# Variant Calling
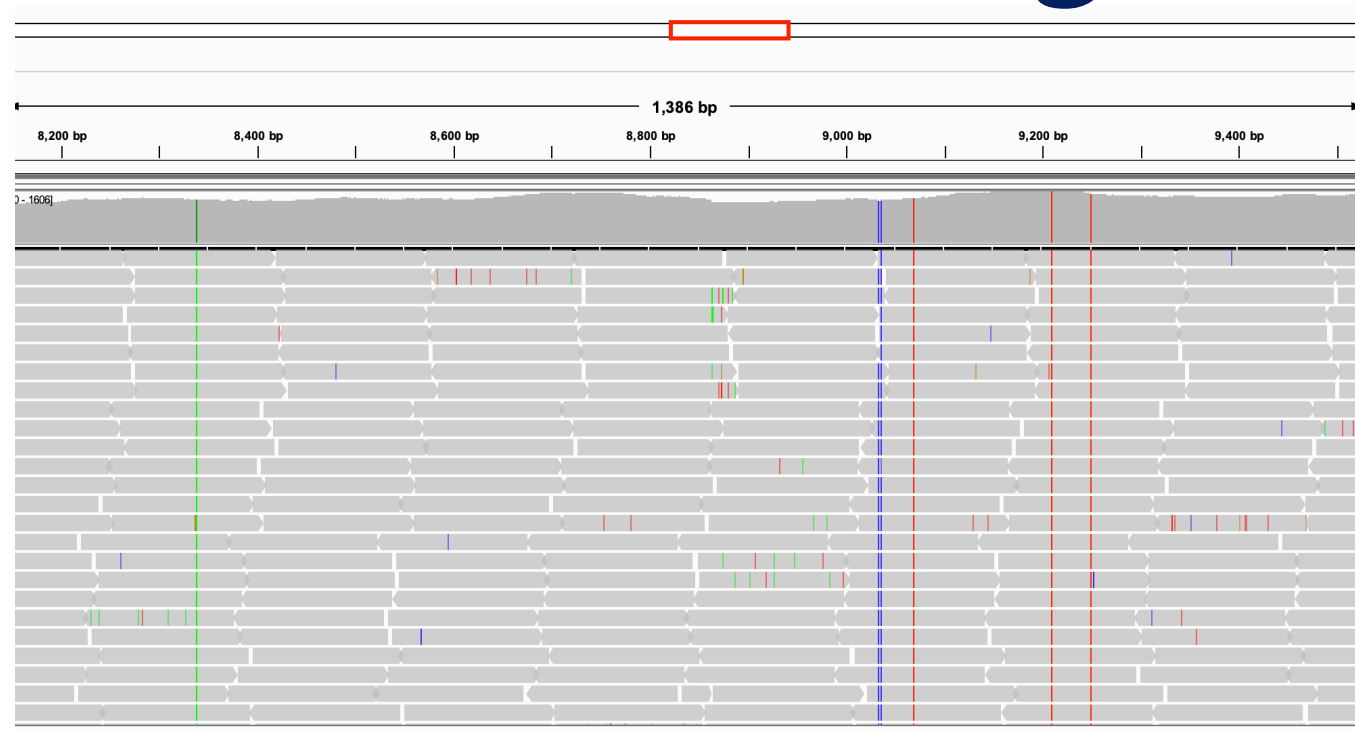


BIOL 435/535: Bioinformatics
April 7th, 2022

# **Mutations –** Inherited changes in nucleotide sequence

| | | | | | | |
|---|---|---|---|---|---|---|
| normal | AUG<br>met | GCC<br>ala | TGC<br>cys | AAA<br>lys | CGC<br>arg | TGG<br>trp |
| silent | AUG<br>met | GCT<br>ala | TGC<br>cys | AAA<br>lys | CGC<br>arg | TGG<br>trp |
| nonsense | AUG<br>met | GCC<br>ala | TGA<br>--- | AAA<br>--- | CGC<br>--- | TGG<br>--- |
| missense | AUG<br>met | GCC<br>ala | GGC<br>arg | AAA<br>lys | CGC<br>arg | TGG<br>trp |
| frameshift<br>(deletion -1) | AUG<br>met | GC-<br>ala | TGC<br>glu | AAA<br>asn | CGC<br>ala | TGG |
| frameshift<br>(insertion +1) | AUG<br>met | GCC<br>ala | C TGC<br>leu | AAA<br>gln | CGC<br>thr | TGG<br>leu |
| insertion +1,<br>deletion -1 | AUG<br>met | GCC<br>ala | C TGC<br>leu | AAA<br>gln | -GC<br>thr | TGG<br>trp |

**synonymous**
**nonsynonymous**

1. How to confidently identify true variants

2. Low-frequency variants

3. Variant annotation

# Class brainstorm:

## What are some important considerations in variant calling?

# How to confidently identify true variants
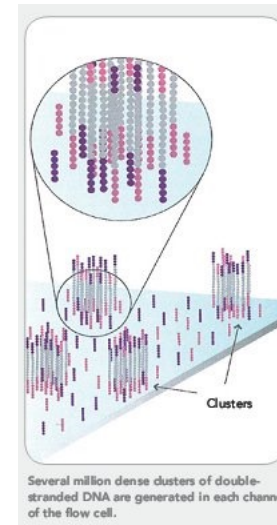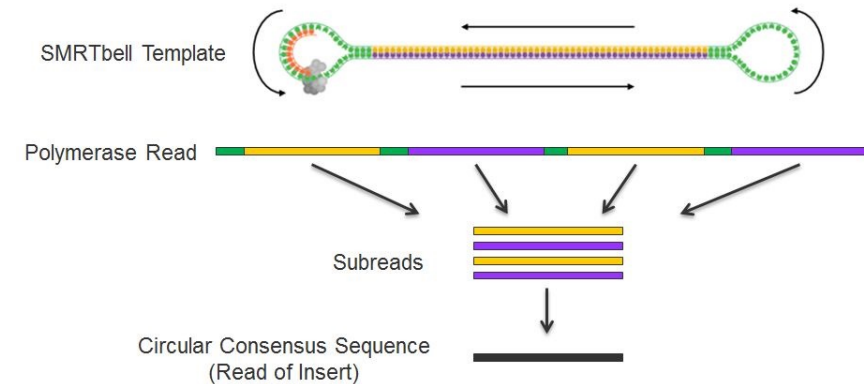
## Primary considerations

- Sequencing method/approach

- Depth of coverage

- Variant type

- Inheritance mode

- Reference quality

- Bioinformatic tool

# How to confidently identify true variants

## Primary considerations

- **Sequencing method/approach**

- Depth of coverage

- Variant type

- Inheritance mode

- Reference quality
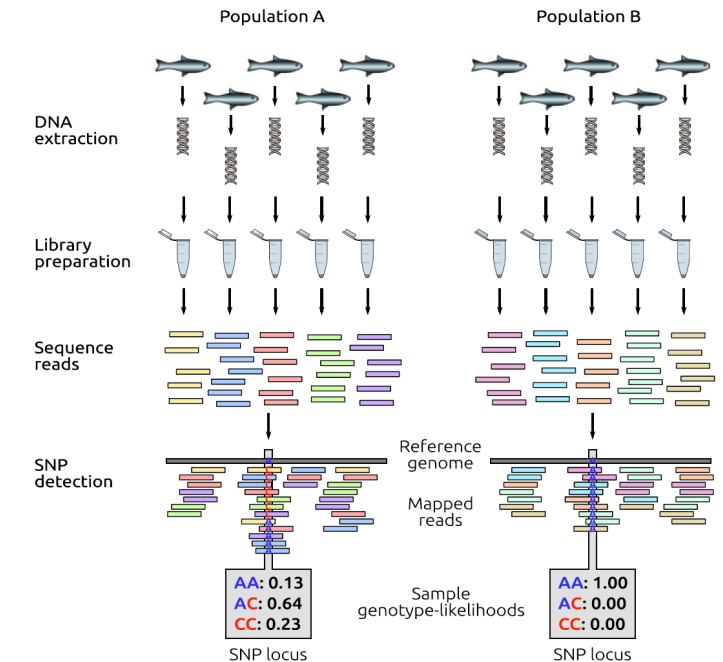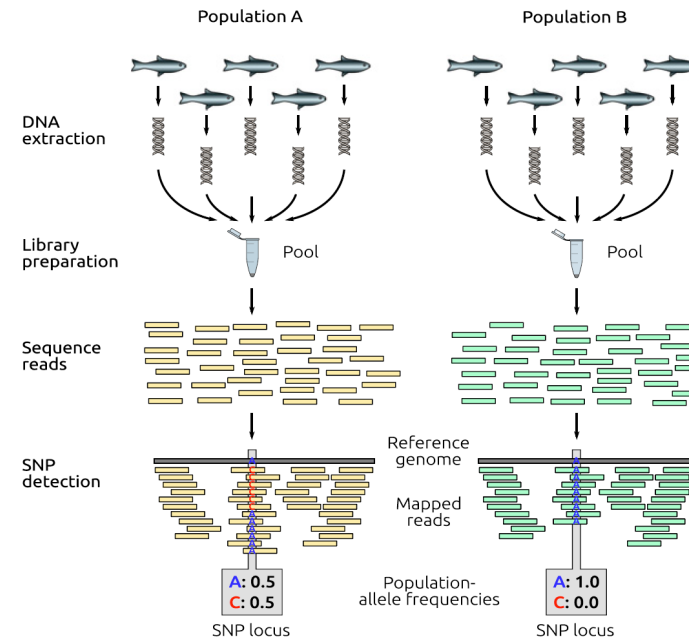
- Bioinformatic tool

## Short vs. long read:



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

Clusters

**Illumina**

SMRTbell Template

Polymerase Read

Subreads

Circular Consensus Sequence
(Read of Insert)

**Pacbio**

# How to confidently identify true variants

## Primary considerations

Individual vs. Multiplex vs. PoolSeq

- **Sequencing method/approach**

- Depth of coverage

- Variant type

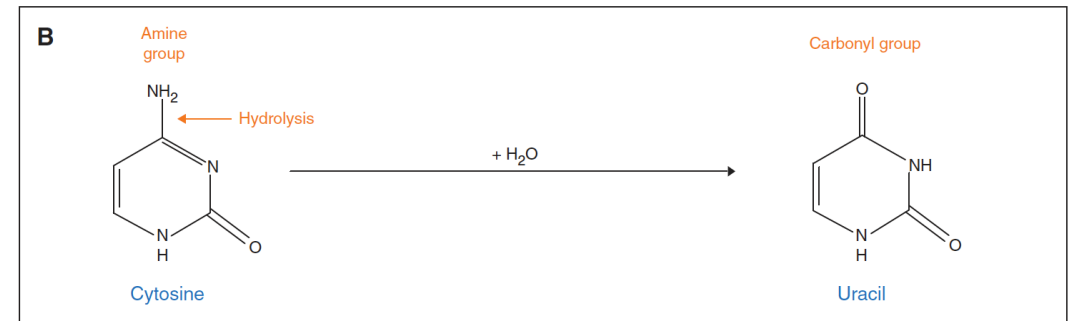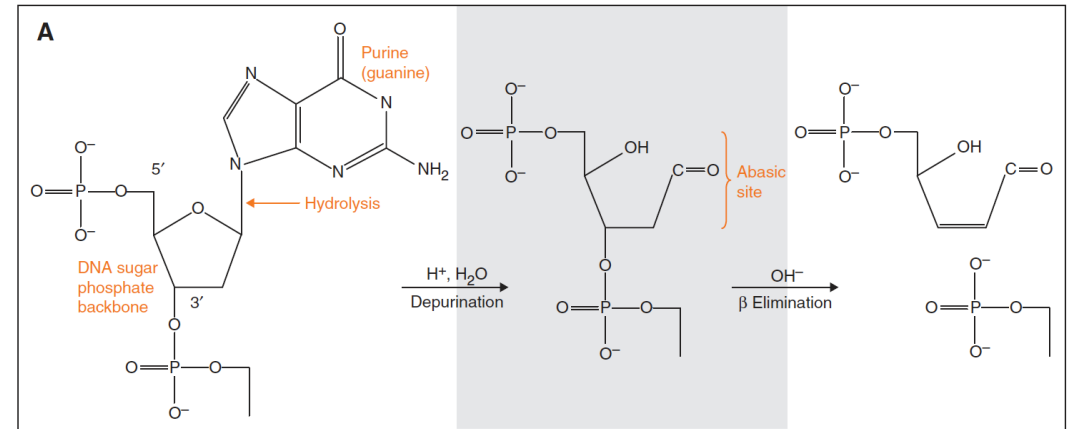- Inheritance mode

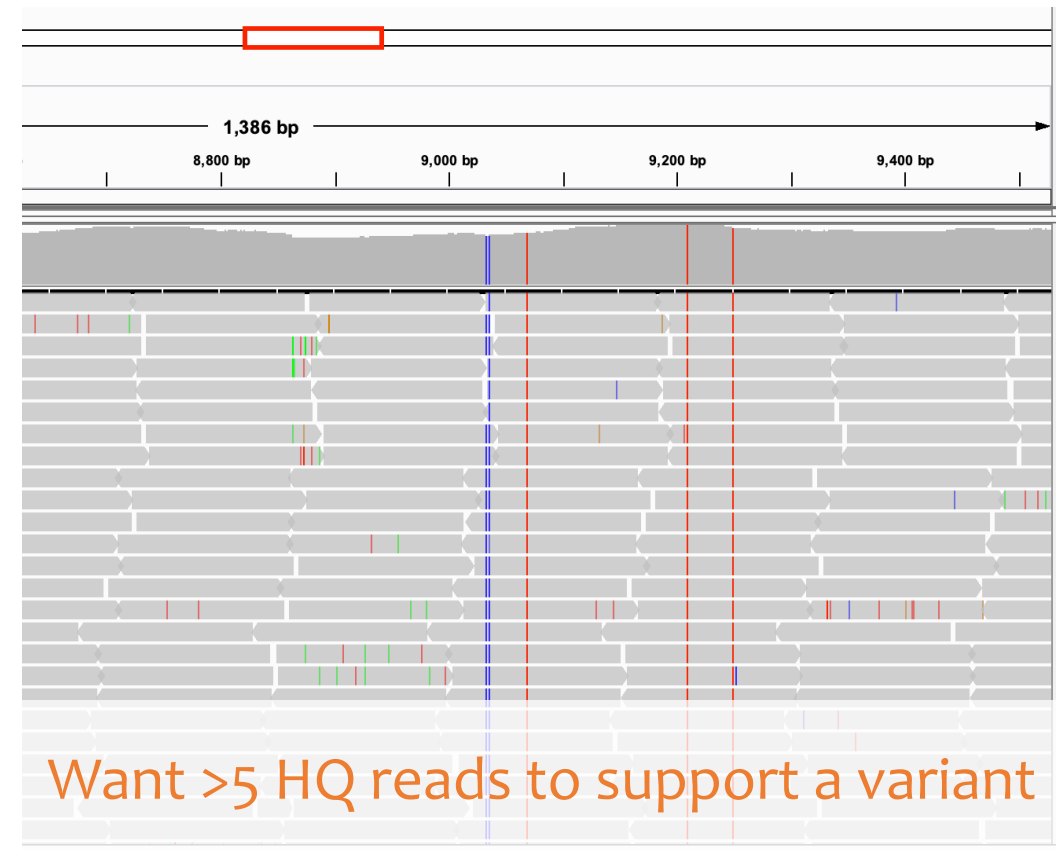- Reference quality

- Bioinformatic tool



Pardo & Ruzzante. 2017, *Mol Ecol*
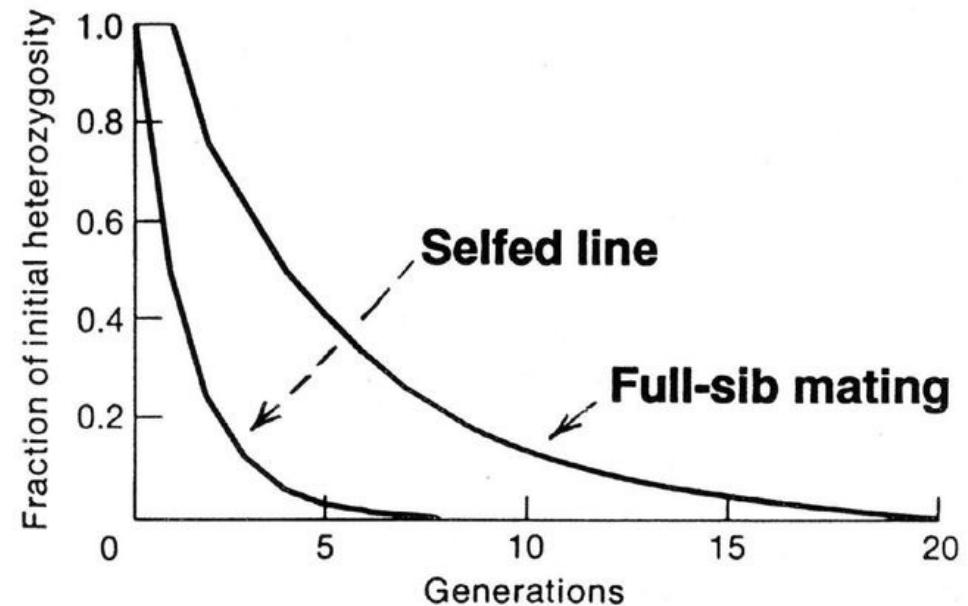
# How to confidently identify true variants

## Primary considerations

- **Sequencing method/approach**

- Depth of coverage

- Variant type

- Inheritance mode

- Reference quality

- Bioinformatic tool

## DNA Damage (e.g., Ancient DNA)



**Most C->T (and A->G) changes are the result of DNA Damage**
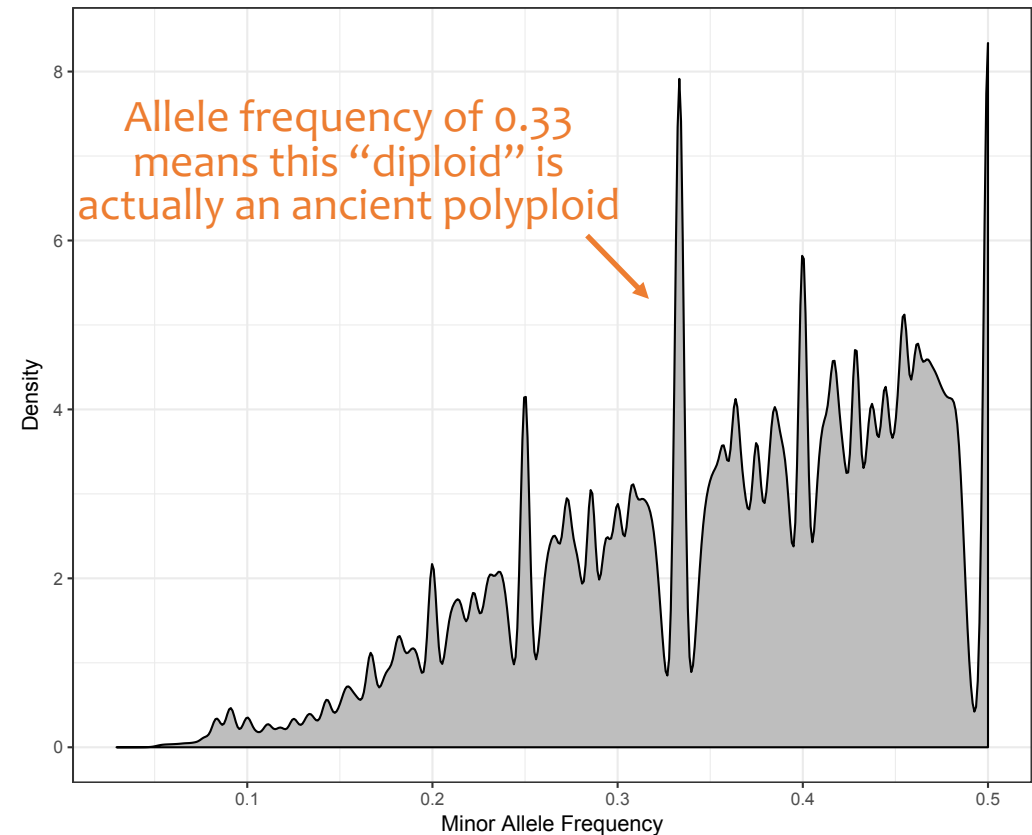
Depardy et al., 2013, *CSHP*

# How to confidently identify true variants

## Primary considerations

- Sequencing method/approach

- **Depth of coverage**

- Variant type

- Inheritance mode

- Reference quality

- Bioinformatic tool

### Goldilocks Principle of Read Depth



1,386 bp

8,800 bp    9,000 bp    9,200 bp    9,400 bp
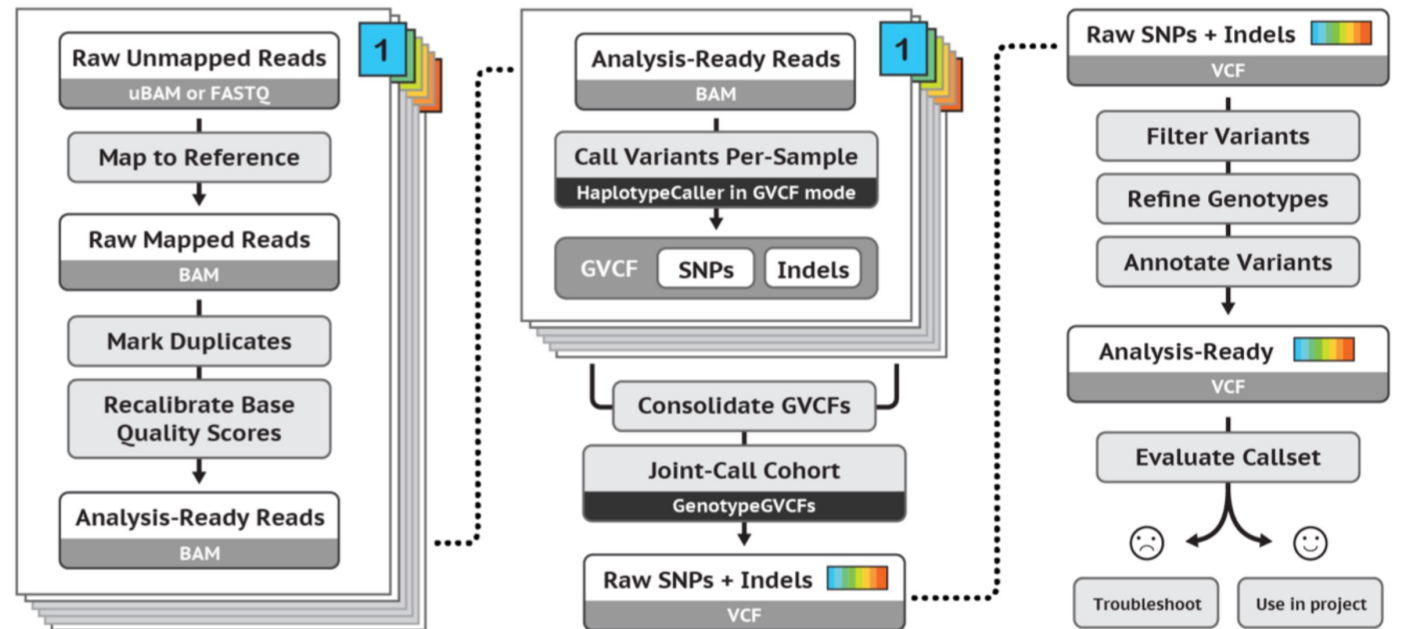
Want >5 HQ reads to support a variant

# How to confidently identify true variants

## Primary considerations

- Sequencing method/approach
- Depth of coverage
- **Variant type**
- Inheritance mode
- Reference quality
- Bioinformatic tool

SNV vs. In/Del vs. SV



Single Nucleotide Variant

Deletion

Insertion

Tandem Duplication

Interspersed Duplication

Inversion

Translocation

Copy Number Variant

**Types of Variants**

# How to confidently identify true variants

## Primary considerations

- Sequencing method/approach

- Depth of coverage

- Variant type

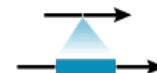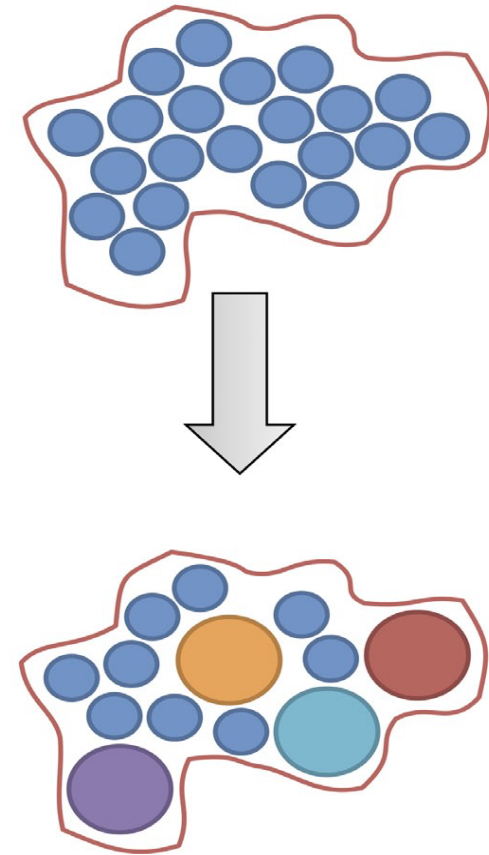- **Inheritance mode**

- Reference quality

- Bioinformatic tool

## Selfing vs. Outcrossing

# How to confidently identify true variants

## Primary considerations

- Sequencing method/approach

- Depth of coverage

- Variant type

- Inheritance mode

- **Reference quality**

- Bioinformatic tool

### Non-model reference genomes are often complex

Allele frequency of 0.33 means this "diploid" is actually an ancient polyploid

# How to confidently identify true variants

## Primary considerations

- Sequencing method/approach

- Depth of coverage

- Variant type

- Inheritance mode

- Reference quality

- **Bioinformatic tool**

BCFTOOLS vs. GATK vs. FreeBayes vs. DeepVariant

# Low-frequency variants:
## Distinguishing between sequencing errors, DNA damage, and true mutations

**Types of low-frequency variants:**

- Somatic mutations

- Heteroplasmy

- Rare polymorphisms in PoolSeq

- DNA Damage

Illumina error rate ~ $10^{-3}$

Coverage > 1000x means that you'll have an erroneous base call at **EVERY SITE**

**How to mitigate?**

# Low-frequency variants:
## Distinguishing between sequencing errors, DNA damage, and true mutations

**Types of low-frequency variants:**

- Somatic mutations

- Heteroplasmy

- Rare polymorphisms in PoolSeq

- DNA Damage



Sloan et al., 2017, *TIBTECH*

# Variant annotation:

## SnpEff, VEP, ANNOVAR

### Effect on function:

- Synonymous vs. Nonsynonymous
- Loss-of-Function

### Genomic location:

- Gene of interest
- Position in gene (e.g., exon/intron/UTR)

# Visualizing with IGV

# Structural Variant Detection



Types of Variants

# Structural variant detection depends upon context, experimental design

**Sample heterogeneity** (↓)

- Individual re-sequencing

- Somatic tissue

- Population-level sampling

- Tumor cells

- Plant mitochondria

Schoenfeld & Fox 2013

# Structural variant detection depends upon context, experimental design

**Sample heterogeneity**

- Individual re-sequencing

- Somatic tissue

- Population-level sampling

- Tumor cells

- **Plant mitochondria**



Unseld et al. 1997

# Three general approaches for structural variant calling



**De novo assembly, alignment**

**Short read mapping**

**Long read mapping**

Mahmoud et al. 2019

# Assembly-to-assembly comparison



**De novo assembly, alignment**

- <u>Genome-to-genome alignment</u>
- Low sample heterogeneity
- Easy to visualize with a dot plot
- **Diploid assemblies**

# Short-read mapping to identify structural rearrangements



**Short read mapping**

- Allow reads to map to multiple locations
- Depth allows for high sample heterogeneity
- Often difficult to visualize, as reads don't often span structure
- **Good for finding breakpoints in heterogeneous samples**

Mahmoud et al. 2019

# Long-read mapping to identify structural rearrangements



**Long read mapping**

- Allow reads to map to multiple locations
- Don't need a ton of depth
- Reads are long enough to span across structural element
- Occasional library artifacts to watch out for
- **Good for homogeneous samples**

Mahmoud et al. 2019

# Long-reads are better than short reads for identifying structural rearrangements



Mahmoud et al. 2019

# Assembly-based methods are more sensitive than mapping-based methods



Mahmoud et al. 2019

# Three general approaches for structural variant calling

## Structural variant callers

Mahmoud et al. 2019

# Reconstructing tandem duplications – problems with consensus



Stellwagen & Burns 2021

# Reconstructing tandem duplications – problems with consensus

# Next up:
# Performing variant calling with GATK



**Next week's homework:**

1. GATK "best practices" pipeline on Stickleback Illumina reads,

2. SV detection w/ PB reads