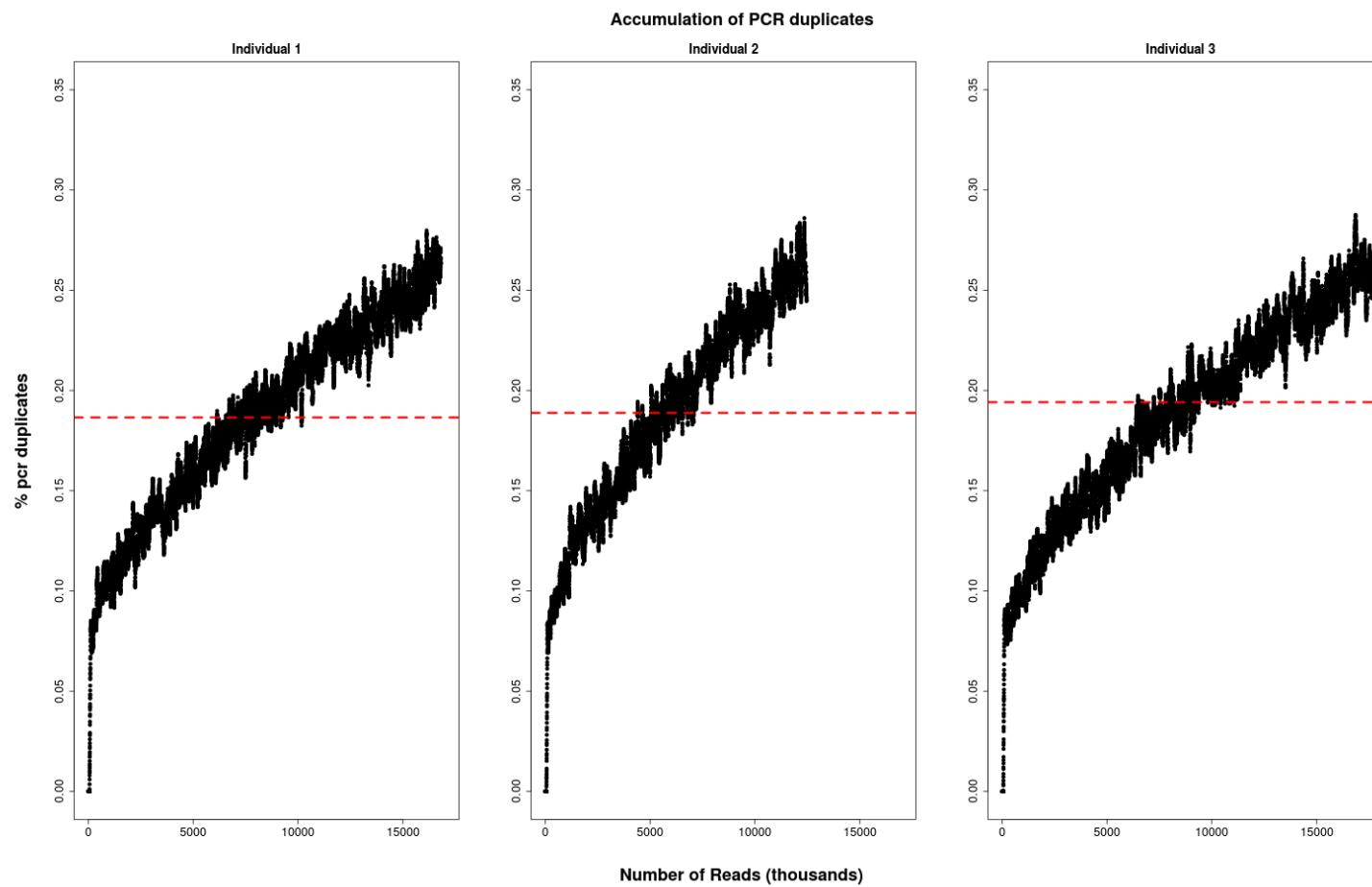
Rad seq



More depth = more duplicates



PCR duplicates accumulate as you process reads



For more info: <https://www.molecularecologist.com/2016/08/25/the-trouble-with-pcr-duplicates/>

Short Lecture: Variant calling

Simple in theory:

AGCTAGCCTAGTGACTTAAGCCTGA
AGCTAGCCTAGTGACGAAGCCTGA

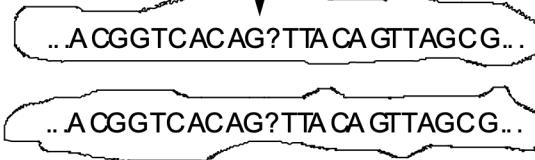
But in practice:

AGCTAGCCTAGTGACTTAAGCCTGA
AGCTAGCCTAGTGACGAAGCCTGA
AGCTAGCCTAGTGACGAAGCCTGA
AGCTAGCCTAGTGACGAAGCCTGA
AGCTAGCCTAGTGACGAAGCCTGA
AGCTAGCCTAGTGACGAAGCCTGA
AGCTAGCCTAGTGACGAAGCCTGA
AGCTAGCCTAGTGACGAAGCCTGA

What is genetic variation, what is an error?

Genotype likelihood

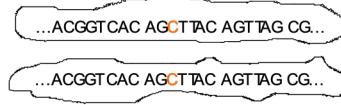
Position of known SNP in
the species, with alleles
T and **C**



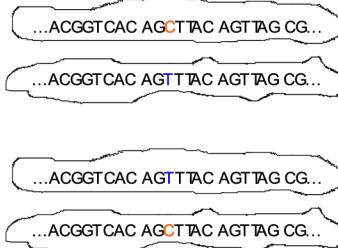
Two homologous
chromosomes within
an individual

The possible genotypes are:

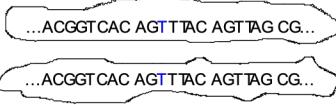
CC



CT or TC



TT



The data are: 4 reads covering that site,
and
the associated base quality scores

Read #	Read	Observed Base	PHRED-scaled base quality score
1	CAG C TTACA	C	32 (A)
2	ACAG C T	C	37 (F)
3	G T TTA	T	35 (D)
4	AG C TTACAG	C	33 (B)

Read #	Read	Observed Base	PHRED-scaled base quality score	3 hypotheses:	
1	CAGCTTACA	C	32 (A)	• H1: CC	What is the evidence in our data for each hypothesis?
2	ACAGCT	C	37 (F)	• H2: CT	
3	GTAA	T	35 (D)	• H3: TT	
4	AGCTTACAG	C	33 (B)		

$$L(H_i | D) = P(D|H)$$

$$L(H_i | D) = L(CC | D)$$

To calculate likelihood → the probability of observing our data (C,C,T,C) given our hypothesis (CC) is true

We assume reads are randomly sampled from each chromosome and there is an error rate

If $e = 0.00063$

Genotype	$P(R1 = C G = CC)$	$1-e$	0.99937
likelihoods if we have only one read that is C	$P(R1 = C G = CT)$	$\frac{1}{2} (1 - e) + \frac{1}{2} (e) = \frac{1}{2}$	0.5
	$P(R1 = C G = TT)$	e	0.00063

Read #	Read	Observed Base	PHRED-scaled base quality score	3 hypotheses:
1	CAGCTTACA	C	32 (A)	
2	ACAGCT	C	37 (F)	
3	GTAA	T	35 (D)	
4	AGCTTACAG	C	33 (B)	<ul style="list-style-type: none"> H1: CC H2: CT H3: TT

If we have only read 3 (T), basically the opposite results:

If $e = 0.00063$

$$P(R3 = T | G = CC) \quad e \quad 0.00063$$

$$P(R3 = T | G = CT) \quad \frac{1}{2}(1 - e) + \frac{1}{2}(e) = \frac{1}{2} \quad 0.5$$

$$P(R3 = T | G = TT) \quad 1 - e \quad 0.99937$$

What if we have one additional C?

Multiply the probabilities together to get total likelihood:

$$L(G=CC | B1=C, B3=T) = 0.99937 * 0.00063 = 0.0006$$

$$L(G=CT | B1=C, B3=T) = 0.5 * 0.5 = 0.25$$

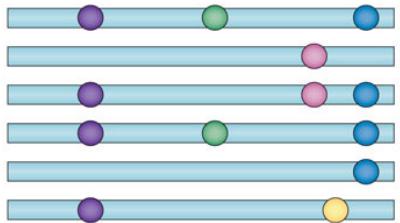
$$L(G=TT | B1=C, B3=T) = 0.99937 * 0.00063 = 0.0006$$

Short Lecture: Selection Scans

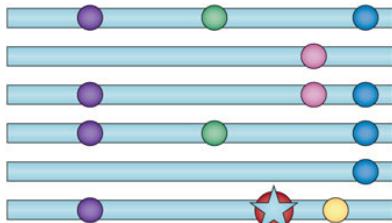
Types of selective sweeps

a Classic selective sweep

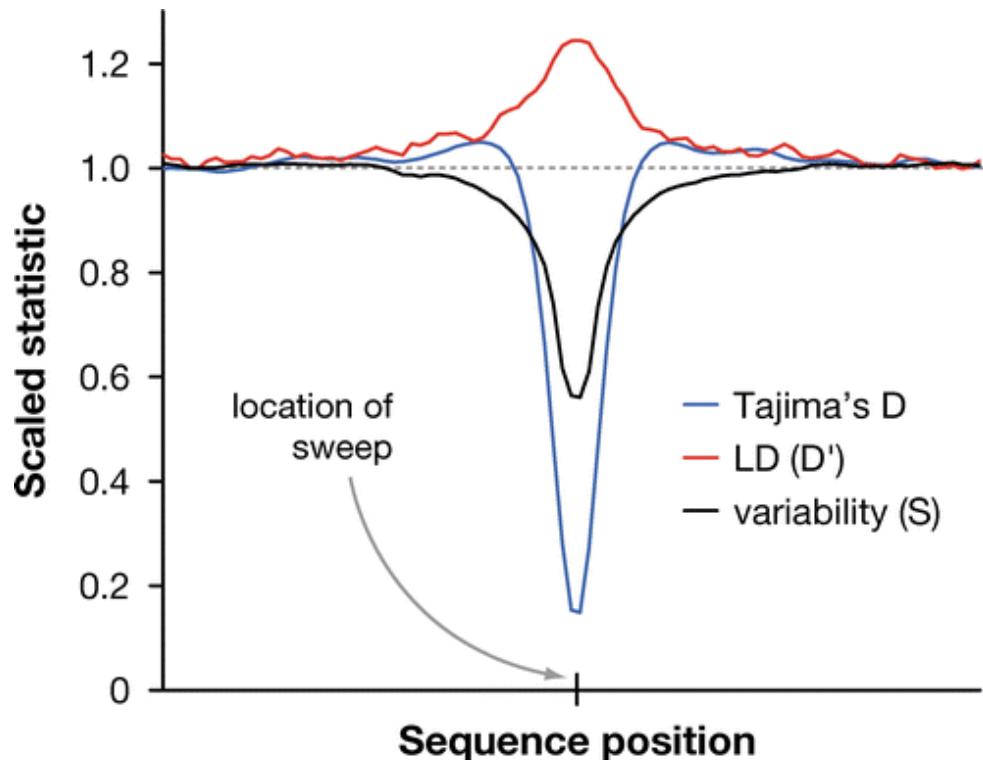
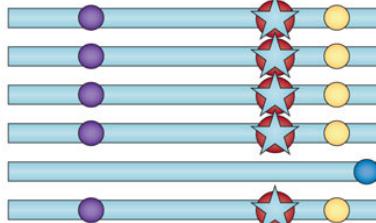
Neutral variation



An advantageous mutation arises

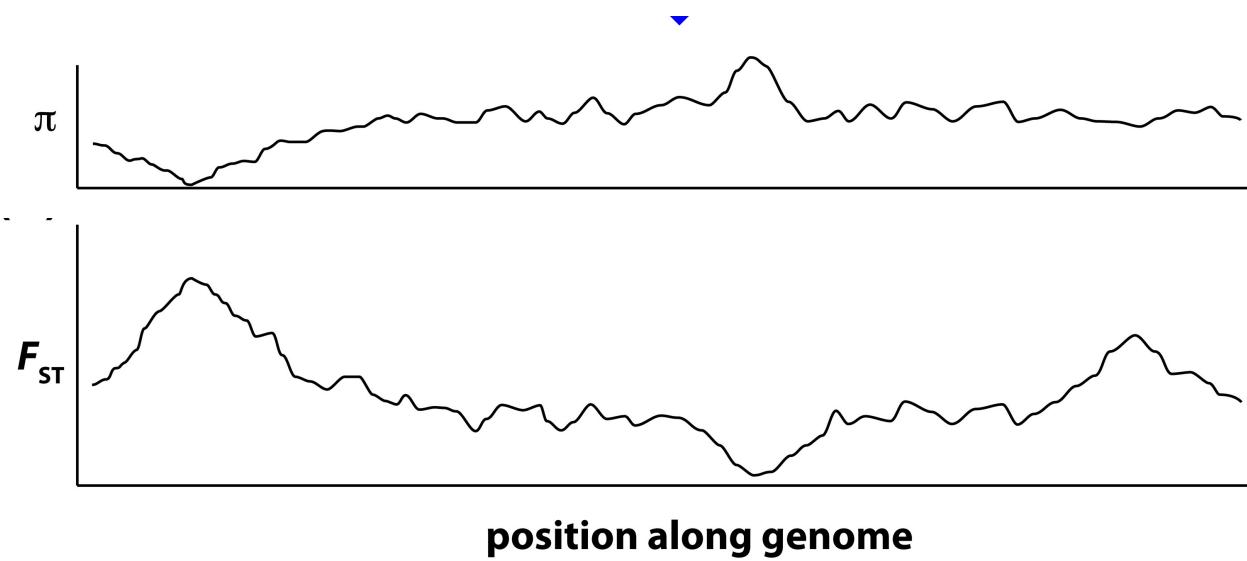


Over time, the advantageous mutation approaches fixation



Elevated: Fst, LD

Decreased: Tajima's D, genetic diversity



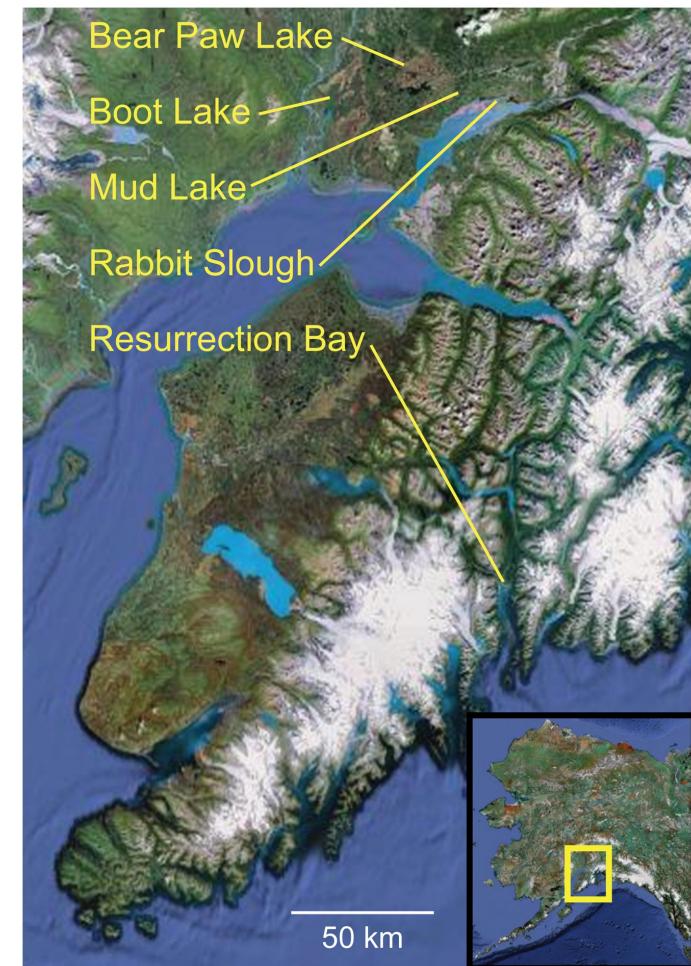
OPEN ACCESS Freely available online

PLOS GENETICS

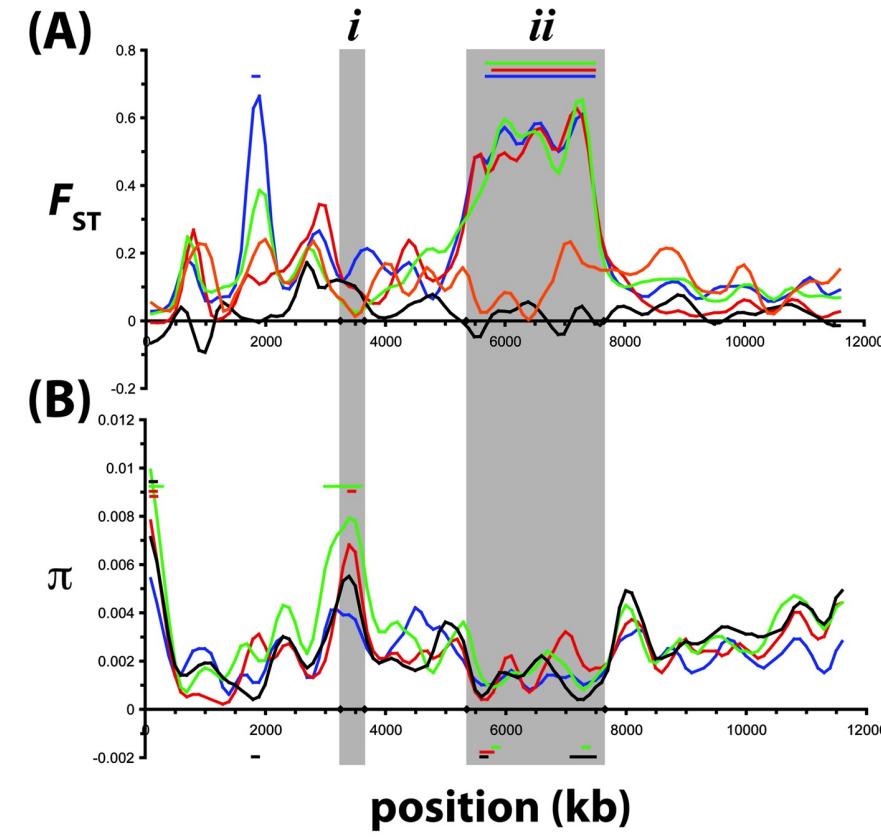
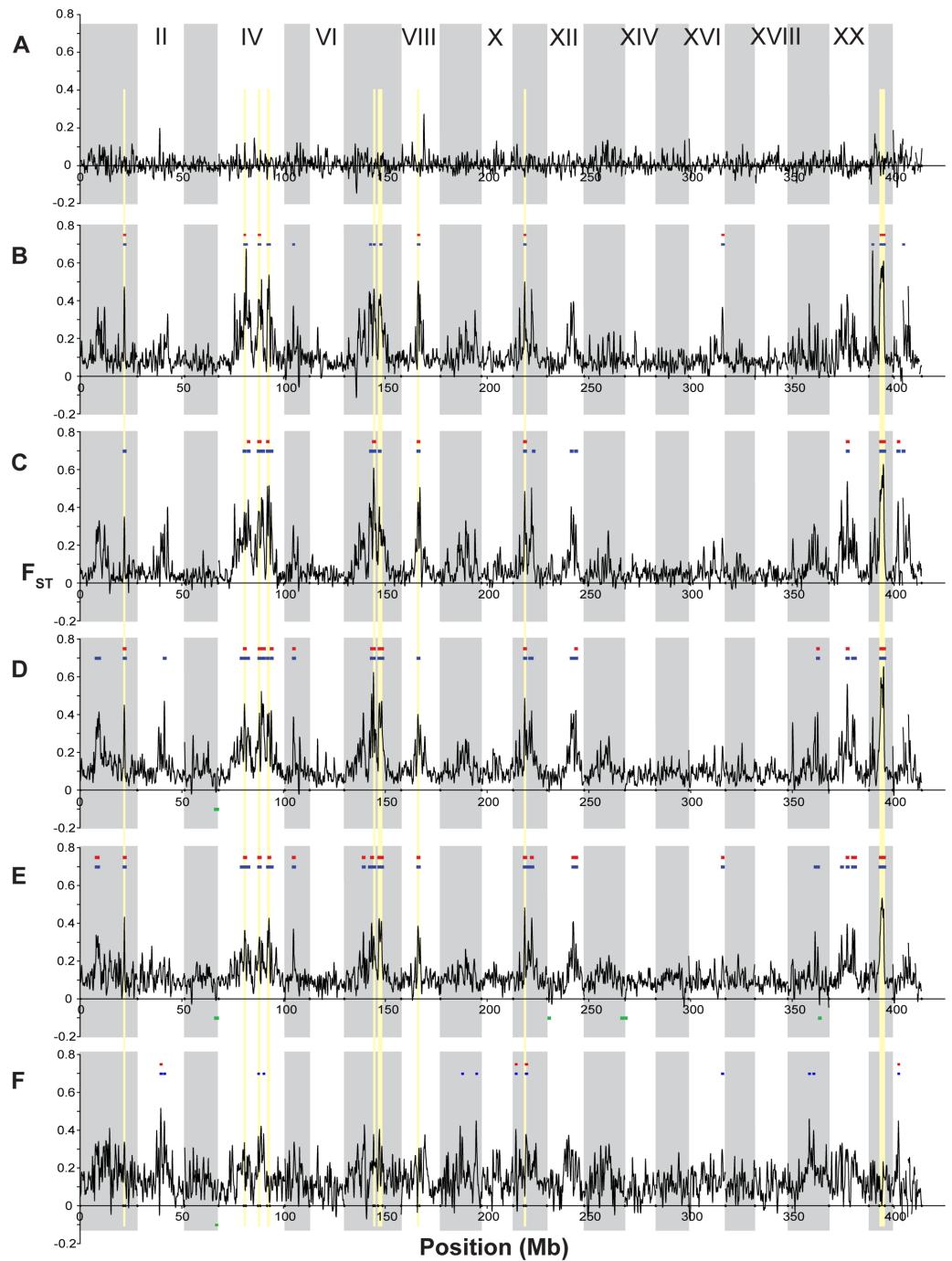
Population Genomics of Parallel Adaptation in Threespine Stickleback using Sequenced RAD Tags

Paul A. Hohenlohe^{1,3}, Susan Bassham^{1,3}, Paul D. Etter², Nicholas Stiffler³, Eric A. Johnson², William A. Cresko^{1*}

¹ Center for Ecology and Evolutionary Biology, University of Oregon, Eugene, Oregon, United States of America, ² Institute of Molecular Biology, University of Oregon, Eugene, Oregon, United States of America, ³ Genomics Core Facility, University of Oregon, Eugene, Oregon, United States of America

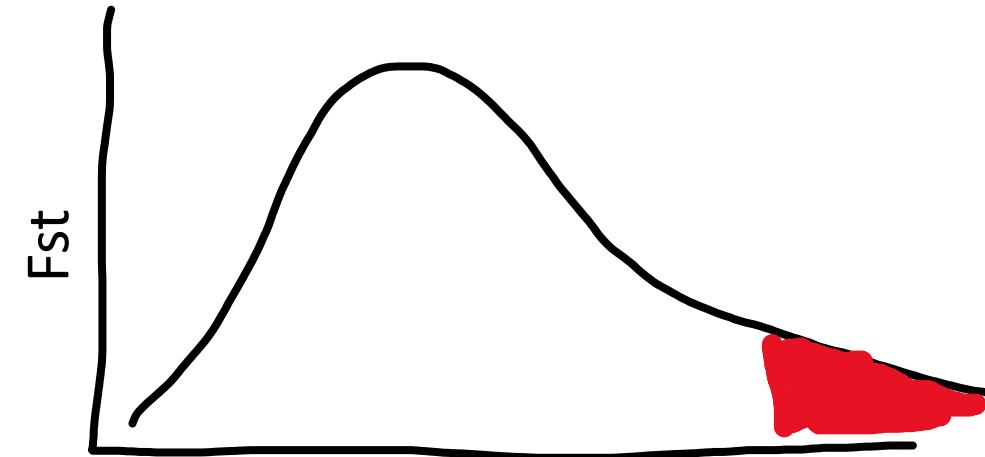
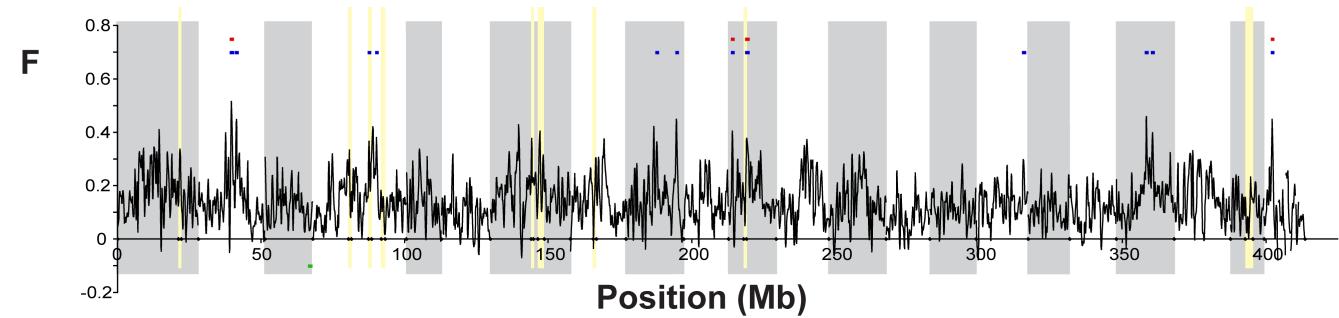


Hohenlohe, Phillips, and Cresko 2010



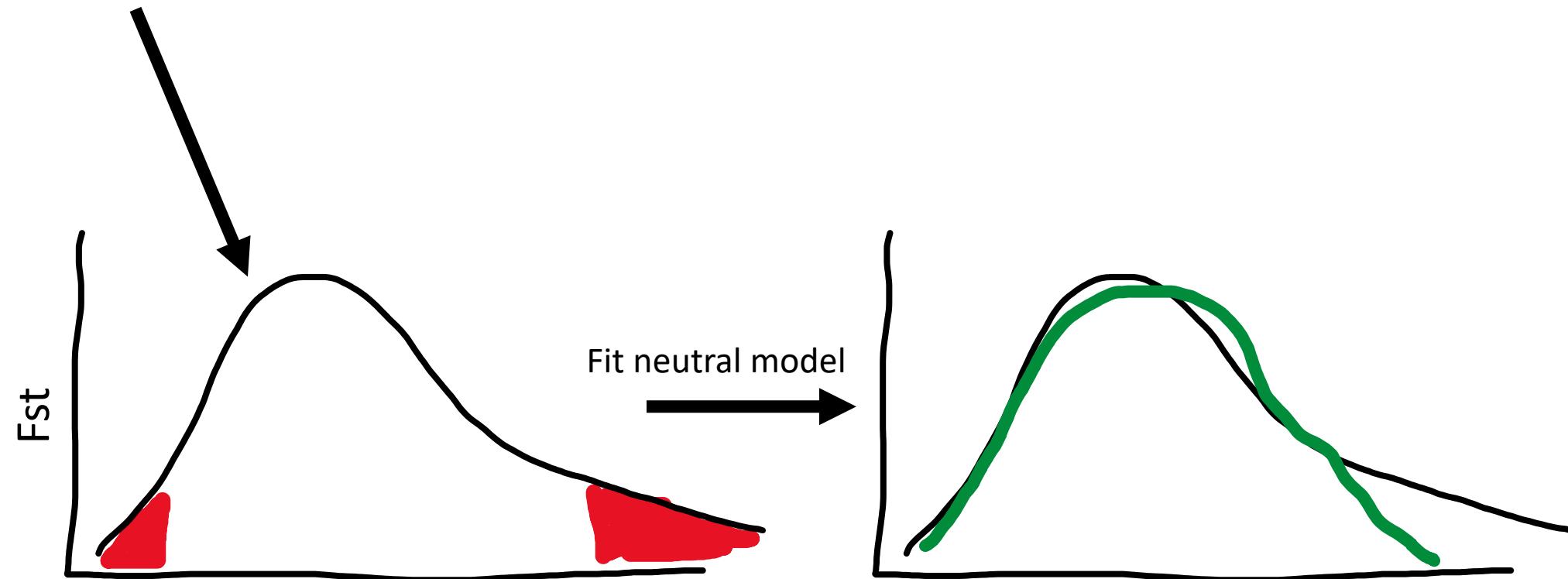
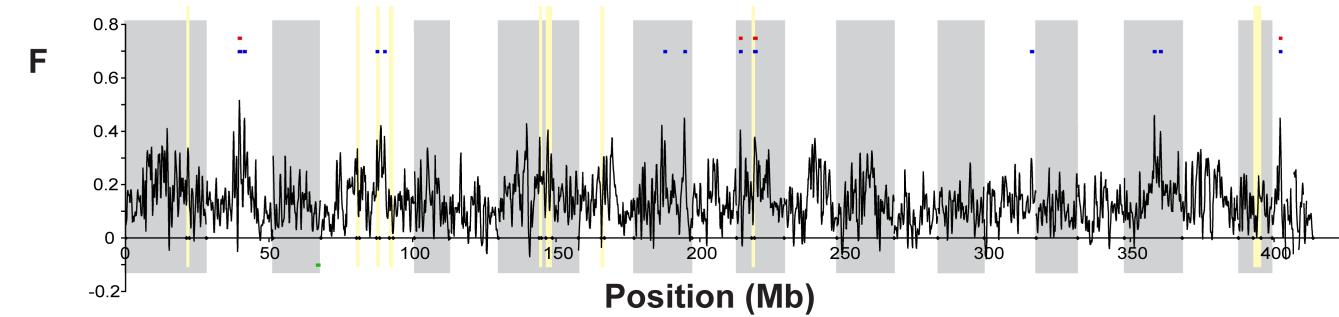
Detecting outlier regions

- What is the null distribution?
 - How to define “outliers”
 - Empirical: Tail of the distribution



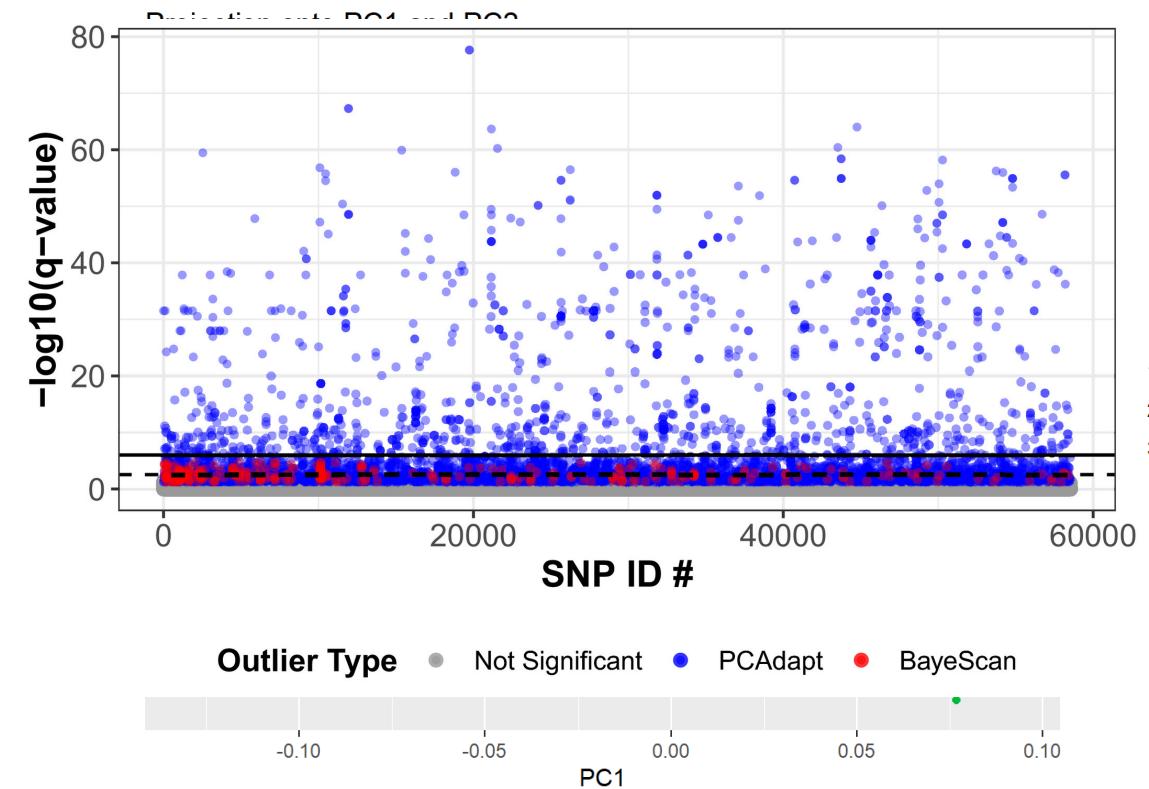
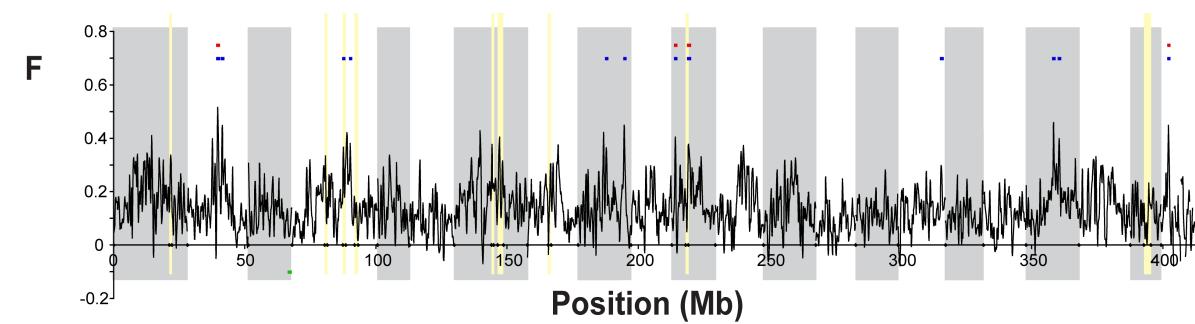
Detecting outlier regions

- What is the null distribution?
 - How to define “outliers”
 - Empirical: Tail of the distribution
 - OutFLANK



Detecting outlier regions

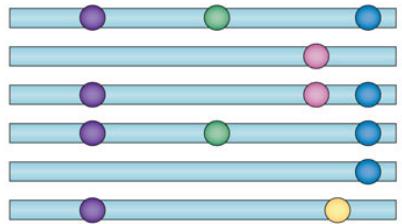
- Estimate covariance/structure among all populations, look for outliers after correction
- Baypass, bayenv2
- PCAdapt
 - PCA based
 - Don't need population assignments
 - markers excessively related to population structure are candidates for local adaptation



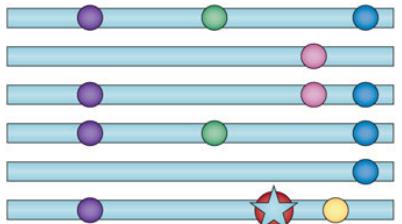
Types of selective sweeps

a Classic selective sweep

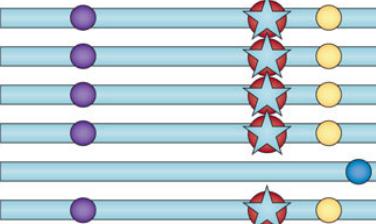
Neutral variation



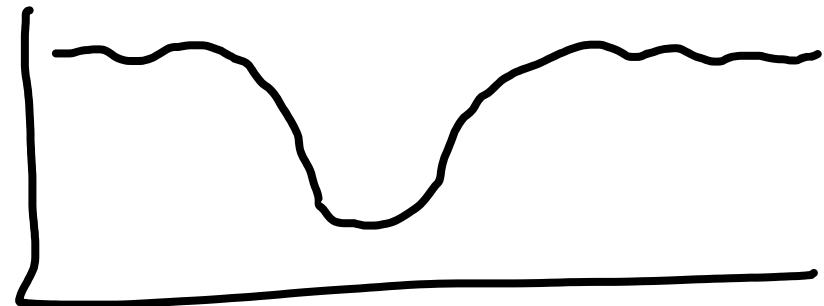
An advantageous mutation arises



Over time, the advantageous mutation approaches fixation

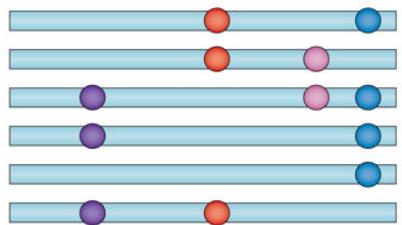


Haplotype diversity

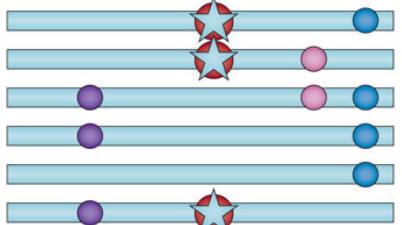


b Selection from standing variation

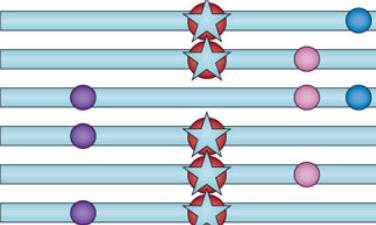
Neutral variation



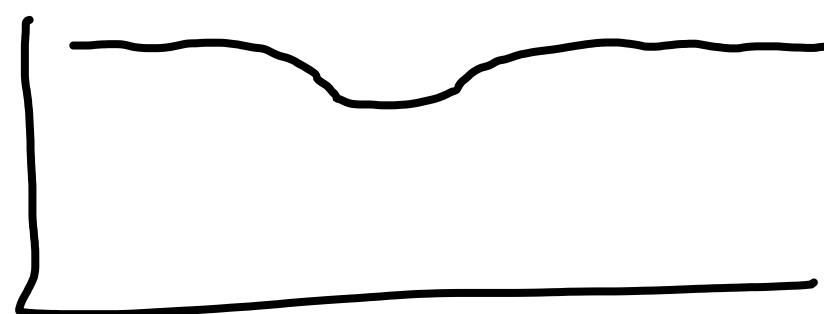
A variant becomes adaptive in a new environment



Over time, the advantageous mutation approaches fixation

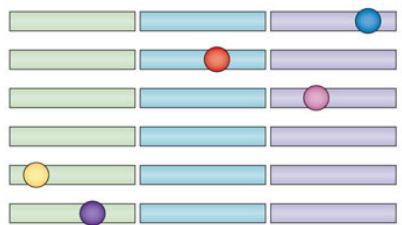


Haplotype diversity

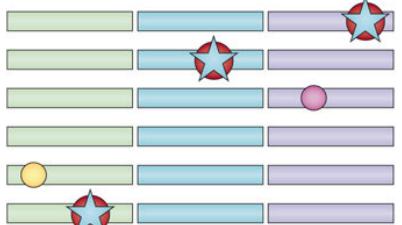


c Selection on a complex trait

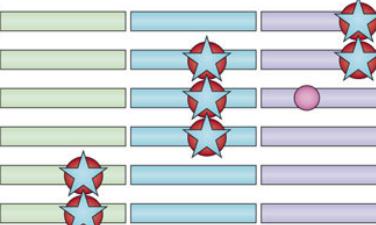
Neutral variation



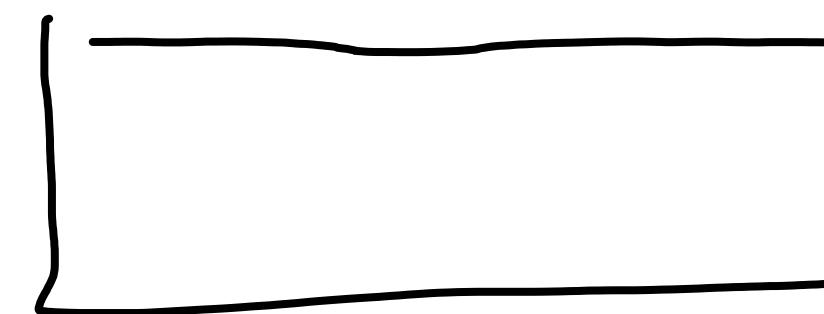
A set of variants becomes adaptive in a new environment



Over time, the set of variants becomes more common

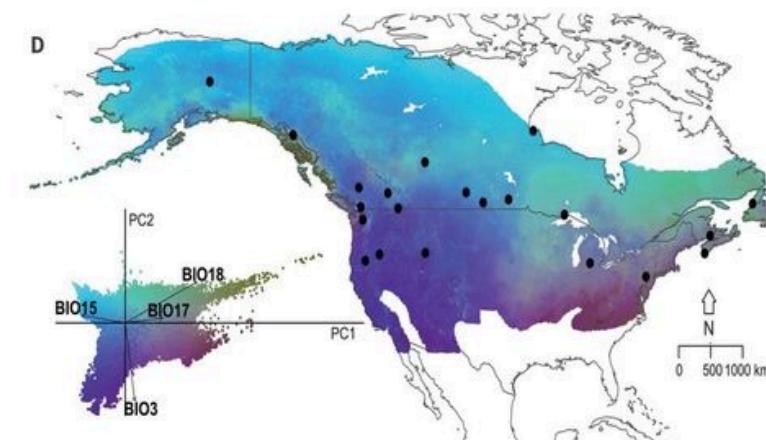


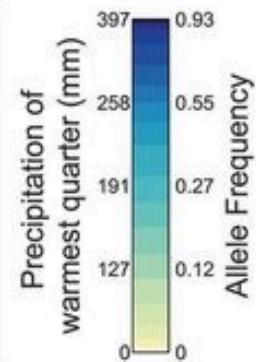
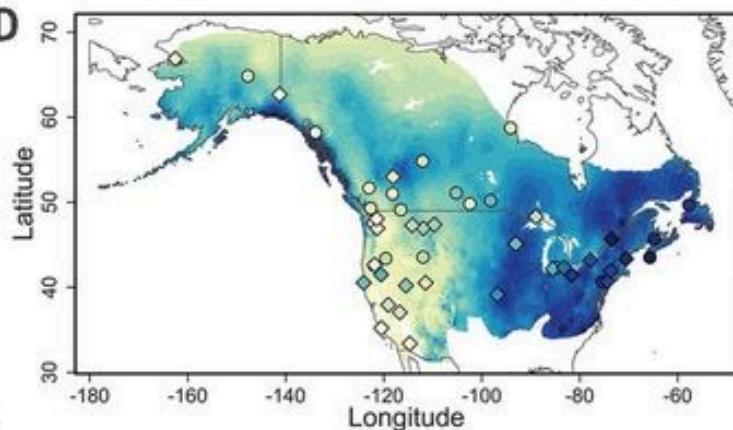
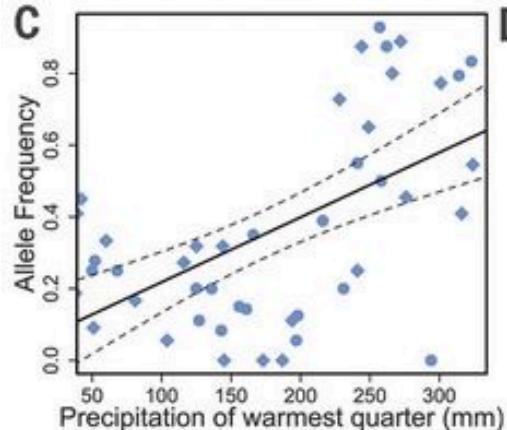
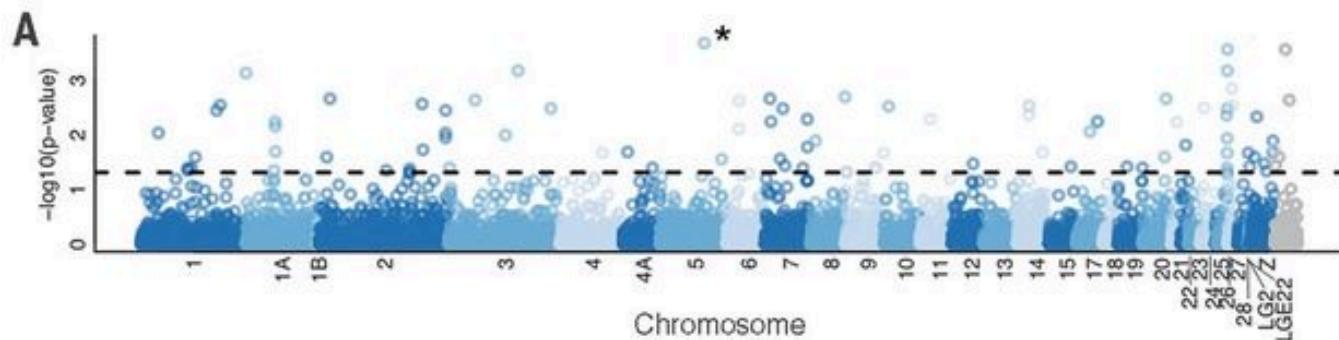
Haplotype diversity



Gene environment associations

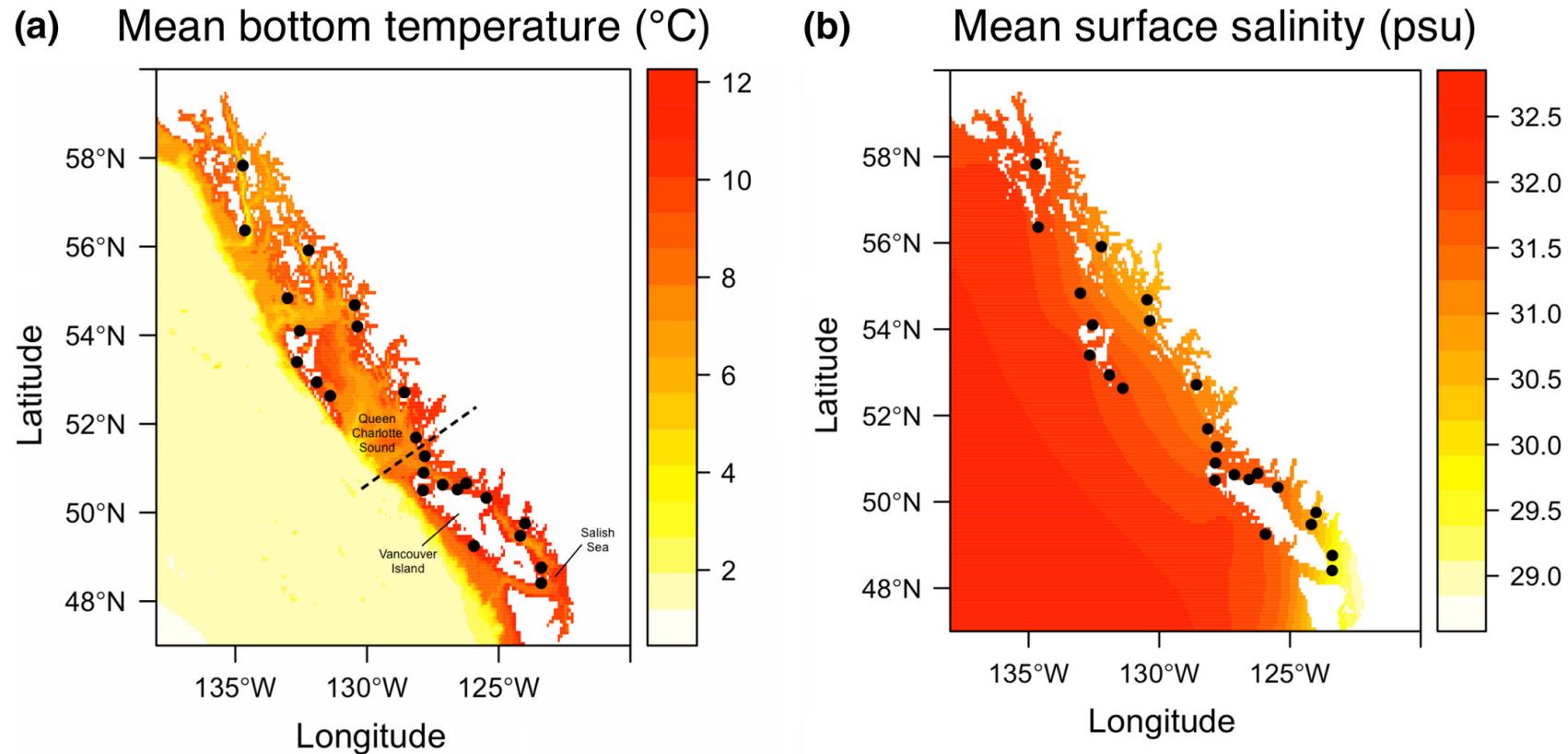
- Where is there a correlation between the environment and genetics?
- Important!
 - Must control for population structure!
- Methods:
 - Bayenv2
 - latent factor mixed models: LFMM2
 - Machine learning
 - Gradient Forests, random forests





Putatively adaptive genetic variation in the giant California sea cucumber (*Parastichopus californicus*) as revealed by environmental association analysis of restriction-site associated DNA sequencing data

Amanda Xuereb¹  | Christopher M. Kimber²  | Janelle M. R. Curtis³ |
Louis Bernatchez⁴  | Marie-Josée Fortin¹ 



Putatively adaptive genetic variation in the giant California sea cucumber (*Parastichopus californicus*) as revealed by environmental association analysis of restriction-site associated DNA sequencing data

Amanda Xuereb¹  | Christopher M. Kimber²  | Janelle M. R. Curtis³ | Louis Bernatchez⁴  | Marie-Josée Fortin¹ 

