

RAD-seq in Roscoff

Matthieu Bruneaux

2015-03-10

Mini-workshop about ddRAD

Introduction about RAD-seq

- ▶ RAD? RAD-seq? ddRAD?
- ▶ Applications
- ▶ Workflow

Practicals

- ▶ One complete project, from raw reads to final results
- ▶ Cherry-picking of some analysis steps
- ▶ Open questions

Objectives

- ▶ Overview of RAD-seq
- ▶ Arouse curiosity
- ▶ Give useful pointers

Disclaimer about the speaker!

- ▶ Not a population geneticist, not a bioinformatician
- ▶ Evolutionary biologist who dropped into a RAD-seq project when he was a small post-doc
- ▶ Some things said here are probably incorrect or plainly wrong!

What are RAD markers?

Miller et al. 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}

Description of RAD markers

- ▶ **Restriction site associated DNA** fragments
- ▶ Used with micro-array systems
- ▶ Similar to RFLP or AFLP, but many more markers

RAD - Miller et al. 2007 (6 steps)

Digest - tag - shear

1) Digest DNA samples



2) Ligate Linkers



3) Shear



RAD - Miller et al. 2007 (6 steps)

Purify - release - type

4) Purify RAD tags



S1



S2

5) Release RAD tags



S1



S2

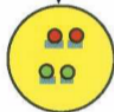
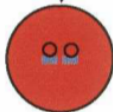
6) Label and hybridize
to identify or type
RAD markers



S1

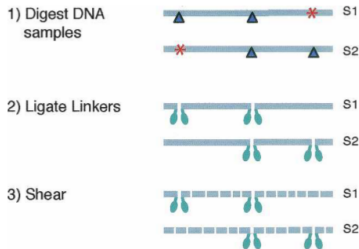


S2

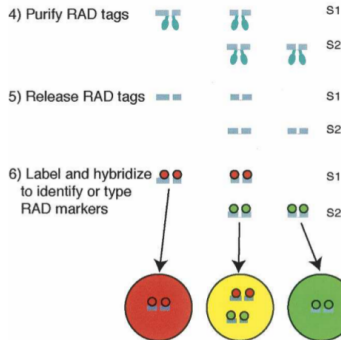


RAD - Miller et al. 2007 (method summary)

Digest - tag - shear



Purify - release - type



Demonstration

- ▶ Mapping breakpoint on a *Drosophila* chromosome
- ▶ Identification of the lateral plate locus in threespine stickleback

Advantage of the method

- ▶ Easy-to-produce genotyping resource for **non-model species**
- ▶ **Moderate cost**
- ▶ **Genetic mapping** possible (if markers location known)
- ▶ **Bulk genotyping** possible

But note that...

- ▶ At this point **the restriction site is the polymorphic marker**
- ▶ **One restriction enzyme** only is used

What is RAD-seq?

Baird et al. 2008

OPEN ACCESS Freely available online



Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers

Nathan A. Baird¹, Paul D. Etter¹, Tressa S. Atwood², Mark C. Currey³, Anthony L. Shiver¹, Zachary A. Lewis¹, Eric U. Selker¹, William A. Cresko³, Eric A. Johnson^{1*}

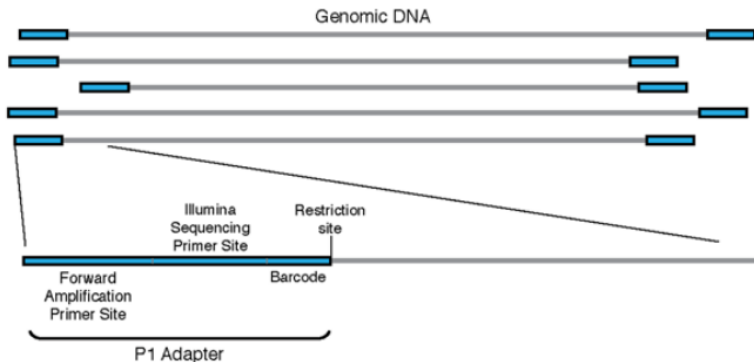
1 Institute of Molecular Biology, University of Oregon, Eugene, Oregon, United States of America, **2** Floragenex, Eugene, Oregon, United States of America, **3** The Center for Ecology and Evolutionary Biology, University of Oregon, Eugene, Oregon, United States of America

RAD-seq

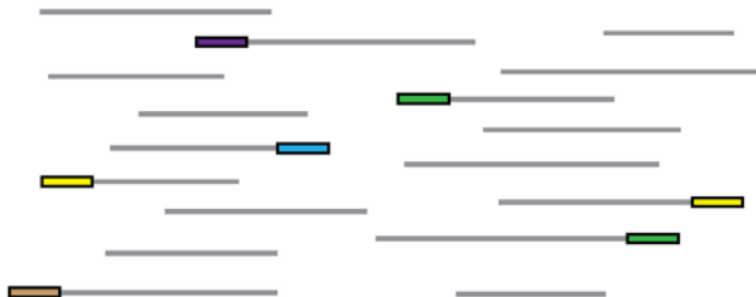
- ▶ RAD fragments with **high-throughput sequencing** (Illumina)
- ▶ SNP identified by **sequence polymorphism** and **site disruption**
- ▶ Can be used **with or without reference genome**

RAD-seq - Baird 2008

A *Ligate P1 Adapter to digested genomic DNA*

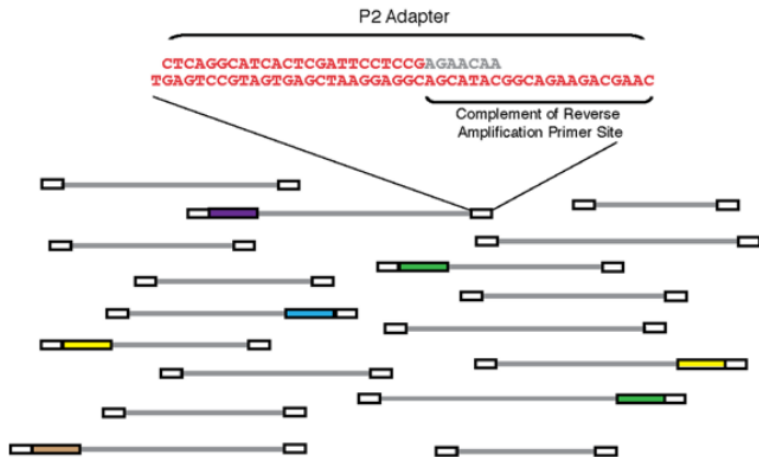


B *Pool barcoded samples and shear*



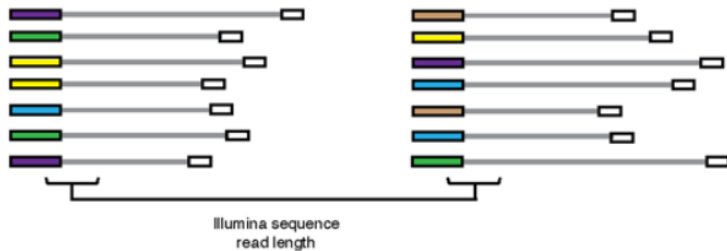
RAD-seq - Baird 2008

C *Ligate P2 Adapter to sheared fragments*



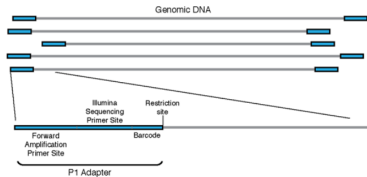
RAD-seq - Baird 2008

D *Selectively amplify RAD tags*

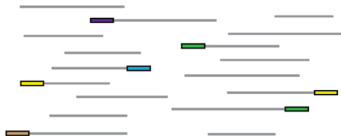


RAD-seq - Baird 2008

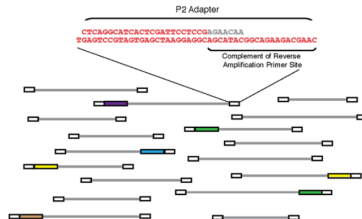
A *Ligate P1 Adapter to digested genomic DNA*



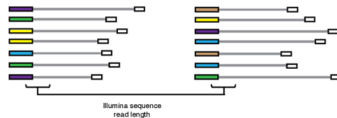
B *Pool barcoded samples and shear*



C *Ligate P2 Adapter to sheared fragments*



D *Selectively amplify RAD tags*

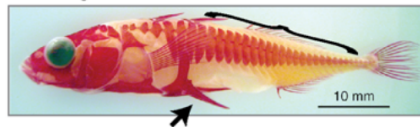


Demonstration

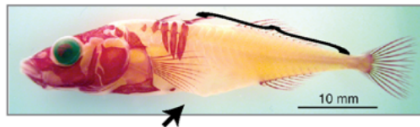
- ▶ Discover **13000 SNP** in threespine stickleback and in *Neurospora*
- ▶ **Barcoding** system for multiplexing
- ▶ **Marker density** can be tuned by the choice of restriction enzyme

Threespine stickleback

Rabbit Slough



Bear Paw



Population genomics of parallel adaptation - Hohenlohe 2010

A major paper

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*}

Method

- ▶ Model: threespine stickleback
- ▶ Comparison of 3 freshwater and 2 marine populations
- ▶ 20 individuals per population, individual barcodes
- ▶ Single reads (not paired ends)

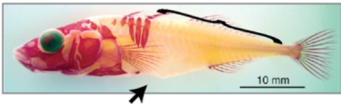
Population genomics of parallel adaptation - Hohenlohe 2010

Gasterosteus aculeatus

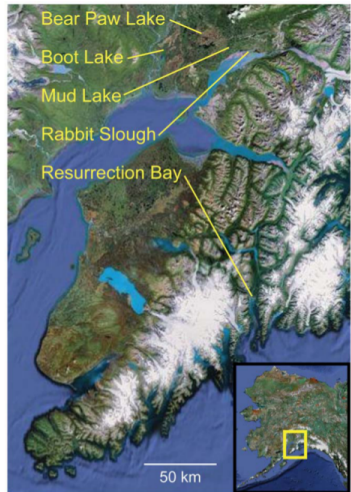
Rabbit Slough

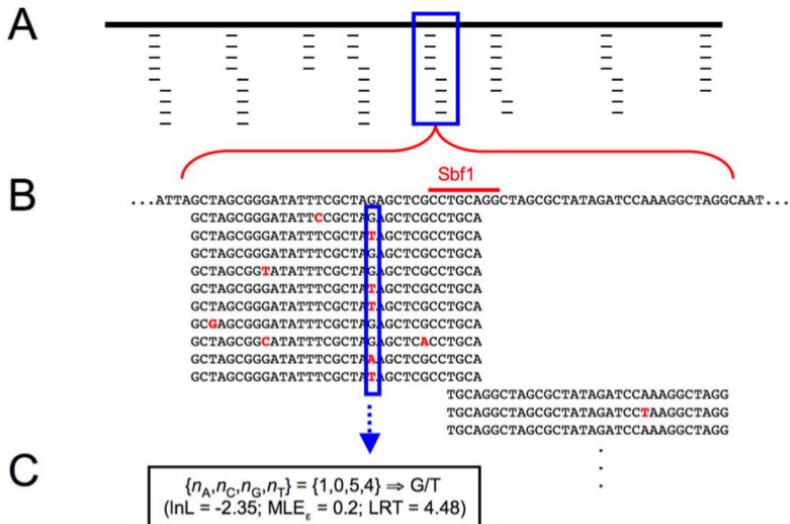


Bear Paw



Locations



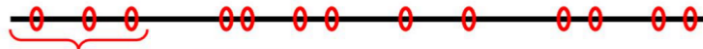


D

```
Ind 1 nnnnnGCTAGCGGGATATTTTCGCTAGAGCTCGCCTGCAGGCTAGCGCTATAGATCCAAAGGCTAGGnnnnn
      nnnnnGCTAGCGGGATATTTTCGCTATAGCTCGCCTGCAGGCTAGCGCTATAGATCCAAAGGCTAGGnnnnn
Ind 2 nnnnnGCTAGCGGGATATTTTCGCTAGAGCTCGCCTGCAGGCTAGCGCTATAGATCCTAAGGCTAGGnnnnn
      nnnnnGCTAGCGGGATATTTTCGCTAGAGCTCGCCTGCAGGCTAGCGCTATAGATCCTAAGGCTAGGnnnnn
Ind 3 nnnnnGCTAGCGGGATATTTTCGCTATAGCTCGCCTGCAGGCTAGCGCTATAGATCCAAAGGCTAGGnnnnn
      nnnnnGCTAGCGGGATATTTTCGCTATAGCTCGCCTGCAGGCTAGCGCTATAGATCCAAAGGCTAGGnnnnn
```

⋮
⋮

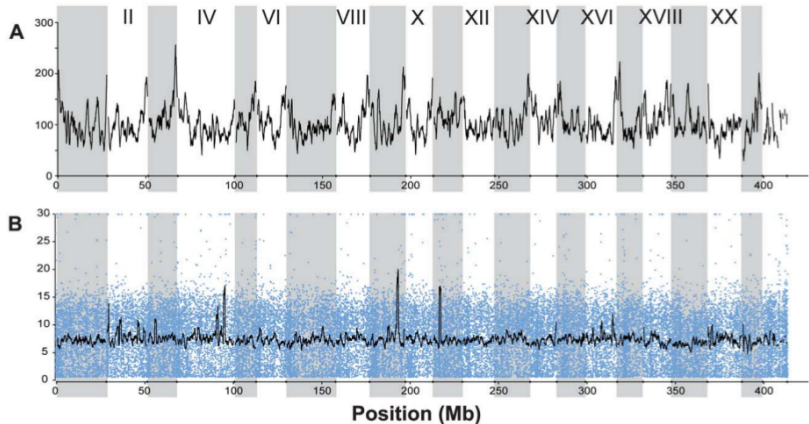
E



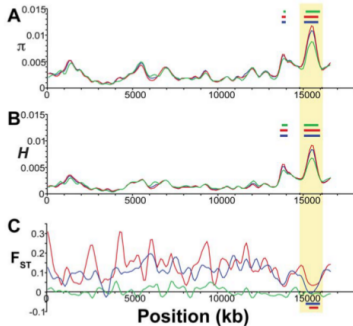
F



Hohenlohe 2010 - Genome profiles



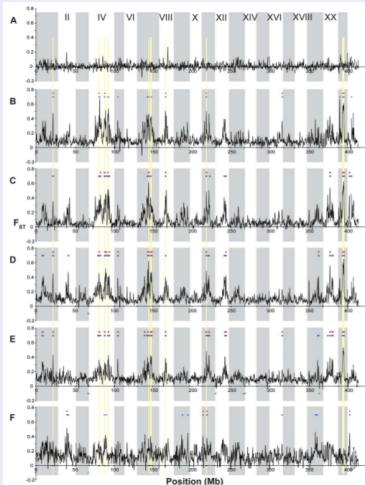
- ▶ A: number of RAD tags per 1Mb
- ▶ B: Coverage per RAD per individual in one run (16 individuals - black line is average)



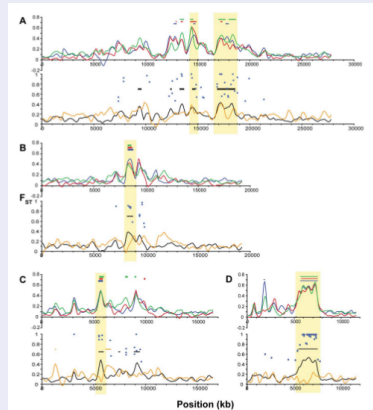
Evidence for balancing selection

- ▶ A: Nucleotide diversity, B: heterozygosity across all five populations (blue), three FW (red) or two SW (green)
- ▶ C: F_{ST} between FW and SW (blue), among FW (red) and among SW (green)
- ▶ Horizontal bars shows regions of significantly elevated or reduced values on the profile

Genome-wide differentiation among populations



Differentiation among SW and FW, zoom on LG



Highlights

- ▶ RAD-seq on **natural populations**, 45000 SNPs in 100 individuals
- ▶ **Barcoded** samples
- ▶ Genome profiling, **kernel smoothing** and **permutation testing**

But note that...

- ▶ Genome available
- ▶ Single reads

What is paired-end RAD-seq?

Etter 2011

OPEN ACCESS Freely available online



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

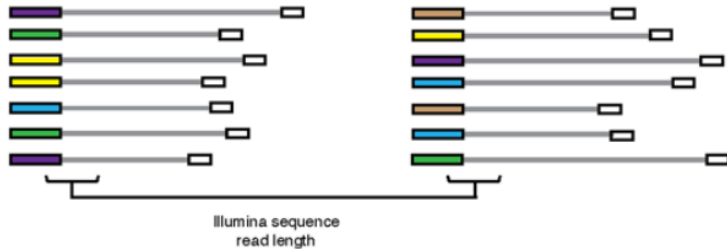
Paul D. Etter¹, Jessica L. Preston¹, Susan Bassham², William A. Cresko², Eric A. Johnson^{1*}

Method

- ▶ Paired-end sequencing of RAD fragments to build contigs on the randomly sheared side
- ▶ Demonstration with threespine and *E. coli* sequencing
- ▶ Up to 5kb contigs with circularization step

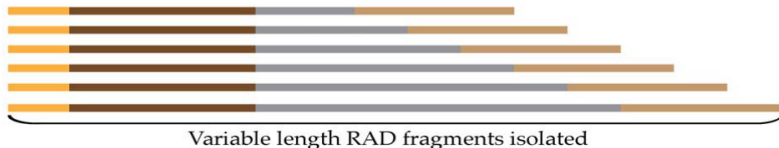
Single-reads RAD-seq

D *Selectively amplify RAD tags*



Paired-ends RAD-seq

B)



C)



Notes

- ▶ The **stacked end** is useful for **high coverage** work (SNP calling, allele frequency estimates)
- ▶ The **echelon end** is useful for contig building, but **base coverage is lower**

What is double-digest RAD-seq?

Peterson et al. 2012

OPEN ACCESS Freely available online



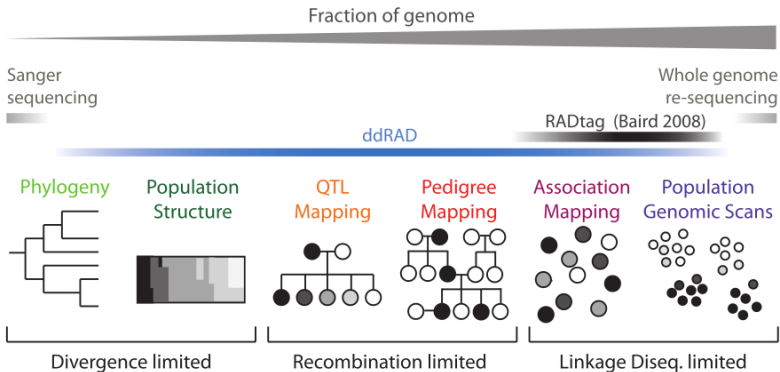
Double Digest RADseq: An Inexpensive Method for *De Novo* SNP Discovery and Genotyping in Model and Non-Model Species

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

Method

- ▶ Two enzyme double digest followed by precise size selection
- ▶ Library contains only fragments close to target size
- ▶ Read counts across regions are expected to be correlated between individuals

Double digest RAD tag



What is paired-end double RAD?

Bruneaux et al. 2013

Molecular evolutionary and population genomic analysis of the nine-spined stickleback using a modified restriction-site-associated DNA tag approach

MATTHIEU BRUNEAUX,^{*1} SUSAN E. JOHNSTON,^{*1} GÁBOR HERCZEG,[†] JUHA MERILÄ,[†]
CRAIG R. PRIMMER^{*} and ANTI VASEMÄGI^{*‡}

Method

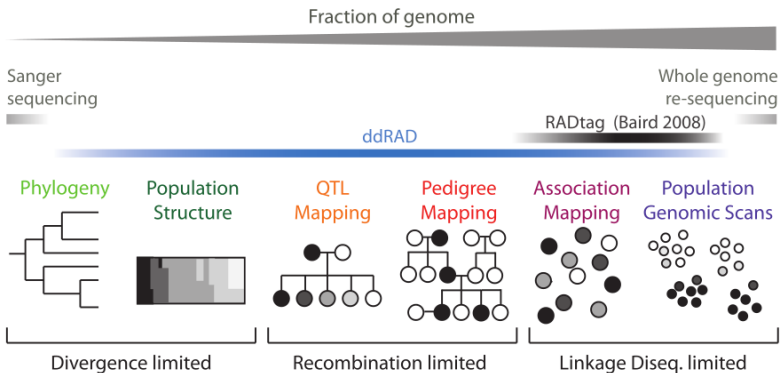
- ▶ Two enzyme double digestion
- ▶ Paired-end sequencing on after size-selection
- ▶ You will hear more about it soon (see practicals)

Paired-end double RAD

Add a picture

Uses of RAD tags

From Peterson 2012



Uses of RAD tags

Population genomics

Hess 2013 - Pacific lamprey

QTL mapping

Houston 2012

Phylogeography

Emerson 2011

Phylogenies

Rubin 2012

There are also some potential issues...

- ▶ PCR-duplicates
- ▶ individual vs pool genotyping for allele frequencies
- ▶ Comparison SNP vs microsat (deFaveri)

Conclusion

In a nutshell

- ▶ **RAD tags**: versatile method of **genome complexity reduction**
- ▶ **RAD-seq**: large scale discovery of SNPs, affordable
- ▶ Useful for both **model** and **non-model** organisms
- ▶ **Just a tool**: the downstream analyses are still **your expertise**

General workflow scheme

Development of pipelines and tools

Rainbow, STACKS, GATK, dDocent

Tools for NGS can be used for RAD

Simple scripts can be used also

This is one thing I want to show during the practical
Get a good grip and a good feeling/understanding about the data with simple, straightforward methods before choosing to apply more complex methods which rely on third-party scripts. It is important to understand what the third party scripts do!

One complete project

Tour of other tools and specific analyses

To illustrate some specific points (e.g. likelihood or bayesian based genotyping or allele frequency estimates or F_{st} calculations, ...)