# Genomic analysis in non-model organisms

2013 Workshop on Genomics
**Český Krumlov**

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

UNIVERSITY OF OREGON

# Outline for today's lecture

Genomic data and non-model organism research

RAD-seq for ecological & evolutionary genomics

Genomically enabling a non-model organism

*Stacks* software pipeline

1850

Evolution

1900

Conditions
of
Existence

Systematics
Ecology

Genetics

Unity
of
Type

Population
Genetics

Experimental
Embryology

1950

Modern
Synthesis

Molecular
Genetics

Evolutionary
Genetics

Developmental
Genetics

2000

functional evolutionary genomics

# Model organism research has been very important

Vertebrate **zygotes** or embryos
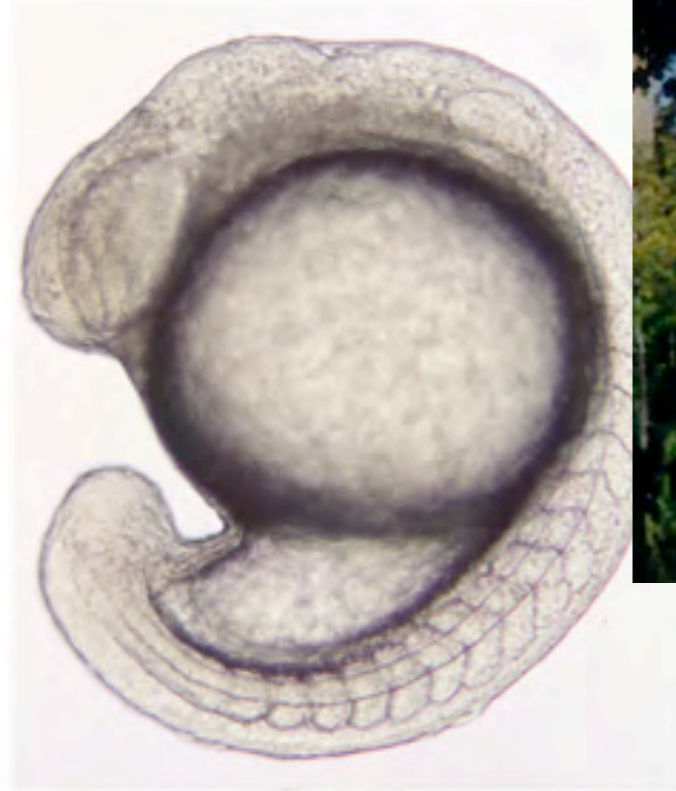


28 day human

19h zebrafish

Video by Don Kane

# Model organism research has been very important

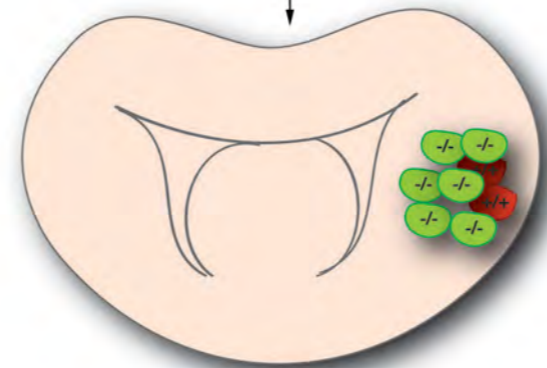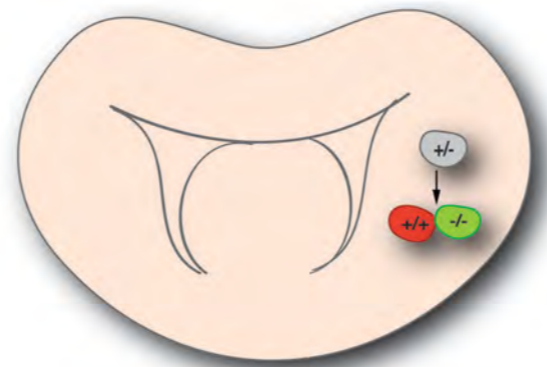Vertebrate **zygotes** or embryos



28 day human

19h zebrafish

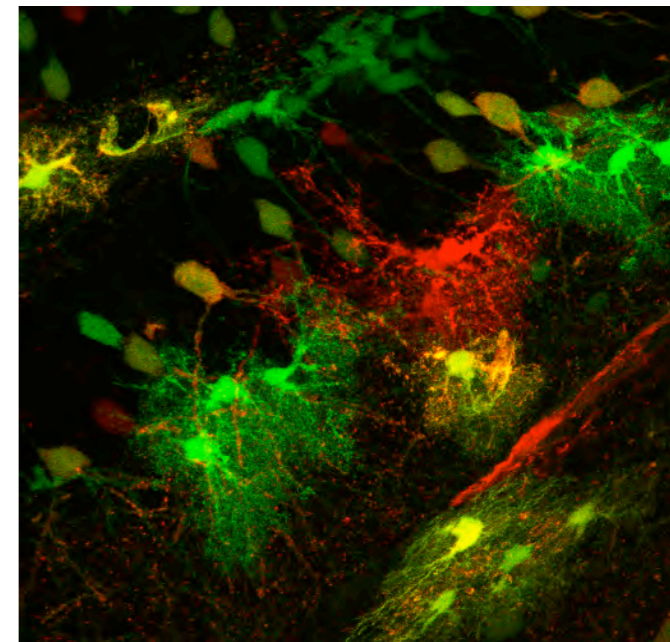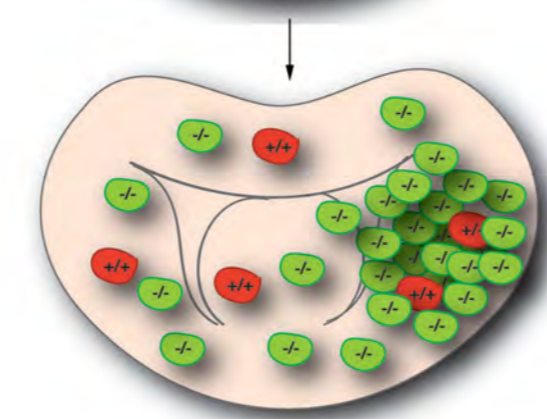# Studying brain cancer using somatic evolutionary genomics in a model organism

pre-cancerous

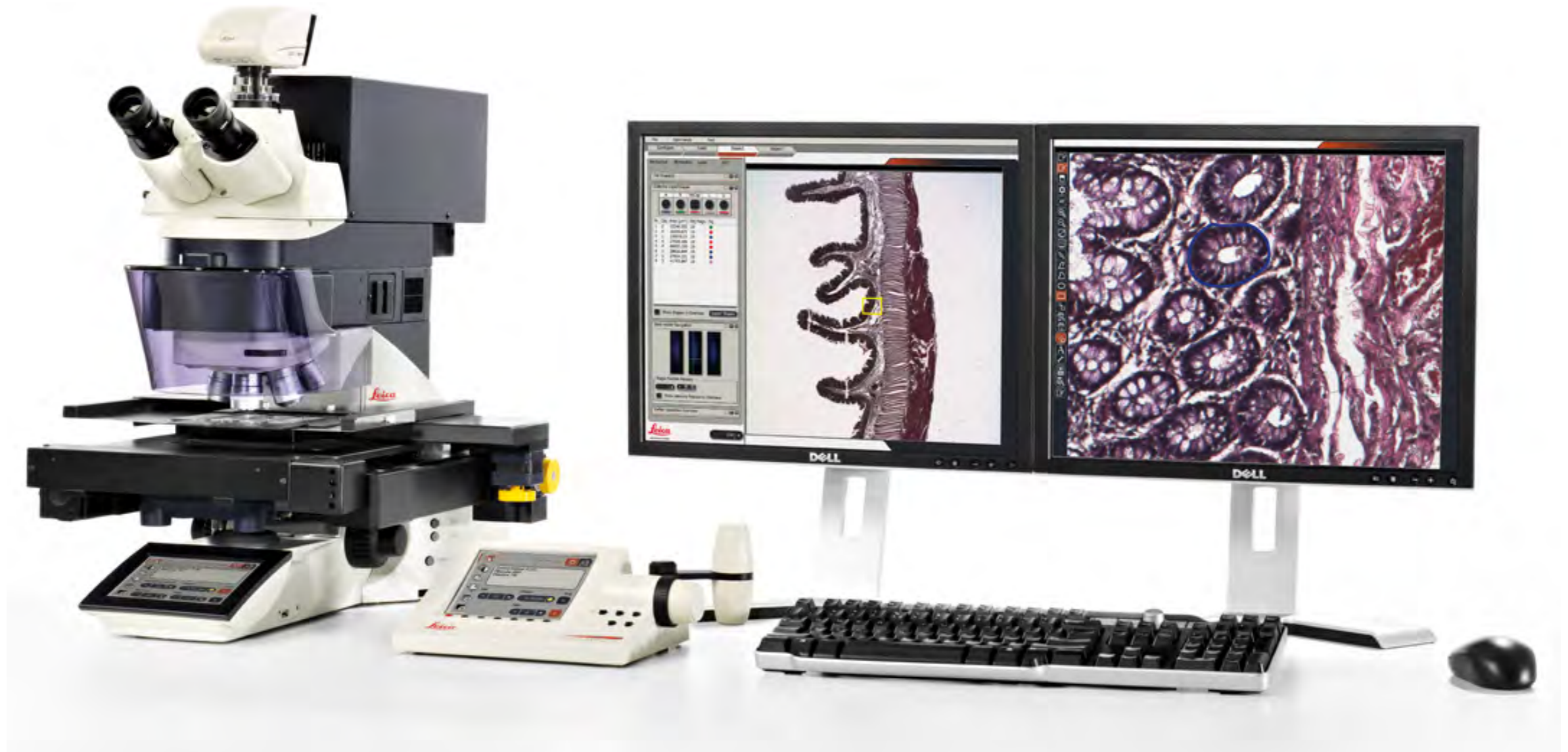tumor



Hui Zong, Rui Galvao, Julian Catchen and Susie Bassham

# Laser Capture Microdissection of cells

# Transcriptomic and genomic analysis of cells
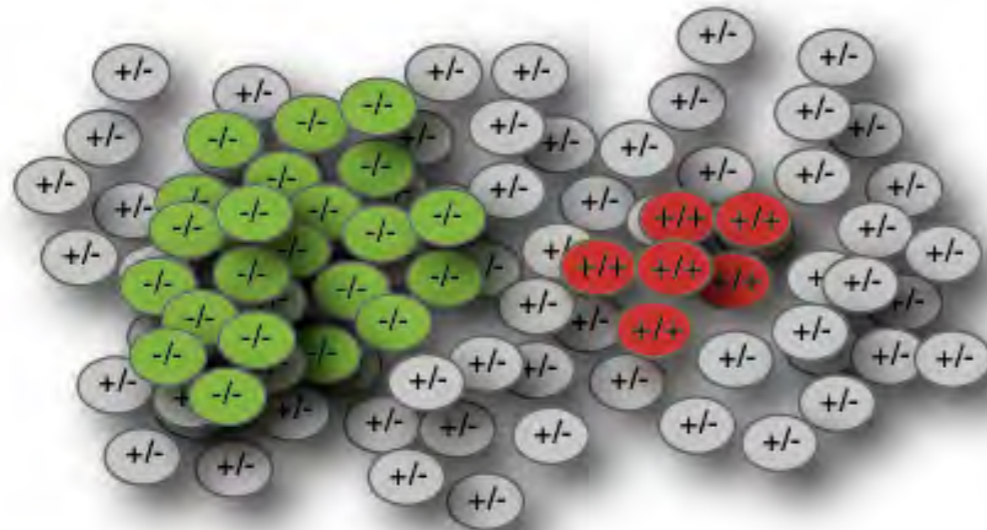


*Sequence cells here...*

*... and here*

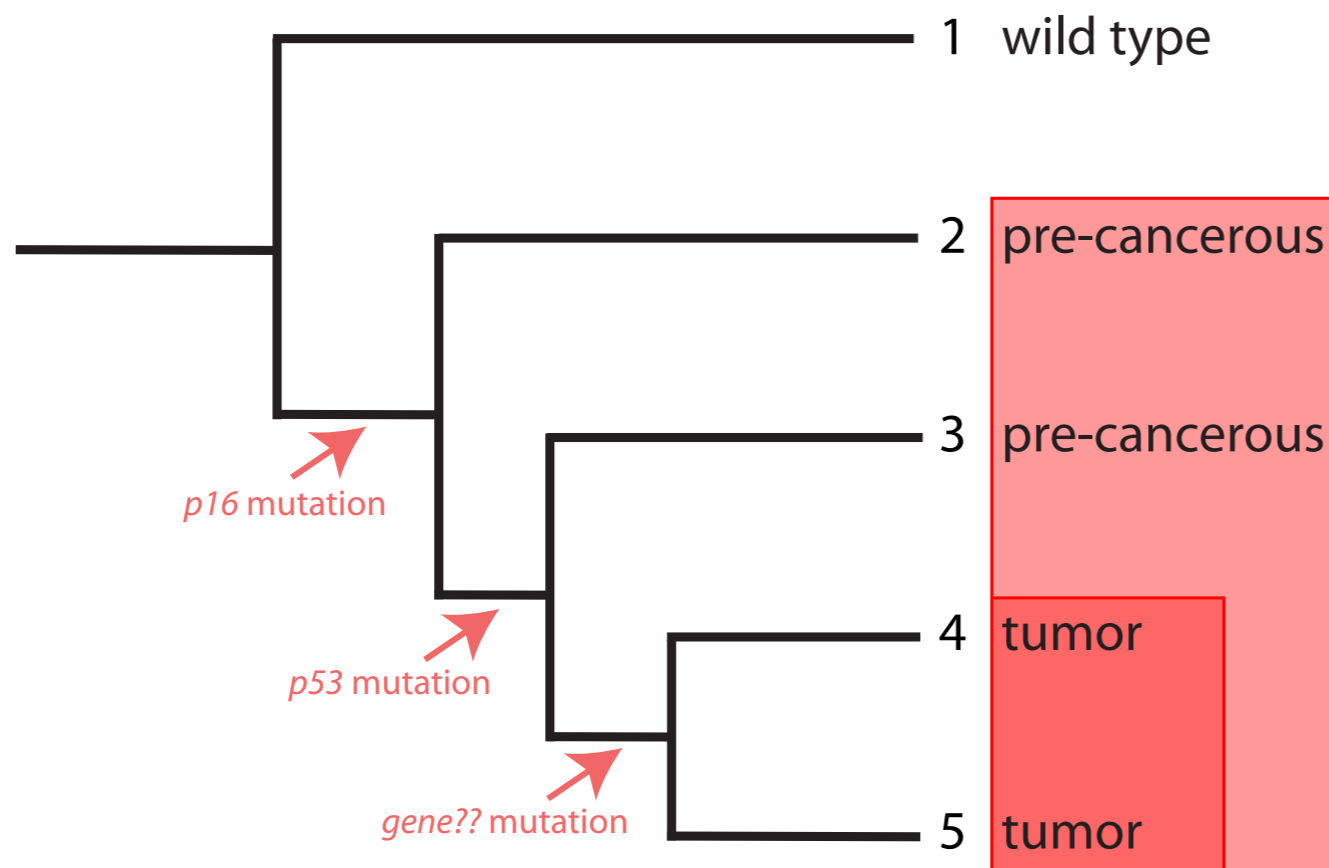# Transcriptomic and genomic analysis of cells



*Sequence cells here...*

*... and here*

# Multiple lines of genetic evidence for causative mutations

## Gene Mutations



| | |
|---|---|
| 1 | wild type |
| 2 | pre-cancerous |
| 3 | pre-cancerous |
| 4 | tumor |
| 5 | tumor |

*p16* mutation

*p53* mutation

*gene??* mutation

## Gene Expression



Expression Level

1

2

3

4

5

WT    *p16*    Gene??

# Genomic rearrangements in cancer cells

How can modern genomics improve studies of non-model organisms??

# How do the major differences among lineages evolve?

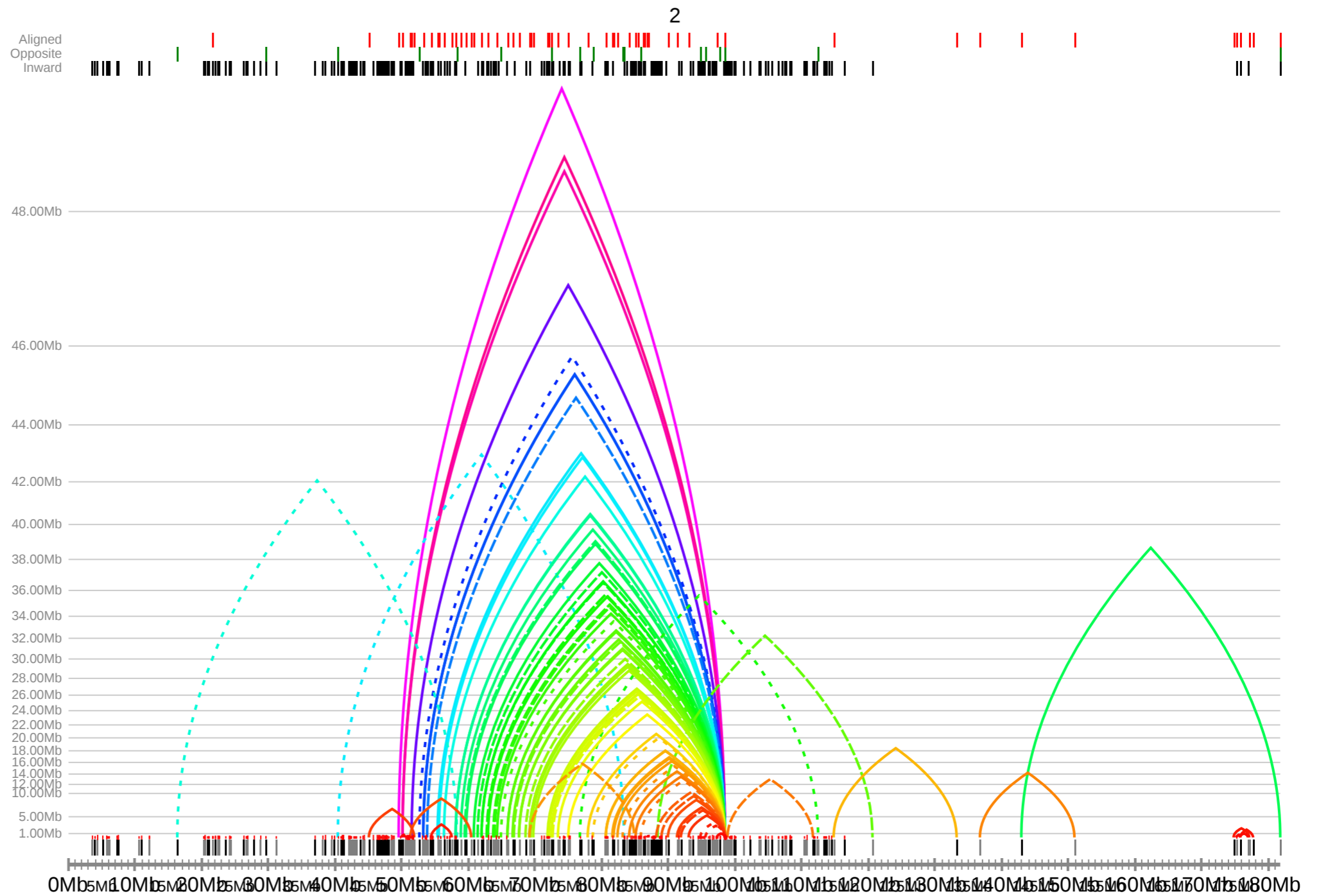# How are organisms related to one another?



Phylogeography of pocket
gophers using mtDNA

Avise, 1979

# 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

# 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 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 unordered markers

# Population genomics of ordered markers

## Genomic architecture:

- Distribution adaptive variation across the genome

- Correlations among genomic regions (linkage disequilibrium)

- Interactions among genomic regions (e.g. epistasis)

- Recombination rates and chromosomal inversions



Hohenlohe et al (2010) *Int J Plant Sci 171:1059*

# How do we 'genomically enable' research on non-model organisms?

1. Genetic Markers & Maps

2. Physical Maps

3. Transcriptomes

4. Gene Expression Analyses

# In the field and in the lab until a few years ago....

# The open source genomics breakthrough

*Next generation sequencing, high performance computing
and new analytical approaches have
fundamentally changed the scope of
studies of non-model organisms*

Should we just sequence everything?

# Why not sequence the entire genome??

- Still prohibitively expensive for many studies
    - Human height GWAS; over 15,000 individuals assayed
    - Identified many new regions contributing to the variation
    - Still only identified a fraction of the heritability
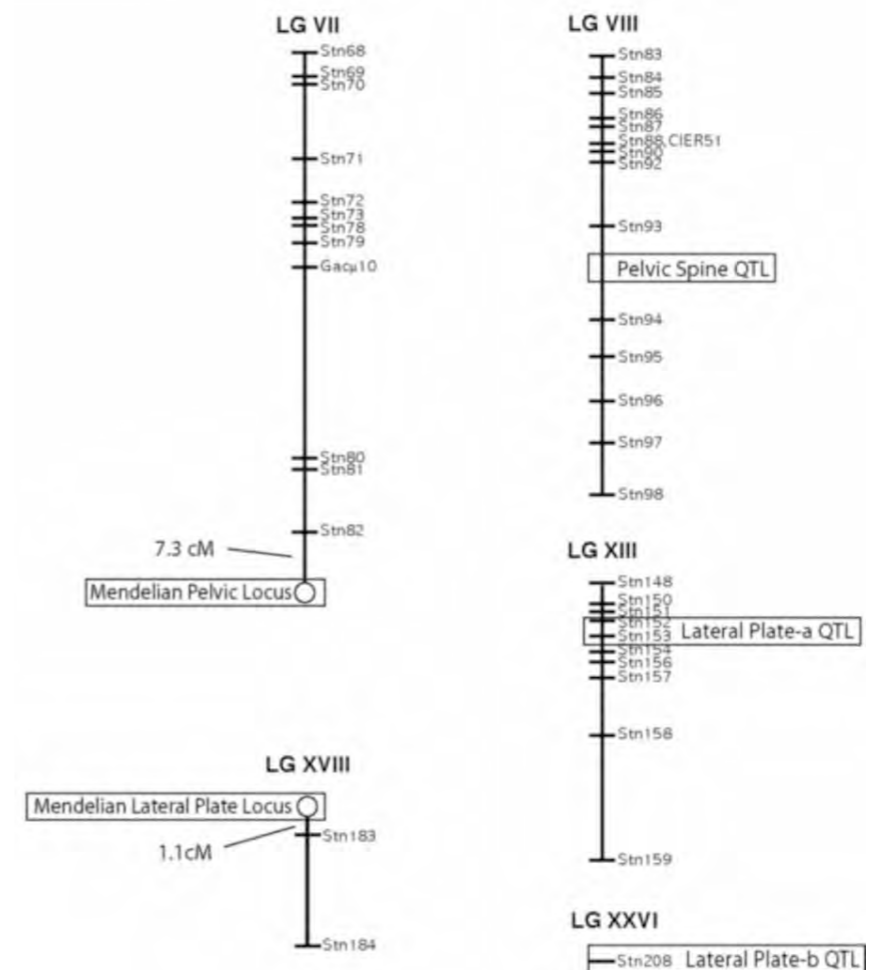
- For many studies a full sequence isn't necessary
    - the genomes of many organisms are organized in linkage blocks
    - well spaced markers will provide the necessary coverage
    - the cost of genotyping will almost always be a fraction of full sequencing

- Genetic maps are very useful in genomic studies
    - a high density genetic map can facilitate genome assembly
    - genomes may be segregating a lot of structural variation

# Alternative approach -
# Reduced representation NGS for genotyping

- Focus the sequencing on a homologous set of tags spread throughout the genome

- Can lead to the simultaneous identification and typing of single nucleotide polymorphisms (SNPs)

- The cost will always be a fraction of the cost of resequencing the genome
  - i.e. 1% genome coverage will be less than 1% the cost
  - often the coverage is more even than whole genome sequencing

- Can allow thousands of genomes to be assayed in just a few weeks

- WHY NOT - some cases complete genomic sequence is necessary
  - when linkage disequilibrium blocks (LD) are very short
  - Inferring patterns of LD may be easiest with full sequences

# Different flavors of Reduced Representation Library (RRL) Sequencing for genotyping

- Common acronyms
    - **RRL** - **R**educed **R**epresentation **L**ibrary
    - **GBS** - **G**enotyping **B**y **S**equencing
    - **CRoPS** - **C**omplexity **R**eduction **o**f **P**olymorphic **S**equences
    - **MSG** - **M**ultiplex **S**hotgun **G**enotyping
    - **RAD** - **R**estriction site **A**ssociated **D**NA

- All rely on restriction enzyme digestion

- RRL, CRoPS, MSG and GBS use one or two restriction enzymes only

- RAD uses an extra shearing step to capture all restriction sites

- Incorporation of barcodes on adaptors for multiplexing

- Aligned against a reference genome or assembled *de novo*

- Statistical issues
    - new level of sampling variation (sequencing in addition to biological)
    - sequencing error and problems for aligning or clustering

# What is RAD-seq?

(Restriction-site Associated DNA)



Illumina

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]

[1]Institute for Molecular Biology, University of Oregon, Eugene, Oregon 97403, USA, [2]Center for Ecology & Evolutionary Biology, University of Oregon, Eugene, Oregon 97403, USA; [3]Institute of Neuroscience, University of Oregon, Eugene, Oregon 97403, USA

2008

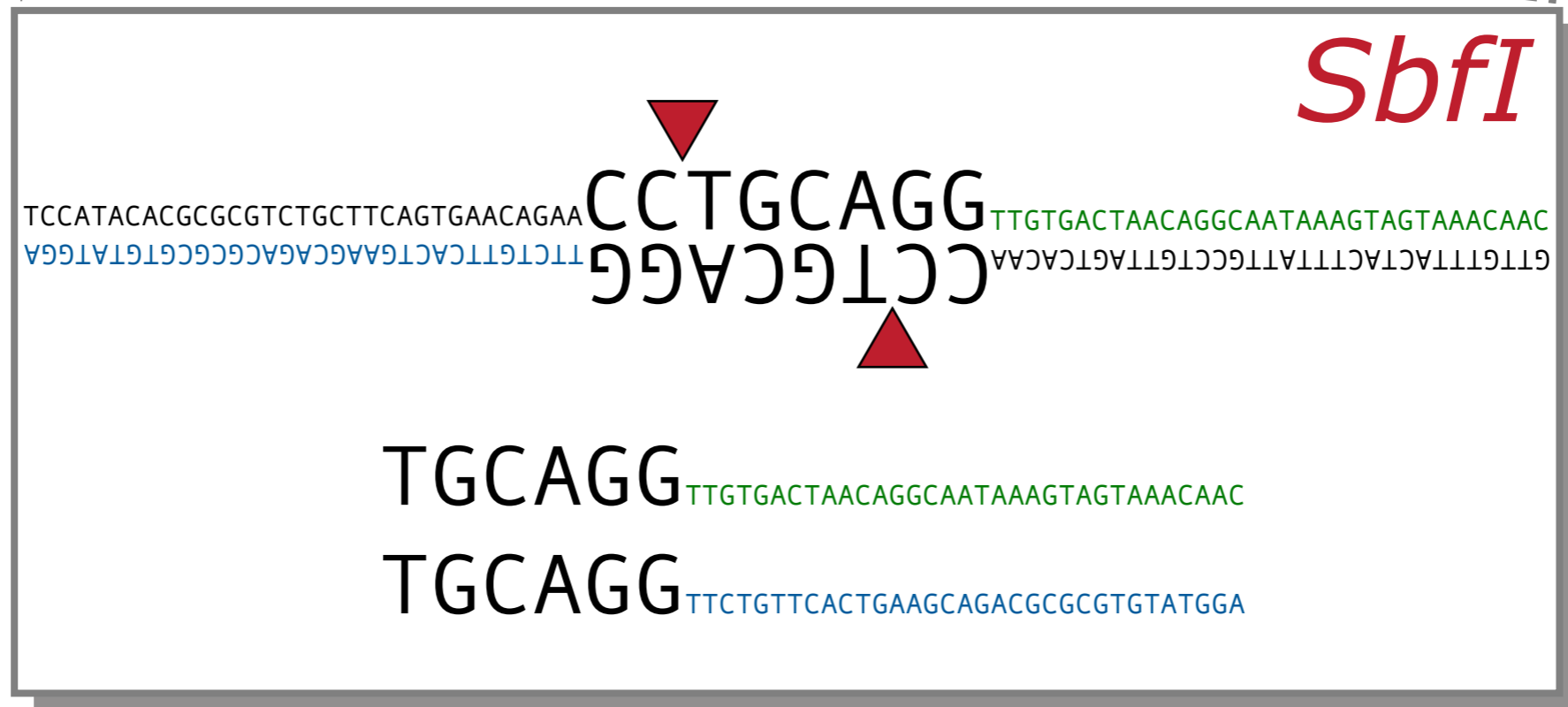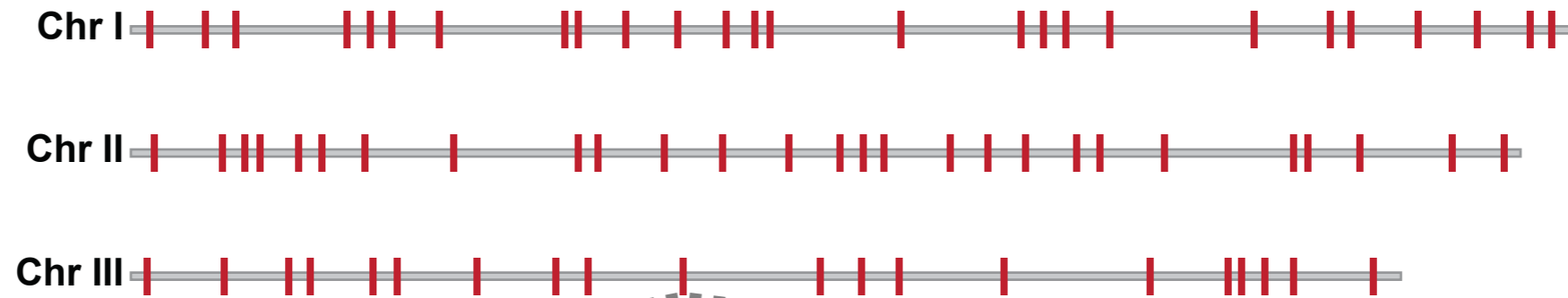OPEN ACCESS Freely available online                                    PLoS one

Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers

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

1 Institute of Molecular Biology, University of Oregon, Eugene, Oregon, United States of America, 2 Phiogenix, Eugene, Oregon, United States of America, 3 The Center for Ecology and Evolutionary Biology, University of Oregon, Eugene, Oregon, United States of America
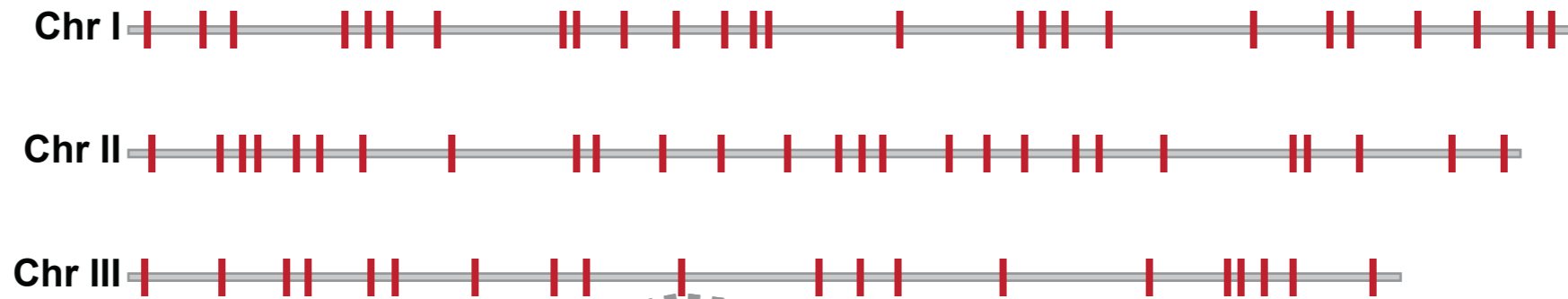
# What is RAD-seq?

(Restriction-site Associated DNA)

# What is RAD-seq?

(Restriction-site Associated DNA)

# What is RAD-seq?

(Restriction-site Associated DNA)

Chr I

22,830 *SbfI* sites in threespine stickleback

~ 45,000 RAD-Tags

HiSeq Illumina Lane:
160 million reads, 96 barcoded individuals

① *SbfI*

TCCATACACGCGCGTCTGCTTCAGTGAACAGAA CCTGCAGG TTGTGACTAACAGGCAATAAAGTAGTAAACAAC
AGGTATGTGCGCGCAGACGAAGTCACTTGTCTT GGACGTCC AACACTGATTGTCCGTTATTTCATCATTTGTTG

TGCAGG TTGTGACTAACAGGCAATAAAGTAGTAAACAAC

TGCAGG TTCTGTTCACTGAAGCAGACGCGCGTGTATGGA

② *SbfI*

TCCATACACGCGCGTCTGCTTCAGTGAACAGAA CCTGCAGG TTGTGACTAACAGGCAATGAAGTAGTAAACAAC
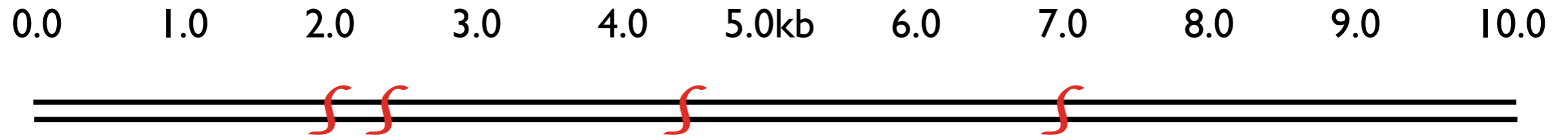AGGTATGTGCGCGCAGACGAAGTCACTTGTCTT GGACGTCC AACACTGATTGTCCGTTACTTCATCATTTGTTG

TGCAGG TTGTGACTAACAGGCAAT G/A AAGTAGTAAACAAC

TGCAGG TTCTGTTCACTGAAGCAGACGCGCGTGTATGGA

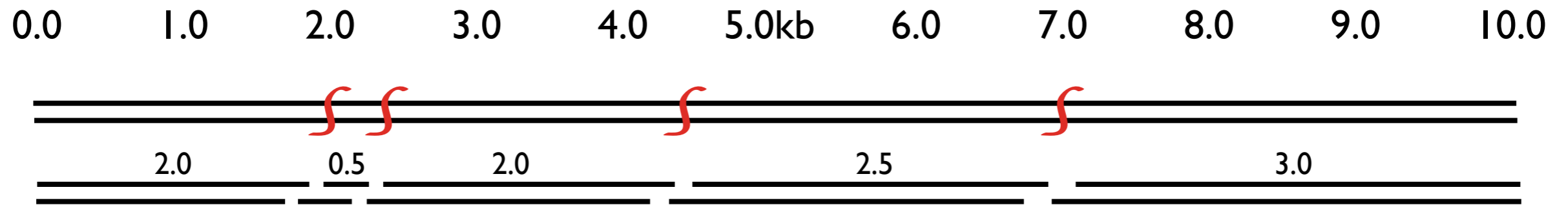# Restriction Enzyme (RE) digestion and first adaptor ligation
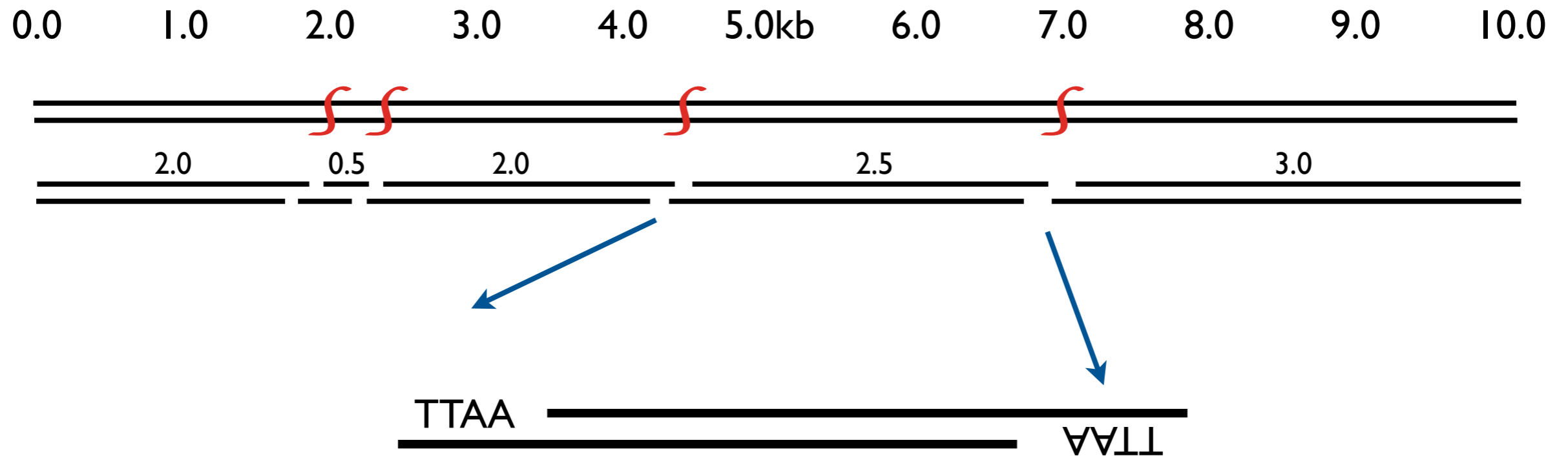
0.0　　1.0　　2.0　　3.0　　4.0　　5.0kb　　6.0　　7.0　　8.0　　9.0　　10.0

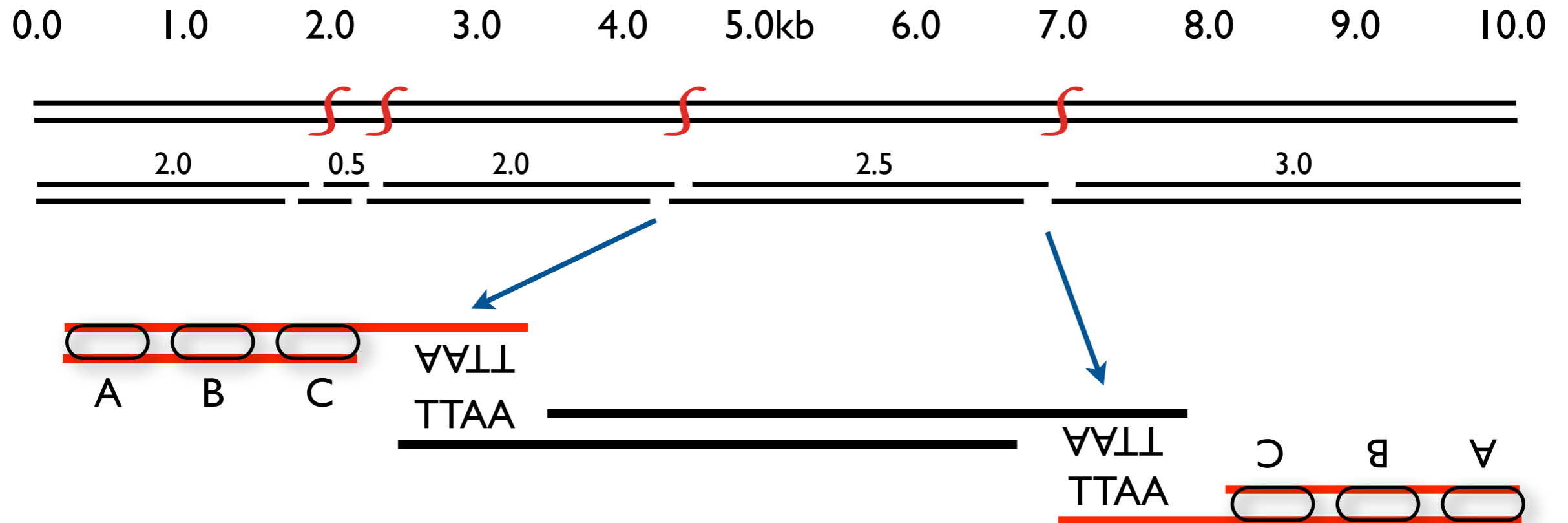# Restriction Enzyme (RE) digestion and first adaptor ligation

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

# Restriction Enzyme (RE) digestion and first adaptor ligation

# Restriction Enzyme (RE) digestion and first adaptor ligation

# Restriction Enzyme (RE) digestion and first adaptor ligation

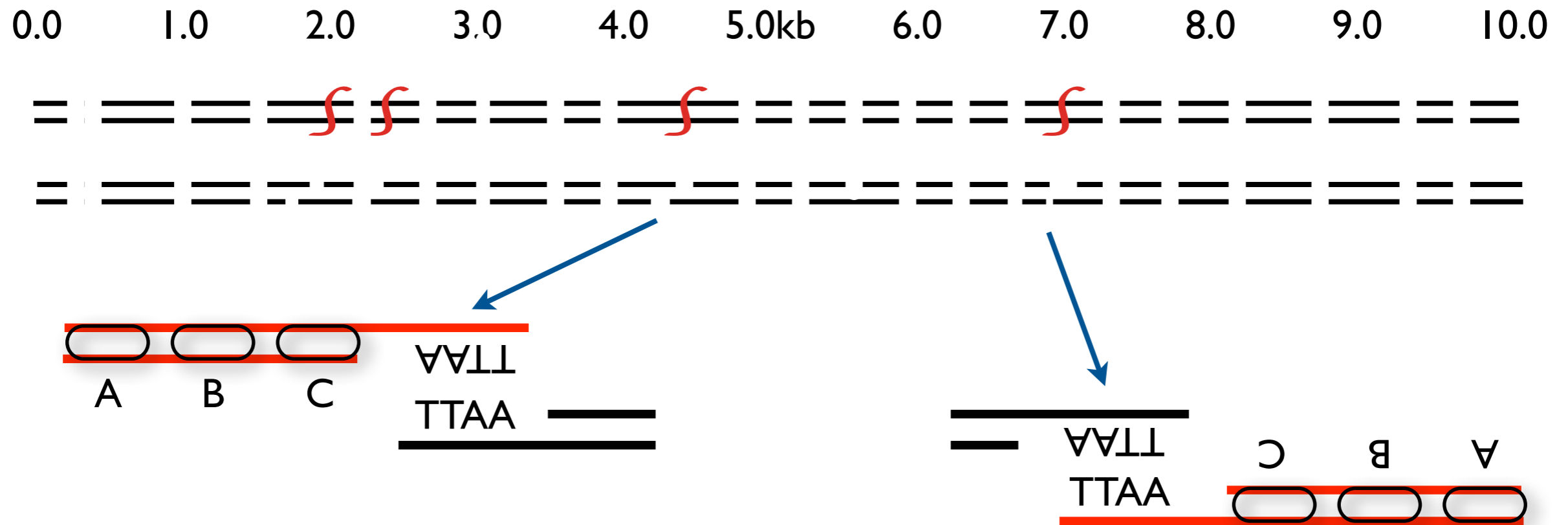0.0      1.0      2.0      3.0      4.0      5.0kb    6.0      7.0      8.0      9.0    10.0

2.0      0.5     2.0          2.5             3.0

A     B     C

ꓯꓯꓕꓕ

TTAA

ꓯꓯꓕꓕ

TTAA

C    B    A

A = Amplification primer
B = Sequencing primer
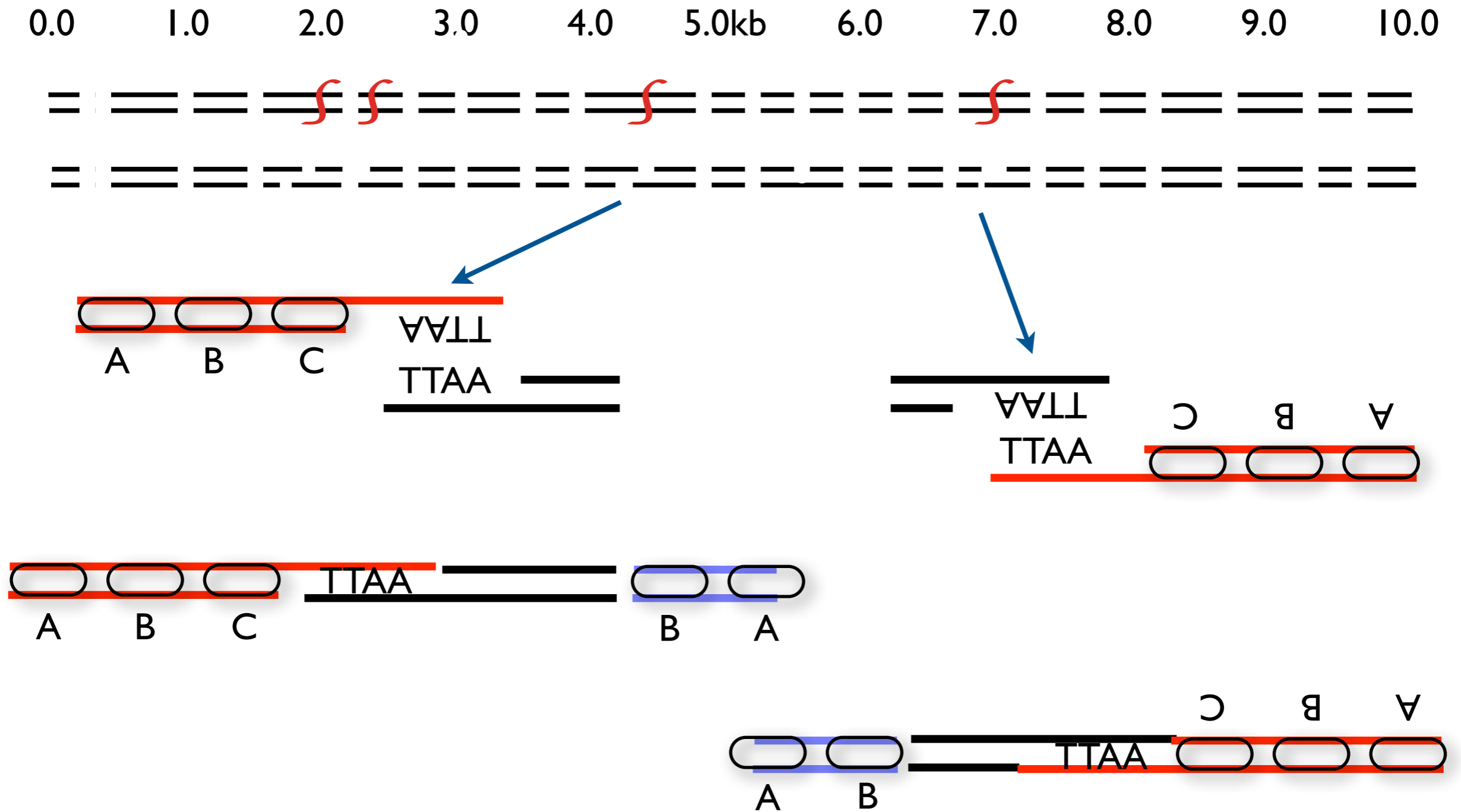C = Barcode

# Shearing and second adaptor ligation



A = Amplification primer
B = Sequencing primer
C = Barcode

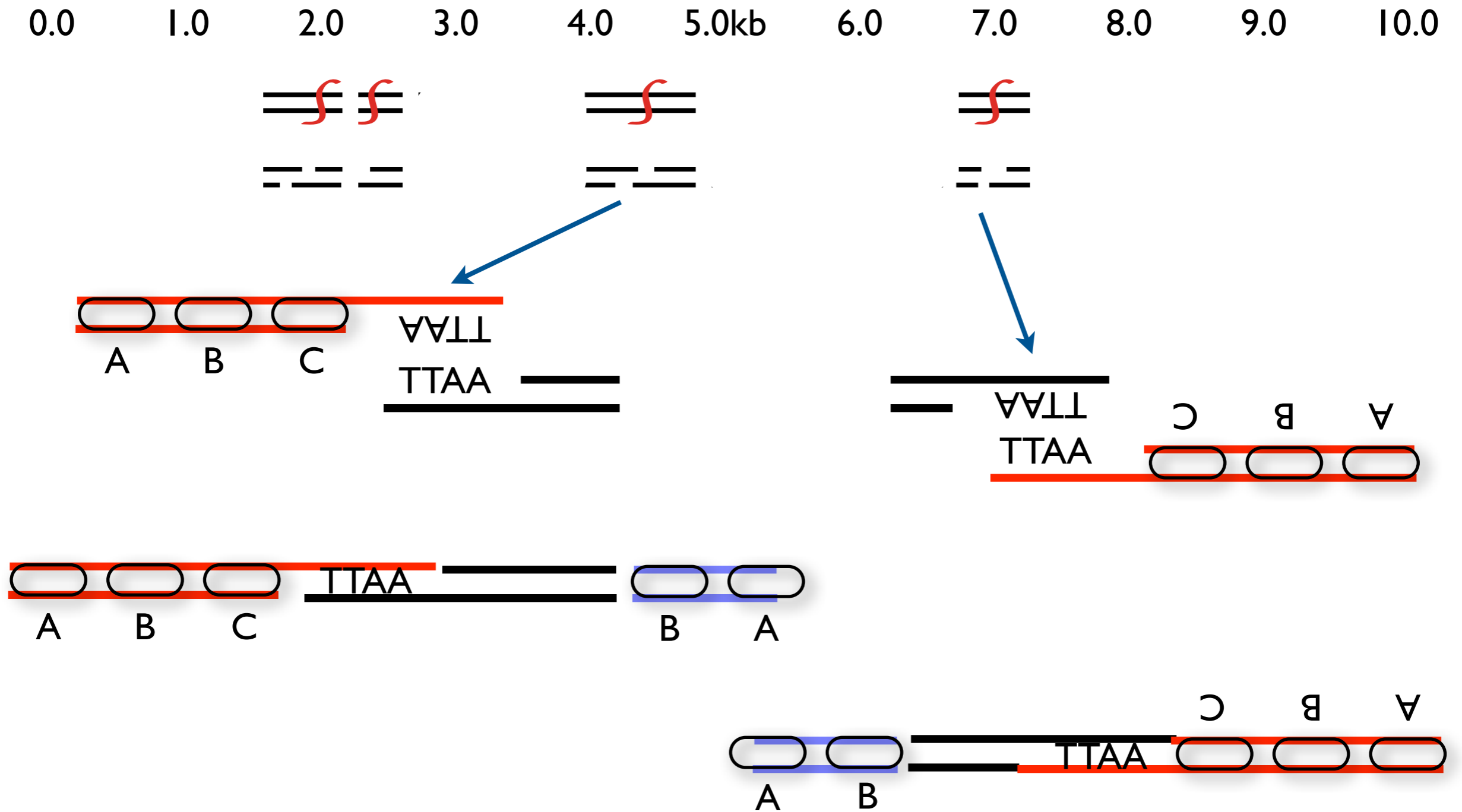# Shearing and second adaptor ligation



A = Amplification primer
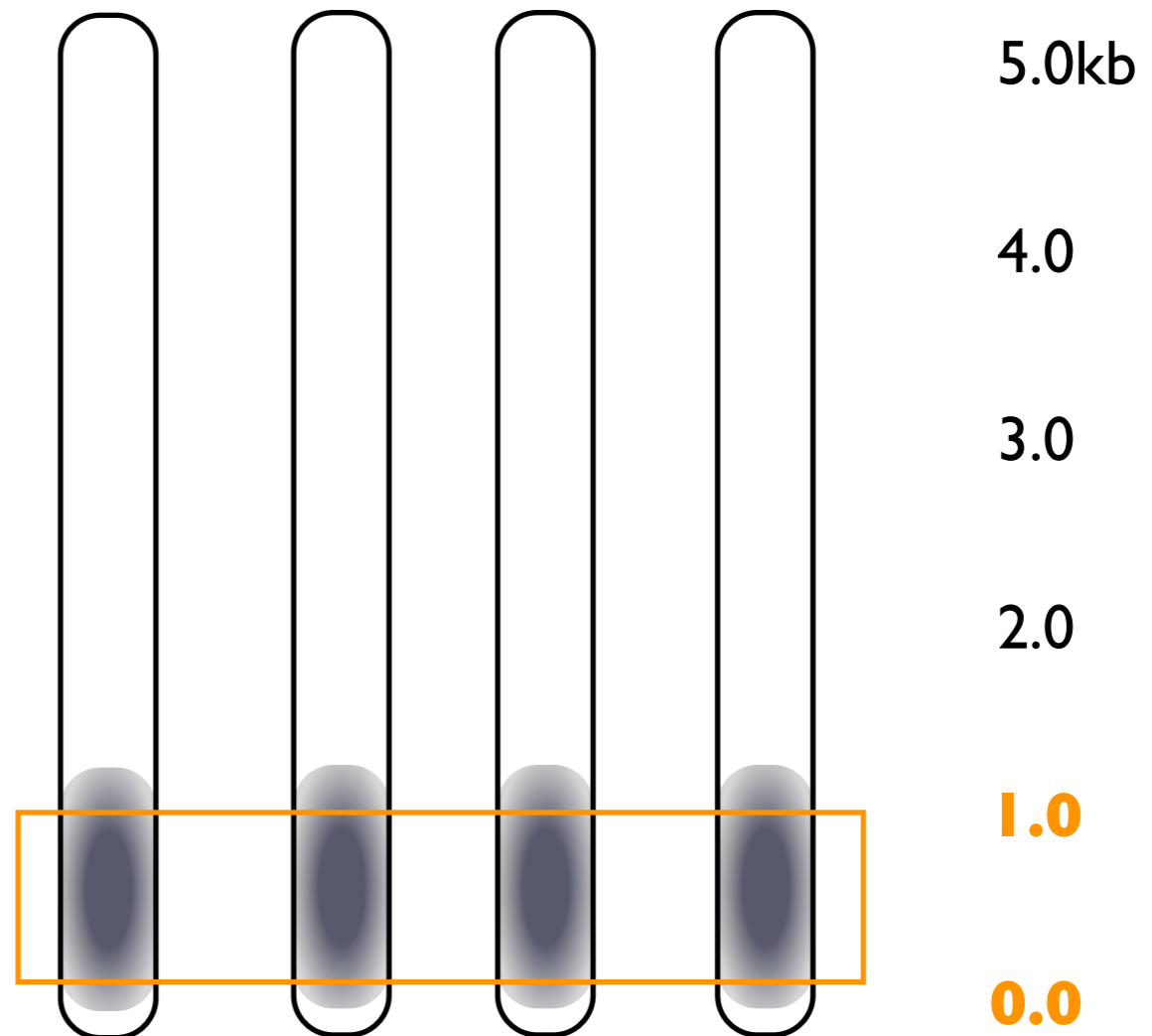B = Sequencing primer
C = Barcode

# Shearing and second adaptor ligation



A = Amplification primer
B = Sequencing primer
C = Barcode

# Shearing makes consistent fragments for sequencing



A = Amplification primer
B = Sequencing primer
C = Barcode

# 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

# 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

A = Amplification primer
B = Sequencing primer
C = Barcode

36bp

# Summary of plusses and minuses of RAD family

|  | Sheared RAD | Single or ddRAD | 2b-RAD |
|---|---|---|---|
| plusses | - Consistent reads<br>- Local assemblies<br>- Identify PCR duplicates | - Fewer steps<br>- Easier marker scaling | - Fewest steps<br>- Easiest marker scaling |
| minuses | - Shearing step<br>- Scaling requires different enzymes | - Multiple enzymes<br>- Poor consistency<br>- PCR duplicates | - Very short reads<br>- PCR duplicates |
|  |  |  |  |

# Benefits of random shearing in RAD



A) Restriction sites in genome

RAD tag sequence read

Sheared-end reads

B) Variable length RAD fragments isolated

C) **200-1200bp in length**

Contigs assembled from the sheared-end reads for each RAD tag

Acquire paired-end sequence

stack

Match to marker catalog

Collate/Assemble PE reads

TGCAGGGGTATTAGCATAA

AACTAATTTTTCACTAGCCATCTTGAATGTGAGTAGCATTTTAAGTAACTATAATTG

Associate markers / PE contigs with ESTs

BLASTn

BLASTn

EST Library

CTGAAGGAGCTGTTACCGGACACCAGCCGGCGCTACGAGAACAAAGCTGGGACCTTCATCACGGGAATCGATGTCACCTCAAAGGAGGCCGTGGAGAAGAAGGAGCAGCGGGCCCGGCGGTTCTCGTGTGG
TTCCACTTCCGCGCCGAGGTCAACCTGGCCCAGAGGAACGTGGTGCTGGACCGGGACAGGATGAGGAAGGCTGTGCCTAAAGCGCGCTTGGAGGCTCTGCACGTGACTGGAGTTGATGAGGGTGTCGTGTA
ATGAGCACTCAGGACGTGTTTGGCTACTTCAAGGAGTATCCCCCGCCCACATCGAGTGGATAGACGACACCTCCTGTAATGTGGTGTGGCTTGATGACGTCACCTCAACCCGAGCCCTCCGGGACGTCCT
CGGCGGAGCCAATACTACATGAAATATGGGAACCCCAACTACGGGGGATGAAGGGAATCCTCAGCAACTCCTGGAAGCGGCGCTACCATTCCCGGCGGATCCAG

ATCAACATGAGCCGCATGCCGGACCCCGCCGCCACCAAGGCCGACCCGGAGGAGAGCGCCCCCACACGGAGCGGGAGGGGACCGGCGAGGAAGGGACTCGGACGAGGAGGCGGAGGAAGGCTGG

GAGGTTGAGGATGATGAGGATGAGAAGAGCAGCACTGGCAGTGCGGGGAAGCCCAGCGACAGTGAGGAAGAGTCGGAGAAGAAGCCCGCTCCTGAAATCACAGAGACTGACGAGCTGTCCAACACCGCCAC
CAGGCCGAGAGAGAGTCGCTCCTCCGAAATGACCTGCGGCCCGCCACCAAACCCTTCAAGGGGAACAAGCTGTTCCTGAGGTTTGCCACTCAGGATGATAAGAAAGAACTGGGCGCCGCCCCCGCCGTACA
ACCAAGAAGAGCCTCATCGGAGACAGCGTGGGGCTGAC

Assign orthology to: markers PE contigs ESTs

BLASTx

BLASTx

BLASTx

Human Genome

Zebrafish Genome

# Random shearing benefits in RAD

Eliminating PCR duplicates:

Eliminating PCR duplicates:

Considerations for RAD-seq 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**

# 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
```

# Differentiating SNPs from error



Restriction enzyme recognition site

Reference genome sequence

sequence reads

# Differentiating SNPs from error

# Differentiating SNPs from error

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

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

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

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

$\{n_A, n_C, n_G, n_T\} = \{1,0,5,4\} \Rightarrow G/T$
$(\ln L = -2.35; \text{MLE}_\varepsilon = 0.2; \text{LRT} = 4.48)$

# 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 $\qquad (0.25)^n \times$ genome size = expected # 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 *C. remanei* (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



40 million clusters per flow cell

20 microns



Agilent Bioanalyzer

[FU]

200

100

0

35    150    300    500   1000    10380   [bp]

# 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

# Threespine stickleback, *Gasterosteus aculeatus*

● *Ancestral Oceanic Populations*

Marine and Anadromous
Old (> 10 million years)
Phenotypically similar

● *Derived Freshwater Populations*

Lake and stream
Young (<15,000 years)
Phenotypically diverse

# Threespine stickleback, *Gasterosteus aculeatus*

🔴 *Ancestral Oceanic Populations*

Marine and Anadromous
Old (> 10 million years)
Phenotypically similar

🔵 *Derived Freshwater Populations*

Lake and stream
Young (<15,000 years)
Phenotypically diverse



Rundle and McKinnon 2002

# Threespine stickleback, *Gasterosteus aculeatus*

🔴 *Ancestral Oceanic Populations*
Marine and Anadromous
Old (> 10 million years)
Phenotypically similar

🔵 *Derived Freshwater Populations*
Lake and stream
Young (<15,000 years)
Phenotypically diverse



Pelvic Structure     Lateral plates

QTL mapping

Cresko et al. 2004. PNAS
Colosimo et al. 2005. Science
Shapiro et al. 2004. Nature
Albert et al. 2008. Evolution
Miller et al. 2007. Cell
Chan et al. 2010. Nature

# Signatures of natural selection across the genome



20 individuals in each of 5 popn's
2 Ocean & 3 Freshwater
45,000 SNPs in each individual

Hohenlohe, Bassham et al. 2010. PLoS Genetics

# What genomic regions are subject to selection during parallel evolution?

# Parallel signatures of selection across the genomes



$F_{st}$

Bear Paw
(mean $F_{ST}$ = 0.121)

Boot
(mean $F_{ST}$ = 0.112)

Mud
(mean $F_{ST}$ = 0.117)

Genomic location (mBases)

# Previously identify quantitative trait loci (QTLs) are under selection



Natural populations

*Eda*  *Enigma*  *Foxi3b*  *FGFR*

Lateral plate major locus
on LGIV (4000 SNPs)

# Extensive LD across the genome



Freshwater

iHH

Ocean

Position (Mb)

Hohenlohe et al. 2012. Philosophical Transactions of the Royal Society of London

# Extensive LD across the genome
# More in oceanic than in freshwater populations



Freshwater

iHH

Ocean

Position (Mb)

Hohenlohe et al. 2012. Philosophical Transactions of the Royal Society of London

# Could genome rearrangements in the stickleback genome be affecting these patterns?

Julian Catchen, Susie Bassham and Kate Ituarte

## Genome Assembly

| | ♂ | ♀ | ♂ | ♀ |
|---|---|---|---|---|
| N50 | 17,417 bp | 18,982 bp | 15,555 bp | 15,534 bp |
| Max | 199,905 bp | 192,283 bp | 238,768 bp | 254,734 bp |
| Total | 488.8 Mb | 472.5 Mb | 456.4 Mb | 473.4 Mb |
| Median Coverage | 24.6x | 26.5x | 24.1x | 25.8x |

# Genome Assembly

# Paired-end Alignments

# Genetic Map Construction

# F1 Pseudo-testcross

# F1 Pseudo-testcross

Male Parent

A₁ B₁ C₁

A₂ B₁ C₂

X

Female Parent

A₁ B₁ C₁

A₁ B₂ C₃

Progeny

A₁ B₁ C₂

A₁ B₁ C₁

A₁ B₁ C₃

A₂ B₁ C₂

Infer

A C

B C

Male map

Female map

Combined map

A B C

♂ X ♀

♂ X ♀

93 progeny

66,071 loci

5,351 markers

93 progeny

45,301 loci

3,927 markers

# 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

G3: Genes, Genomes, Genetics

# Stacks

*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

Aligned
Opposite
Inward

5.00Mb

1.00Mb

LG2|

70cM
60cM
50cM
40cM
30cM
20cM
10cM

1Mb 2Mb 3Mb 4Mb 5Mb 6Mb 7Mb 8Mb 9Mb 10Mb 11Mb

$F_{st}$

0.8
0.7
0.6
0.5
0.4
0.3
0.2
0.1
0
-0.1
-0.2

5000 10000

Linkage Group XXI

RS (Marine)    Boot

Like
Bear Paw

Inverted

# How quickly does stickleback evolution occur?

Montague Is.

Danger Is.

Middleton Is.

25 km

photo: eoimages.gsfc.nasa.gov

Montague Is.

Danger Is.

Middleton Is.

25 km

photo: eoimages.gsfc.nasa.gov

1955

8 km

1964

1967

2002

N

photos: BLM and Aero-Metric, Inc., Anchorage, AK.

743 fish
27,878 RAD loci
110,000 SNPs

photo: e-Terra, LLC.

Marine x Marine  $F_{ST}$ 0.001

15
23
7
6
14
28
8
16
13
17
12
22
11

photo: e-Terra, LLC.

Marine x Marine $F_{ST}$ 0.001
Fresh x Fresh $F_{ST}$ 0.052

photo: e-Terra, LLC.

# NJ F$_{ST}$ tree



Mid15  Mid6  Mid11
Mid28  Mid8
Mid16
Mid7

Mid12

15

23

Mid14

7

Mid22

6

Mid13  Mid23
Mid17

28
8
16
17  12  13

14

22

11

Mont36 Mont35
Mont37

Danger4

0.02

photo: © Terra, JL4

Smoothed Fst

Linkage Group XXI

Fresh vs Marine

7

17

Rabbit SI x Rabbit SI

Linkage Group XXI

Fresh vs Marine

7

17

Bear Paw Lk

Boot Lk   vs Marine

Mud Lk

(Hohenlohe, Bassham et al. 2010)

Linkage Group XXI

17 (Marine)   23 (Marine)   8 (Both)

Like
Bear Paw

Inverted

# Other recent uses of RAD-seq

Quantitative Trait Loci (QTL) mapping

Population genomics & Genome Wide Association Studies (GWAS)

Phylogenetics and phylogeography

Genetic mapping, comparative genomics

*de novo* Genome assembly

Identifying signatures of selection in natural populations

Inferring parentage and pedigrees in the wild

Quantitative genetics in outbred populations

Allele specific transcriptional profiling using RNA-seq

What if you don't have a genome sequence?

Genomically enabling very non-model organisms

Spotted green pufferfish
*Tetraodon nigroviridis*

Japanese pufferfish
*Takifugu rubripes*

Three-spined stickleback
*Gasterosteus aculeatus*

Medaka
*Oryzias latipes*

Platyfish
*Xiphophorus maculatus*

Zebrafish
*Danio rerio*

Goldeye
*Hiodon alosoides*

Bowfin
*Amia calva*

Spotted gar
*Lepisosteus oculatus*

amphibians

birds

lizard

mammals

shark

lamprey

700  600  500  400  300  200  100  0 MYA

Andrew Nishida, Julian Catchen, Susie Bassham, Clay Small and Adam Jones

# Seahorses, sea dragons and pipefishes

# Gasterosteidae and Syngnathidae are historically considered to be closely related



Seahorses

Pipefish

Seadragons

Stickleback

0.1 substitutions / site

Wilson et al. 2003

# Gulf Pipefish
# Syngnathus scovel

- 160 mm (6.3")
- reversed sex roles
- sexual dimorphism
- specialized suction feeding
- no sequences in international databases

# We're really interested in the head and body axis

# *Solution: 'genomically enable' pipefish*

1) A high quality transcriptome

2) Very dense RAD genetic map

3) Deep coverage shotgun sequencing of genome

4) Order genomic and transcriptomic contigs against the RAD reference map

# Pipefish Transcriptome

# Building an EST database in pipefish



Pipefish embryonic mRNA

↓

Illumina sequencing:
100 nt, paired-end

↓

200 million reads (two lanes)

↓

Assembly of transcripts

# Transcriptome



30,000 solid contigs

Mean depth of coverage = 24X

Nearly all of the expected genes in the genome

Number of contigs (y-axis)

Contig length (kb) (x-axis)

# Transcriptome

# We could use these genes right away
## *Dlx2a* and *Dlx5a* expression in pipefish

# Pipefish Genetic Map

# Genetic map workflow

Generated an F1 family of 103 individuals

RAD sequenced the parents and offspring

Analyzed the data using *Stacks*

Paired end local assemblies

Output to JoinMap format

Created Linkage map

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

# Pipefish Genome Project

# Genome workflow

Generated DNA from a single individual

Random Illumina shotgun sequencing

Removed highly repetitive kmers

Produced *several* different genome assemblies

# Illumina genomic libraries for pipefish genome

paired end 101bp

500–700bp

25x

mate pair

4500–7500bp

2x

overlapping

40x

150–250bp

paired end RAD

ACTCTC

500–1200bp

15–25x of 3% of the genome

# Pipefish genome assembly version 0.99
## Nearly the whole genome is covered

| Coverage | Scaffolds | Contigs | Scaffold N50 | Contig N50 |
|---|---|---|---|---|
| All (66.6x) | 33,911 | 307,317 | 26,109 | 1,840 |

| Max | Average Length | Total Length | Gap Length | % |
|---|---|---|---|---|
| 198,155 | 9,916.35 | 336,273,415 | 38,303,839 | (11.39%) |

# Bringing it all together; the spotted gar



Amores, Catchen et al. 2011. Genetics

94 Individuals
15,076 Markers
8,046 Mapped
974 In Genes

94 Individuals
15,076 Markers
8,046 Mapped
974 In Genes

| Organism | Markers |
|---|---|
| Silver carp | 483 |
| Guppy | 790 |
| Barramundi | 240 |
| Catfish | 331 |
| Sea bass | 368 |
| Cichlid | 204 |
| Platyfish | 290 |
| Halibut | 604 |
| Sea bream | 204 |

Physical Contig Size Distribution

Physical Contig Size Distribution

Ordered Contigs ✕
Unordered Contigs +

Physical Contig Size Distribution

- multiple RAD sites per segregating marker means that more of the genome can be tiled

- Mis-assemblies are easily identified

*Danio rerio* (vertical axis, 1–25)

*Lepisosteus oculatus* (horizontal axis, 1–29)

Hsa17    Loc10    Dre12

1:1      1:2

Dre3

*Homo sapiens*      *Lepisosteus oculatus*      *Danio rerio*

*Lepisosteus oculatus*

Genomics can be a tool for enabling new ecology and evolution research
- documenting patterns of genetic variation
- identifying the molecular genetic basis of important phenotypic variation
- assessing how ecological processes structure this genetic variation in genomes
- RAD-seq is a powerful tool for SNP identification and genotyping
- analytical and computational approaches are challenging but

Not your father'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 is good enough - it doesn't have to be completely closed

*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 Python is essential**

G3: Genes, Genomes, Genetics

*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* workflow



PROCESS_RADTAGS

# *Stacks* workflow

# *Stacks* workflow

# *Stacks* workflow

# *Stacks* workflow

1  (1 tags)                                                                          tags per page  10

| Id | SNP | Consensus | Matching Parents | Progeny | Marker | Ratio | Genotypes |
|---|---|---|---|---|---|---|---|
| ˅ 103 annotate | Yes [2nuc] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC | 2 | 92 / 91 | ab/ac | aa: 25 (27.5%)<br>ab: 24 (26.4%)<br>ac: 18 (19.8%)<br>bc: 24 (26.4%) | 91 |

**SNPs**

Column: 52; G/A
Column: 70; T/G

**Alleles**

a : GT
b : GG
c : AG

**Matching Samples**

View: ☑ Haplotypes ☐ Allele Depths ☐ Genotypes

| Male | Female | Progeny 1 | Progeny 2 | Progeny 3 | Progeny 4 | Progeny 5 | Progeny 6 | Progeny 7 | Progeny 8 |
|---|---|---|---|---|---|---|---|---|---|
| GT / GG | AG / GT | GT | AG / GG | GG / AG | GG / GT | GG / AG | AG | GT / GG | AG / GT |
| Progeny 9 | Progeny 10 | Progeny 11 | Progeny 12 | Progeny 13 | Progeny 14 | Progeny 15 | Progeny 16 | Progeny 17 | Progeny 18 |
| GT | GT | GG / GT | GT / AG | GG / AG | GT / AG | GT / GG | GG / GT | GG / AG | GT |
| Progeny 19 | Progeny 20 | Progeny 21 | Progeny 22 | Progeny 23 | Progeny 24 | Progeny 25 | Progeny 26 | Progeny 27 | Progeny 28 |
| GT / AG | AG / GG | GT / AG | AG / GT | GG / AG | GG / AG | GT | GG / GT | GG / AG | GG / GT |
| Progeny 29 | Progeny 31 | Progeny 32 | Progeny 33 | Progeny 34 | Progeny 35 | Progeny 36 | Progeny 37 | Progeny 38 | Progeny 39 |
| GT / GG | GT | GT | GT | GT | GT / GG | GT | GT / AG | GT | AG / GT |
| Progeny 40 | Progeny 41 | Progeny 42 | Progeny 43 | Progeny 44 | Progeny 45 | Progeny 46 | Progeny 47 | Progeny 48 | Progeny 49 |
| GT | GT | GT | GT / GG | GG / GT | GT | GG / GT | GG / AG | GT | GT / GG |
| Progeny 50 | Progeny 51 | Progeny 52 | Progeny 53 | Progeny 54 | Progeny 55 | Progeny 56 | Progeny 57 | Progeny 58 | Progeny 59 |
| GT | GT | GT / AG | GG / GT | GT / GG | AG / GG | GT | AG / GT | GT / AG | GG / GT |
| Progeny 60 | Progeny 61 | Progeny 62 | Progeny 63 | Progeny 64 | Progeny 65 | Progeny 66 | Progeny 67 | Progeny 68 | Progeny 70 |
| GT / GG | GT / GG | GT / AG | GG / AG | GG / GT | GT | GT | GG / GT | GT | GG / AG |
| Progeny 71 | Progeny 72 | Progeny 73 | Progeny 74 | Progeny 75 | Progeny 76 | Progeny 77 | Progeny 78 | Progeny 79 | Progeny 80 |
| GG / AG | AG / GG | GT | GG / AG | GT / GG | GT | GG / AG | GG / AG | GT / GG | GT |
| Progeny 81 | Progeny 82 | Progeny 83 | Progeny 84 | Progeny 85 | Progeny 86 | Progeny 87 | Progeny 88 | Progeny 89 | Progeny 90 |
| GT / AG | GT / AG | GG / AG | GT | GT / GG | GT / GG | GT | GG / AG | GT | GG / AG |
| Progeny 91 | Progeny 92 | Progeny 93 | Progeny 94 | | | | | | |
| AG / GG | GT / AG | AG / GG | GG / AG | | | | | | |

1 (1 tags)　　　　　　　　　　　　　　　　　　　tags per page 10

| Id | SNP | Consensus | Matching Parents | Progeny | Marker | Ratio | Genotypes |
|---|---|---|---|---|---|---|---|
| ˅ 103 annotate | Yes [2nuc] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC | 2 | 92 / 91 | ab/ac | aa: 25 (27.5%) ab: 24 (26.4%) ac: 18 (19.8%) bc: 24 (26.4%) | 91 |

**SNPs**
Column: 52; G/A
Column: 70; T/G

**Alleles**
a : GT
b : GG
c : AG

**Matching Samples**

View: ☑ Haplotypes ☑ Allele Depths ☐ Genotypes

| Male | Female | Progeny 1 | Progeny 2 | Progeny 3 | Progeny 4 | Progeny 5 | Progeny 6 | Progeny 7 | Progeny 8 |
|---|---|---|---|---|---|---|---|---|---|
| GT / GG 34 / 13 | AG / GT 12 / 14 | GT 7 | AG / GG 8 / 16 | GG / AG 26 / 14 | GG / GT 15 / 11 | GG / AG 14 / 8 | AG 29 | GT / GG 22 / 11 | AG / GT 12 / 5 |

| Progeny 9 | Progeny 10 | Progeny 11 | Progeny 12 | Progeny 13 | Progeny 14 | Progeny 15 | Progeny 16 | Progeny 17 | Progeny 18 |
|---|---|---|---|---|---|---|---|---|---|
| GT 25 | GT 23 | GG / GT 32 / 14 | GT / AG 22 / 7 | GG / AG 7 / 8 | GT / AG 7 / 8 | GT / GG 2 / 3 | GG / GT 19 / 14 | GG / AG 9 / 4 | GT 15 |

| Progeny 19 | Progeny 20 | Progeny 21 | Progeny 22 | Progeny 23 | Progeny 24 | Progeny 25 | Progeny 26 | Progeny 27 | Progeny 28 |
|---|---|---|---|---|---|---|---|---|---|
| GT / AG 6 / 3 | AG / GG 6 / 9 | GT / AG 18 / 9 | AG / GT 4 / 5 | GG / AG 7 / 6 | GG / AG 8 / 10 | GT 7 | GG / GT 10 / 16 | GG / AG 3 / 3 | GG / GT 4 / 5 |

| Progeny 29 | Progeny 31 | Progeny 32 | Progeny 33 | Progeny 34 | Progeny 35 | Progeny 36 | Progeny 37 | Progeny 38 | Progeny 39 |
|---|---|---|---|---|---|---|---|---|---|
| GT / GG 8 / 5 | GT 11 | GT 10 | GT 17 | GT 20 | GT / GG 7 / 3 | GT 8 | GT / AG 12 / 4 | GT 9 | AG / GT 12 / 7 |

| Progeny 40 | Progeny 41 | Progeny 42 | Progeny 43 | Progeny 44 | Progeny 45 | Progeny 46 | Progeny 47 | Progeny 48 | Progeny 49 |
|---|---|---|---|---|---|---|---|---|---|
| GT 9 | GT 5 | GT 9 | GT / GG 9 / 12 | GG / GT 3 / 6 | GT 6 | GG / GT 4 / 11 | GG / AG 3 / 7 | GT 18 | GT / GG 5 / 6 |

| Progeny 50 | Progeny 51 | Progeny 52 | Progeny 53 | Progeny 54 | Progeny 55 | Progeny 56 | Progeny 57 | Progeny 58 | Progeny 59 |
|---|---|---|---|---|---|---|---|---|---|
| GT 18 | GT 9 | GT / AG 8 / 5 | GG / GT 10 / 8 | GT / GG 5 / 6 | AG / GG 8 / 10 | GT 22 | AG / GT 17 / 16 | GT / AG 23 / 24 | GG / GT 25 / 13 |

| Progeny 60 | Progeny 61 | Progeny 62 | Progeny 63 | Progeny 64 | Progeny 65 | Progeny 66 | Progeny 67 | Progeny 68 | Progeny 70 |
|---|---|---|---|---|---|---|---|---|---|
| GT / GG 12 / 18 | GT / GG 22 / 29 | GT / AG 7 / 23 | GG / AG 15 / 11 | GG / GT 13 / 20 | GT 44 | GT 27 | GG / GT 23 / 17 | GT 30 | GG / AG 14 / 13 |

| Progeny 71 | Progeny 72 | Progeny 73 | Progeny 74 | Progeny 75 | Progeny 76 | Progeny 77 | Progeny 78 | Progeny 79 | Progeny 80 |
|---|---|---|---|---|---|---|---|---|---|
| GG / AG 15 / 7 | AG / GG 9 / 6 | GT 42 | GG / AG 31 / 29 | GT / GG 15 / 22 | GT 41 | GG / AG 14 / 17 | GG / AG 25 / 17 | GT / GG 29 / 14 | GT 34 |

| Progeny 81 | Progeny 82 | Progeny 83 | Progeny 84 | Progeny 85 | Progeny 86 | Progeny 87 | Progeny 88 | Progeny 89 | Progeny 90 |
|---|---|---|---|---|---|---|---|---|---|
| GT / AG 17 / 29 | GT / AG 29 / 24 | GG / AG 16 / 25 | GT 41 | GT / GG 14 / 24 | GT / GG 6 / 4 | GT 15 | GG / AG 5 / 11 | GT 18 | GG / AG 5 / 17 |

| Progeny 91 | Progeny 92 | Progeny 93 | Progeny 94 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| AG / GG 14 / 13 | GT / AG 12 / 6 | AG / GG 7 / 7 | GG / AG 3 / 2 | | | | | | |

1 (1 tags)　　　　　　　　　　　　　　　　　　　tags per page 10

http://genome.uoregon.edu/stacks/catalog.php?id=1&db=gartut_radtags&p=1&pp=10&filter_type[]=cata&filter_cata=103&filter_alle_l=1&

1  (1 tags)

tags per page  10

| Id | SNP | Consensus | Matching Parents | Progeny | Marker | Ratio | Genotypes |
|---|---|---|---|---|---|---|---|
| ˅ 103 annotate | Yes [2nuc] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC | 2 | 92 / 91 | ab/ac | aa: 25 (27.5%) ab: 24 (26.4%) ac: 18 (19.8%) bc: 24 (26.4%) | 91 |

**SNPs**
Column: 52; G/A
Column: 70; T/G

**Alleles**
a : GT
b : GG
c : AG

**Matching Samples**

View: ☑ Haplotypes ☑ Allele Depths ☑ Genotypes

| Male | Female | Progeny 1 | Progeny 2 | Progeny 3 | Progeny 4 | Progeny 5 | Progeny 6 | Progeny 7 | Progeny 8 |
|---|---|---|---|---|---|---|---|---|---|
| GT / GG 34 / 13 | AG / GT 12 / 14 | GT 7 aa | AG / GG 8 / 16 bc | GG / AG 26 / 14 bc | GG / GT 15 / 11 ab | GG / AG 14 / 8 bc | AG 29 AC | GT / GG 22 / 11 ab | AG / GT 12 / 5 ac |

| Progeny 9 | Progeny 10 | Progeny 11 | Progeny 12 | Progeny 13 | Progeny 14 | Progeny 15 | Progeny 16 | Progeny 17 | Progeny 18 |
|---|---|---|---|---|---|---|---|---|---|
| GT 25 aa | GT 23 aa | GG / GT 32 / 14 ab | GT / AG 22 / 7 ac | GG / AG 7 / 8 bc | GT / AG 7 / 8 ac | GT / GG 2 / 3 ab | GG / GT 19 / 14 ab | GG / AG 9 / 4 bc | GT 15 aa |

| Progeny 19 | Progeny 20 | Progeny 21 | Progeny 22 | Progeny 23 | Progeny 24 | Progeny 25 | Progeny 26 | Progeny 27 | Progeny 28 |
|---|---|---|---|---|---|---|---|---|---|
| GT / AG 6 / 3 ac | AG / GG 6 / 9 bc | GT / AG 18 / 9 ac | AG / GT 4 / 5 ac | GG / AG 7 / 6 bc | GG / AG 8 / 10 bc | GT 7 AC | GG / GT 10 / 16 ab | GG / AG 3 / 3 bc | GG / GT 4 / 5 ab |

| Progeny 29 | Progeny 31 | Progeny 32 | Progeny 33 | Progeny 34 | Progeny 35 | Progeny 36 | Progeny 37 | Progeny 38 | Progeny 39 |
|---|---|---|---|---|---|---|---|---|---|
| GT / GG 8 / 5 ab | GT 11 aa | GT 10 aa | GT 17 aa | GT 20 aa | GT / GG 7 / 3 ab | GT 8 aa | GT / AG 12 / 4 ac | GT 9 aa | AG / GT 12 / 7 ac |

| Progeny 40 | Progeny 41 | Progeny 42 | Progeny 43 | Progeny 44 | Progeny 45 | Progeny 46 | Progeny 47 | Progeny 48 | Progeny 49 |
|---|---|---|---|---|---|---|---|---|---|
| GT 9 aa | GT 5 aa | GT 9 aa | GT / GG 9 / 12 ab | GG / GT 3 / 6 ab | GT 6 AC | GG / GT 4 / 11 ab | GG / AG 3 / 7 bc | GT 18 aa | GT / GG 5 / 6 ab |

| Progeny 50 | Progeny 51 | Progeny 52 | Progeny 53 | Progeny 54 | Progeny 55 | Progeny 56 | Progeny 57 | Progeny 58 | Progeny 59 |
|---|---|---|---|---|---|---|---|---|---|
| GT 18 aa | GT 9 aa | GT / AG 8 / 5 ac | GG / GT 10 / 8 ab | GT / GG 5 / 6 ab | AG / GG 8 / 10 bc | GT 22 aa | AG / GT 17 / 16 ac | GT / AG 23 / 24 ac | GG / GT 25 / 13 ab |

| Progeny 60 | Progeny 61 | Progeny 62 | Progeny 63 | Progeny 64 | Progeny 65 | Progeny 66 | Progeny 67 | Progeny 68 | Progeny 70 |
|---|---|---|---|---|---|---|---|---|---|
| GT / GG 12 / 18 ab | GT / GG 22 / 29 ab | GT / AG 7 / 23 ac | GG / AG 15 / 11 bc | GG / GT 13 / 20 ab | GT 44 aa | GT 27 aa | GG / GT 23 / 17 ab | GT 30 aa | GG / AG 14 / 13 bc |

| Progeny 71 | Progeny 72 | Progeny 73 | Progeny 74 | Progeny 75 | Progeny 76 | Progeny 77 | Progeny 78 | Progeny 79 | Progeny 80 |
|---|---|---|---|---|---|---|---|---|---|
| GG / AG 15 / 7 bc | AG / GG 9 / 6 bc | GT 42 aa | GG / AG 31 / 29 bc | GT / GG 15 / 22 ab | GT 41 aa | GG / AG 14 / 17 bc | GG / AG 25 / 17 bc | GT / GG 29 / 14 ab | GT 34 aa |

| Progeny 81 | Progeny 82 | Progeny 83 | Progeny 84 | Progeny 85 | Progeny 86 | Progeny 87 | Progeny 88 | Progeny 89 | Progeny 90 |
|---|---|---|---|---|---|---|---|---|---|
| GT / AG 17 / 29 | GT / AG 20 / 24 | GG / AG 16 / 25 | GT 41 | GT / GG 14 / 24 | GT / GG 5 / 4 | GT 15 | GG / AG 5 / 11 | GT 18 | GG / AG 5 / 17 |

http://genome.uoregon.edu/stacks/tag.php?db=gartut_radtags&batch_id=1&sample_id=2&tag_id=73

# Stacks

version 0.998

## Batch #1 [2011-08-10; 80bp Lepisosteus oculatus F1 Genetic Map RAD-Tag Samples]

## RAD-Tag Sample #2 [female]

### ˅ Sequence #73

| Catalog ID | Depth | SNPs | | Alleles | | Deleveraged? | Lumberjackstack? | Blacklisted? |
|---|---|---|---|---|---|---|---|---|
| #103 | 26x | Column: 52 | G/A | AG | 46.15% | False | False | False |
| | | Column: 70 | T/G | GT | 53.85% | | | |

| | Relationship | Seq ID | Sequence |
|---|---|---|---|
| | consensus | | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| | model | | OOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOOEOOOOOOOOOOOOOOOOOOOEOOOO |
| 1 | primary | CAGTC_2_0018_768_1365_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACAAAGCAACACTTCACAGTCCC |
| 2 | primary | CAGTC_2_0029_1628_1751_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACAAAGCAACACTTCACAGTCCC |
| 3 | primary | CAGTC_2_0053_1692_1388_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACAAAGCAACACTTCACAGTCCC |
| 4 | primary | CAGTC_2_0058_1588_1038_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACAAAGCAACACTTCACAGTCCC |
| 5 | primary | CAGTC_2_0059_1524_1186_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACAAAGCAACACTTCACAGTCCC |
| 6 | primary | CAGTC_2_0094_1356_1854_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACAAAGCAACACTTCACAGTCCC |
| 7 | primary | CAGTC_2_0096_1791_1246_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACAAAGCAACACTTCACAGTCCC |
| 8 | primary | CAGTC_2_0021_877_296_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 9 | primary | CAGTC_2_0024_307_735_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 10 | primary | CAGTC_2_0025_108_810_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 11 | primary | CAGTC_2_0039_1252_1764_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 12 | primary | CAGTC_2_0061_596_159_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 13 | primary | CAGTC_2_0068_1310_997_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 14 | primary | CAGTC_2_0070_644_2040_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 15 | primary | CAGTC_2_0074_328_659_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 16 | primary | CAGTC_2_0075_1668_1862_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 17 | primary | CAGTC_2_0079_1481_505_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 18 | primary | CAGTC_2_0084_805_1974_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 19 | primary | CAGTC_2_0100_481_1043_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATTCCC |
| 20 | secondary | CAGTC_2_0014_728_1008_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACAAAGCAACACTTCACAGACCC |
| 21 | secondary | CAGTC_2_0016_86_1022_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCAACACTTCACATACCC |
| 22 | secondary | CAGTC_2_0042_426_1001_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACCAAGCAACACTTCACAGTCCC |
| 23 | secondary | CAGTC_2_0052_867_1387_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCCGTGGACCGAGAGCACAAAGCAACACTTCACAGTCCC |
| 24 | secondary | CAGTC_2_0012_221_1043_1[35245] | TGCAGGAGCCCTCCCACTAGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACAAAGCAACACTTCACAGTCCC |
| 25 | secondary | CAGTC_2_0095_120_1067_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCGCAAAGCCACACTTCACATCCCC |
| 26 | secondary | CAGTC_2_0077_1003_356_1[35245] | TGCAGGAGCCCTCCCACTCGCTGATGGCCACTCCATTCAGTGGACCGAGAGCACAAAGCAACACCTCACAGTCCC |

last updated: Sun Jan 8 08:59:35 PST 2012