

Tackling Genomic Data:

The Advantages and Applications of Bioinformatics Pipelines

Brad Reil

What's the plan...?

- What is a Bioinformatics Pipeline?
- Why do we need bioinformatics pipelines?
- How do we generate data for a pipeline?
- How does a pipeline operate (STACKS example)?
- What types of analyses do pipelines facilitate?
- What technical know-how is required?

What is a Bioinformatics Pipeline?

Bioinformatics – a Definition¹

(Molecular) **bio – informatics**: bioinformatics is conceptualising biology in terms of molecules (in the sense of Physical chemistry) and applying “**informatics techniques**” (derived from disciplines such as applied maths, computer science and statistics) to **understand** and **organise** the **information** associated with these molecules, on a **large scale**. In short, bioinformatics is a management information system for molecular biology and has many **practical applications**.

¹As submitted to the Oxford English Dictionary.



Pipeline

A set of **data processing elements** connected in **series**, where the output of one element is the input of the next one. - Wikipedia

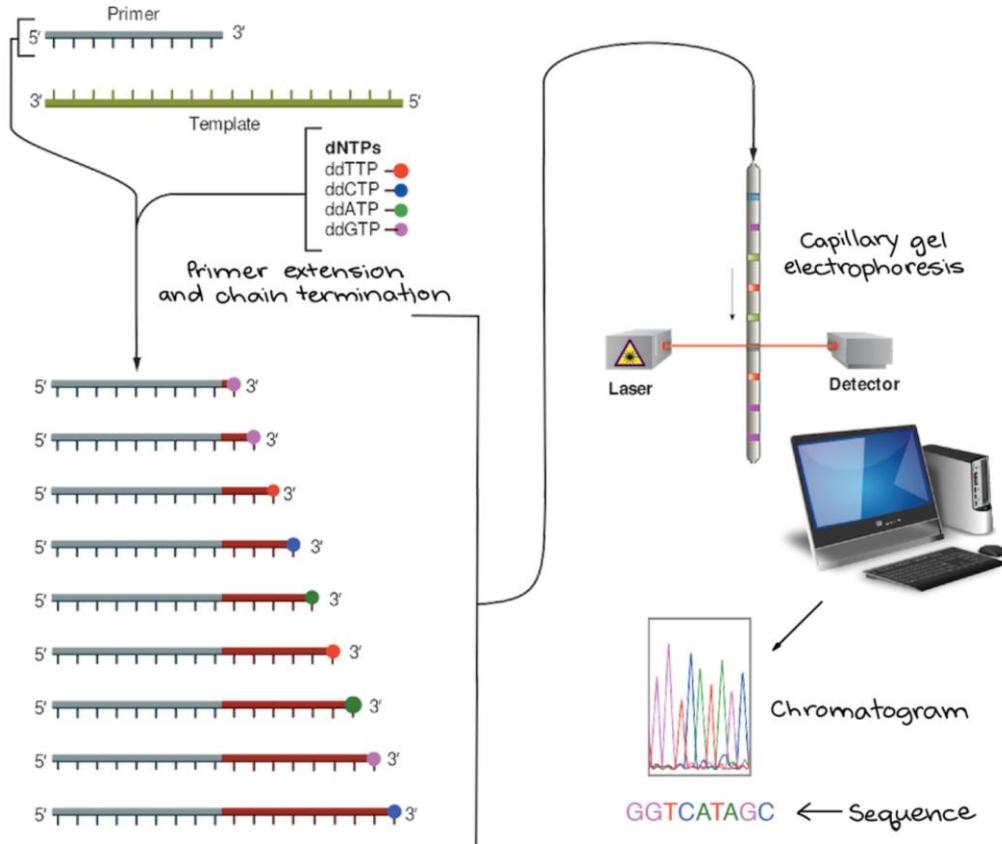


Bainformatics Pipeline

A set of **data processing elements** connected in **series**, used to understand and organize the information associated with large scale data sets derived from biological molecules.

Why do we need these pipelines?

Well, we started with...



Sanger Sequencing

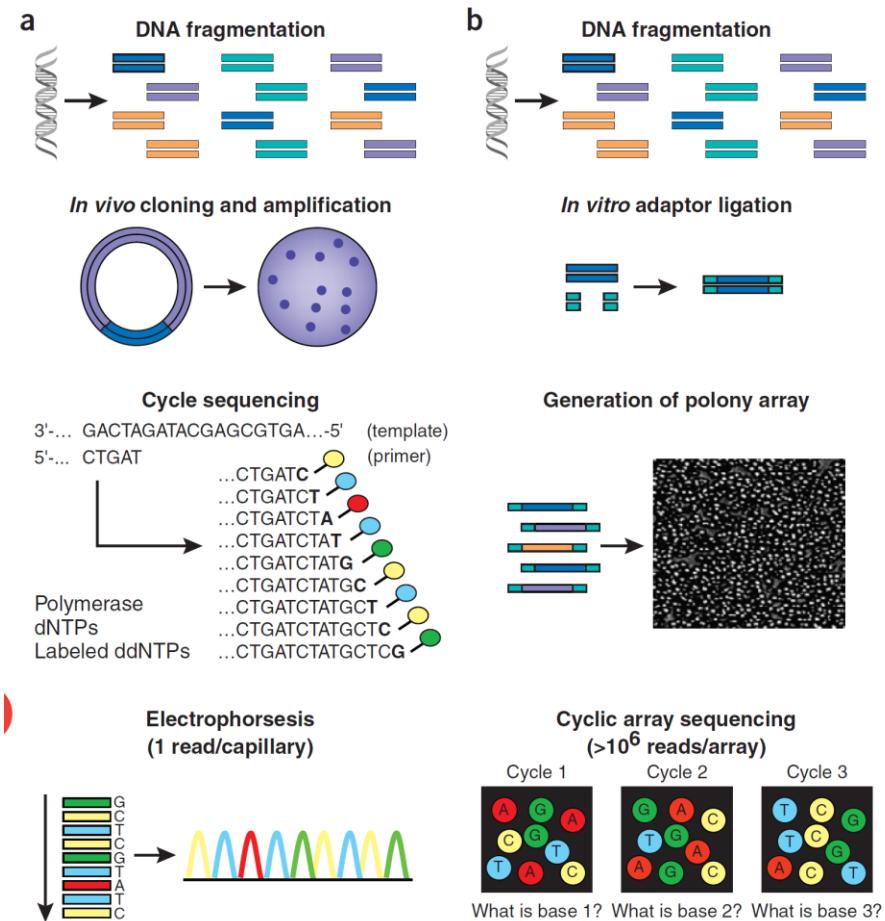
Image modified from "Sanger sequencing," by Estevezj (CC BY-SA 3.0). The modified image is licensed under a (CC BY-SA 3.0) license.

Why do we need these pipelines?

Enter Next-Generation Sequencing Techniques...

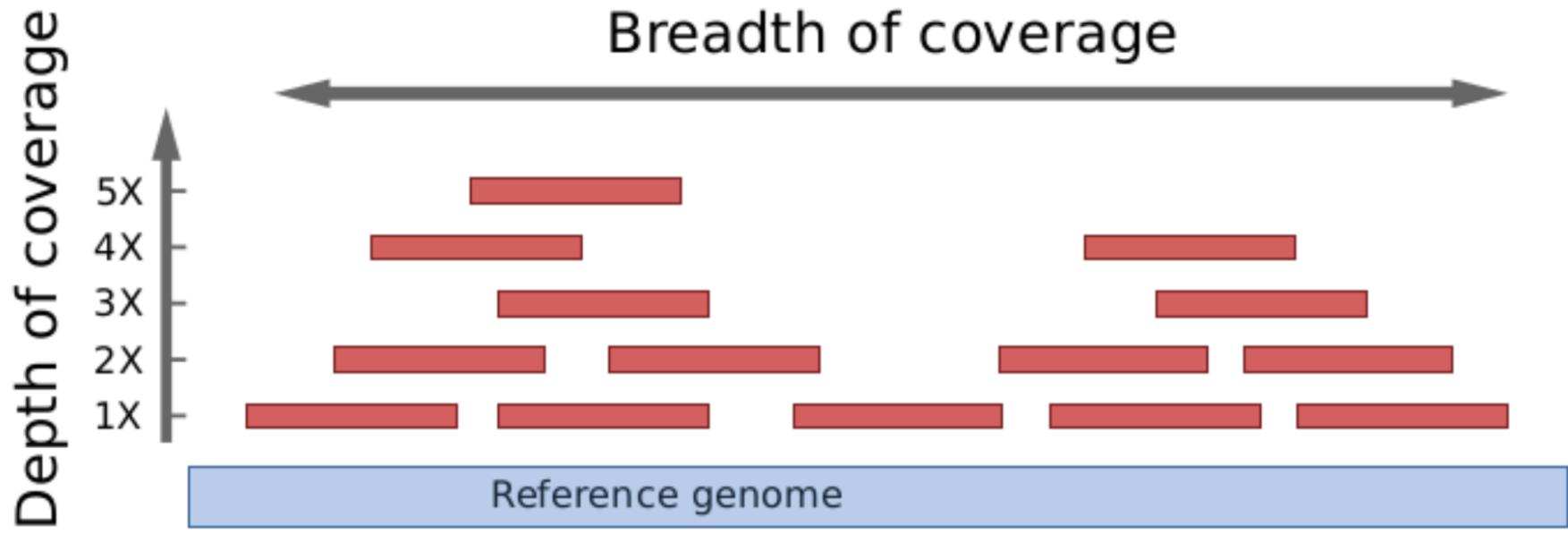
Next-Generation Sequencing

Term used to describe a collection of molecular sequencing techniques capable of producing reads for large arrays of different DNA fragments simultaneously.



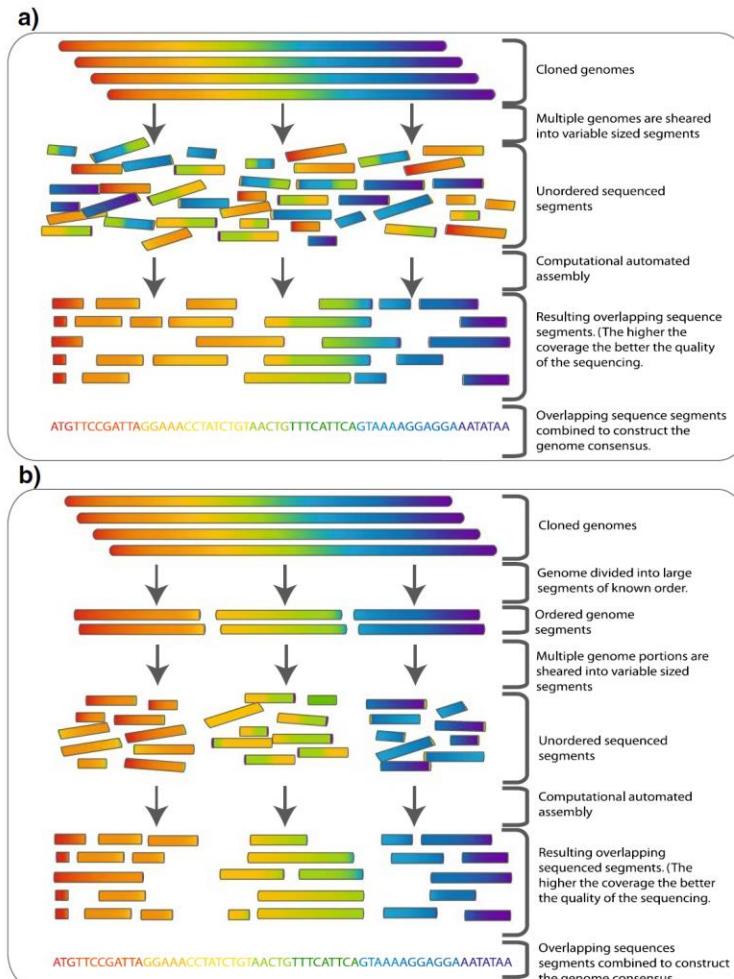
Sampling the Genome - DNA Fragmentation

Some quick terminology...



Sampling the Genome - DNA Fragmentation

So, how do we **fragment** and sample the genome...



Shotgun Sequencing

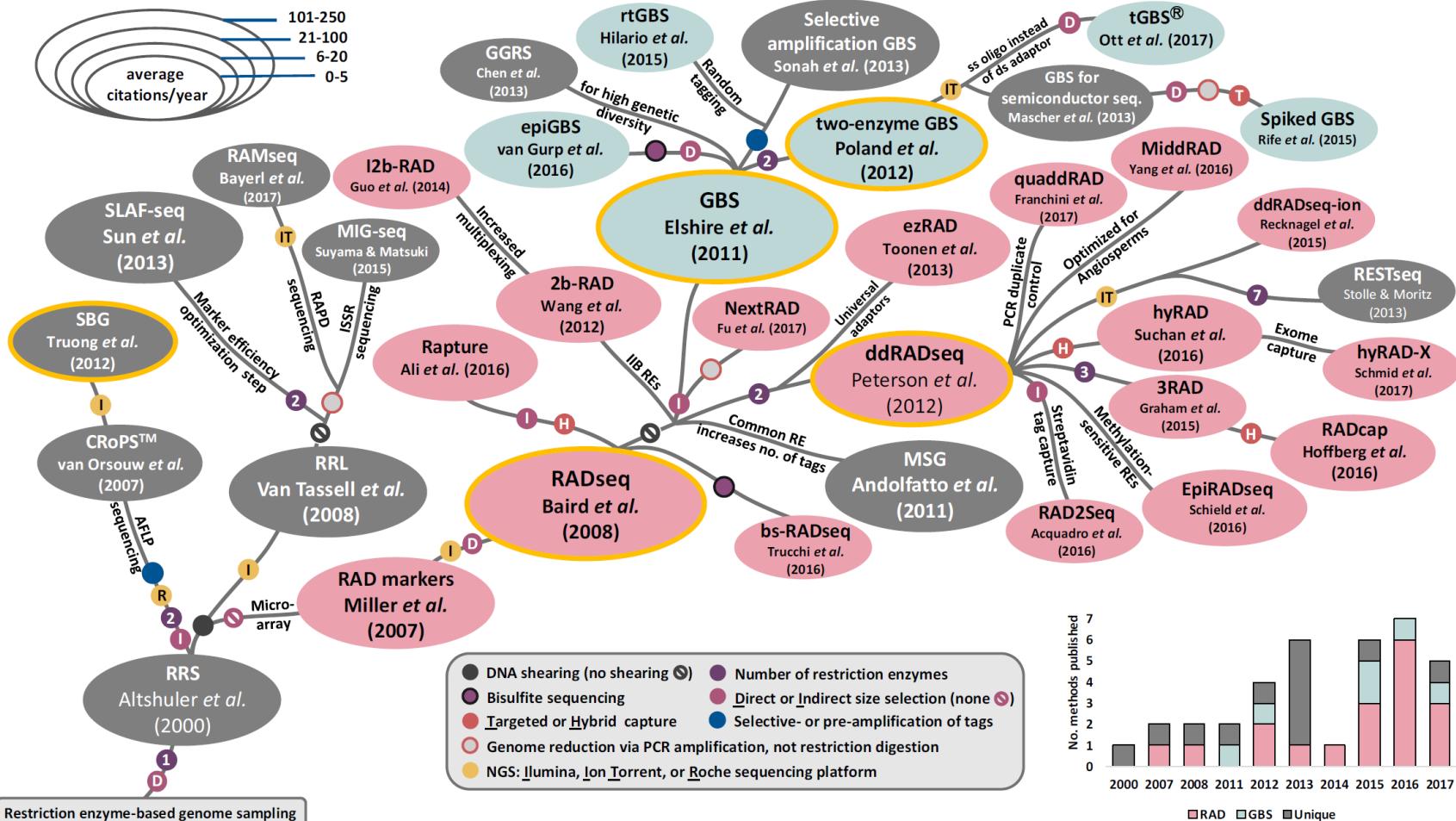
DNA Randomly Sheered By:

- Vortexing
- Glass Beads
- Sonication
- Radiation

The fragments in a specific size range are then isolated and sequenced

DNA Fragmentation by Restriction Enzymes

So, how do we **fragment** and sample the genome...

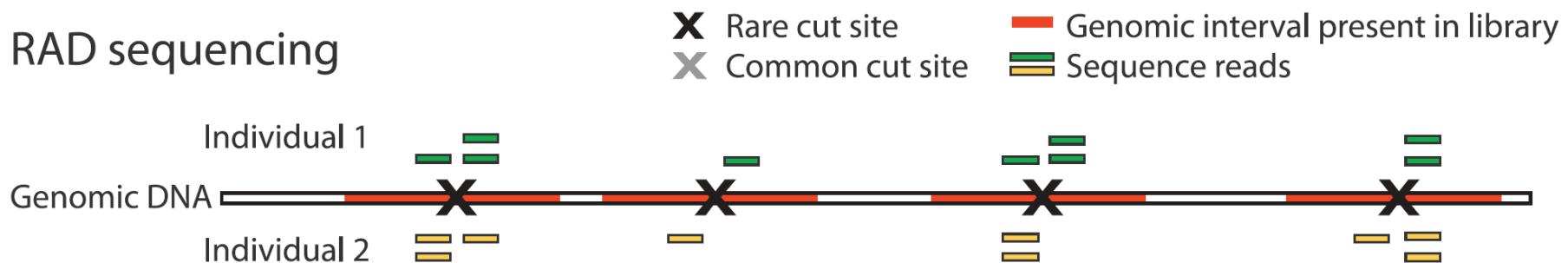


DNA Fragmentation by Restriction Enzymes

So, how do we **fragment** and sample the genome...

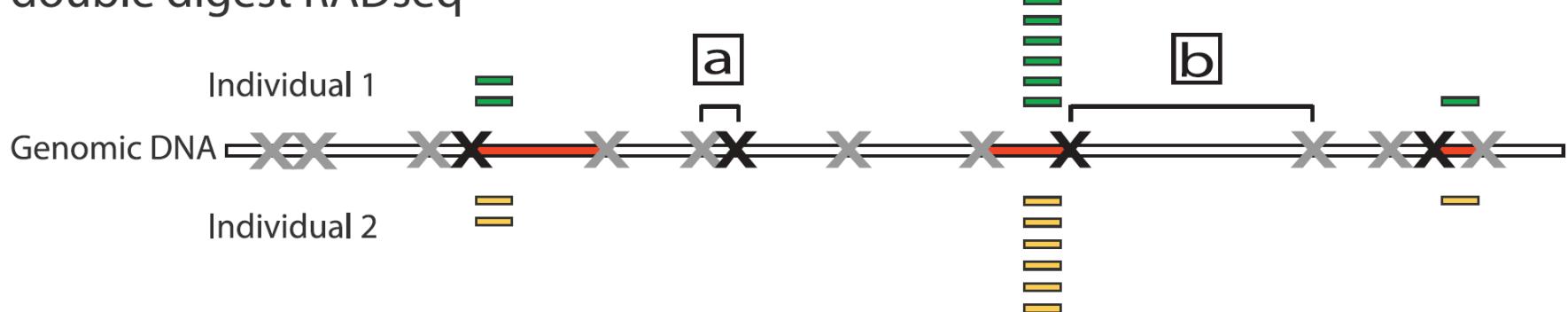
A

RAD sequencing



B

double digest RADseq

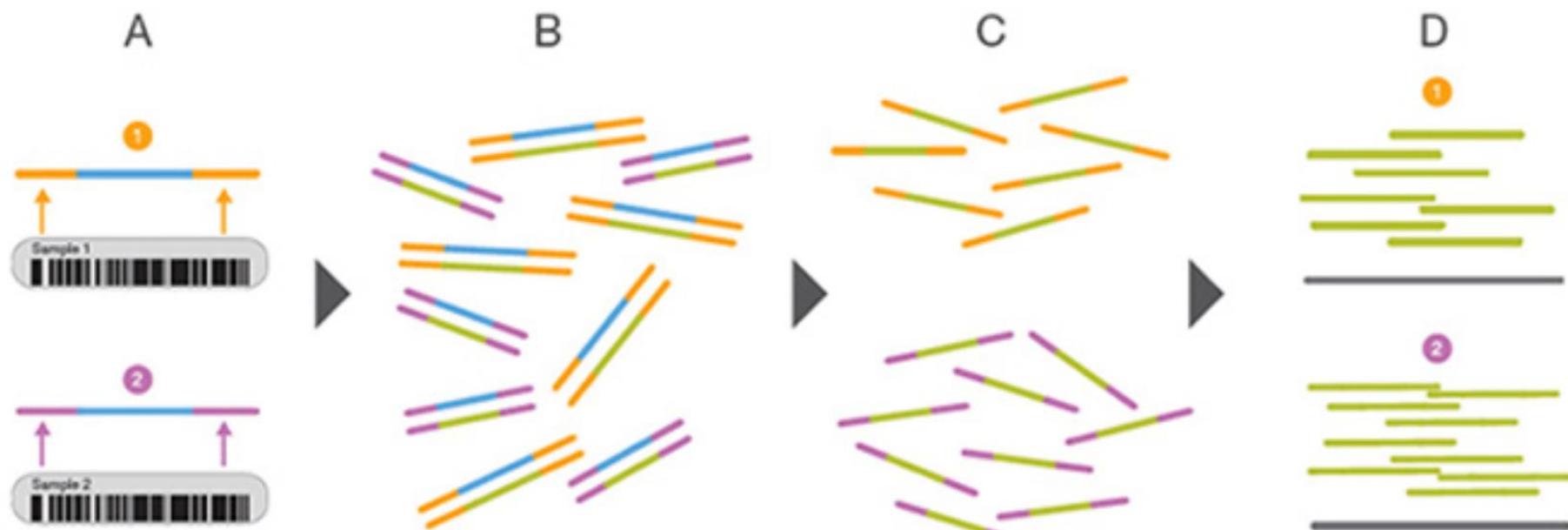


DNA Multiplexing (Combinatorial Indexing)

And how do we **organize** large amounts of DNA...

DNA Multiplexing

The process of reducing repetitive sequencing steps by the mixing together of different DNA samples. This is achieved by tagging the samples isolated from a particular sample with a unique barcoding sequence via PCR.

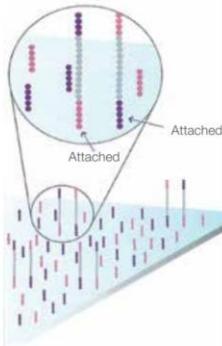


Next-Gen Sequencing Platforms

Now the genome is sampled, organized, and ready to **sequence**

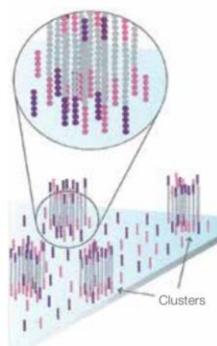
Illumina Sequencing

Figure 6: Denature the Double-Standed Molecules



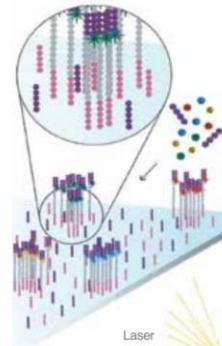
Denaturation leaves single-stranded templates anchored to the substrate.

Figure 7: Complete Amplification



Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

Figure 8: Determine First Base



The first sequencing cycle begins by adding four labeled reversible terminators, primers, and DNA polymerase.

Figure 9: Image First Base



After laser excitation, the emitted fluorescence from each cluster is captured and the first base is identified.

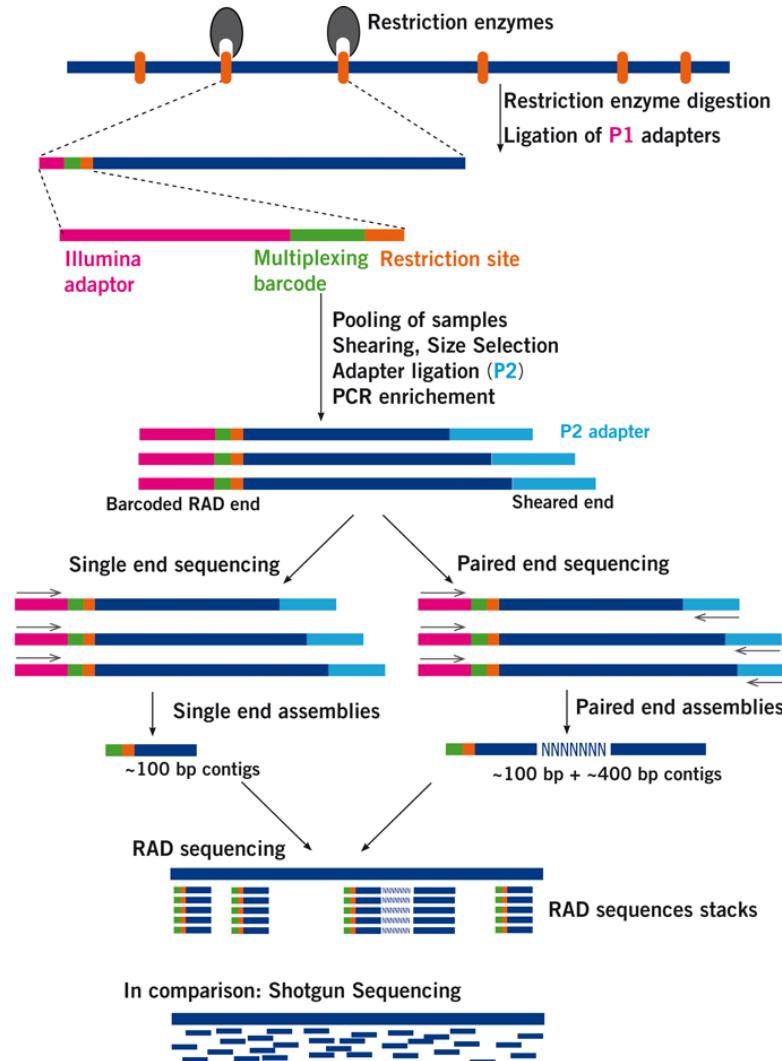
Next-Gen Sequencing Platforms

Now the genome is sampled, organized, and ready to **sequence**

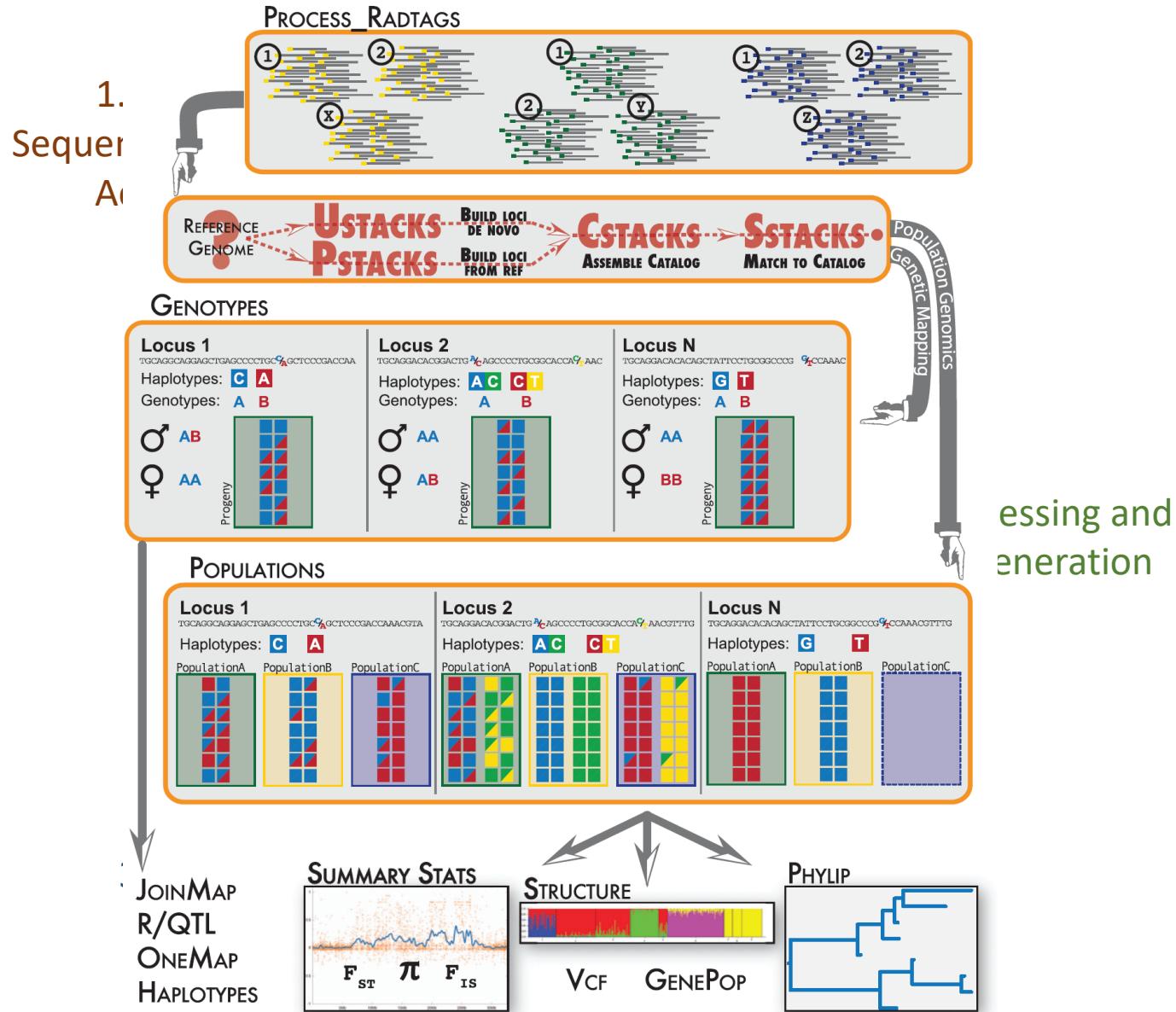
Table 2 Comparison of sequencing instruments, sorted by cost/Mb, with expected performance by mid 2011

| Instrument | Run time ^a | Millions of reads/run | Bases/read ^b | Yield Mb/run | Reagent cost/run ^c | Reagent cost/Mb | Minimum unit cost (% run) ^d |
|--|-----------------------|-----------------------|-------------------------|--------------|-------------------------------|-----------------|--|
| 3730xl (capillary) | 2 h | 0.000096 | 650 | 0.06 | \$96 | \$1500 | \$6 (1%) |
| Ion Torrent – ‘314’chip | 2 h | 0.10 | 100 | >10 | \$500 | <\$50 | ~\$750 (100%) |
| 454 GS Jr. Titanium | 10 h | 0.10 | 400 | 50 | \$1100 | \$22 | \$1500 (100%) |
| Starlight* | † | ~0.01 | >1000 | † | † | † | † |
| PacBio RS | 0.5–2 h | 0.01 | 860–1100 | 5–10 | \$110–900 | \$11–180 | † |
| 454 FLX Titanium | 10 h | 1 | 400 | 500 | \$6200 | \$12.4 | \$2000 (10%) |
| 454 FLX+ ^e | 18–20 h | 1 | 700 | 900 | \$6200 | \$7 | \$2000 (10%) |
| Ion Torrent – ‘316’chip* | 2 h | 1 | >100 | >100 | \$750 | <\$7.5 | ~\$1000 (100%) |
| Helicos ^f | N/A | 800 | 35 | 28 000 | N/A | NA | \$1100 (2%) |
| Ion Torrent – ‘318’chip* | 2 h | 4–8 | >100 | >1000 | ~\$925 | ~\$0.93 | ~\$1200 (100%) |
| Illumina MiSeq* | 26 h | 3.4 | 150 + 150 | 1020 | \$750 | \$0.74 | ~\$1000 (100%) |
| Illumina iScanSQ | 8 days | 250 | 100 + 100 | 50 000 | \$10 220 | \$0.20 | \$3000 (14%) |
| Illumina GAIIX | 14 days | 320 | 150 + 150 | 96 000 | \$11 524 | \$0.12 | \$3200 (14%) |
| SOLiD – 4 | 12 days | >840 ^g | 50 + 35 | 71 400 | \$8128 | <\$0.11 | \$2500 (12%) |
| Illumina HiSeq 1000 | 8 days | 500 | 100 + 100 | 100 000 | \$10 220 | \$0.10 | \$3000 (12%) |
| Illumina HiSeq 2000 | 8 days | 1000 | 100 + 100 | 200 000 | \$20 120 ^h | \$0.10 | \$3000 (6%) |
| SOLiD – 5500 (PI)* | 8 days | >700 ^g | 75 + 35 | 77 000 | \$6101 | <\$0.08 | \$2000 (12%) |
| SOLiD – 5500xl (4hq)* | 8 days | >1410 ^g | 75 + 35 | 155 100 | \$10 503 ^h | <\$0.07 | \$2000 (12%) |
| Illumina HiSeq 2000 – v3 ^{i*} | 10 days | ≤3000 | 100 + 100 | ≤600 000 | \$23 470 ^h | ≥\$0.04 | ~\$3500 (6%) |

Alright so where are we...



And now we are ready for our pipeline...



The STACKS pipeline breakdown

Stacks

Stacks is a software pipeline for building loci from short-read sequences, such as those generated on the Illumina platform. Stacks was developed to work with restriction enzyme-based data, such as RAD-seq, for the purpose of building genetic maps and conducting population genomics and phylogeography.

Raw Reads

process_radtags
process_shortreads
clone_filter
kmer_filter

Core

ustacks
pstacks
cstacks
sstacks
genotypes
populations
rxstacks

3

Execution control

denovo_map.pl
ref_map.pl
load_radtags.pl

2

Utilities

index_radtags.pl
export_sql.pl
sort_read_pairs.pl
exec_velvet.pl

The STACKS pipeline breakdown

Some more quick terminology...

Stack

A series of identical reads from a given sample that have been anchored to a homologous location in the genome (reference genome or genome assembled de novo) and arranged thus to form a 'stack'; a series of stacks is used to form a loci.

Loci

The location of a gene or mutation along the genome

Allele

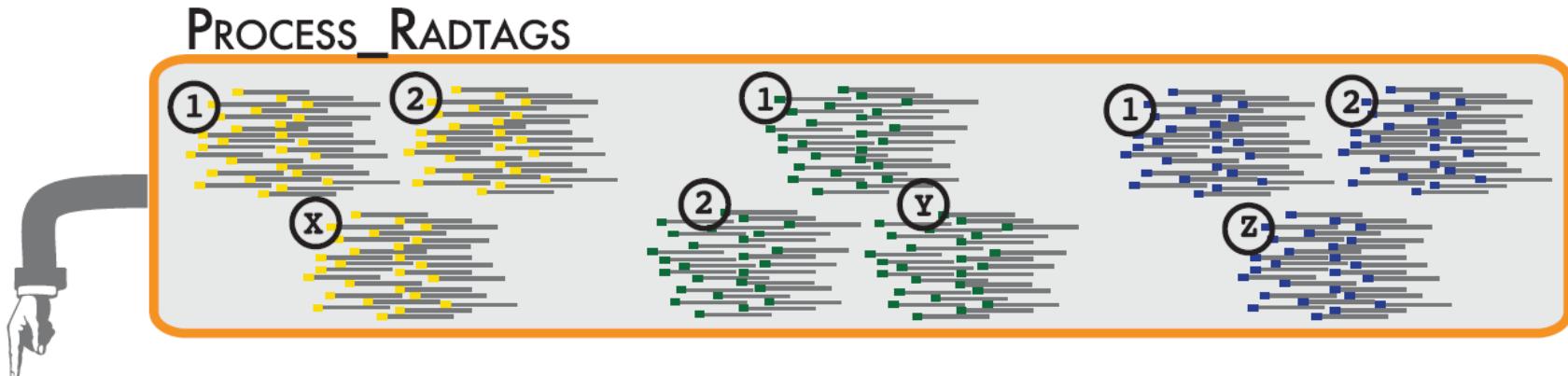
The particular form of a gene or mutation found at a given loci along the genome

STACKS pipeline 1: process_radtags

process_radtags

A function to de-multiplex and clean raw sequence ready

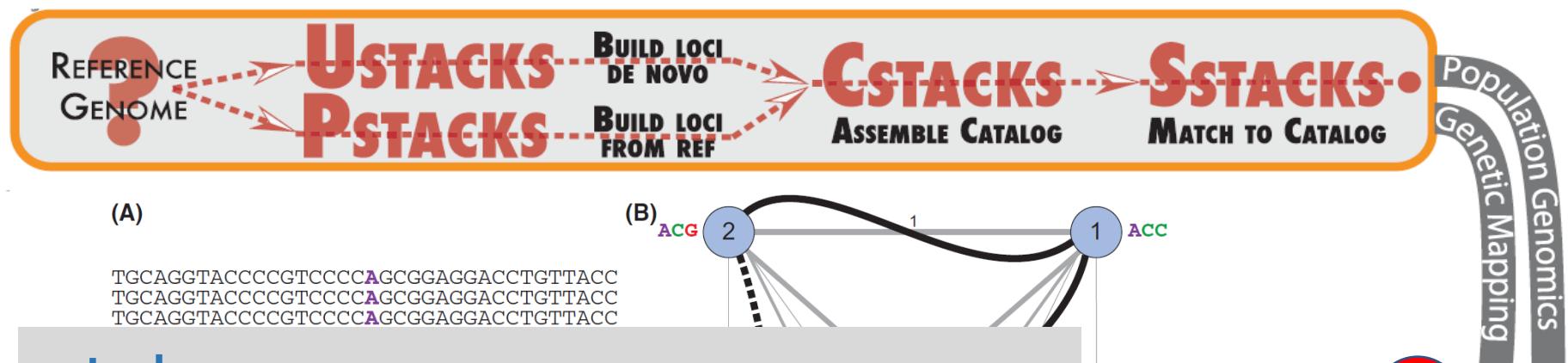
```
process_radtags -f XXXXX.fastq.gz -i gzip -o ./demultiplexed_fastqs/ -b  
XXXXX.txt -e nlaIII -c -q -r
```



STACKS pipeline 2: denovo_map.pl

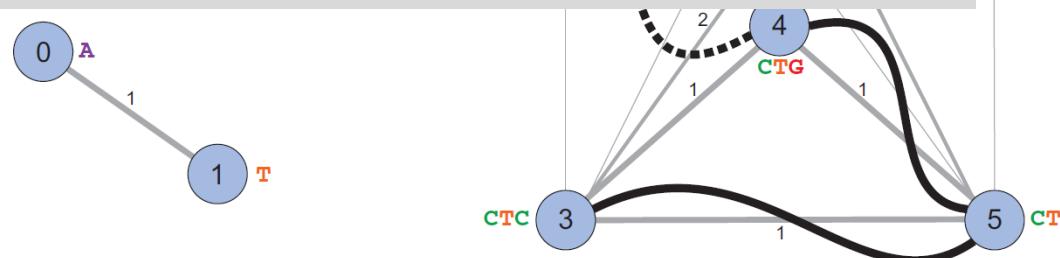
Denovo_map.pl

A wrapper function that runs all three main programs (ustacks, cstacks, and sstacks) comprised within the *de novo* loci assembly version of the stacks pipeline



ustacks

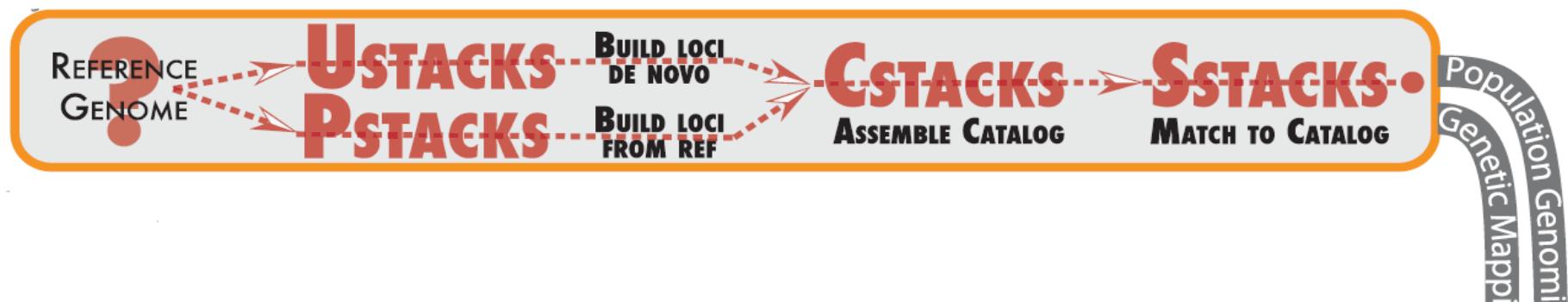
Data from each individual are grouped into loci and polymorphic sites are identified



STACKS pipeline 2: denovo_map.pl

Denovo_map.pl

A wrapper function that runs all three main programs (ustacks, cstacks, and sstacks) comprised within the *de novo* loci assembly version of the stacks pipeline



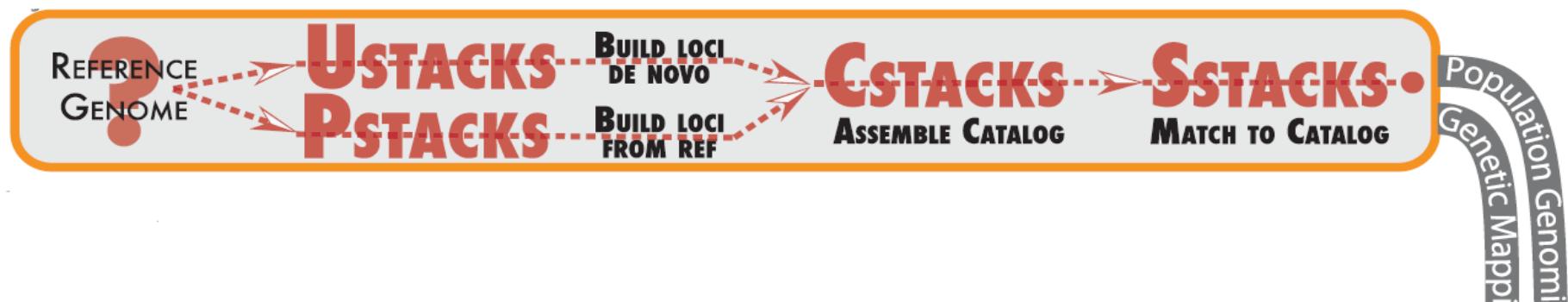
cstacks

Loci are grouped together across individuals and a catalogue is written

STACKS pipeline 2: denovo_map.pl

Denovo_map.pl

A wrapper function that runs all three main programs (ustacks, cstacks, and sstacks) comprised within the *de novo* loci assembly version of the stacks pipeline



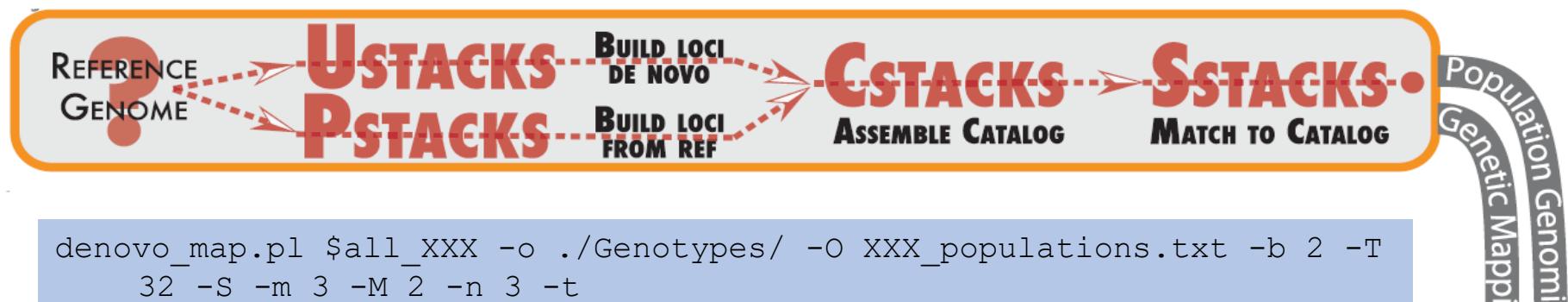
sstacks

Loci from each individual are matched against the catalogue to determine the allelic state at each locus in each individual

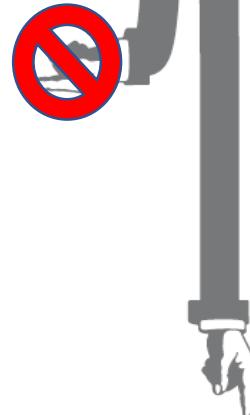
STACKS pipeline 2: denovo_map.pl

Denovo_map.pl

A wrapper function that runs all three main programs (ustacks, cstacks, and sstacks) comprised within the *de novo* loci assembly version of the stacks pipeline



```
denovo_map.pl $all_XXX -o ./Genotypes/ -O XXX_populations.txt -b 2 -T  
32 -S -m 3 -M 2 -n 3 -t
```



STACKS pipeline 3: denovo_map.pl

populations

tabulates the state of loci within and among populations, calculates population genetics statistics and exports to a number of additional, useful formats

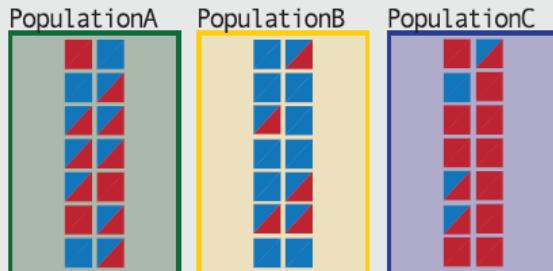
```
populations -b 2 -t 32 -P ./Genotypes/ -M ./XXX_populations.txt --renz nlaIII  
-m 10 -p 1 -r 0.5 -k --fstats -f p_value --genomic --ordered_export --  
write_random_snp --fasta --vcf --plink --genepop --structure --phylip
```

POPULATIONS

Locus 1

TGCAGGCAGGAGCTGAGCCCCCTGC**C**A GCTCCCGACCAAACGTA

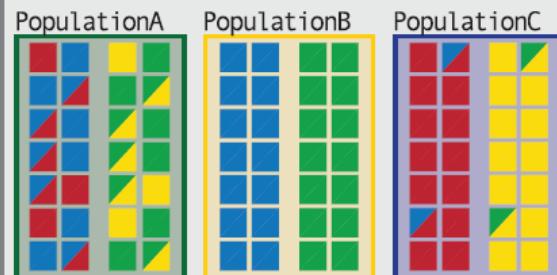
Haplotypes: **C** **A**



Locus 2

TGCAGGACACGGACTG**A****C** AGCCCCCTGGGGCACCA**C****T** AACGTTTG

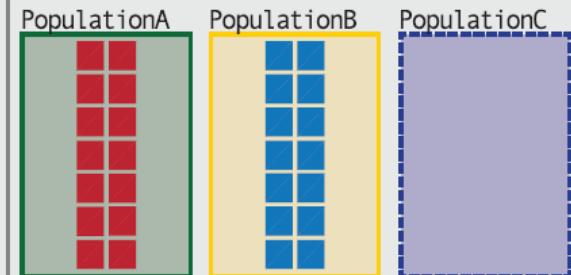
Haplotypes: **AC** **CT**



Locus N

TGCAGGACACACAGCTATTCTCTGCAGGCCCG**G****T** CCAAACGTTTG

Haplotypes: **G** **T**



Ok so where are we now...

So we sent our raw sequence data through the pipeline and have...

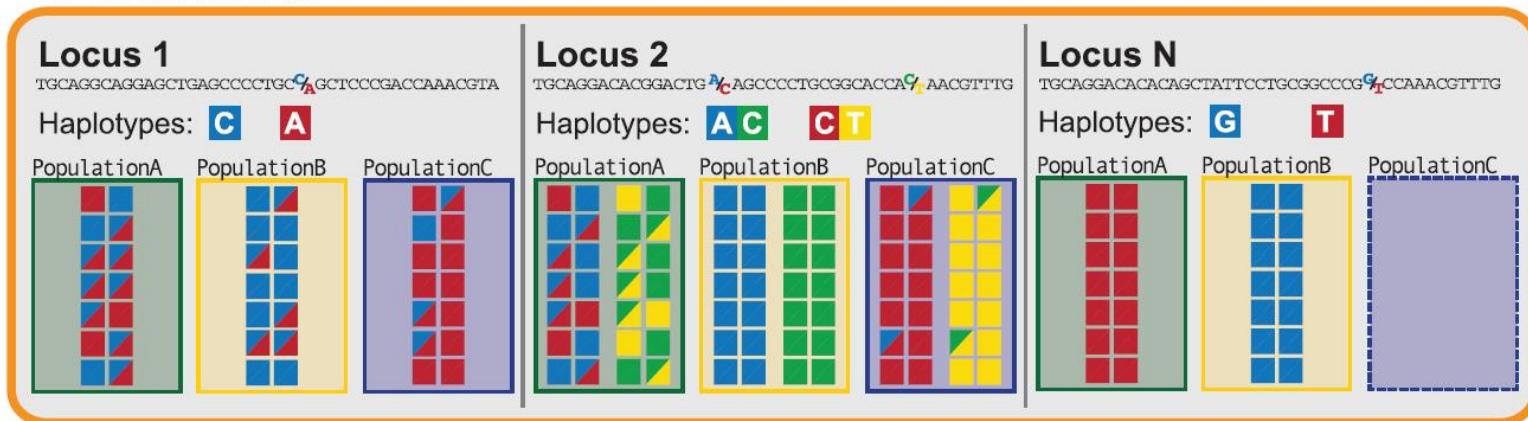
Data Matrices With...

- Series of SNP's From Across the Genome
- Unlinked
- Primarily Nuclear
- Coding and Noncoding
- Partially Filtered
- In VCF or other desired output format

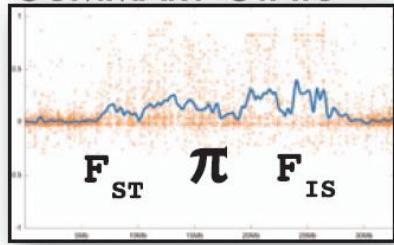
Ok so where are we now...

And now we can use this data to perform various analyses

POPULATIONS



SUMMARY STATS



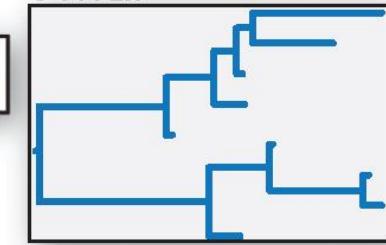
STRUCTURE



VCF

GENEPOP

PHYLIP



NOTE: Methods are best for exploring recent evolutionary history (no more than ~10 mya)

Before that a few more tools

Some tools for additional filtering and formatting

Welcome to VCFtools

VCFtools is a program package designed for working with VCF files, such as those generated by the 1000 Genomes Project. The aim of VCFtools is to provide easily accessible methods for working with complex genetic variation data in the form of VCF files.

This toolset can be used to perform the following operations on VCF files:

- Filter out specific variants
- Compare files
- Summarize variants
- Convert to different file types
- Validate and merge files
- Create intersections and subsets of variants



PGDSpider version 2.1.1.2 (July 2017)

An automated data conversion tool for connecting population genetics and genomics programs

Ok now, on to the analyses...

And by way of example, here is a recent paper

ORIGINAL ARTICLE

WILEY [MOLECULAR ECOLOGY](#)

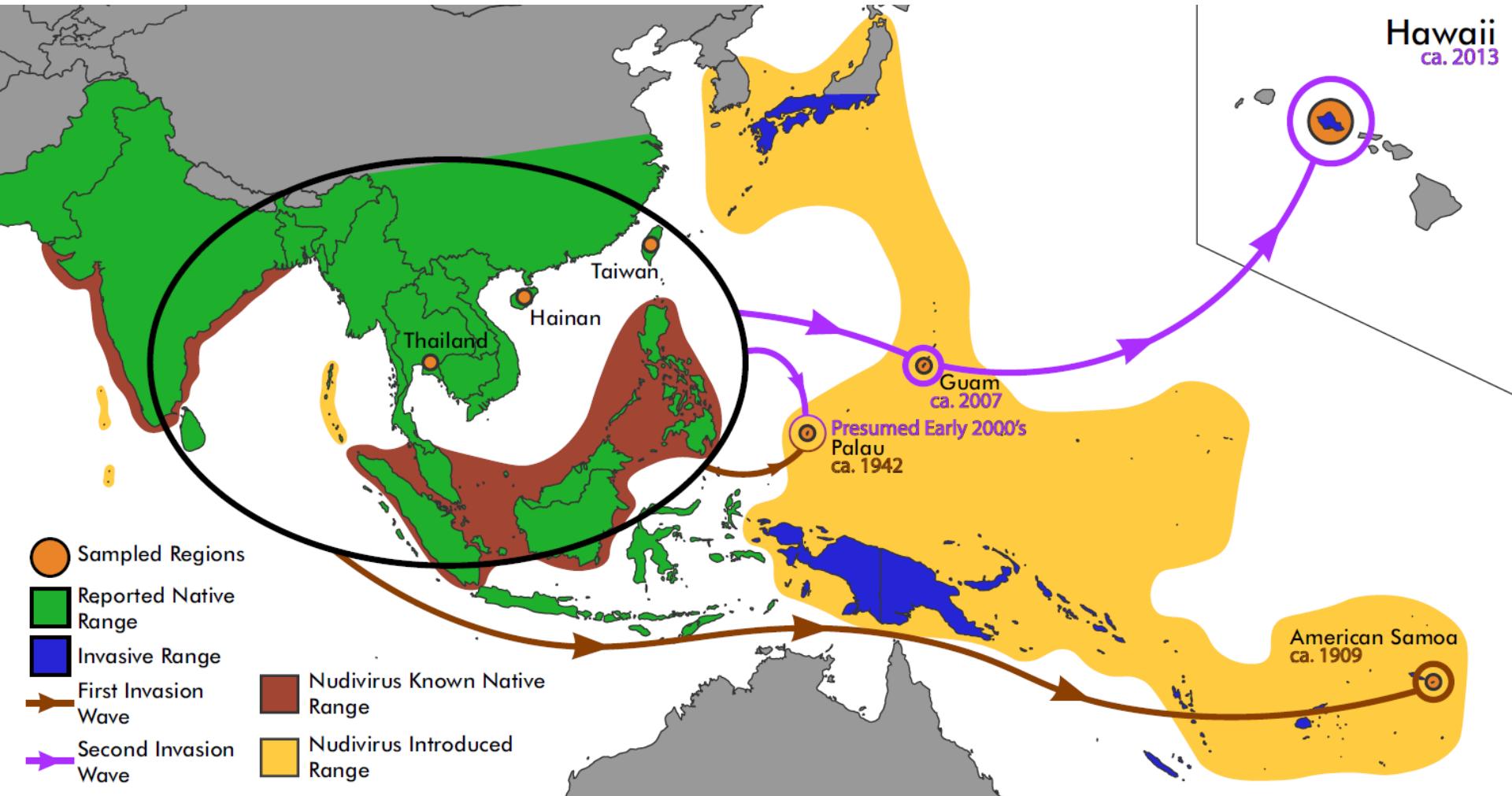
Transpacific coalescent pathways of coconut rhinoceros beetle biotypes: Resistance to biological control catalyses resurgence of an old pest

Jonathan Bradley Reil¹  | Camiel Doorenweerd¹  | Michael San Jose¹  |
Sheina B. Sim^{1,2}  | Scott M. Geib²  | Daniel Rubinoff¹



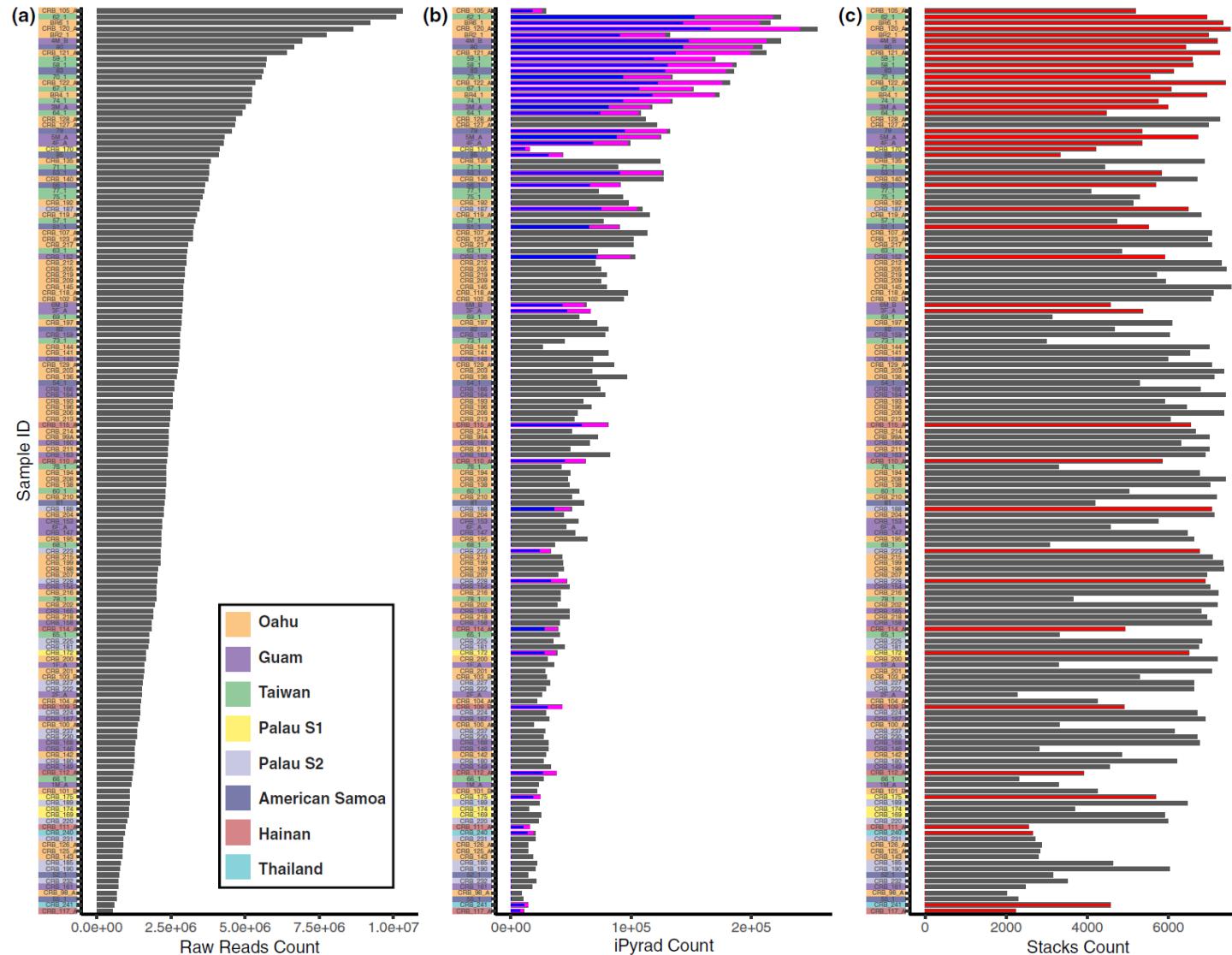
A working example: *Oryctes rhinoceros*

So what is the story...



A working example: *Oryctes rhinoceros*

Data obtained from 172 individuals (151 post quality control) using ddRAD and Illumina Seq.



A working example: *Dryctes rhinoceros*

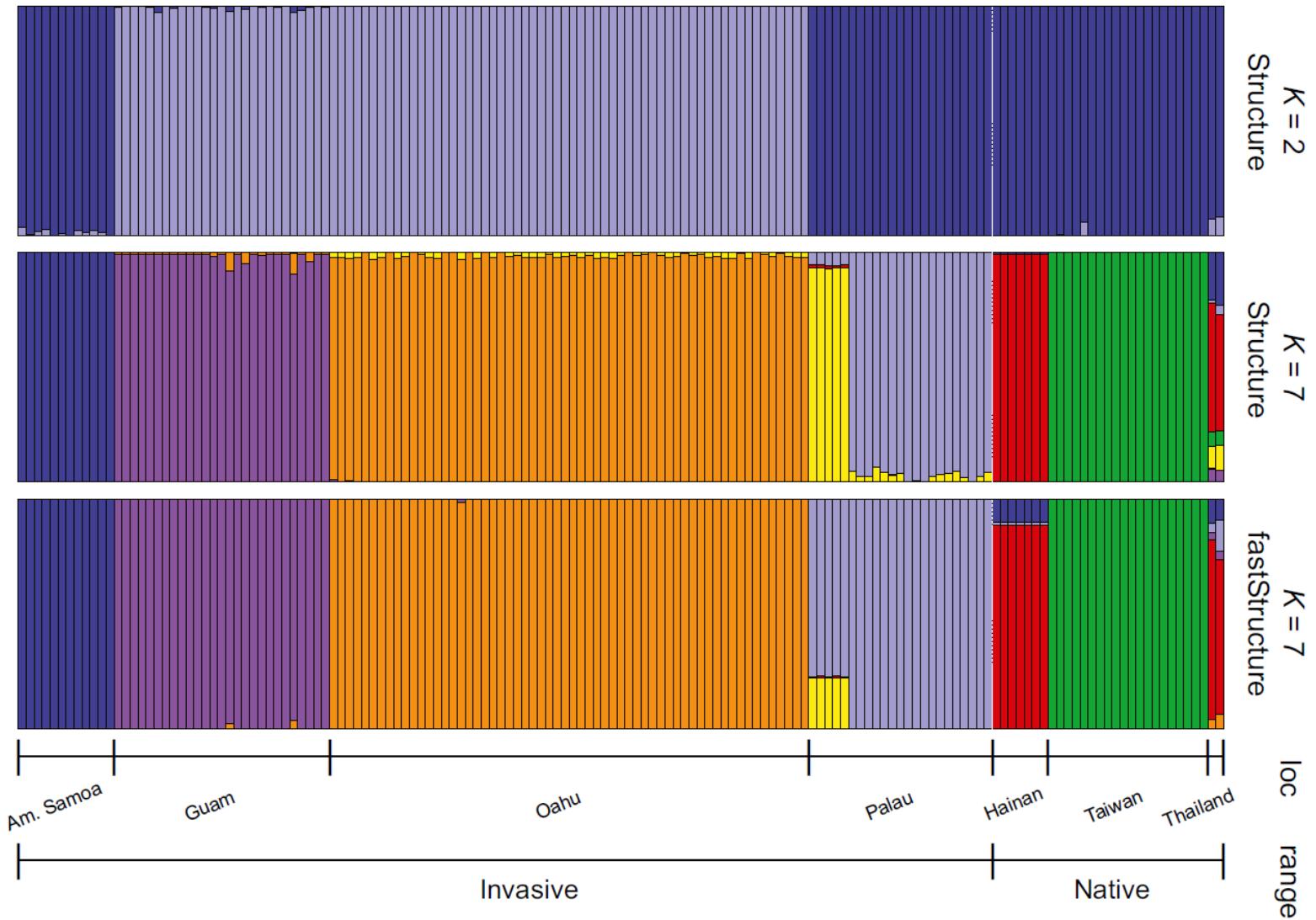
Began our analysis with some population genetics using STRUCTURE



STRUCTURE analyses differences in the distribution of genetic variants amongst populations with a Bayesian iterative algorithm by placing samples into groups whose members share similar patterns of variation. *STRUCTURE* both identifies populations from the data and assigns individuals to that population representing the best fit for the variation patterns found. Typically *STRUCTURE* is the first step in examining population structures that emerge from the sample set to provide a preamble to further genetic analysis or to infer the origins of individuals with unknown population characteristics, especially when population admixture has occurred. As *STRUCTURE* uses the core Bayesian principle of comparing likelihoods, prior information about study samples can be supplied to further shape the analysis.

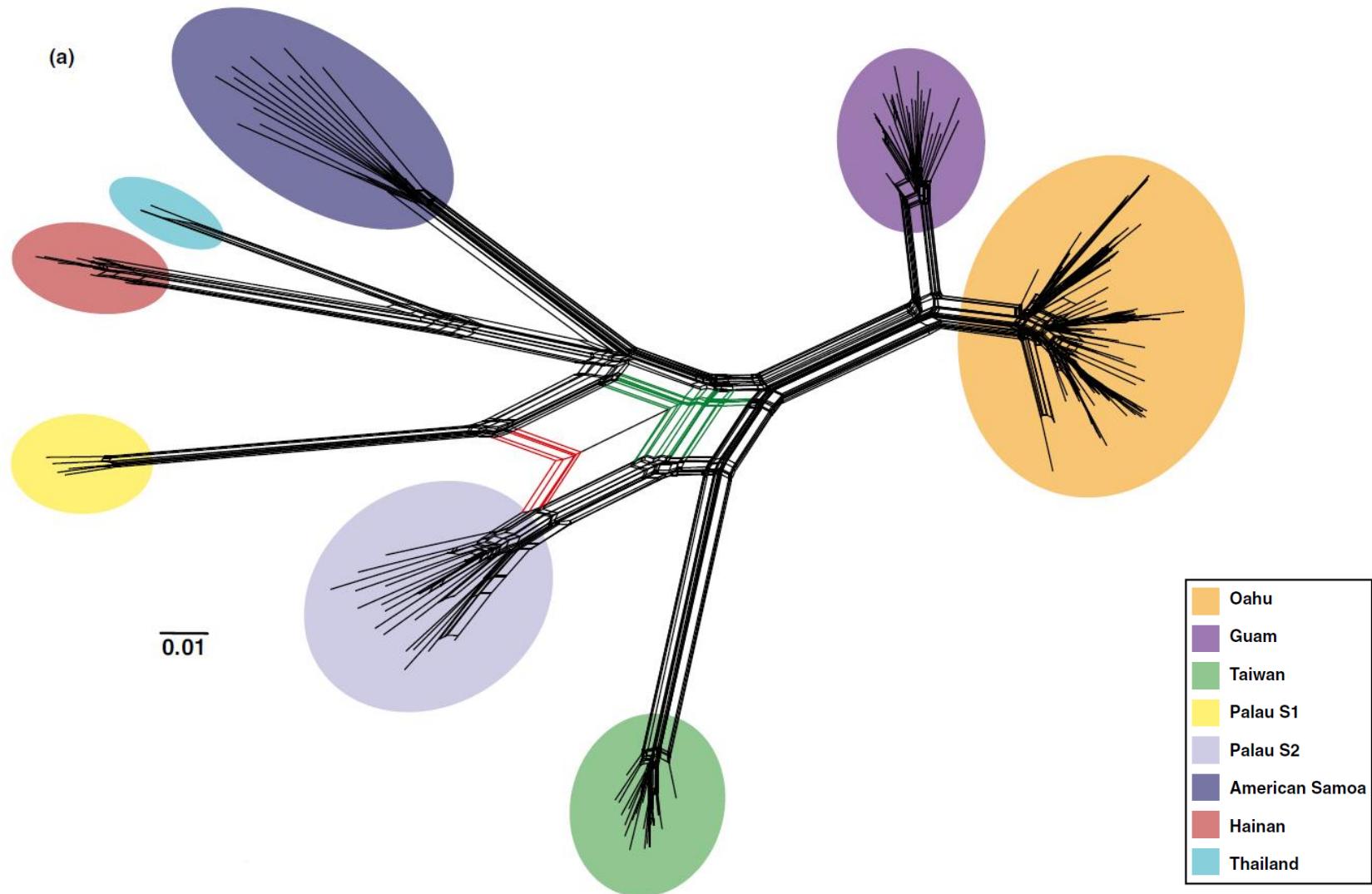
A working example: *Dryctes rhinoceros*

Began our analysis with some population genetics using STRUCTURE



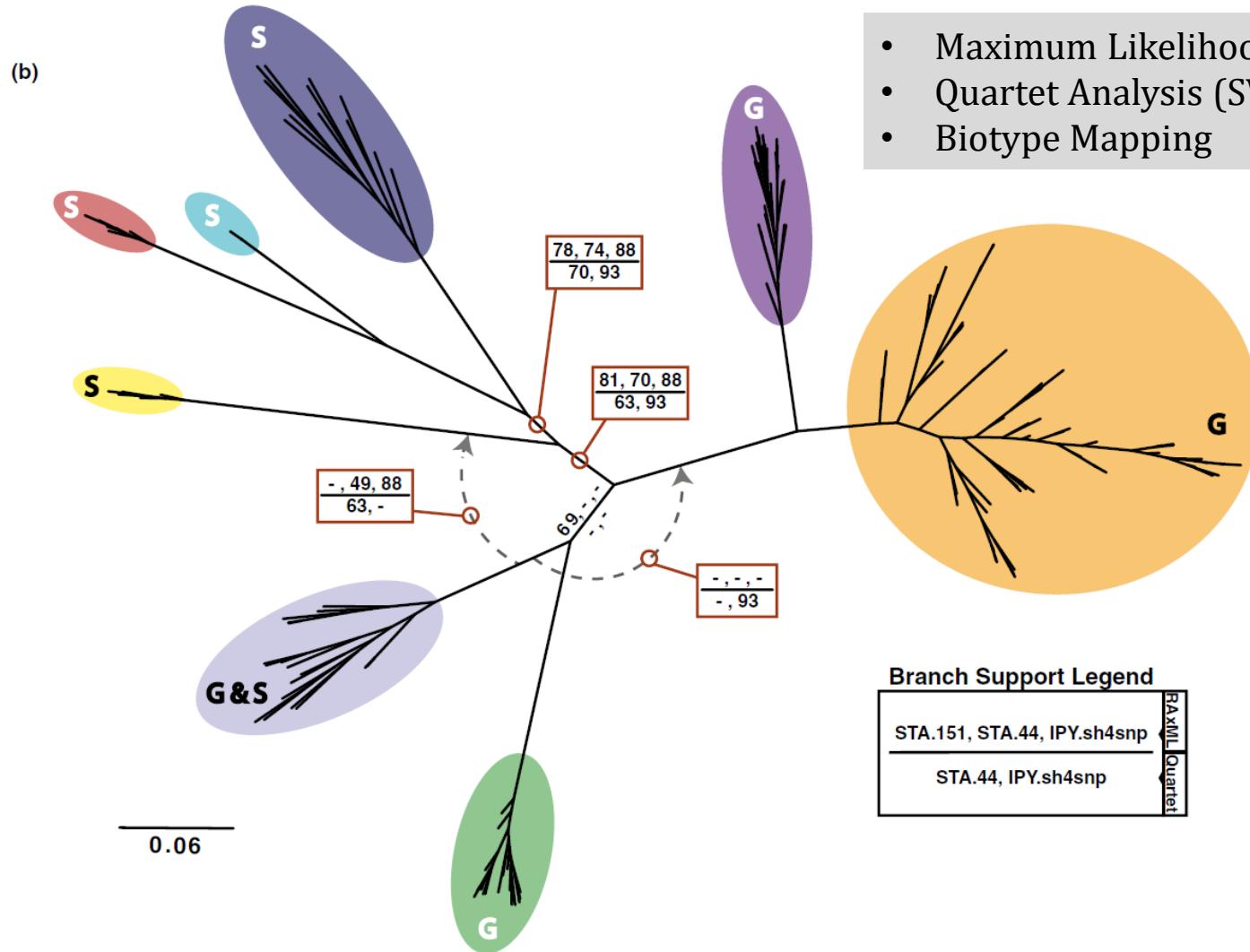
A working example: *Oryctes rhinoceros*

We then proceeded to begin working out Phylogenetic relationships



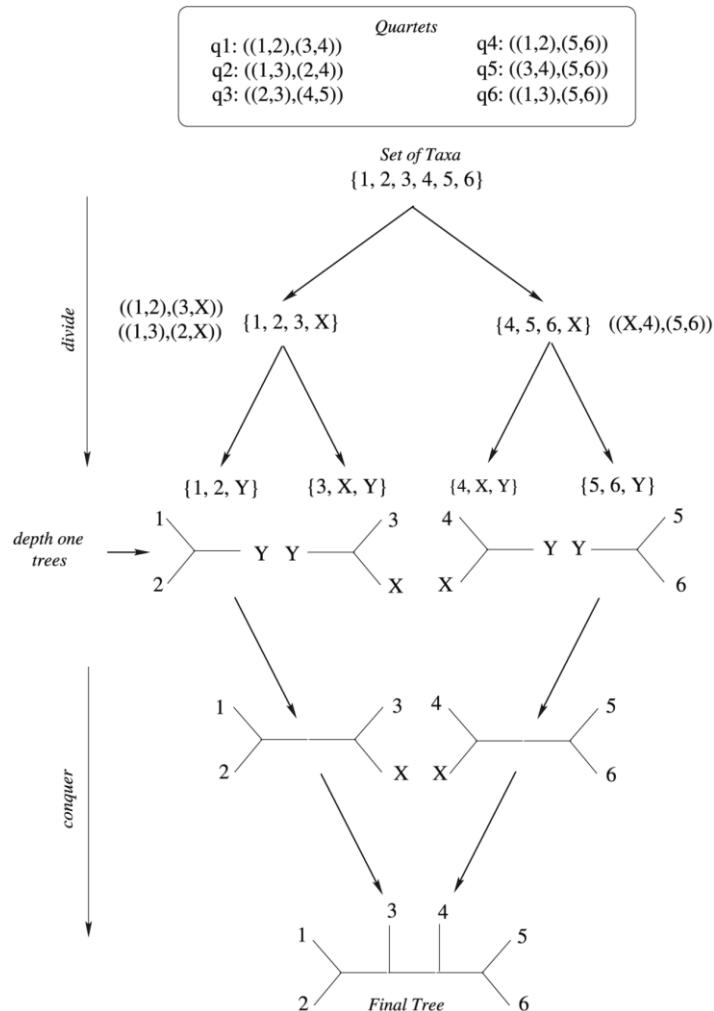
A working example: *Dryctes rhinoceros*

Ok something interesting is going on, lets dig deeper



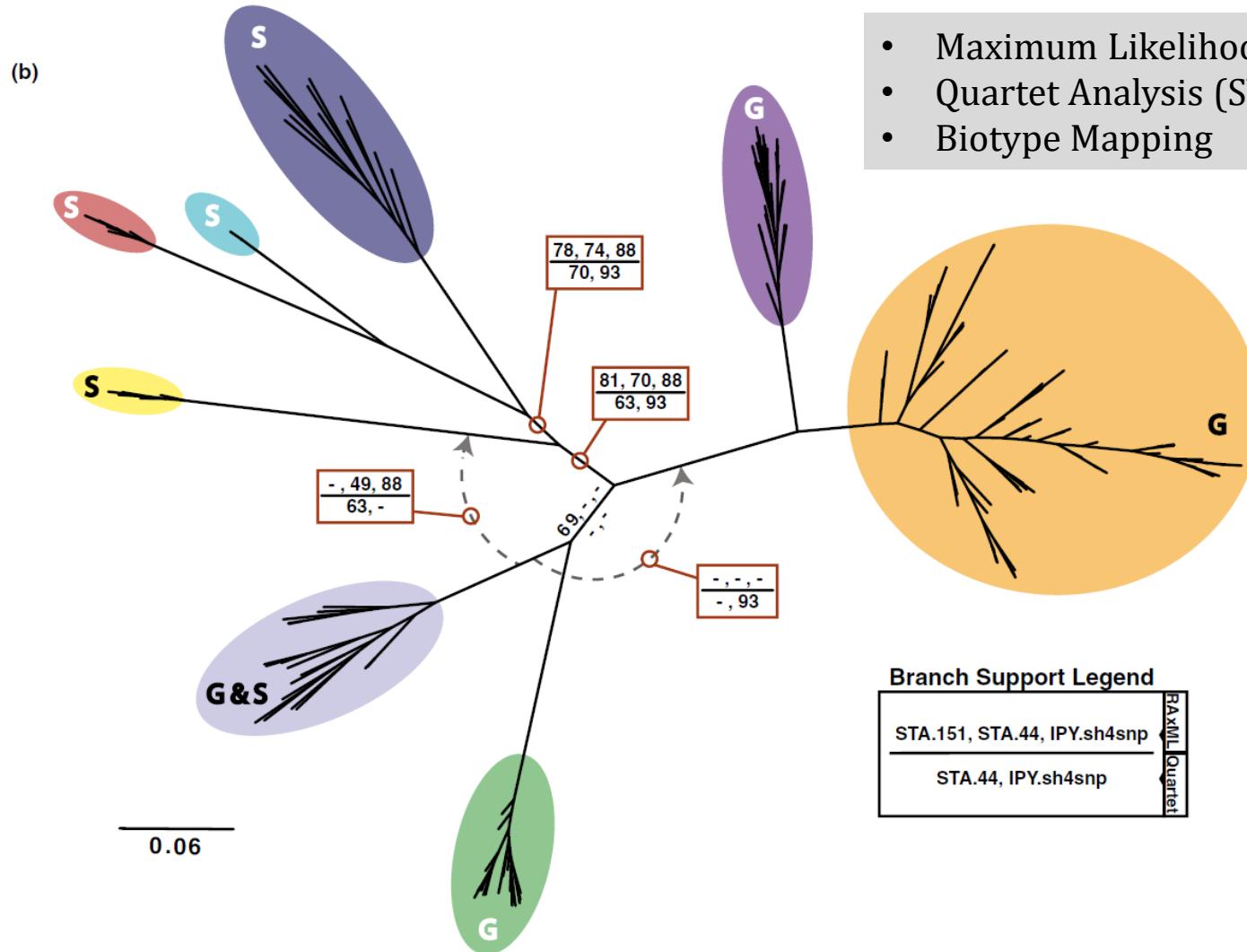
A working example: *Dryctes rhinoceros*

A bit on SVDQuartets and quartet analysis...



A working example: *Dryctes rhinoceros*

Ok something interesting is going on, lets dig deeper

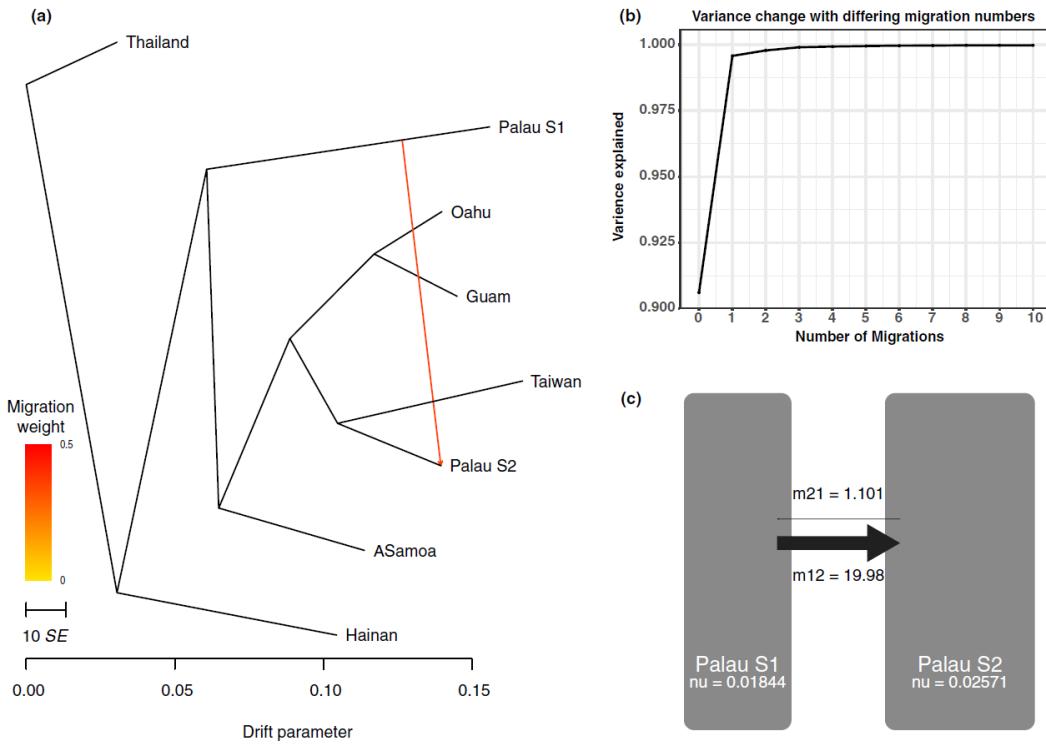


A working example: *Oryctes rhinoceros*

And so we concluded..

Conclusions

- Second Invasion catalyzed by the emergency of biotype G
- Both biotypes present in Palau, leading to gene tree variation
- Biotypes exchanging genetic material
- Beetles in Hawaii, likely originated in Guam
- Biotype G origin unknown, but Taiwan is suspect



So, what does it really take to do this

Who can reasonably take advantage of these techniques?

```
###STACKS PIPELINE SCRIPT (flags varied for different data matrices)#####
##STEP 1: Process Radtags (De-multiplexing): Run this for each set of
##fastq.gz and combine the outputs into a single directory (Ex.
##demultiplexed_all)
process_radtags -f XXXXX.fastq.gz -i gzfastq -o ./demultiplexed_fastqs/ -b
XXXXX.txt -e nlaIII -c -q --rr

##STEP 2: Pass this portion with denovo_map.pl: It parses your
##individuals.txt (file with sample ID's) and files in your combined directory
##into a list usable by denovo_map.pl
#all_XXX=""
#for x in `cat XXX_individuals.txt`;
#do
#all_XXX+="`./demultiplexed_all/$x.fq.gz "
#done

## Denovo assembly of loci, catalog generation, SNP flagging
denovo_map.pl $all_XXX -o ./Genotypes/ -O XXX_populations.txt -b 2 -T 32 -S -
m 3 -M 2 -n 3 -t

##STEP 3: Optional things like SNP/Loci filtering, basic pop. gen. stats, etc
populations -b 2 -t 32 -P ./Genotypes/ -M ./XXX_populations.txt --renz nlaIII
-m 10 -p 1 -r 0.5 -k --fstats -f p_value --genomic --ordered_export --
write_random_snp --fasta --vcf --plink --genepop --structure --phylip
```