# Mapping and Calling Variants

Eve198

Week 4: April 23rd
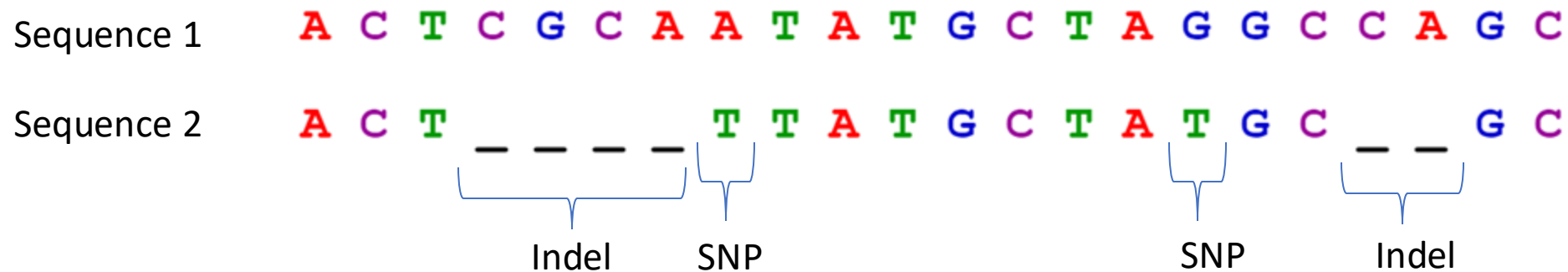
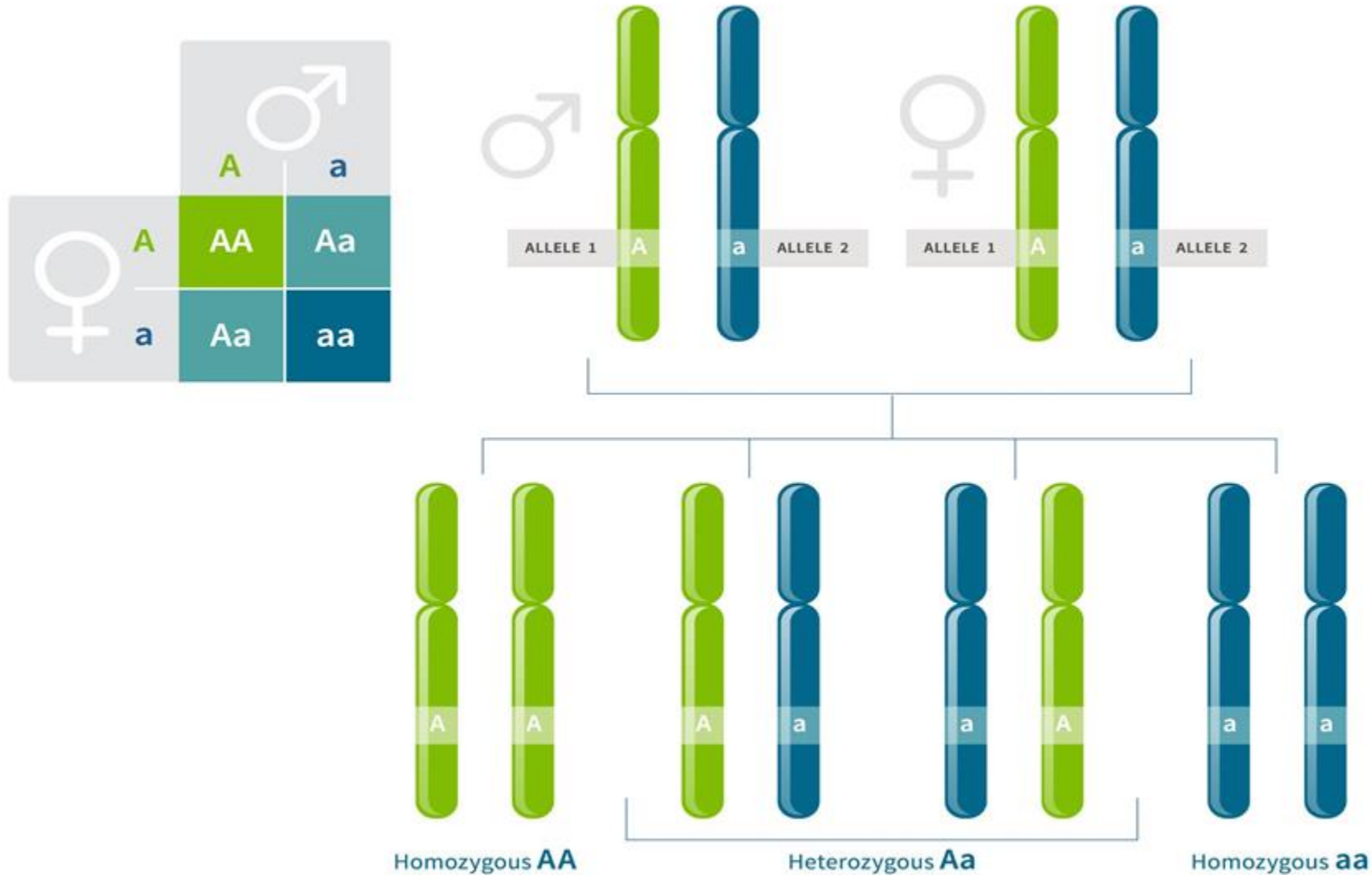Maddie Armstrong & Rachael Bay

# What is a genetic variant?

A region of the genome that differs from the reference (or another genome)

Signifies a mutation and can be a single base-pair, or larger insertion and/or deletion of several base-pairs.
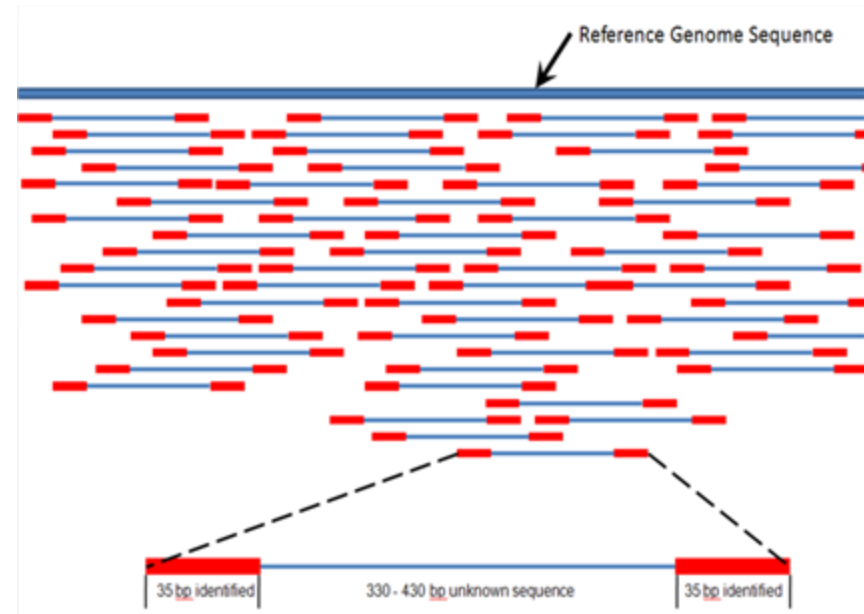
Sequence 1    A C T C G C A A T A T G C T A G G C C A G C

Sequence 2    A C T _ _ _ _ _ T T A T G C T A T G C _ _ G C

Indel    SNP                SNP    Indel

# What is a genotype?

# How do we find a variant?

Map and align sequences from other individuals to a reference genome

- Does It matter what your reference genome is?
    - Is it the same or different species?
    - Is it from the same population?

- Short answer: Yes, it matters!



Reference Genome Sequence

35 bp identified | 330 - 430 bp unknown sequence | 35 bp identified

Genomes are continually being improved & sequenced all the time!

American Genetic Association

OXFORD

## Genome Resources

# A genome assembly of the American black bear, *Ursus americanus*, from California

Megan A. Supple[1,2,*,†] iD, Merly Escalona[3,†] iD, Jillian Adkins[4], Michael R. Buchalski[5] iD, Nicolas Alexandre[1,2] iD, Ruta M. Sahasrabudhe[6] iD, Oanh Nguyen[6] iD, Samuel Sacco[1] iD, Colin Fairbairn[1] iD, Eric Beraut[1] iD, William Seligmann[1] iD, Richard E. Green[3] iD, Erin Meredith[4,‡], Beth Shapiro[1,2,‡] iD

[1]Department of Ecology and Evolutionary Biology, University of California, Santa Cruz, CA, United States,
[2]Howard Hughes Medical Institute, University of California, Santa Cruz, CA, United States,
[3]Department of Biomolecular Engineering, University of California, Santa Cruz, CA, United States,
[4]Wildlife Forensic Lab, Law Enforcement Division, California Department of Fish and Wildlife, Sacramento, CA, United States,
[5]Wildlife Genetics Research Unit, Wildlife Health Laboratory, California Department of Fish and Wildlife, Sacramento, CA, United States,
[6]DNA Technologies and Expression Analysis Core Laboratory, Genome Center, University of California, Davis, CA, United States
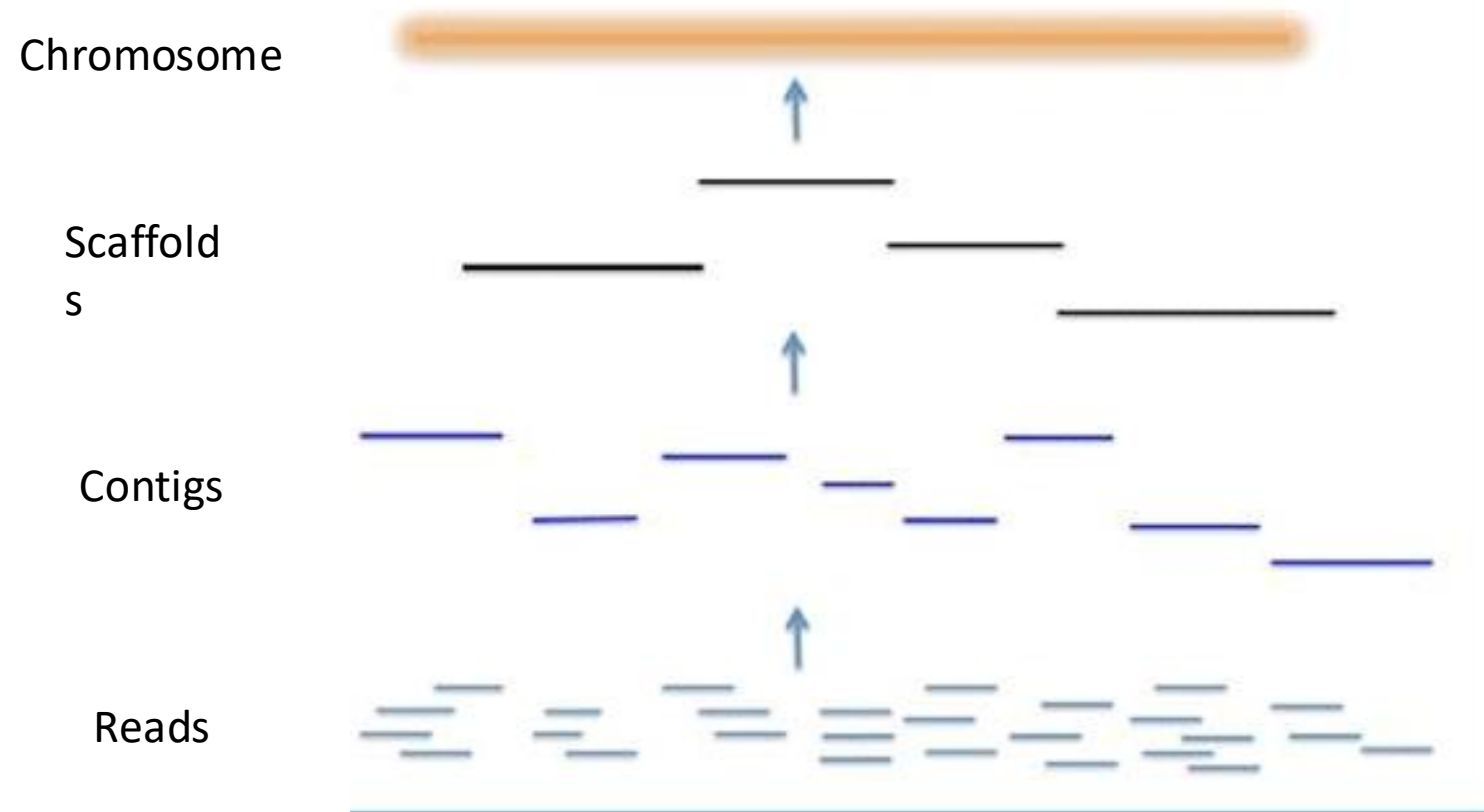
[†‡]These authors contributed equally to this work.
[*]Corresponding author: Email: megan.a.supple@gmail.com

Corresponding Editor: Klaus-Peter Koepfli

The American black bear, *Ursus americanus*, is a widespread and ecologically important species in North America. In California, the black bear plays an important role in a variety of ecosystems and serves as an important species for recreational hunting. While research suggests that the populations in California are currently healthy, continued monitoring is critical, with genomic analyses providing an important surveillance tool. Here we report a high-quality, near chromosome-level genome assembly from a *U. americanus* sample from California. The primary assembly has a total length of 2.5 Gb contained in 316 scaffolds, a contig N50 of 58.9 Mb, a scaffold N50 of 67.6 Mb, and a BUSCO completeness score of 96%. This *U. americanus* genome assembly will provide an important resource for the targeted management of black bear populations in California, with the goal of achieving an appropriate balance between the recreational value of black bears and the maintenance of viable populations. The high quality of this genome assembly will also make it a valuable resource for comparative genomic analyses among black bear populations and among bear species.

https://academic.oup.com/jhered/article/115/5/498/7713838

# Finding variants – some terminology



Chromosome

Scaffolds

Contigs

Reads
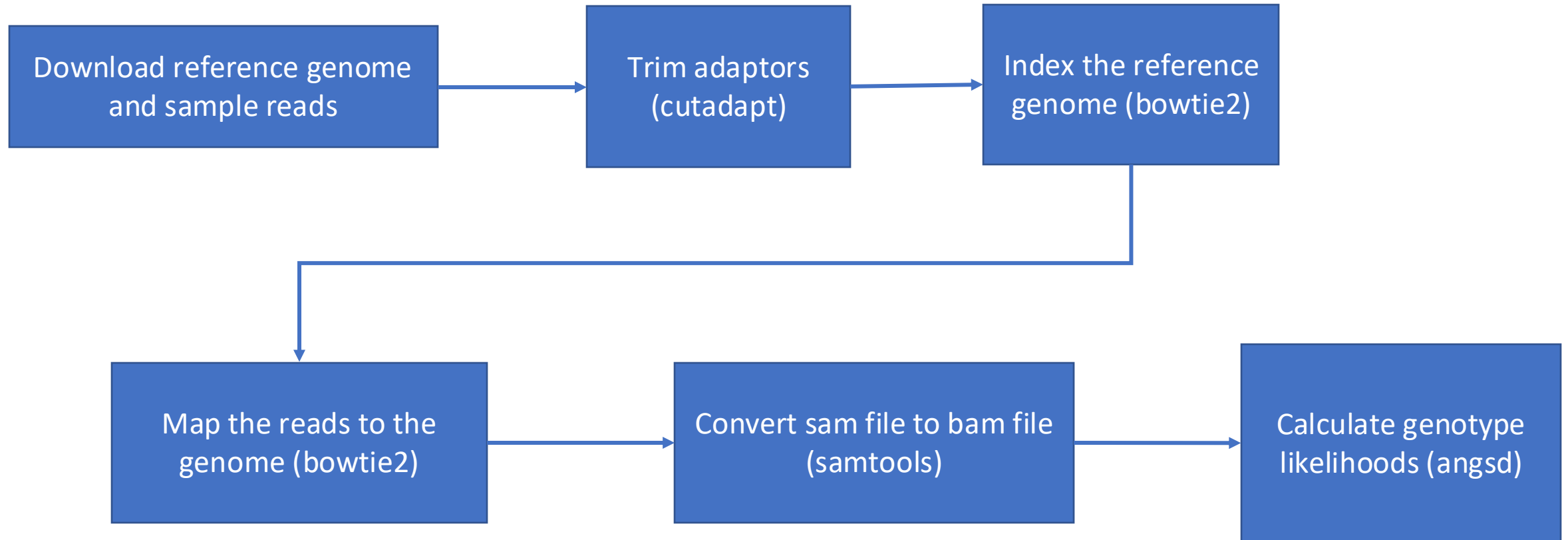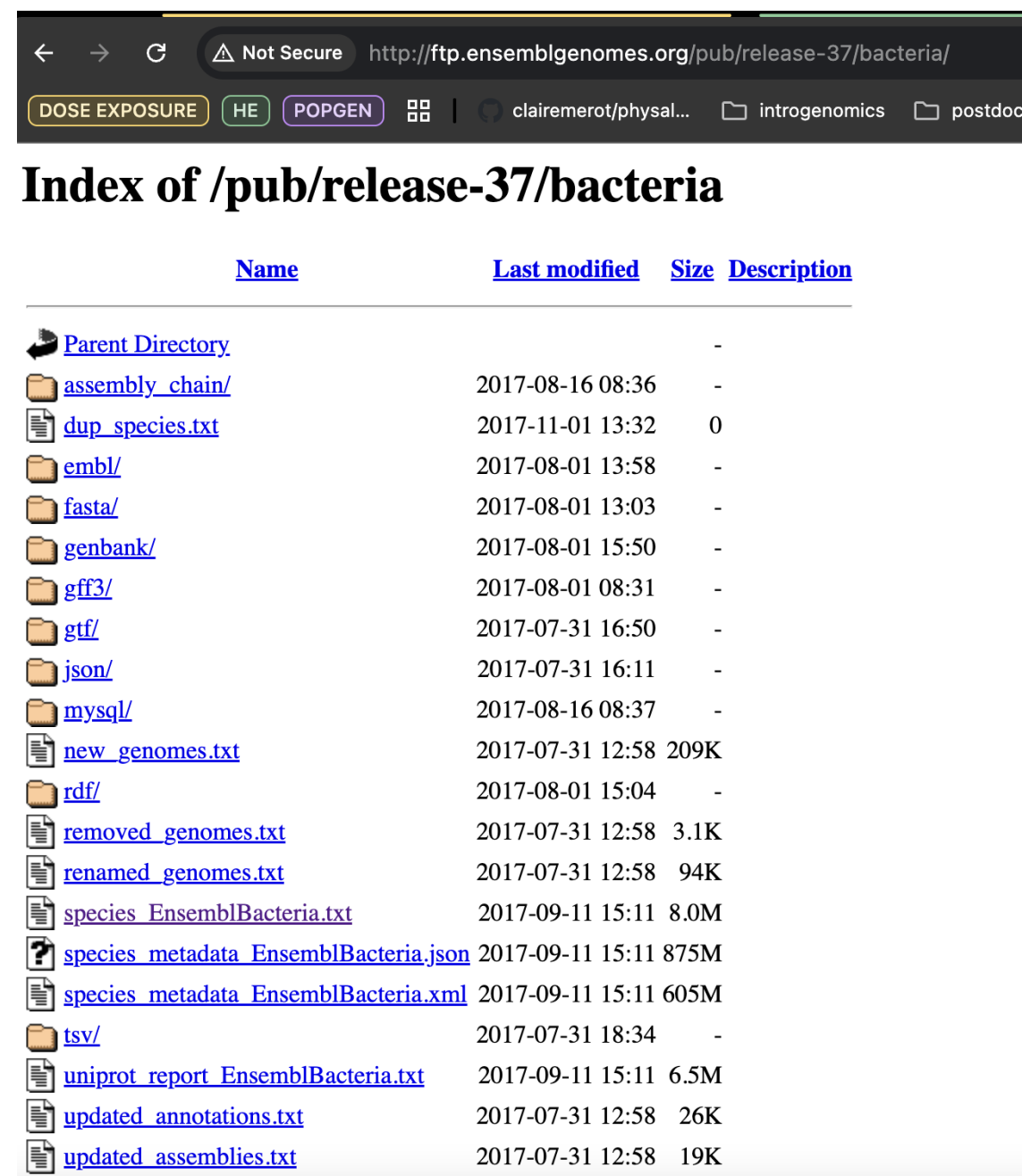
A reference genome is a collection of contigs

# Finding variants – some terminology

```
>KN893585.1 Parastichopus parvimensis isolate Sea Cucumber 01 unplaced genomic scaffold Scaf
fold11424, whole genome shotgun sequence
CATATATGTGAGAGAAAAGTGCATTGACCTGGCTTTAACTTGACAACAAGACTGTTTCCCGCTCGTTACGAATAATCTTCTATATCCTAATA
ATCGCTATGCAAGGCTATAAGcaatatcacataatatcacACCAGCTTGTAAGGTTACATTTAAACATCAGGTGGTTATTTCCAATTAGACA
GTGTTGAAATCCATCAACGCTTCCTGGAAAGTAAAATCCAAACCGTTAAATCATAACCCATCCATATGGTATTGGCTTCCTTGTATATGCTA
CCACCTATATACATGTGAGCCTACAGCAATAATGATTCCTTCCATACACACCACCAAGGAAACCCAACTGCTGGTTTTTGACACATGCCGTA
GGAGTTACAGCCTGTCTTTCATTGCCAACAACATGAGATGACATGATGTTTTCTCTGTCACATTTTGGTGTGAATTTTTCTCGTTTGCTATA
ACCACCGATTTTTACTGTAGGGTTTTAATTCTCAAATTTATCAATAAGTTTGGGTAAGACAAACATGCATTAAACTAAAGTTAGTTTCTCTG
ATCCTCCATTTTGTTCCCAGTCATTGAGGATTATTAAAAGTGACAAAAGGTCTTAGTGGTTAATACTAACTTTTAGAGAGGCAAGAAAATGA
CTTGAAATTTCAGTTTGGGTGACCATCATTTGAGTTAAGGTTCACACAGTTTTAAAGATGCATAGGAATGAgacaaaaggggaaaaaagctT
ACTCCGCGTGGAAATTCAATGACACAACTTCCTGTTCTATGTGATGGACATAACCCCTGTAAGATTTATCTCCTCTTCCGCTTGAATGTGTC
GCATAGAGATGATCTCCTCTGAGTACAGAAGGACGATTCTCGGCTAACCCGGGGACCTGTAAATGAAGAGTTTTACACGTGAGCTAGCGAGA
GGGGGAAGATCGACCACAATTGCAATTATAGTCCGACACAACTGTAATTGCCAAACATACCTGCAGCAACATACTCTTTGGATcccacgttt
tttttattaacaaatgAAATTCTAGACTTTTTGAAGACCAAAACACGTCTTATGGTTTACTATATGAAGCCTACACACTAATGATGTCCTA
AGGTTATGTTACCTATGATAGGCATTATCAATTGTAACTCTTGCAAAATATACACTAACTAACCCCTTGTGTTAATTTTTGGTGAAGGGGTA
TTCAATAGGCCATGAGTGCCAAACATGACATGCTATAGCTATTTTTTTTTCCCACCTAAGTGTGACATTAACTTTATCTCACACTTCTTCAAA
CCTTGCTAGCATAAGCCATATCATTTAGGAAGAAGTGTTAAAATGAGGATGTTTCCATCCTTTACAGACTCCAATCGAAAATTCAAAGACTT
TCAAGATCTAGAAAAGAGGGTTTTCCTTTTCCCTAGGTTTCCCCTCCCTGCCCCATTTTGCAGATCATGAGGGAACATGCATACattagtta
attaaaatatgaaaaacattgttaatGAGGGATGaatgaattttgacaaaaaagaaGAGTAAAGATGACTGGATTTGAATATTTAGaaagct
tttaattttaattcttaaACATTTGAGAATATGCTAAAATTATTGTTTCAAATCGTCAATTAGTACTCTGGCAACATACCTTCAGATTCAGT
AAACCCACATGAGTCCTGCTTTTGGACATTTCAGCTGCTTTCTTGTCGTAGCGGCGAATGTCAACTTTCATTTGATGTTCTTCTGCGTAGAG
GAGATTCTCAAACTTTTGAGAGTAGTTCTCTTCCGAGAGAGGATCCTGCAAGCTTTCGATTCTCCTagtaattaaagaaaatgaaaaagttt
```

A reference genome is a collection of contigs

Typically, in fasta format

# Finding variants - pipeline

# Step 1: Download the data!



wget ftp://ftp.ensemblgenomes.org/pub/release-37/bacteria/species_EnsemblBacteria.txt

# Step 1.2: Quality Control with fastqc



## FastQC Report
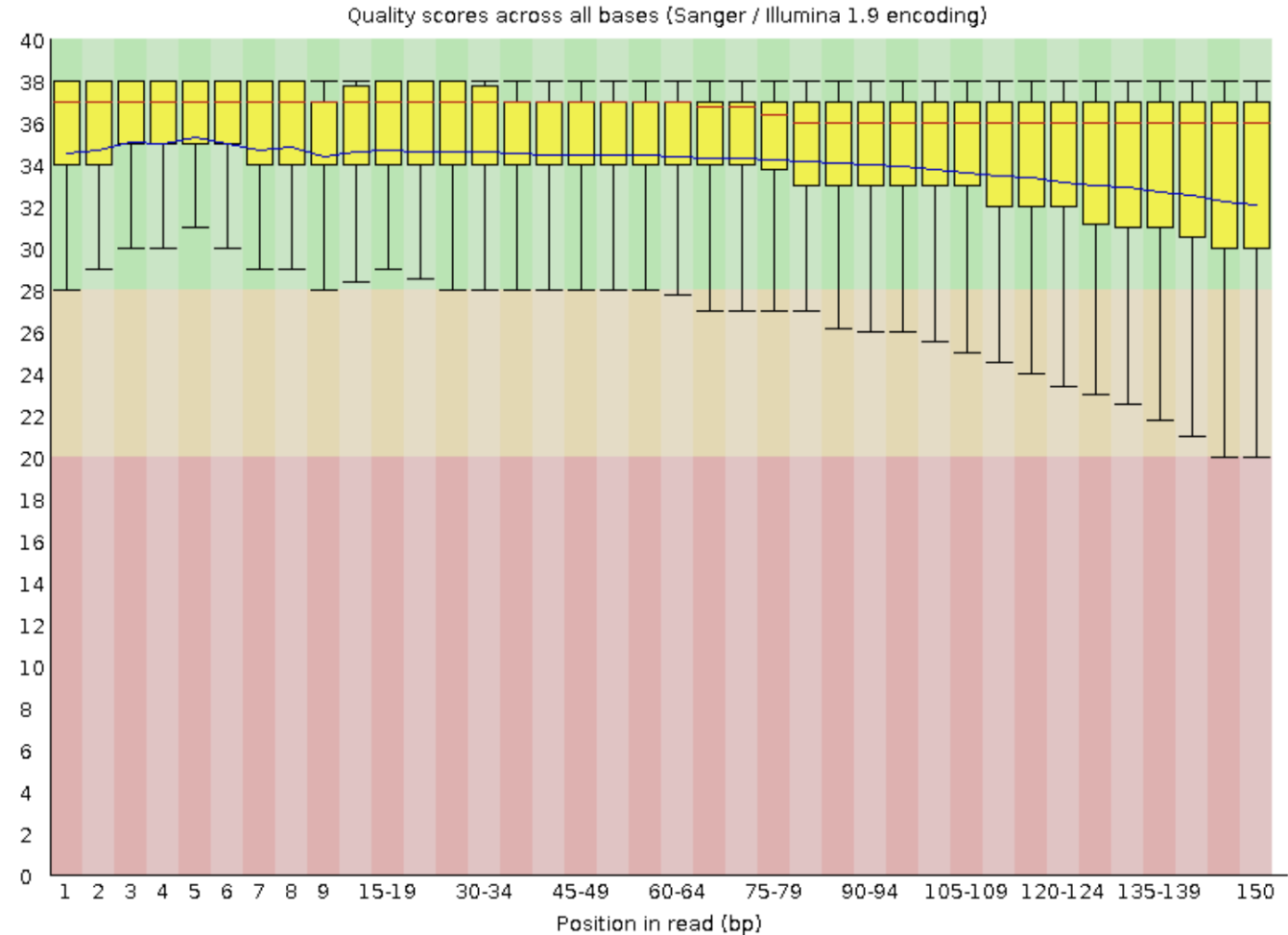
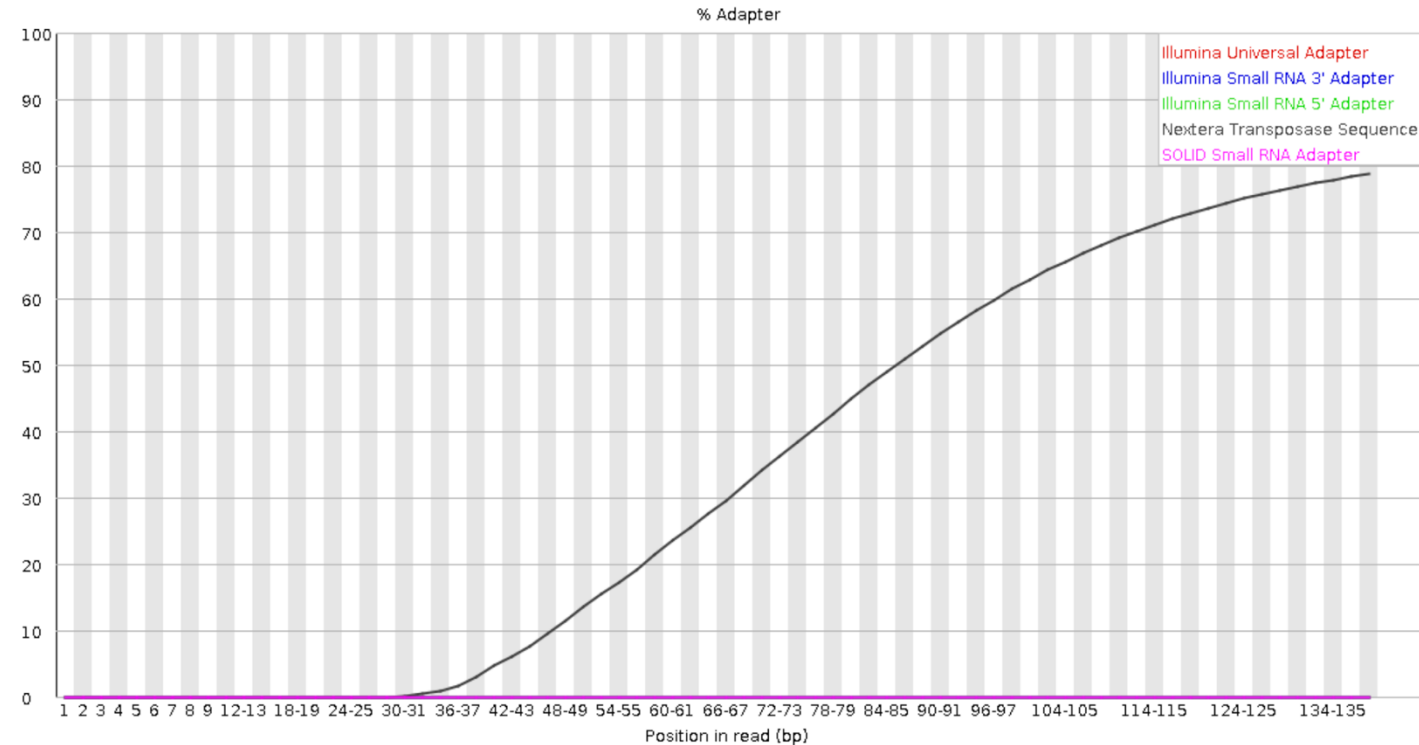### Summary

✅ Basic Statistics
✅ Per base sequence quality
✅ Per sequence quality scores
❌ Per base sequence content
⚠️ Per sequence GC content
✅ Per base N content
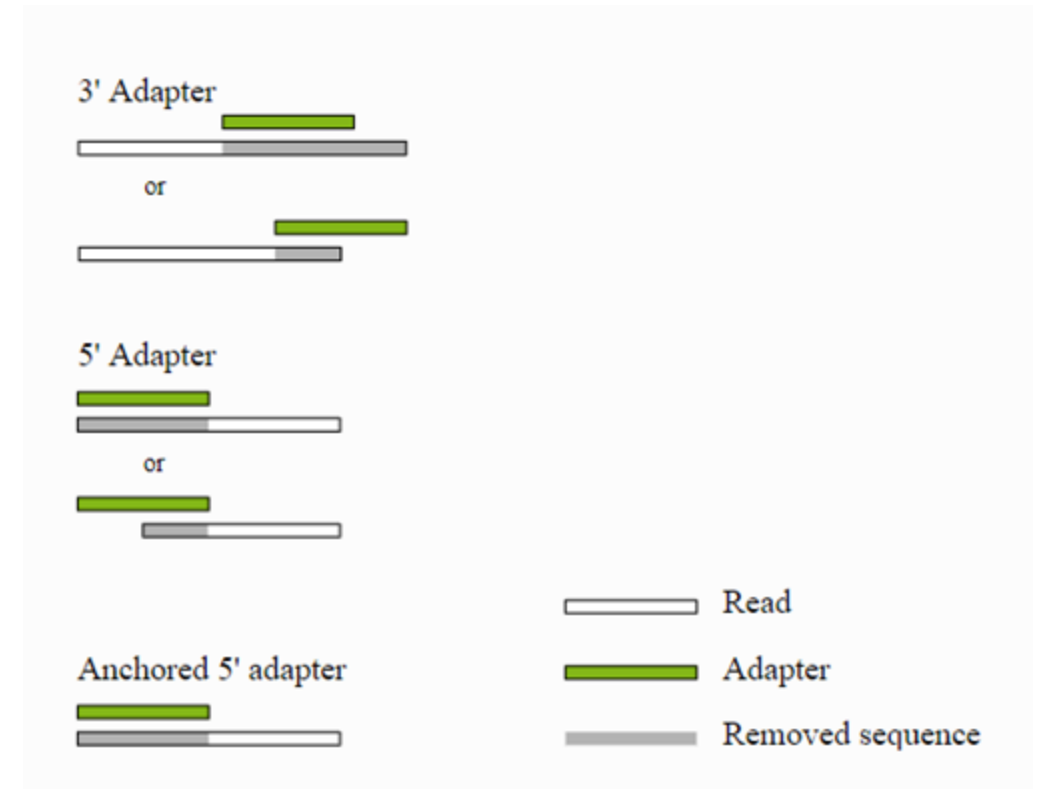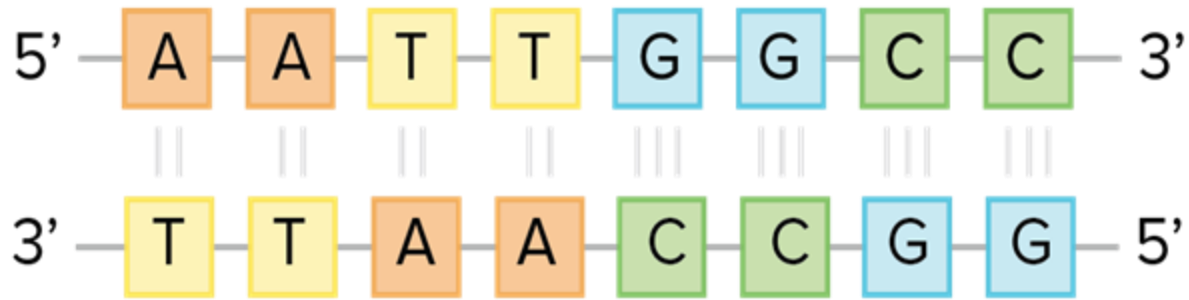✅ Sequence Length Distribution
✅ Sequence Duplication Levels

✅ **Basic Statistics**

| Measure | Value |
|---------|-------|
| Filename | Bir8_1.fq.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 2080506 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 150 |
| %GC | 46 |

✅ **Per base sequence quality**

Quality scores across all bases (Sanger / Illumina 1.9 encoding)

# Step 1.2: Quality Control with fastqc

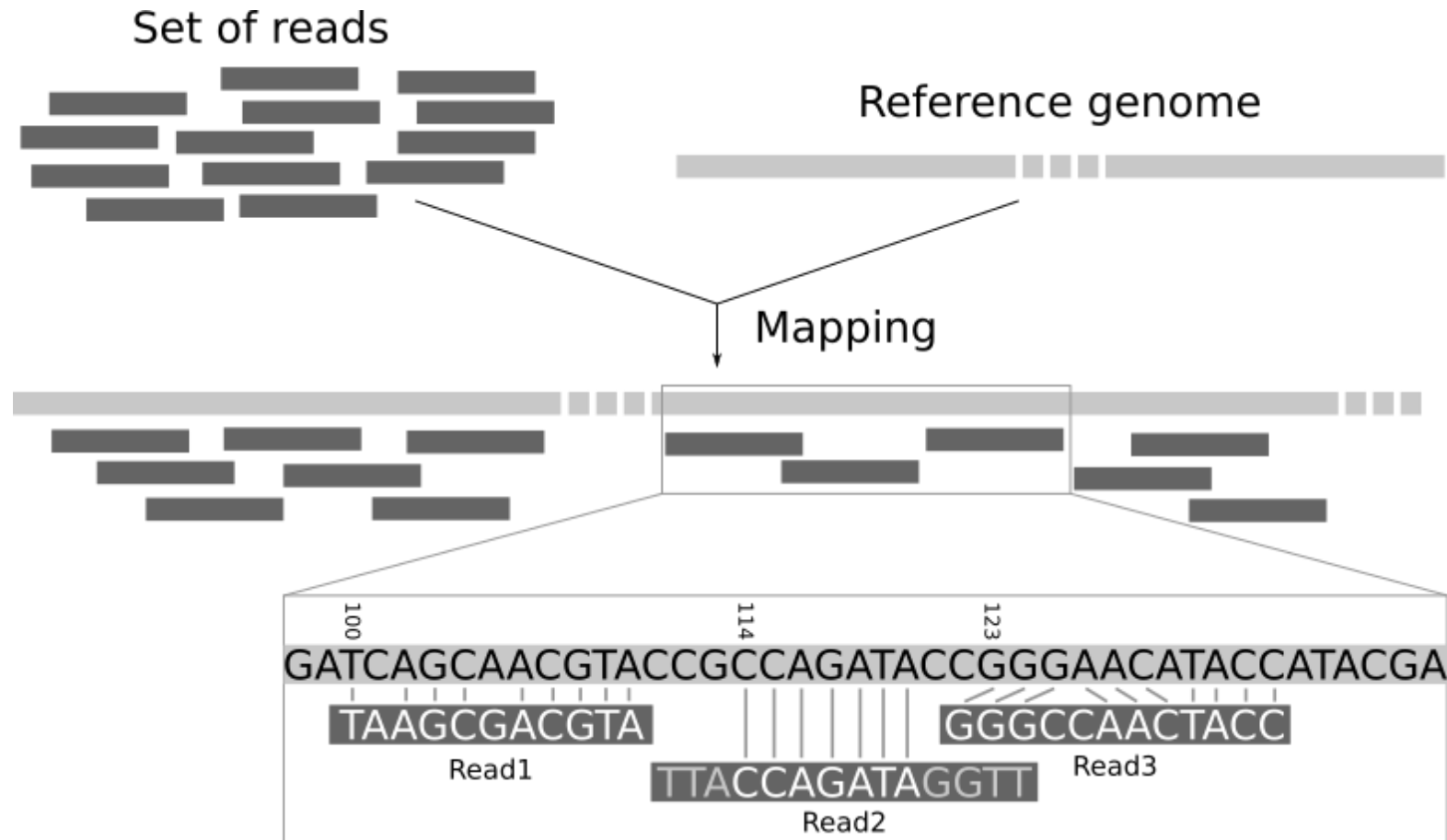# Step 2: Trimming adaptors from reads

# Step 2: Trimming adaptors from reads

# Step 3: Index to Reference Genome

# Step 3: Index to Reference Genome

# Step 4: Mapping to the Reference Genome

# Step 5: SAM and BAM file formats

Sequence Alignment Map, Binary Alignment Map



Head of .sam file



Tail of .sam file

# Step 5: SAM and BAM file formats

Sequence Alignment Map, Binary Alignment Map



Name of read

Name of contig where read aligns

Position on contig where 5' end starts

Alignment information "cigar string"

80M = contiguous match of 80bp

Tail of .sam file

# Step 6: "Calling" a genotype

Two main ways to estimate a genotype

1. Hard genotype calling
2. Genotype likelihoods

Depends on the type of data that you have. If you have reads with a high degree of coverage (many copies of the same read) you can do hard genotype calling. If you have variable coverage or low coverage you can do genotype likelihoods, to account for some uncertainty in the genotype.

# Step 6: "Calling" a genotype: alignment

# Step 6: Genotype likelihoods

In ANGSD    http://www.popgen.dk/angsd/index.php/Genotype_Likelihoods

Accounts for some uncertainty in the genotype estimation

## Theory

Genotype likelihoods are in this context the likelihood the data given a genotype. This is to be understood as we take all the information from our data for a specific position for a single individual, and we use this information to calculate the likelihood for our different genotypes. Since we assume diploid individuals it follows that we have 10 different genotypes.

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|----|----|----|----|----|----|----|----|----|----|
| AA | AC | AG | AT | CC | CG | CT | GG | GT | TT |

And we write the genotype likelihood as

$$L(G = \{A_1, A_2\}|D) \propto Pr(D|G = A_1, A_2), \qquad A_1, A_2 \in \{A, C, G, T\}.$$