

Population genomics

RAD-seq

Eyal Privman

Department of Ecology and Evolution, University of Lausanne

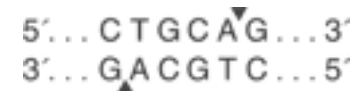
PhD Summer School
June 2012



Swiss Institute of
Bioinformatics

Fire ant population RAD-seq data

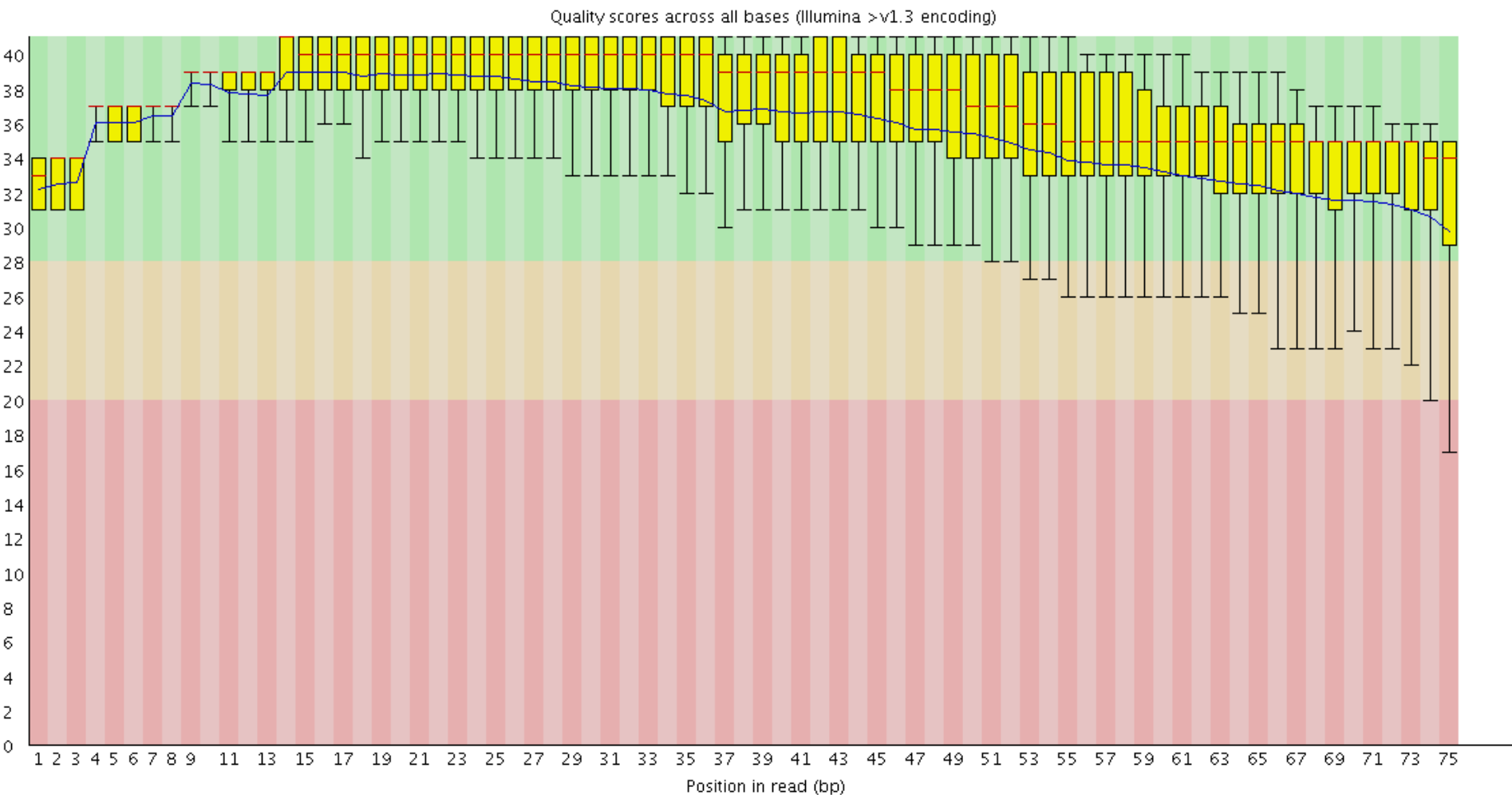
- 94 diploid workers (one per nest in sampled region)
- RAD protocol (Baird *et al.* Plos One 2008)
PstI restriction enzyme:
6-cutter, expect 1 site per 10Kbp
- Multiplexed Illumina 75bp single read sequencing
- **Surprise:** Not one simple population!



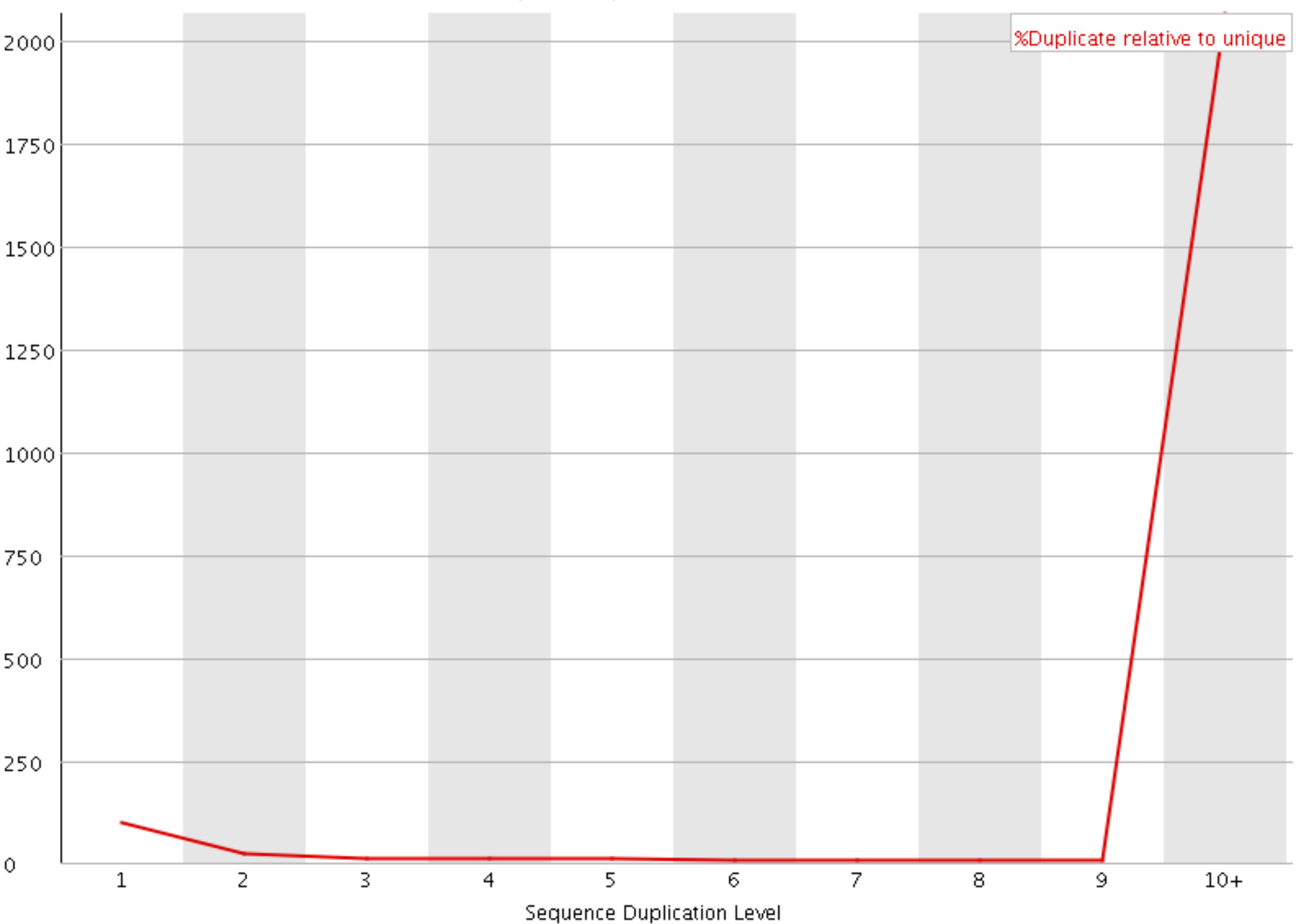
Stacks

- <http://creskolab.uoregon.edu/stacks/>
- Process RAD-seq data from either
 - Parents and progeny samples → linkage map
 - Population samples → population genetics
- Can work with / without reference genome
- Can assemble mini-contigs of second-reads for paired-end sequencing

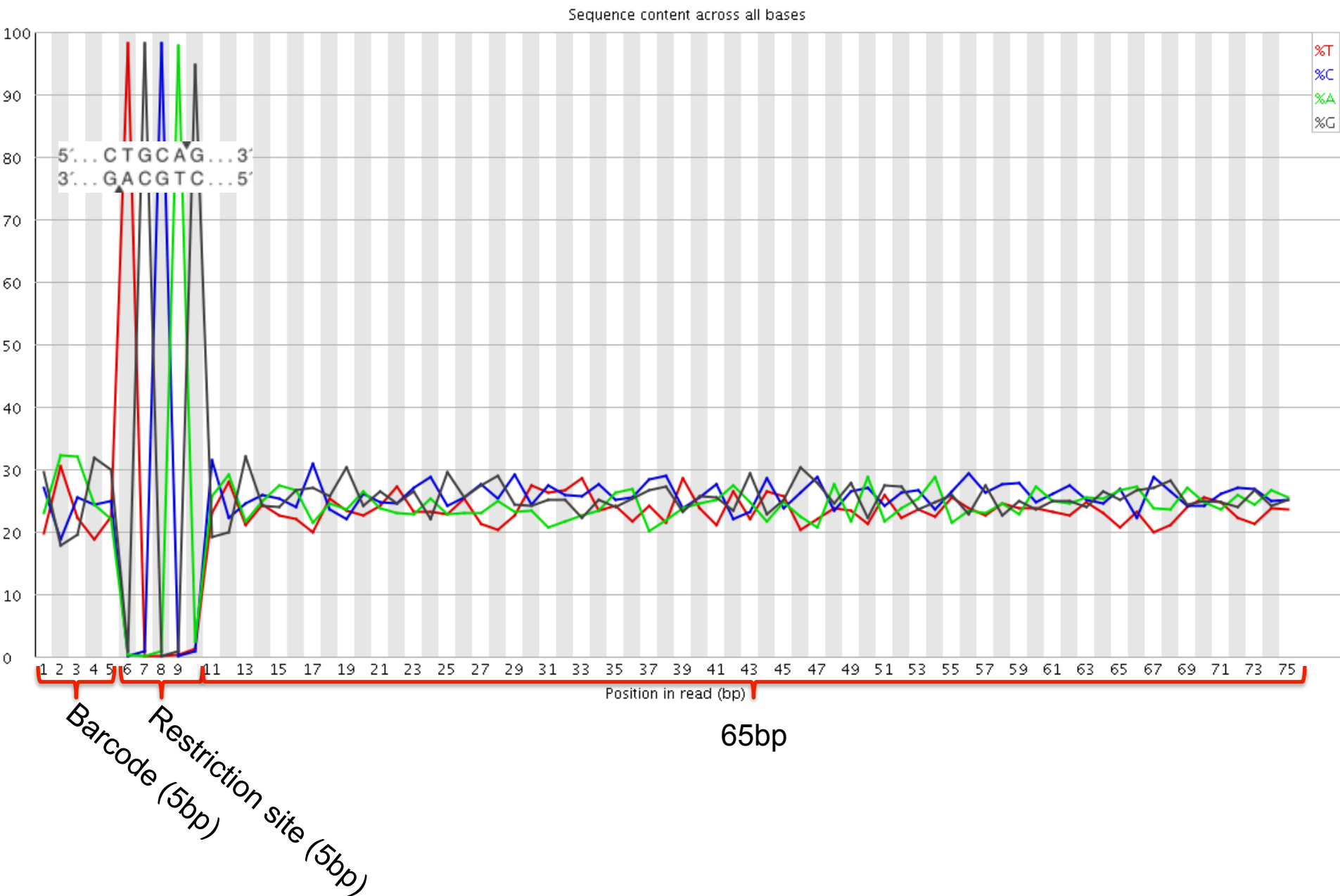
Quality



Sequence Duplication Level $\geq 96.22\%$

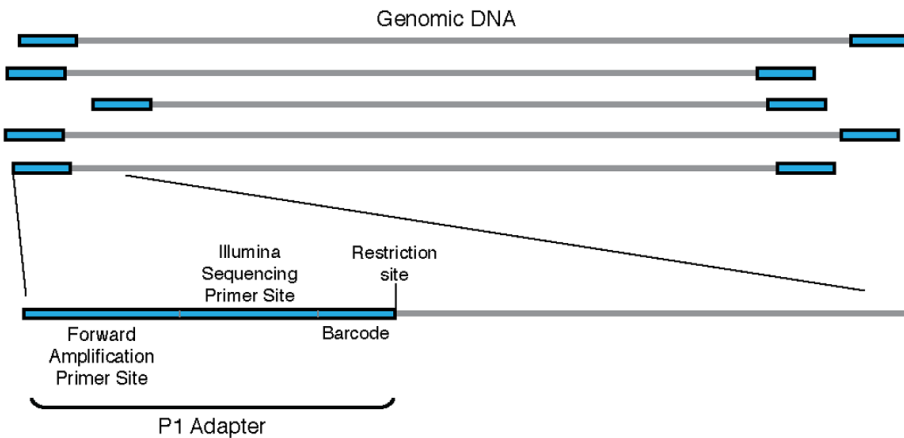


Base content



RAD sequencing

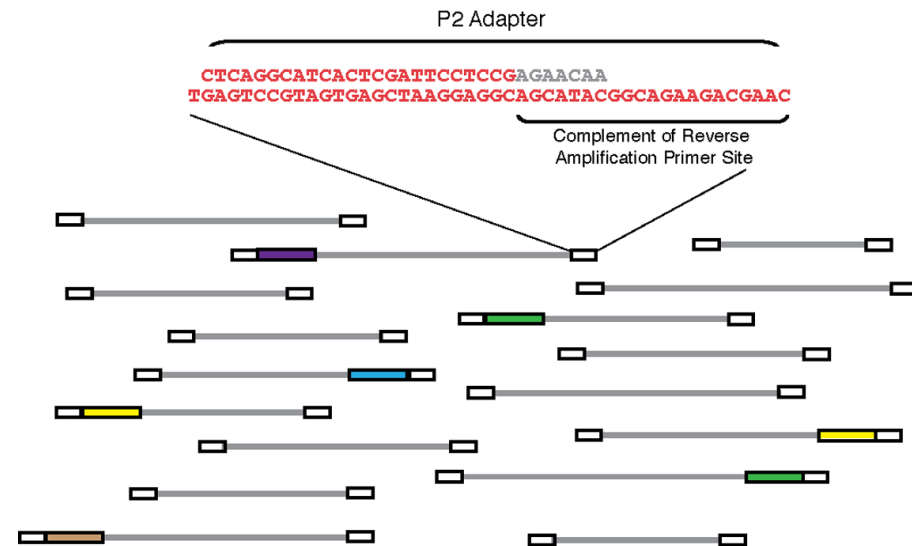
A *Ligate P1 Adapter to digested genomic DNA*



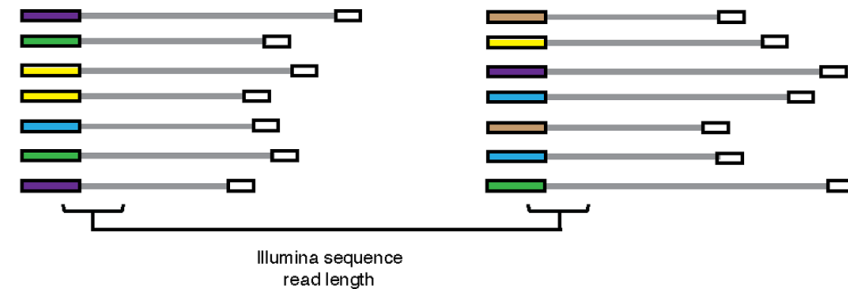
B *Pool barcoded samples and shear*



C *Ligate P2 Adapter to sheared fragments*



D *Selectively amplify RAD tags*



Baird *et al.* Plos One 2008

Pre-processing

process_radtags processes raw Illumina input:

- De-multiplexing and removing barcodes
- Sliding-window quality filter
- Truncate reads

(alternatively – FASTX)

```
$ process_radtags -f in.fastq -o out_dir -i fastq -b  
barcodes -e sbfI -c -q -r -D
```

`f` – path to the input file if processing single-end sequences.

`i` – input file type, either 'bustard' for the Illumina BUSTARD output files, or 'fastq' (default 'fastq').

`o` – path to output the processed files.

`b` – a list of barcodes for this run.

`e` – specify the restriction enzyme to look for (either 'sbfI', 'pstI', 'ecoRI', or 'sgrAI').

`c` – clean data, remove any read with an uncalled base.

`q` – discard reads with low quality scores.

`r` – rescue barcodes and RAD-Tags.

`D` – capture discarded reads to a file.

Mapping

Use **bowtie** to map RAD tags to the genome

- Allow a few mismatches
 - *limits the number of SNPs per RAD tag!*
- Do not allow indels
 - *polymorphic indels will not be found!*

Stacks pipeline

Use RAD tags to the genome

- **ustacks** (unique stacks): Builds loci de novo and detects haplotypes in one individual
- **cstacks** (catalog stacks): Merges loci from multiple individuals to form a catalog
- **sstacks** (search stacks): Matches loci from an individual against a catalog
- **pstacks** (population stacks): Takes cleaned reads aligned to a reference genome, builds stacks based on the genomic locations of the reads, and detects haplotypes in one individual

