# Introduction to short read NGS:

*Library construction, UCE capture and ddRADseq*

The Natural History Museum, London

Autumn 2021

Instructor: Jeff Streicher

j.streicher@nhm.ac.uk

*Litoria iris,* Papua New Guinea

# Unit 4: Double digest restriction-site associated DNA sequencing (ddRADseq)



https://github.com/nhm-herpetology/museum-NGS-training

# Unit 3 Review

- Reduced representation genome sequencing

- Targeted sequence capture and UCEs
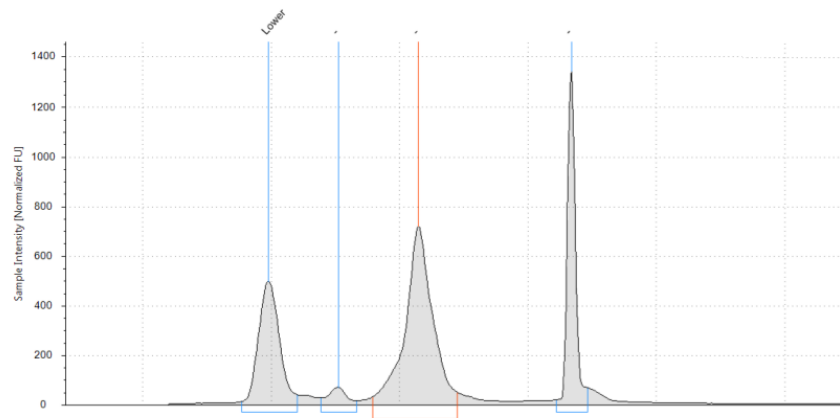
Bioinformatics Lab

- How to download and process UCE data

Molecular Lab

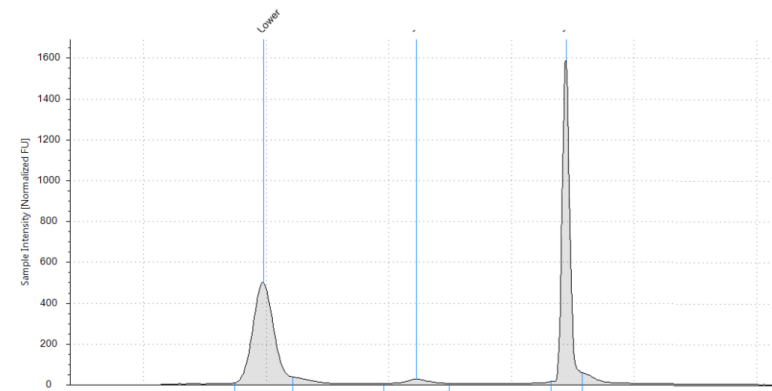- Targeted sequence capture of UCEs from shotgun libraries made in Unit 2

# Tapestation results



Unit 2

Unit 3

# Unit 4 Overview

Lecture

- Restriction Enzymes
- Restriction-site associated DNA sequencing (RADseq)
- ddRADseq

Bioinformatics Lab

- How to process ddRADseq data

Molecular Lab

- Restriction digestion and adapter design (Tomorrow)
- PCR and TapeStation (Monday)

# Reduced-representation NGS sequencing

- Genomes can be large!

- We might want to compare multiple individuals/species

- Targeted Sequence Capture (TSC) [Unit 3]

- Restriction site associated DNA sequencing (RADseq)

# Restriction endonucleases

EcoR1

Restriction enzymes cleave DNA into fragments at or near specific recognition sites within molecules known as restriction sites. These enzymes are found in bacteria and archaea and provide a defense mechanism against invading viruses.
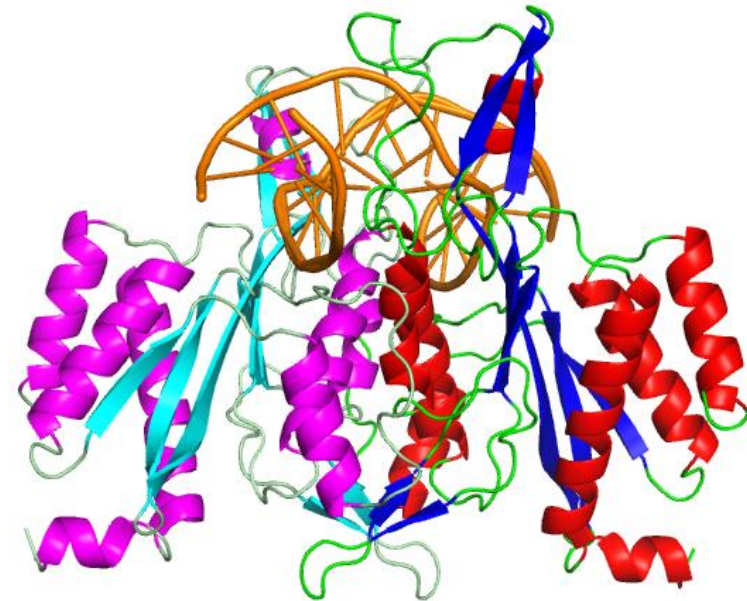
Text mostly from Wikipedia ☺

Image by A2-33 with CCL

# Type II restriction enzymes

They form homodimers, with recognition sites that are usually undivided and palindromic and 4–8 nucleotides in length. They recognize and cleave DNA at the same site throughout the genome and can either cleave at the center of both strands to yield a **blunt end**, or at a staggered position leaving overhangs called **sticky (protruding) ends**.
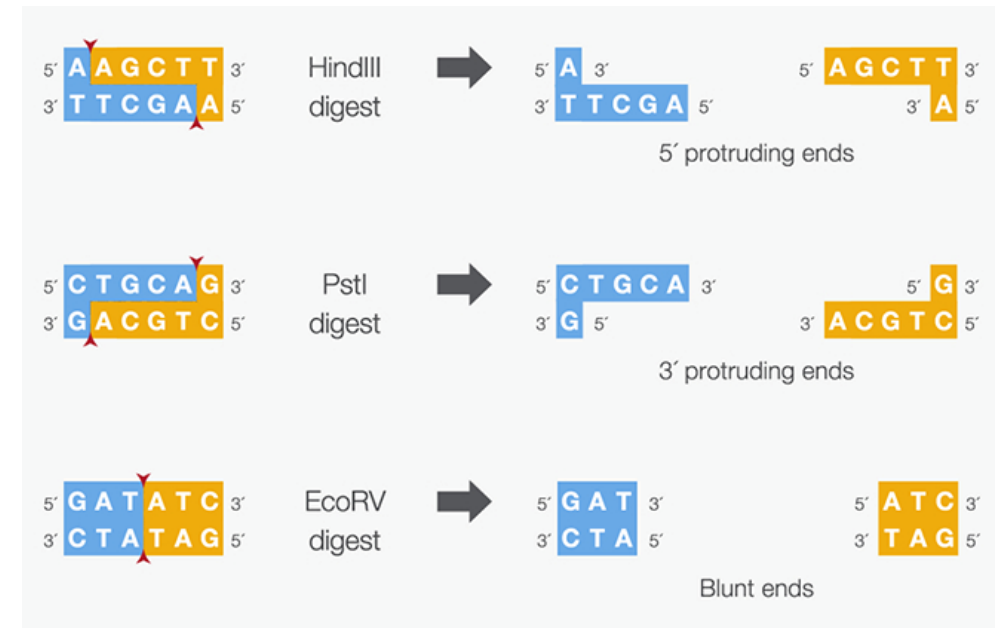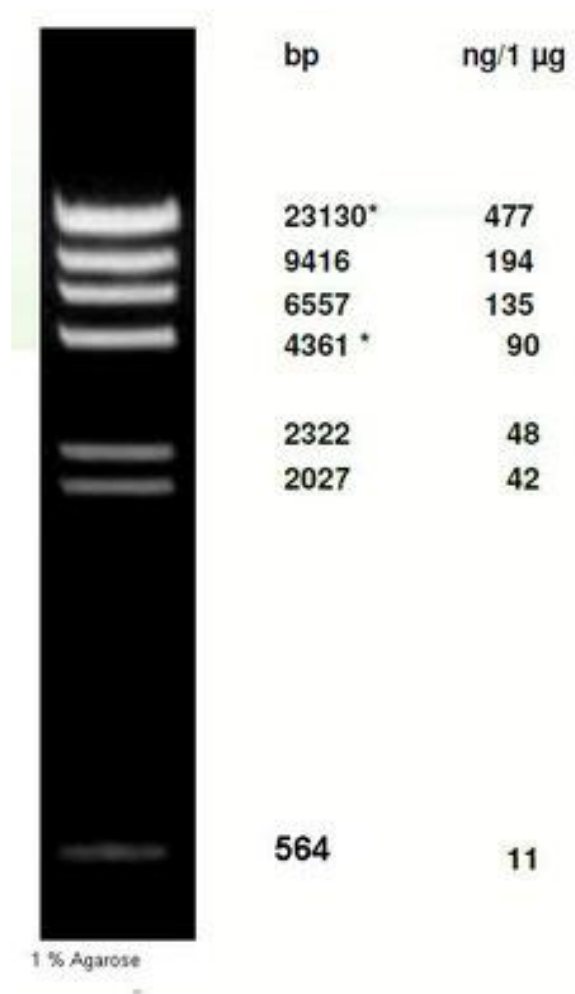
Text mostly from Wikipedia ☺



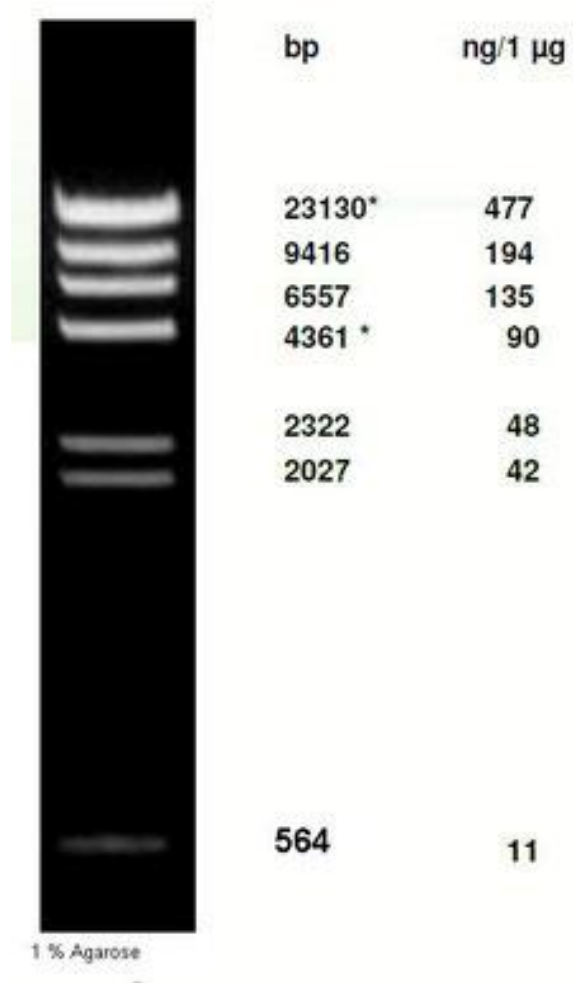Image by Thermo Fisher Scientific

# Restriction digestion

- DNA Ladders (size standards)
- Restriction fragment length polymorphism (RFLP)
- Amplified fragment length polymorphism (AFLP)

| bp | ng/1 µg |
|---|---|
| 23130* | 477 |
| 9416 | 194 |
| 6557 | 135 |
| 4361 * | 90 |
| 2322 | 48 |
| 2027 | 42 |
| 564 | 11 |

1 % Agarose

DNA Ladder (Lambda phage genome) cut with Hind III restriction enzyme

| bp | ng/1 µg |
|---|---|
| 23130* | 477 |
| 9416 | 194 |
| 6557 | 135 |
| 4361 * | 90 |
| 2322 | 48 |
| 2027 | 42 |
| 564 | 11 |

1 % Agarose

## Hind III Recognition Site

5´...A AGCTT...3´
3´...TTCGA A...5´

Image by New England BioLabs

λ phage

head containing DNA
collar
tail
base plate
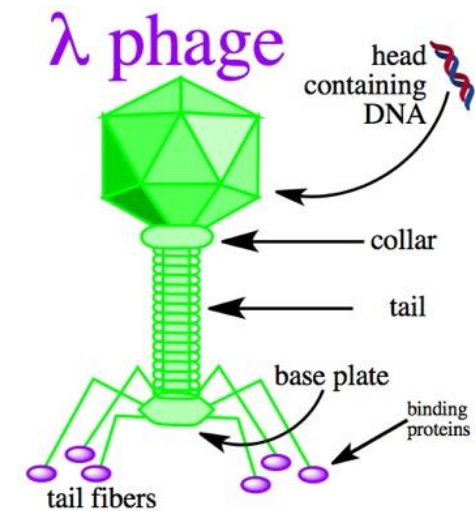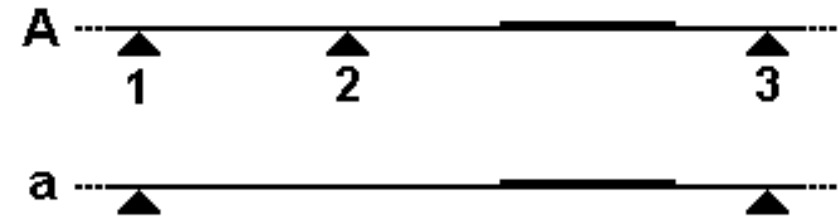binding proteins
tail fibers

Image by Lizanne Koch/Public Domain

DNA Ladder (Lambda phage genome) cut with Hind III restriction enzyme

# Restriction fragment length polymorphism (RFLP)

DNA in a diploid organism

- DNA is digested using a restriction enzyme.

- DNA fragments produced by the digest are then separated by length through agarose gel electrophoresis and transferred to a membrane via the Southern blot procedure.

- Hybridization of the membrane to a labeled DNA probe then determines the length of the fragments which are complementary to the probe.

Text mostly from Wikipedia ☺

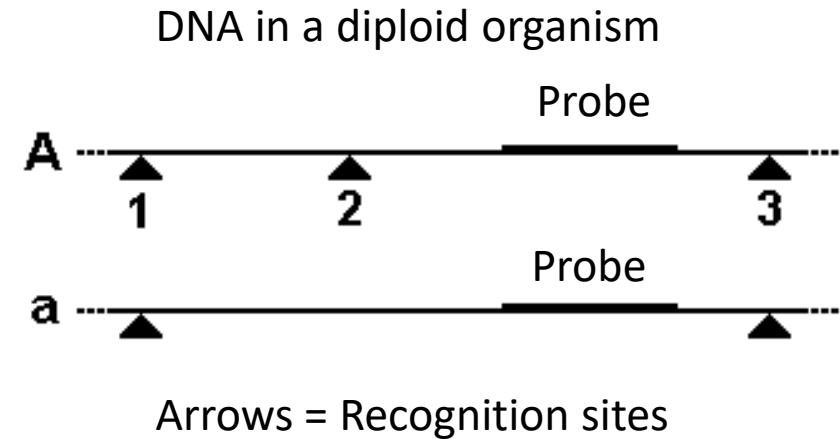# Restriction fragment length polymorphism (RFLP)

- DNA is digested using a restriction enzyme.

- DNA fragments produced by the digest are then separated by length through agarose gel electrophoresis and transferred to a membrane via the Southern blot procedure.

- Hybridization of the membrane to a labeled DNA probe then determines the length of the fragments which are complementary to the probe.

Text mostly from Wikipedia ☺

DNA in a diploid organism

Probe

A

1    2                    3

Probe

a

Arrows = Recognition sites
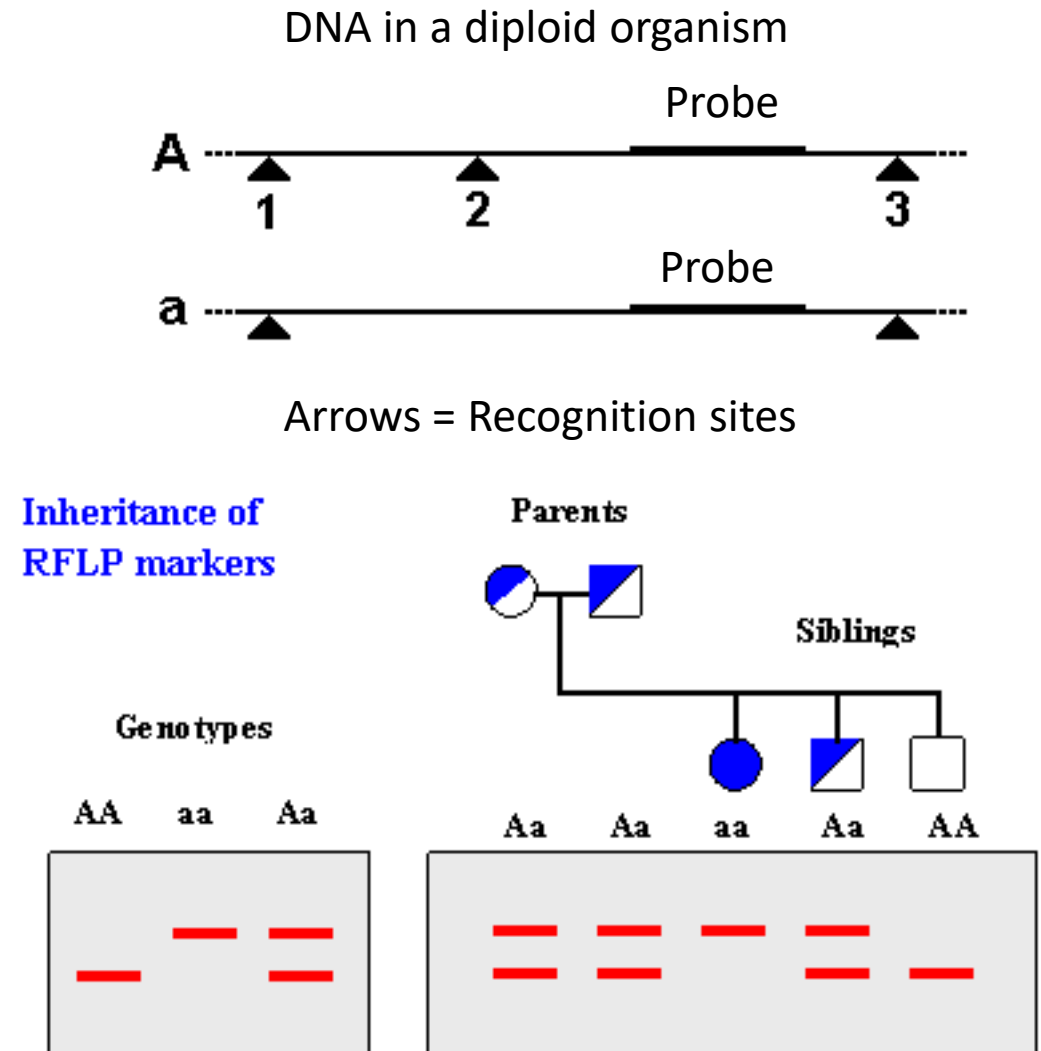
# Restriction fragment length polymorphism (RFLP)

- DNA is digested using a restriction enzyme.

- DNA fragments produced by the digest are then separated by length through agarose gel electrophoresis and transferred to a membrane via the Southern blot procedure.

- Hybridization of the membrane to a labeled DNA probe then determines the length of the fragments which are complementary to the probe.
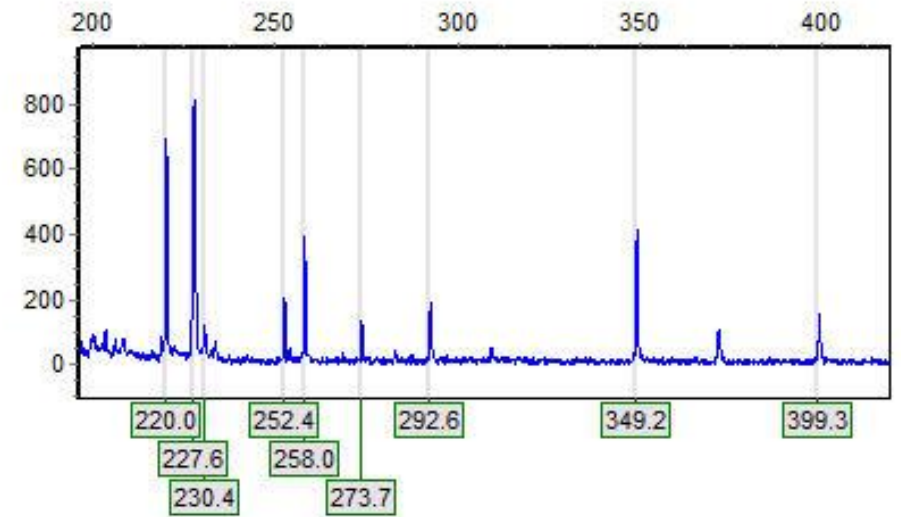
Text mostly from Wikipedia ☺

DNA in a diploid organism

Probe

A ···· ▲     ▲     ▲ ····
   1     2     3

Probe

a ···· ▲     ▲ ····

Arrows = Recognition sites

**Inheritance of RFLP markers**

Parents

Siblings

**Genotypes**

AA   aa   Aa

Aa  Aa  aa  Aa  AA

# Amplified fragment length polymorphism (AFLP)

- AFLP uses restriction enzymes to digest genomic DNA, followed by ligation of adaptors to the sticky ends of the restriction fragments. A subset of the restriction fragments is then selected to be amplified.

- The amplified fragments are separated and visualized on denaturing on agarose gel electrophoresis , either through autoradiography or fluorescence methodologies, or via automated capillary sequencing instruments.
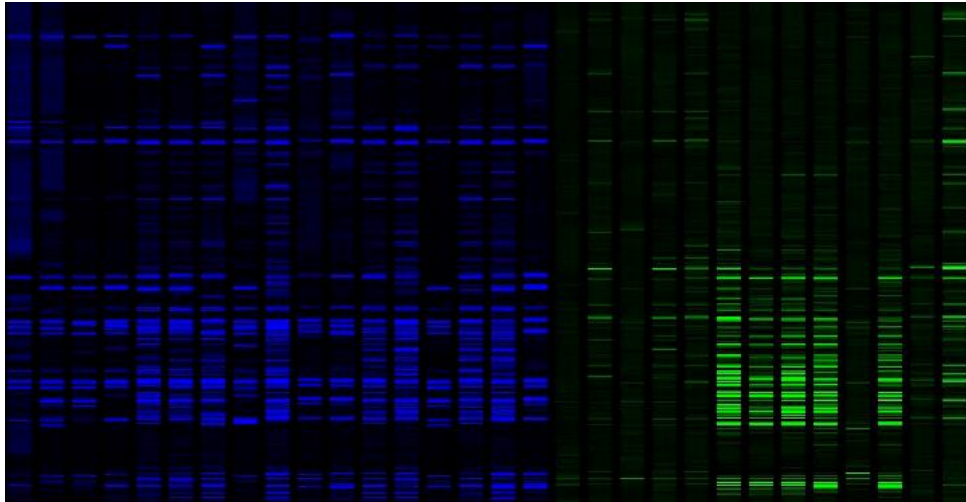
Text mostly from Wikipedia ☺



Public Domain

AFLP data, scoring different size digest fragments

# Amplified fragment length polymorphism (AFLP)



Anthony Genova

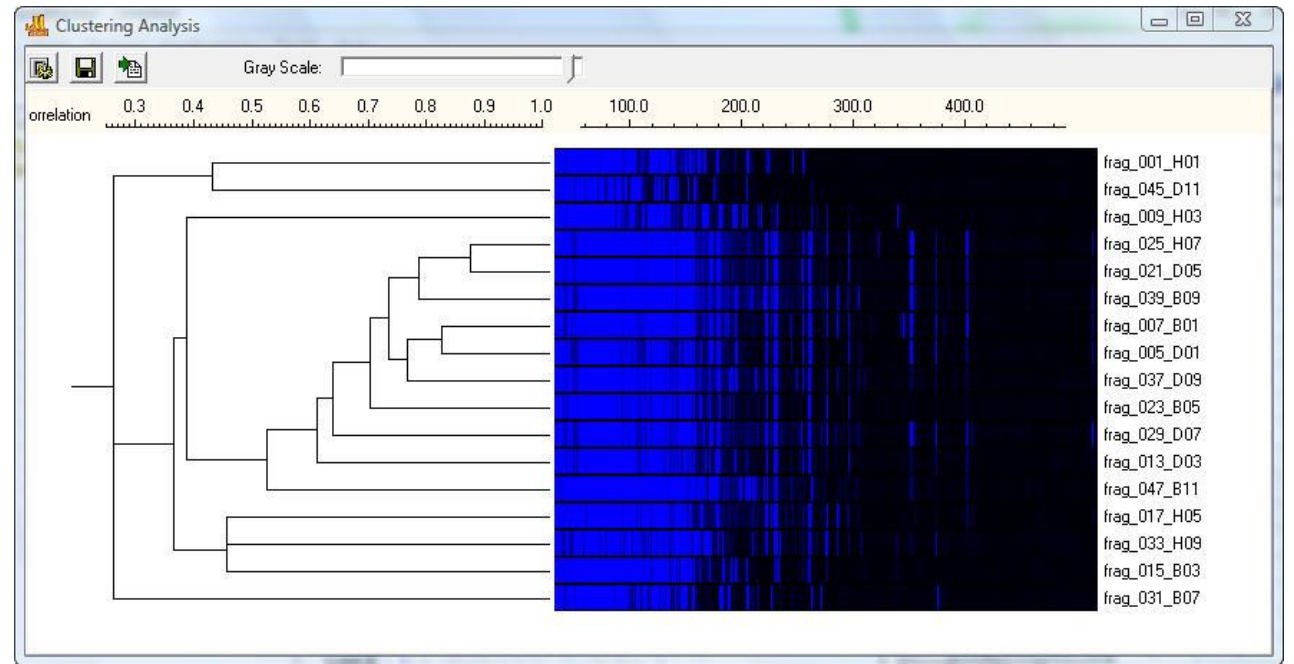Amplified Fragment Length Polymorphisms (*Anolis*)

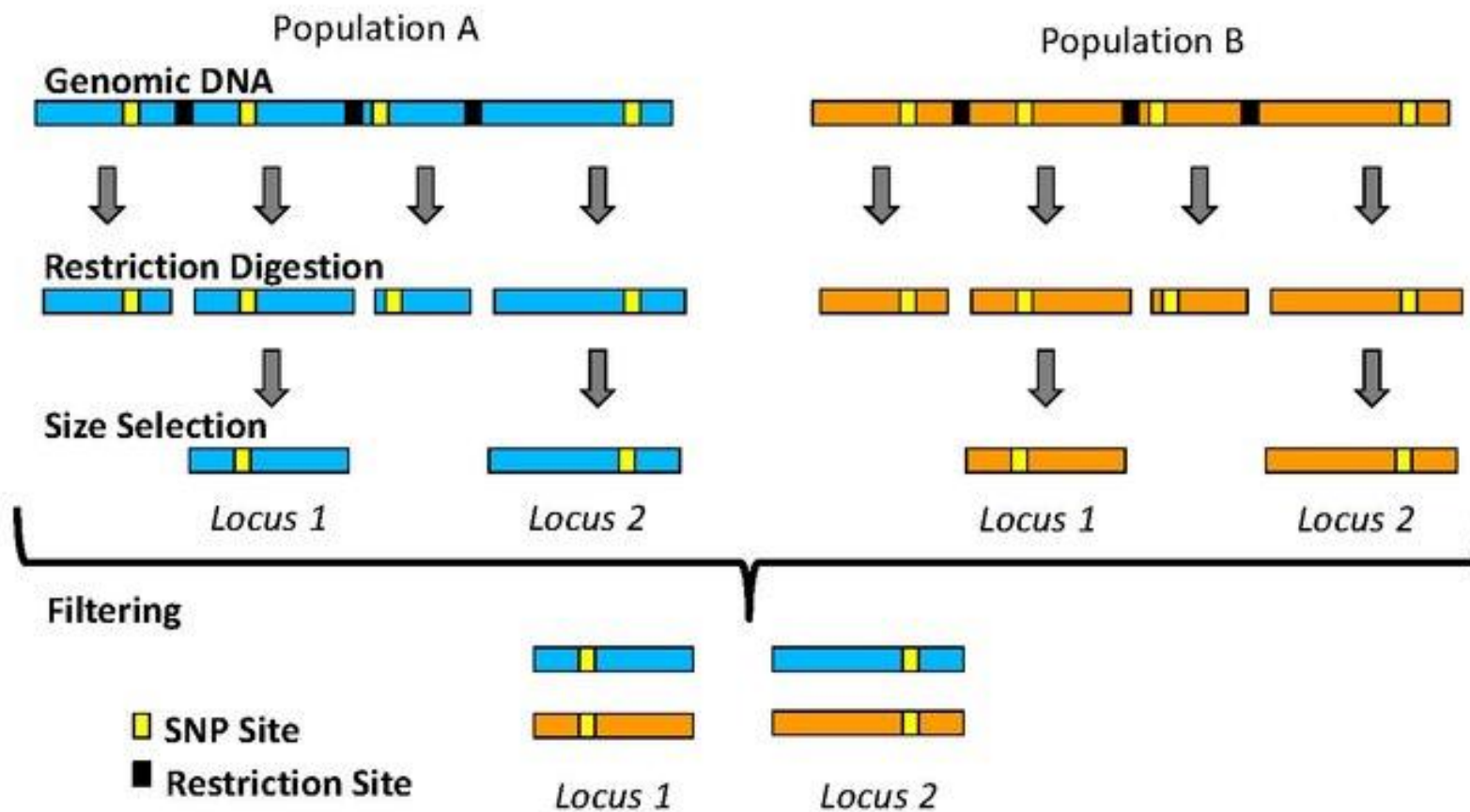Phylogenetic analysis



Image by Gene Marker

# Restriction Enzymes
# +
# NGS short read sequencing

# Reduced-representation NGS sequencing

- Restriction enzymes can be used to identify homologous blocks of the genome

- Valuable for comparative research in genomics, population genetics etc.

- 'RADseq'

Restriction-site Associate DNA Sequencing (RADseq)

Image by Jonathan Clark CCAS 4.0

# 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

Restriction site associated DNA (RAD) tags are a genome-wide representation of every site of a particular restriction enzyme by short DNA tags. Most organisms segregate large numbers of DNA sequence polymorphisms that disrupt restriction sites, which allows RAD tags to serve as genetic markers spread at a high density throughout the genome. Here, we demonstrate the applicability of RAD markers for both individual and bulk-segregant genotyping. First, we show that these markers can be identified and typed on pre-existing microarray formats. Second, we present a method that uses RAD marker DNA to rapidly produce a low-cost microarray genotyping resource that can be used to efficiently identify and type thousands of RAD markers. We demonstrate the utility of the former approach by using a tiling path array for the fruit fly to map a recombination breakpoint, and the latter approach by creating and using an enriched RAD marker array for the threespine stickleback. The high number of RAD markers enabled localization of a previously identified region, as well as a second region also associated with the lateral plate phenotype. Taken together, our results demonstrate that RAD markers, and the method to develop a RAD marker microarray resource, allow high-throughput, high-resolution genotyping in both model and nonmodel systems.

[Supplemental material is available online at www.genome.org.]

Miller et al. 2007, Gen. Res.

# Double Digest RADseq: An Inexpensive Method for *De Novo* SNP Discovery and Genotyping in Model and Non-Model Species

Brant K. Peterson*, Jesse N. Weber, Emily H. Kay, Heidi S. Fisher, Hopi E. Hoekstra

Department of Organismic & Evolutionary Biology, Department of Molecular & Cellular Biology, Museum of Comparative Zoology, Harvard University, Cambridge, Massachusetts, United States of America

**Abstract**

The ability to efficiently and accurately determine genotypes is a keystone technology in modern genetics, crucial to studies ranging from clinical diagnostics, to genotype-phenotype association, to reconstruction of ancestry and the detection of selection. To date, high capacity, low cost genotyping has been largely achieved via "SNP chip" microarray-based platforms which require substantial prior knowledge of both genome sequence and variability, and once designed are suitable only for those targeted variable nucleotide sites. This method introduces substantial ascertainment bias and inherently precludes detection of rare or population-specific variants, a major source of information for both population history and genotype-phenotype association. Recent developments in reduced-representation genome sequencing experiments on massively parallel sequencers (commonly referred to as RAD-tag or RADseq) have brought direct sequencing to the problem of population genotyping, but increased cost and procedural and analytical complexity have limited their widespread adoption. Here, we describe a complete laboratory protocol, including a custom combinatorial indexing method, and accompanying software tools to facilitate genotyping across large numbers (hundreds or more) of individuals for a range of markers (hundreds to hundreds of thousands). Our method requires no prior genomic knowledge and achieves per-site and per-individual costs below that of current SNP chip technology, while requiring similar hands-on time investment, comparable amounts of input DNA, and downstream analysis times on the order of hours. Finally, we provide empirical results from the application of this method to both genotyping in a laboratory cross and in wild populations. Because of its flexibility, this modified RADseq approach promises to be applicable to a diversity of biological questions in a wide range of organisms.

Peterson et al. 2012, PLoS ONE

A

RAD sequencing

X Rare cut site     ▬ Genomic interval present in library
X Common cut site     ▬ Sequence reads

Individual 1

Genomic DNA

Individual 2

B

double digest RADseq

a

b

Individual 1

Genomic DNA

Individual 2

# Examples of cut-sites

5′... C C T G C A▼G G ...3′
3′... G G▲A C G T C C ...5′

Sbfl – 'rare cutter' 6 nucleotide recognition sequence

5′...C▼C G G ...3′
3′...G G C▲C ...5′

MspI – 'common cutter' 3 nucleotide recognition sequence

Fraction of genome

Sanger sequencing

Whole genome re-sequencing

RADtag (Baird 2008)

ddRAD

Phylogeny    Population Structure    QTL Mapping    Pedigree Mapping    Association Mapping    Population Genomic Scans

Divergence limited    Recombination limited    Linkage Diseq. limited

Peterson et al. 2012, PLoS ONE 7: e37135

# ddRADseq Laboratory methods

- DNA extraction
- Restriction digestion with two Type II enzymes
- Ligation of adapters
- Size-selection
- PCR amplification
- Pooling and Illumina sequencing

# DNA extraction



Tissue sampling
(Muscle, liver, etc.)



DNA extraction
(Qiagen kit, Phenol-chloroform, salt extraction etc.)

# Quantification of dsDNA

- Like other genomic library construction protocols (e.g. Unit 2) ddRADseq requires specific starting concentrations of DNA.
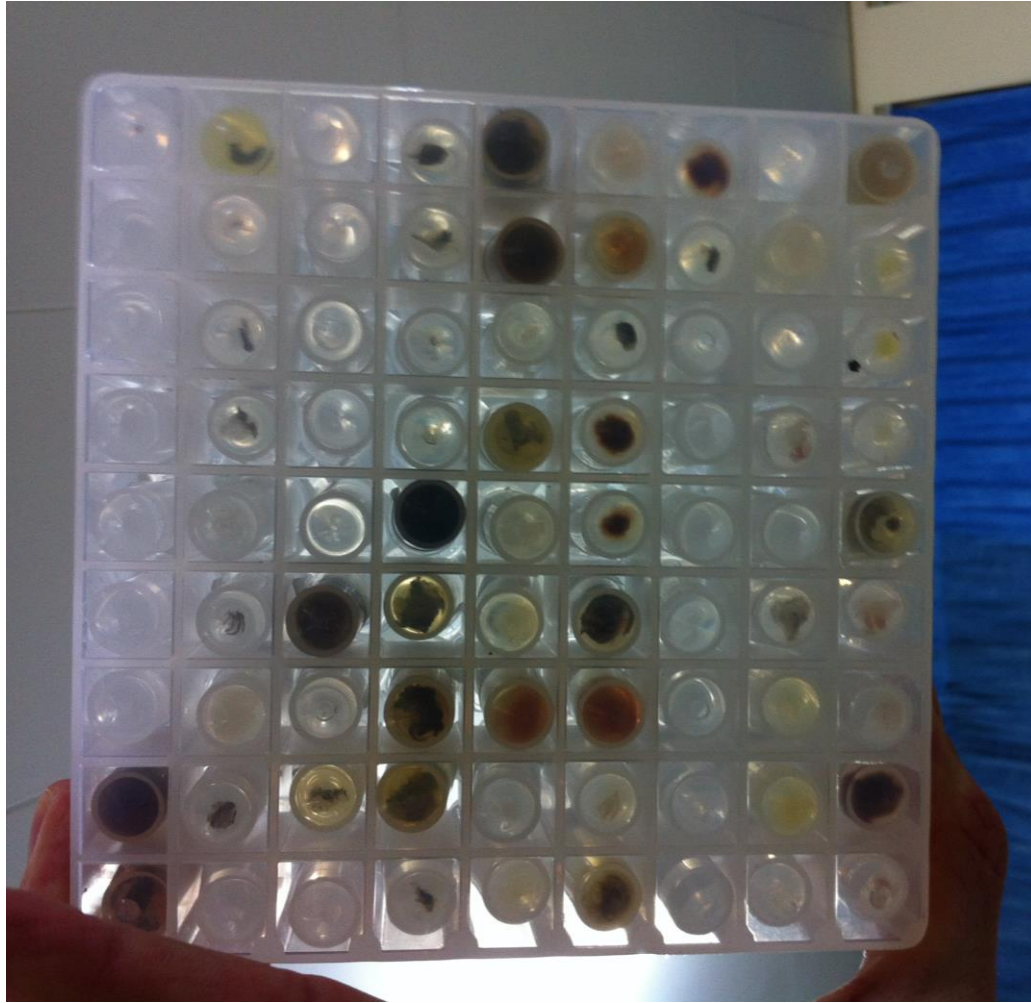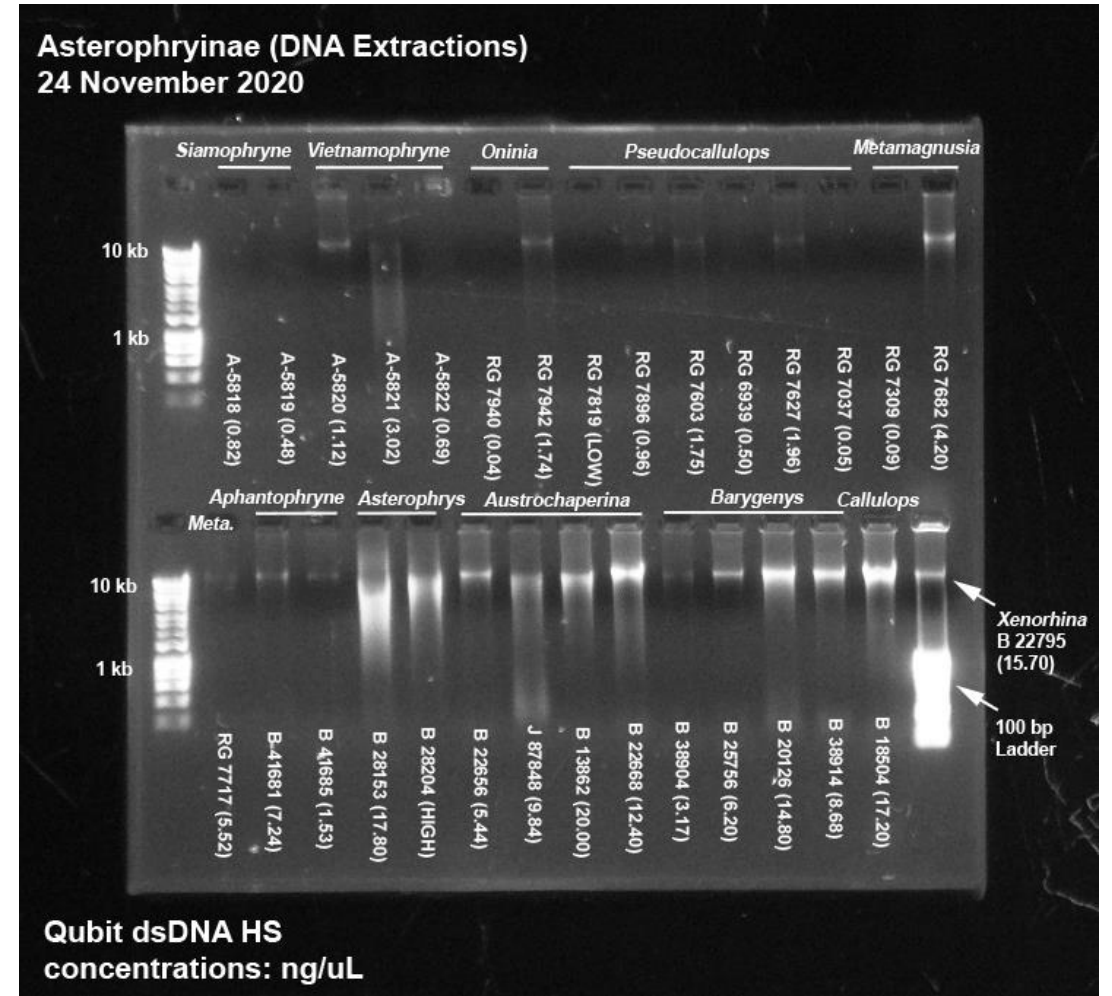
- We need to determine the concentration of double-stranded DNA (dsDNA) before or after the fragmentation.

- One of the most effective ways to do this (IMO) is with fluorometry.

- We will cover this during the molecular labs tomorrow and next week.

| Sample ID | Qubit concentration (ng/uL) | uL needed for 500 ng | uL of water to add |
|-----------|------------------------------|----------------------|---------------------|
| Sample 1  | 10.0                         | 50.0                 | 10.0                |
| Sample 2  | 18.5                         | 27.0                 | 33.0                |
| Sample 3  | 33.2                         | 15.1                 | 44.9                |
| Sample 4  | 80.0                         | 6.3                  | 53.7                |

Table from Unit 2 Molecular Lab Protocol
https://github.com/nhm-herpetology/museum-NGS-training/tree/main/Unit_02/Molecular_Lab



Qubit 2.0 Fluorometer

# Digestion



Digest at optimal temperature
e.g. 37 C for 6 hours



Streicher et al. 2014, Mol. Ecol.

# Ligation of ddRADseq adapters

- Designed specifically for RADseq protocol to match 'sticky' cut sites

- Similar in some ways to the adapters we talked about in Unit 1

- Let's review…

# Illumina Adapter Design

"Stubby, Y-Yoked Adapters"

- One oligo with terminal thymine (Required)

- One oligo with phosphorylated terminal nucleotide (Required)

- Illumina P5 and P7 recognition sequences (Required)

- Read 1 and Read 2 priming sequences (Required)

- Unique Index (for multiplexing; Required)

- Second Index (for multiplexing; Optional)

- Unique Molecular Identifier (UMI; Optional)

# Illumina Adapter Design

"Stubby, Y-Yoked Adapters"

P5 (i5) Illumina sequence

**AATGATACGGCGACCACCGAGATCT**ACACTCTTTCCCTACACGACGCTCTTCCGATC **T**

CGCTCTTCCGATCGTGTGCTCTTCCGATC***PHOS**

Read 1 priming sequence

Terminators

P7 (i7) Illumina sequence

**CAAGCAGAAGACGGCATACGAGAT**CGTGATGTGACTGGAGTTCAGACGTGTGCTCTTCCGATC***PHOS**

Sample Index

Read 2 priming sequence

# Illumina Adapter Design

## "Stubby, Y-Yoked Adapters"

**P5 (i5) Illumina sequence**

**AATGATACGGCGACCACCGAGATCT**ACACTCTTTCCCTACACGACGCTCTTCCGATC **T**

Read 1 priming sequence

GTGACTGGAGTTCAGACGTGTGCTCTTCCGATC**\*PHOS**

Terminators

**P7 (i7) Illumina sequence**

**CAAGCAGAAGACGGCATACGAGAT**CGTGATGTGACTGGAGTTCAGACGTGTGCTCTTCCGATC

Sample Index

Read 2 priming sequence



**AATGATACGGCGACCACCGAGATCT**ACACTCTTTCCCTACACGACGCTCTTCCGATC **T**

GTGCTCTTCCGATCA

CGTGATGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCA

**CAAGCAGAAGACGGCATACGAGAT**CGTGATGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCA

5'        3'

Genomic DNA

3'        5'

A CTAGCCT
T CTAGCCT

# Illumina Adapter Design

"Stubby, Y-Yoked Adapters"

- One oligo with terminal thymine (Required)

- One oligo with phosphorylated terminal nucleotide (Required)

- Illumina P5 and P7 recognition sequences (Required)

- Read 1 and Read 2 priming sequences (Required)

- Unique Index (for multiplexing; Required)

- Second Index (for multiplexing; Optional)

- Unique Molecular Identifier (UMI; Optional)
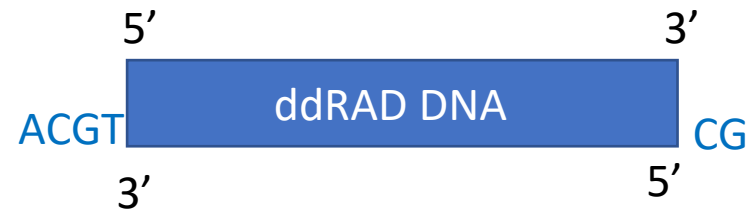
# ~~Illumina~~ Adapter Design

ddRADseq

~~"Stubby, Y-Yoked Adapters"~~

- ~~One oligo with terminal thymine (Required)~~

- One oligo with phosphorylated terminal nucleotide (Required)

- Illumina P5 and P7 recognition sequences (Required)

- Read 1 and Read 2 priming sequences (Required)

- Unique Index (for multiplexing; Required)

- Second Index (for multiplexing; Required)

- Unique Molecular Identifier (UMI; Optional)

- One oligo that matches the rare cutter site

- One oligo that matches the common cutter site

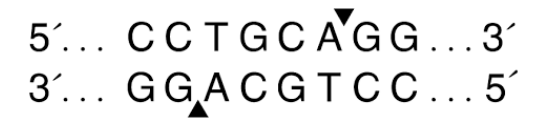# Standard Illumina Library Prep

5'                                3'

A | Genomic DNA | A

3'                                5'

# ddRADseq Illumina Library Prep

5'                                3'

ACGT | ddRAD DNA | CG

3'                                5'

# Standard Illumina Library Prep



# ddRADseq Illumina Library Prep



**SbfI**

5´... C C T G C A G G ...3´
3´... G G A C G T C C ...5´

**MspI**

5´...C C G G ...3´
3´...G G C C ...5´

**P1 Adapter**

```
ACACTCTTTCCCTACACGACGCTCTTCCGATCTNNNNNNNNACTAGGTGCA
TGTGAGAAAGGGATGTGCTGCGAGAAGGCTAGANNNNNNNNTGATCC -5'PHOS
```

Read 1 priming sequence      Unique Molecular Identifier      Sample index     SbfI cutsite remnant

**P1**
ACACTCTTTCCCTACACGACGCTCTTCCGATCT**NNNNNNNNNACTAGGTGCA**

**Adapter**
TGTGAGAAAGGGATGTGCTGCGAGAAGGCTAGA**NNNNNNNNNTGATCC** –5' PHOS

Read 1 priming sequence     Unique Molecular Identifier     Sample index     SbfI cutsite remnant

**P1 Adapter**

ACACTCTTTCCCTACACGACGCTCTTCCGATCT**NNNNNNNNNACTAGGTGCA**

TGTGAGAAAGGGATGTGCTGCGAGAAGGCTAGA**NNNNNNNNNTGATCC** –5'PHOS

**P2 Adapter**

GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT

AACAAGAGCGAGAAGGCTAGAGC – 5'PHOS

Read 1 priming sequence   Unique Molecular Identifier   Sample index   SbfI cutsite remnant

**P1
Adapter**

ACACTCTTTCCCTACACGACGCTCTTCCGATCT**NNNNNNNNACTAGGTGCA**
TGTGAGAAAGGGATGTGCTGCGAGAAGGCTAGA**NNNNNNNNTGATCC** –5'PHOS

**P2
Adapter**

Read 2 priming sequence

GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT

AACAAGAG CGAGAAGGCTAGA**GC** – 5'PHOS

MspI cutsite remnant

Forked adapter permits PCR only from P1 side in first cycle

Unique Molecular Identifier

Read 1 priming sequence    Sample index    SbfI cutsite remnant

**P1 Adapter**

ACACTCTTTCCCTACACGACGCTCTTCCGATCT**NNNNNNNNACTAGGTGCA**
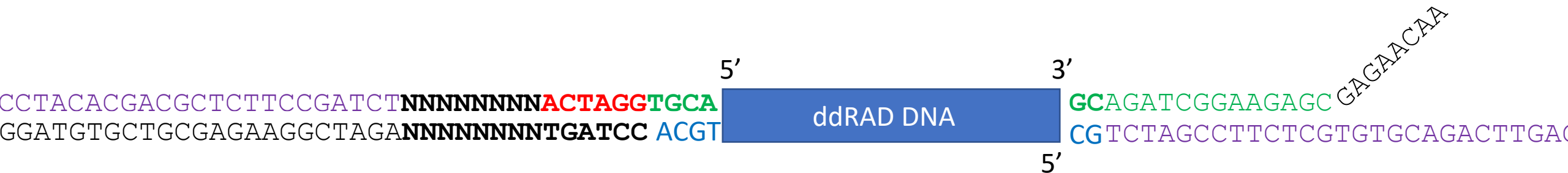TGTGAGAAAGGGATGTGCTGCGAGAAGGCTAGA**NNNNNNNNTGATCC** –5′ PHOS

**P2 Adapter**

Read 2 priming sequence

GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT
AACAAGAG CGAGAAGGCTAGA**GC** – 5′ PHOS

MspI cutsite remnant

Forked adapter permits PCR only from P1 side in first cycle

5′    3′

GAGAACAA

CCTACACGACGCTCTTCCGATCT**NNNNNNNNACTAGGTGCA**    | ddRAD DNA |    **GC**AGATCGGAAGAGC
GGATGTGCTGCGAGAAGGCTAGA**NNNNNNNNTGATCC** ACGT                        CG TCTAGCCTTCTCGTGTGCAGACTTGA

5′

# Size Selection

Illumina sequencers can only sequence DNA fragments >600 nucleotides in size, so making sure that the mean size of fragments in your libraries are smaller is critical.
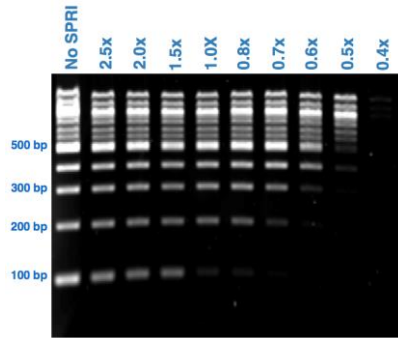


Image from Enseqlopedia
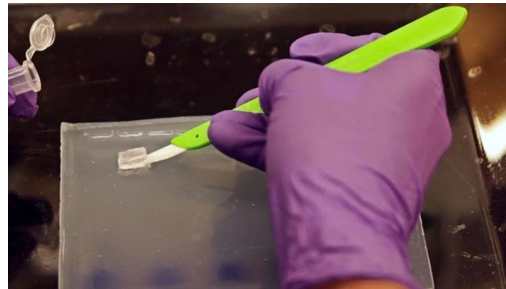
**Bead-based size selection**
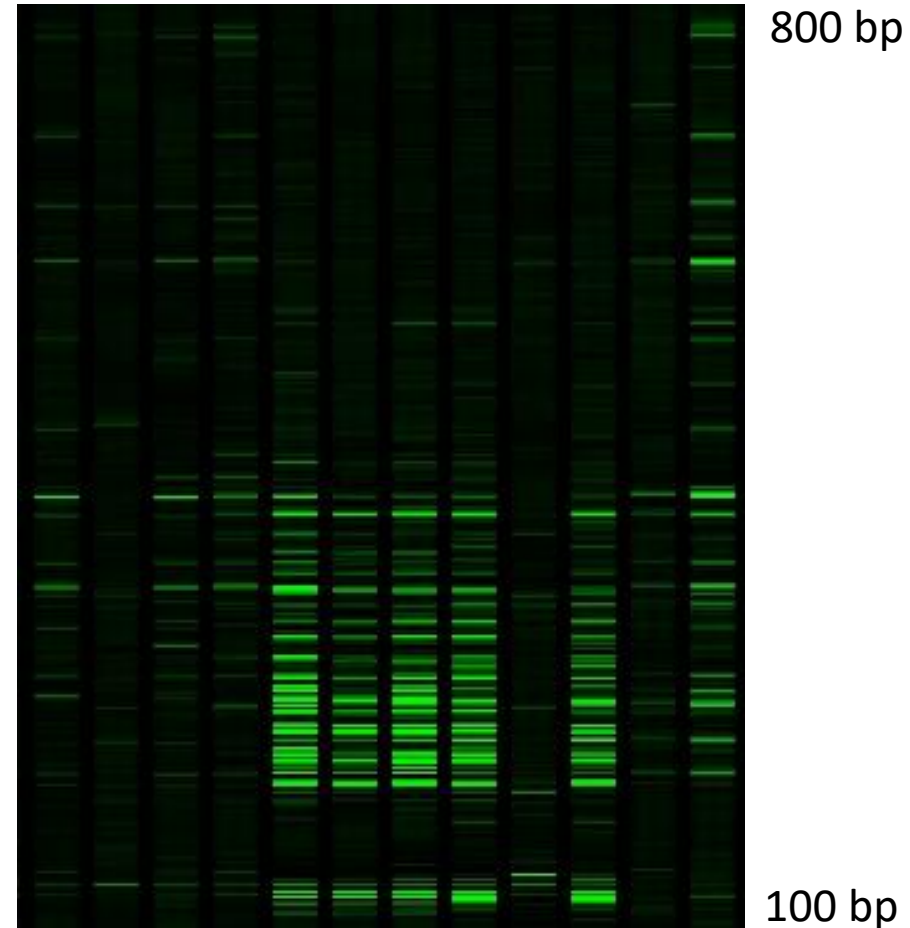


Image from NEB

**Gel-extraction size selection**



Blue Pippin (Sage Science)

**Automated Size Selection**

# ddRADseq – size selection determines which loci you will sequence



800 bp

Amplified Fragment Length Polymorphisms (*Anolis*)
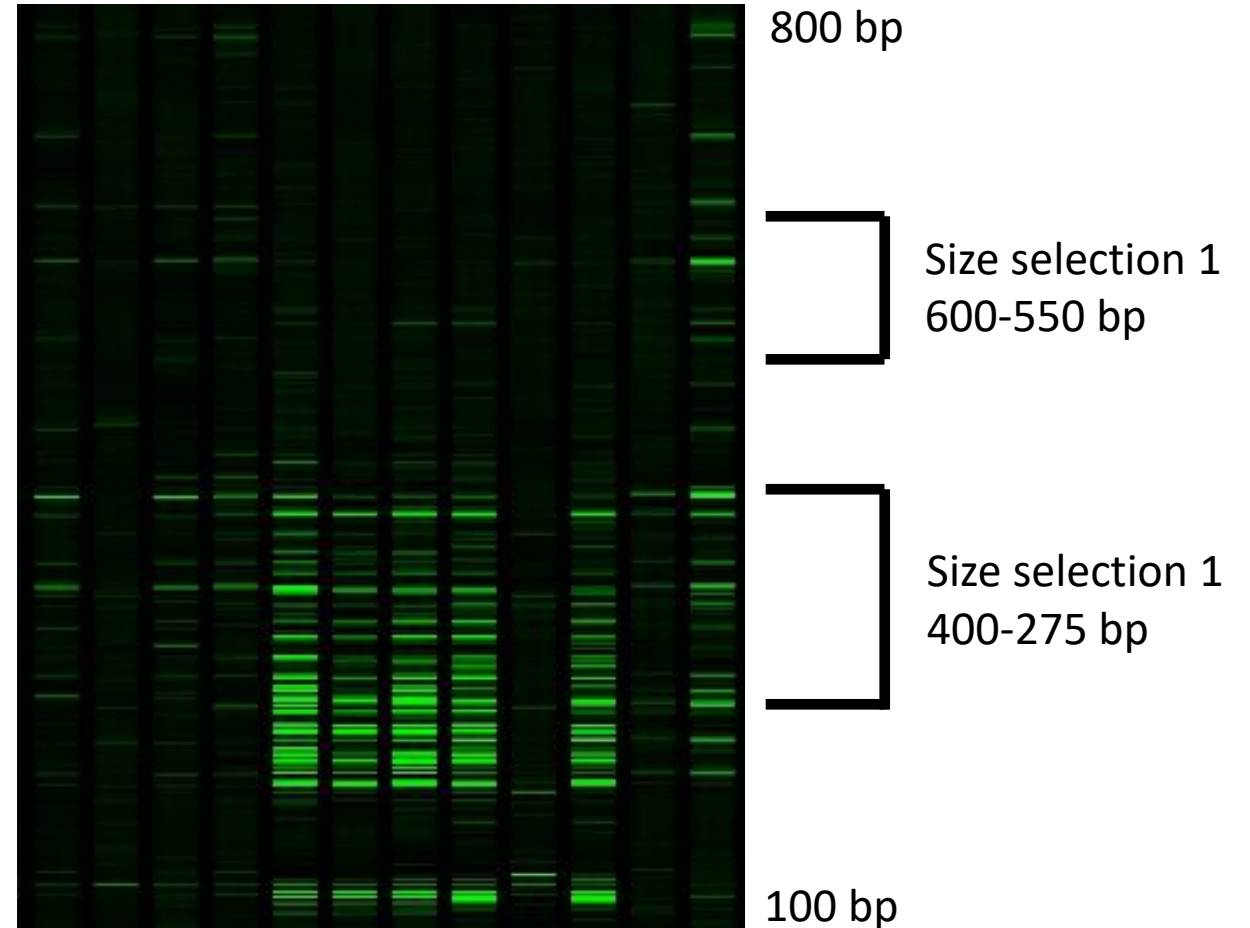
100 bp

Anthony Genova

DISCLAIMER: I made up the AFLP gel sizes for the purpose of demonstration

# ddRADseq – size selection determines which loci you will sequence



800 bp

Size selection 1
600-550 bp

Size selection 1
400-275 bp

100 bp

Anthony Genova

Amplified Fragment Length Polymorphisms (*Anolis*)

DISCLAIMER: I made up the AFLP gel sizes for the purpose of demonstration
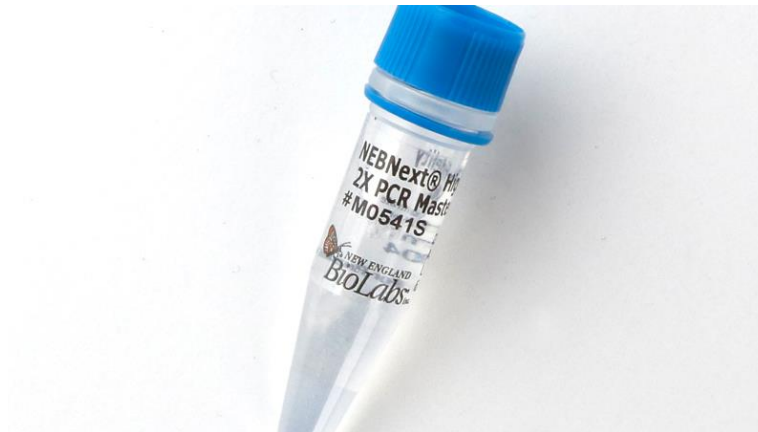
# Limited PCR Amplification

PCR usually of 8-12 cycles

**PCR Primers Standard Illumina Library Prep**
TruSeq P5: AATGATACGGCGACCACCGAGA
TruSeq P7: CAAGCAGAAGACGGCATACGAG

**Hi-Fidelity Polymerase**

# Limited PCR Amplification

PCR usually of 8-12 cycles

**PCR Primers Standard Illumina Library Prep**
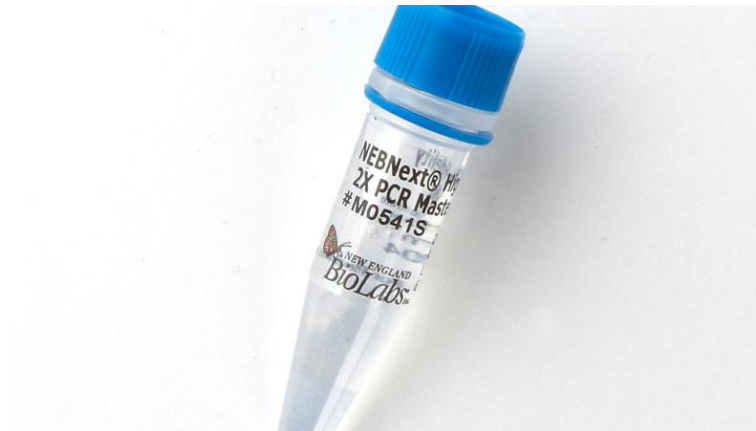TruSeq P5: AATGATACGGCGACCACCGAGA
TruSeq P7: CAAGCAGAAGACGGCATACGAG

**PCR Primers ddRADseq Illumina Library Prep**
PCR 1: AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACG
PCR 2: CAAGCAGAAGACGGCATACGAGATNNNNNNGTGACTGGAGTTCAGACGTGTGC

**Hi-Fidelity Polymerase**

# Why the long primers?

# Why the long primers?

# Why the long primers?

AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACG

PCR 1 primer

ACACTCTTTCCCTACACGACGCTCTTCCGATCT**NNNNNNNNNACTAGGTGCA**

TGTGAGAAGGGATGTGCTGCGAGAAGGCTAGA**NNNNNNNNTGATCC** ACGT

5'

ddRAD DNA

3'

5'                                                                  3'

**GCA**                    GAGAACAA

**ACGT**                    **GC**AGATCGGAAGAGC

ddRAD DNA            **CG** TCTAGCCTTCTCGTGTGCAGACTTGAGGTCAGTG

3'                          5'            CGTGTGCAGACTTGAGGTCACTGTAGTGCTAGAGCATACGGGAGAAGACGAAC

PCR 2 primer

# Why the long primers?

Illumina P5 flow cell

**AATGATACGGCGACCACCGA**GATCTACACTCTTTCCCTACACGACG

PCR 1 primer

ACACTCTTTCCCTACACGACGCTCTTCCGATCT**NNNNNNNNN**<span style="color:red">**ACTAGG**</span><span style="color:green">**TGCA**</span>

TGTGAGAAGGGATGTGCTGCGAGAAGGCTAGA**NNNNNNNNNTGATCC** <span style="color:blue">ACGT</span>

5'

ddRAD DNA

3'

---

5'                                                    3'              <span style="color:green">**GC**</span>AGATCGGAAGAGC    GAGAACAA

<span style="color:green">**GCA**</span>

<span style="color:blue">ACGT</span>        ddRAD DNA                         <span style="color:blue">CG</span>TCTAGCCTTCTCGTGTGCAGACTTGAGGTCAGTG                    <span style="color:#3ba9d6">*Second index*</span>    Illumina P7 flow cell

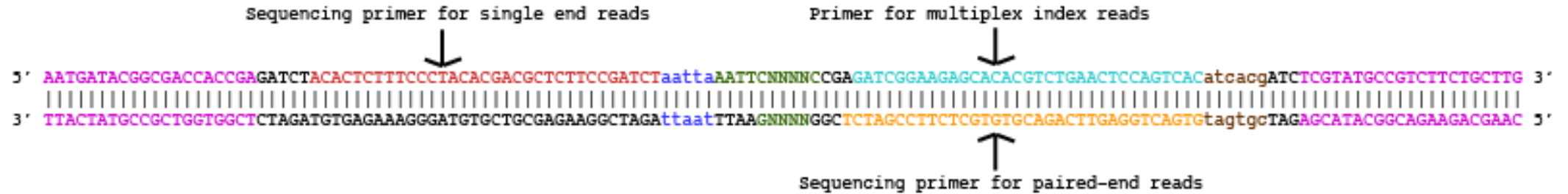3'                                                    5'                   CGTGTGCAGACTTGAGGTCACTG<span style="color:#3ba9d6">**TAGTGC**</span>TAG**AGCATACGGGAGAAGACGAAC**
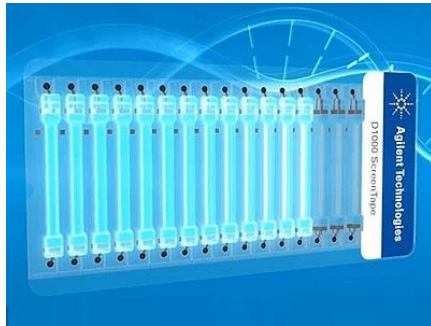
PCR 2 primer

Final sequencing library

Sequencing primer for single end reads

Primer for multiplex index reads

5' AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCTaattaAATTCNNNNCCGAGATCGGAAGAGCACACGTCTGAACTCCAGTCACatcacgATCTCGTATGCCGTCTTCTGCTTG 3'
3' TTACTATGCCGCTGGTGGCTCTAGATGTGAGAAAGGGATGTGCTGCGAGAAGGCTAGAttaatTTAAGNNNNGGCTCTAGCCTTCTCGTGTGCAGACTTGAGGTCAGTGtagtgcTAGAGCATACGGCAGAAGACGAAC 5'

Sequencing primer for paired-end reads

DNA Sequence Legend

READ 1 primer
READ 2 primer
MULTIPLEX READ primer
genomic DNA
barcode (aatta) - inline
index (atcacg) - multiplex
flowcell annealing

Peterson et al. 2012, PLoS ONE 7: e37135

# Quantification of genomic DNA libraries

- Reasonably precise estimates of DNA concentration are needed for Illumina sequencer input



D1000 Screentape (Agilent)

TapeStation 2200 (Agilent)

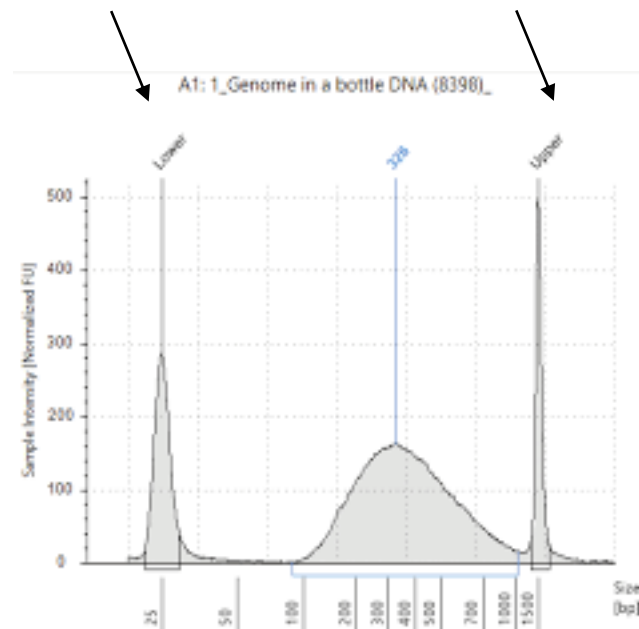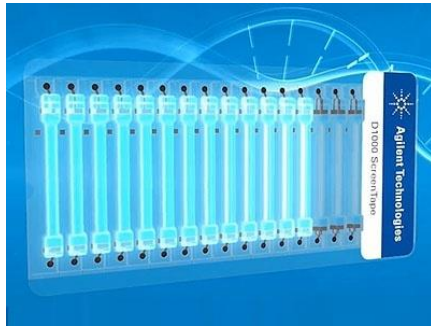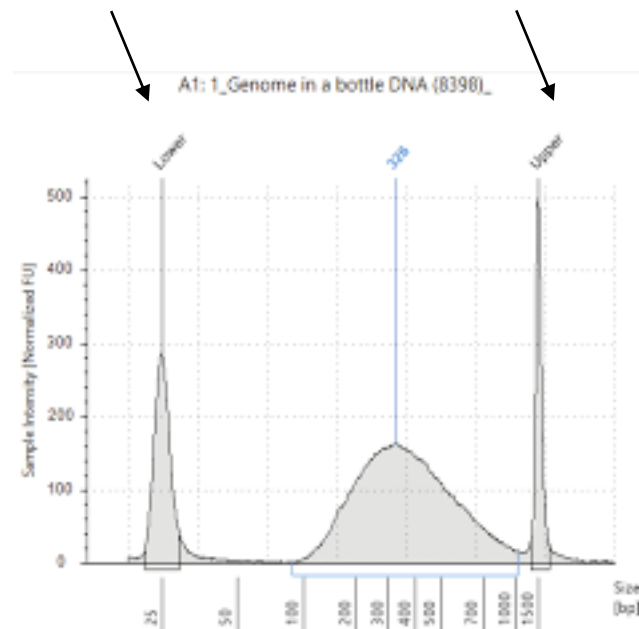Lower size standard

Upper size standard

Illumina ddRADseq library

Image from Agilent

# Quantification of ddRADseq libraries

- Reasonably precise estimates of DNA concentration are needed for Illumina sequencer input



D1000 Screentape (Agilent)

TapeStation 2200 (Agilent)

Lower size standard

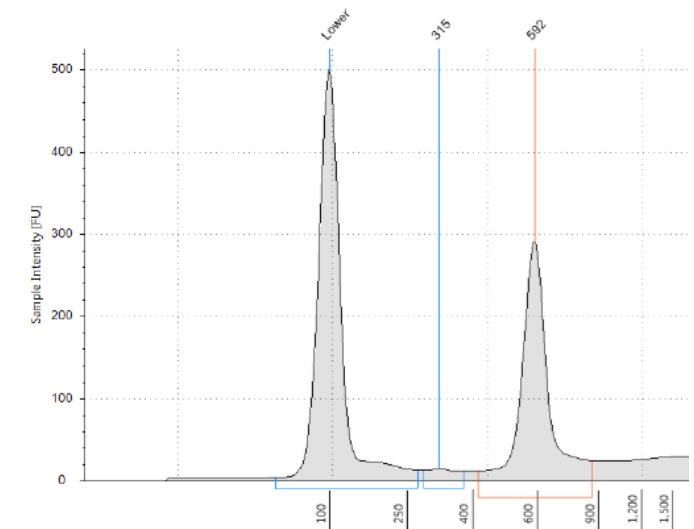Upper size standard

A1: 1_Genome in a bottle DNA (8398)_

Image from Agilent

Illumina ddRADseq library

B1: Agama atra_FINAL LIBRARY

Image from K. Alujevic

Tighter size selection for smaller number of RAD loci

# Analyzing ddRADseq data

- Stacks (Catchen et al.) provides a convenient way to demultiplex and sort data

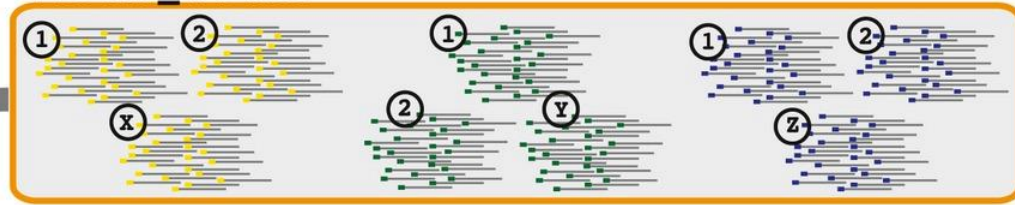- We will use this software in the bioinformatics lab today

- A brief overview…

## Stacks: an analysis tool set for population genomics

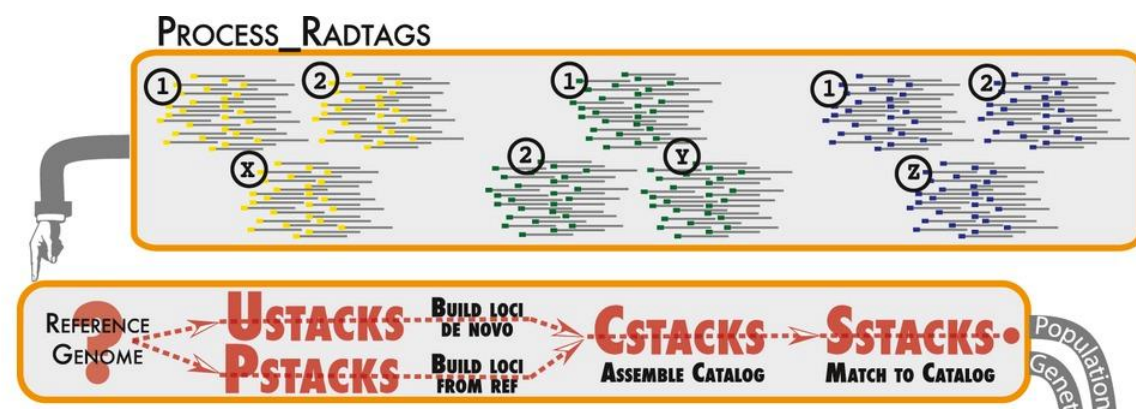Julian Catchen, Paul A. Hohenlohe, Susan Bassham, Angel Amores, William A. Cresko ✉

Original Article | 🔓 Full Access

# Stacks: an analysis tool set for population genomics

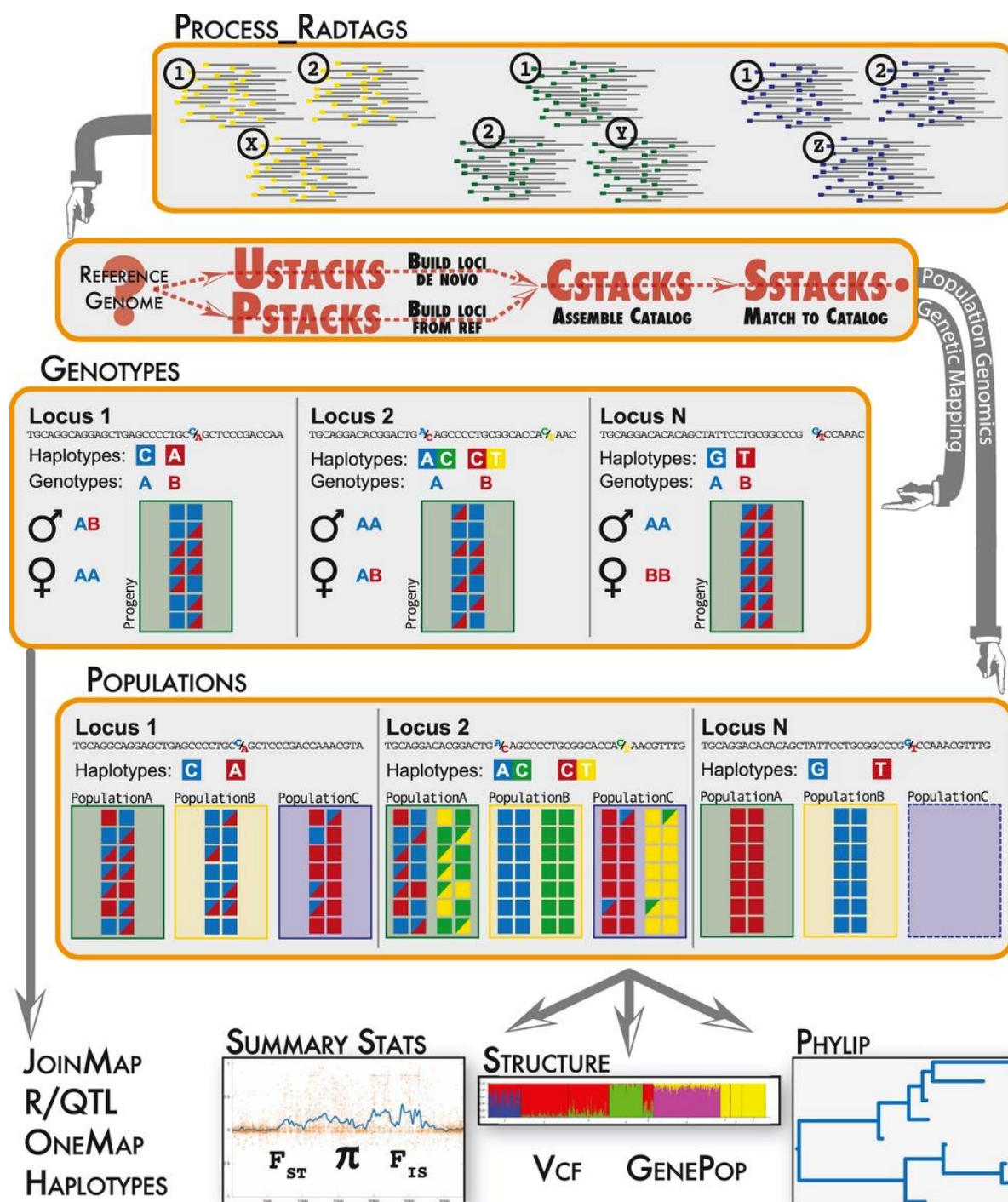Julian Catchen, Paul A. Hohenlohe, Susan Bassham, Angel Amores, William A. Cresko ✉

# Stacks: an analysis tool set for population genomics

Julian Catchen, Paul A. Hohenlohe, Susan Bassham, Angel Amores, William A. Cresko ✉

# Unit 4: Double digest restriction-site associated DNA sequencing (ddRADseq)

Bioinformatics Lab



https://github.com/nhm-herpetology/museum-NGS-training