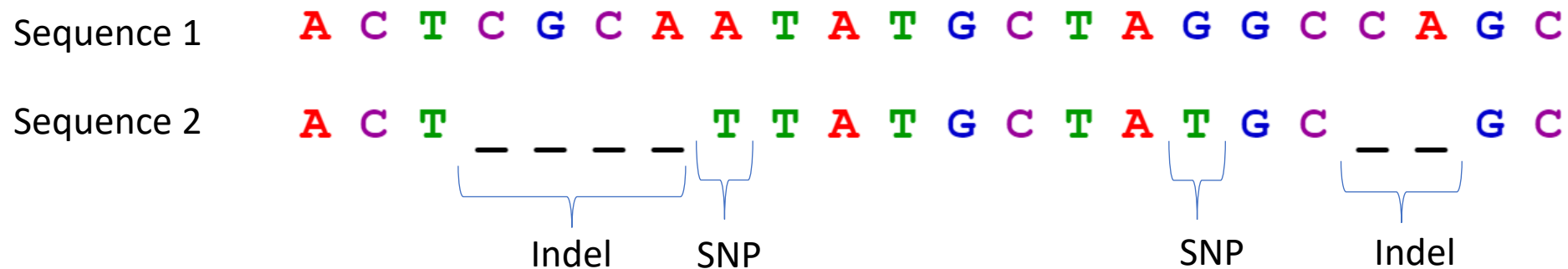# Marine Genomics

April 13th 2021

Mapping and calling variants

# What is a genetic variant?

A region of the genome that differs from the reference (or another genome)
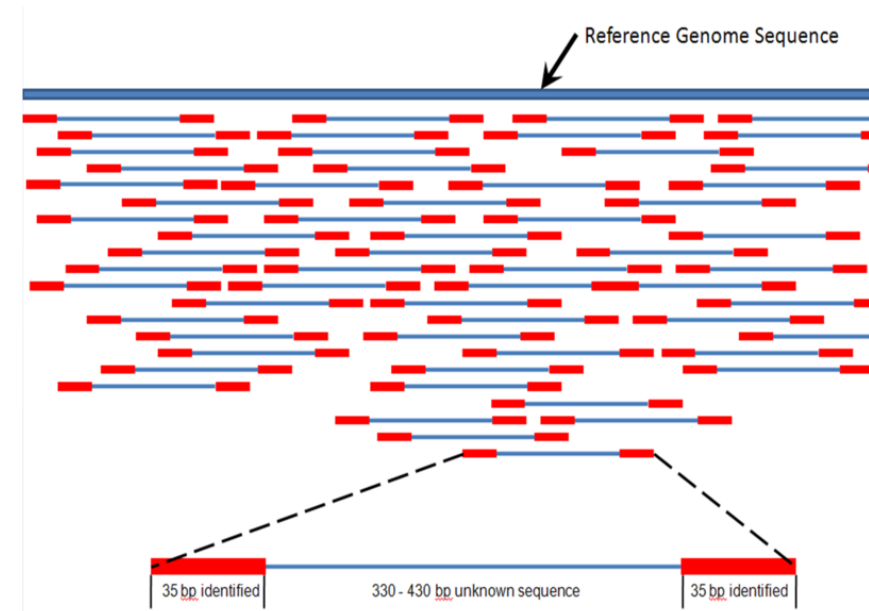
Signifies a mutation and can be a single base-pair, or larger insertion and/or deletion of several base-pairs.

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

# 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

- More genomes are being sequenced all the time

- Many marine organisms don't yet have a genome sequence available

## The structure, function and evolution of a complete human chromosome 8

Glennis A. Logsdon, Mitchell R. Vollger, […] Evan E. Eichler ✉

### Abstract

The complete assembly of each human chromosome is essential for understanding human biology and evolution[1,2]. Here we use complementary long-read sequencing technologies to complete the linear assembly of human chromosome 8. Our assembly resolves the sequence of five previously long-standing gaps, including a 2.08-Mb centromeric α-satellite array, a 644-kb copy number polymorphism in the β-defensin gene cluster that is important for disease risk, and an 863-kb variable number tandem repeat at chromosome 8q21.2 that can function as a neocentromere. We show that the centromeric α-satellite array is generally methylated except for a 73-kb hypomethylated region of diverse higher-

# Finding variants – some terminology

Chromosome

Scaffolds

Contigs

Reads

A reference genome is a collection of contigs

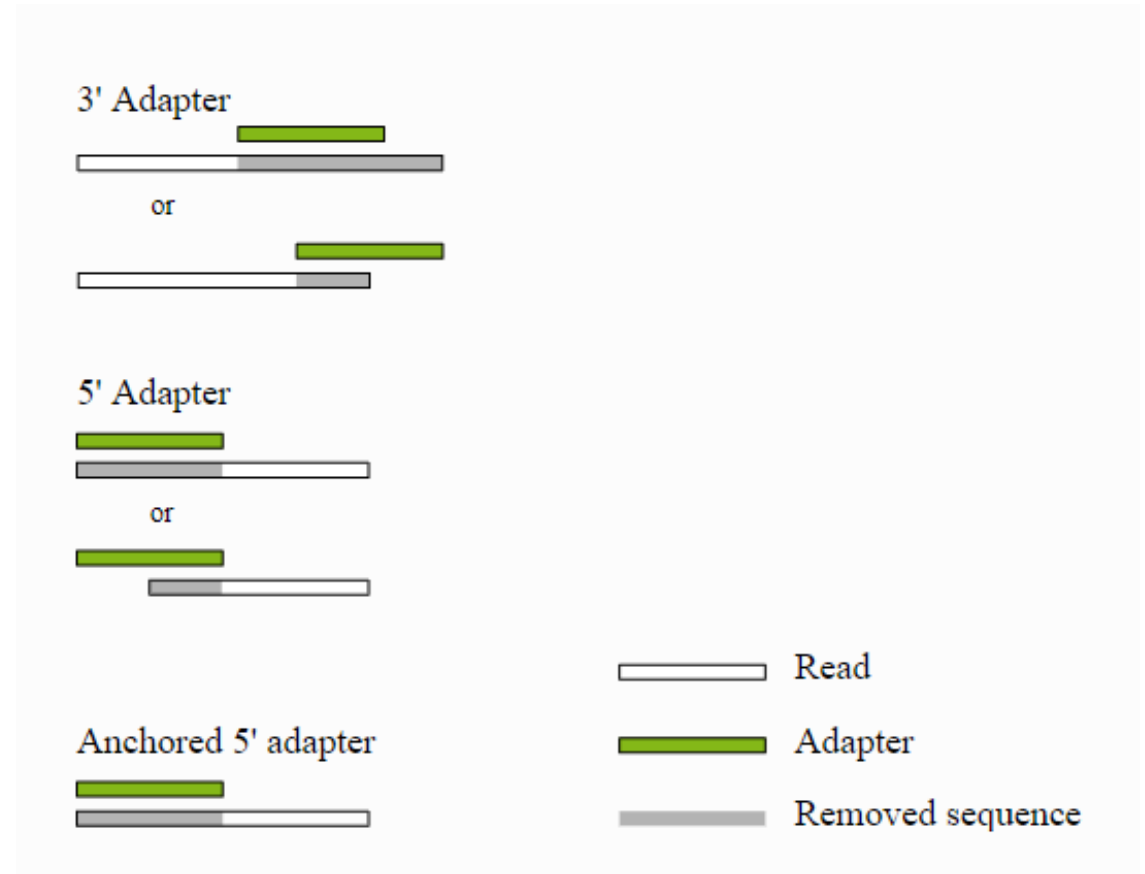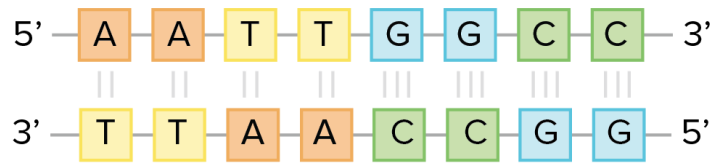# Finding variants – some terminology



A reference genome is a collection of contigs

Typically in fasta format

# Finding variants - pipeline

- Get reads and genome (download from git hub)

- Trim adapters off of reads (cutadapt)

- Index genome (bowtie2)

- Map reads to genome -> generate a sam file -> convert to bam file (bowtie2, samtools)

- Calculate genotype likelihoods (angsd, samtools)

- Happy dance

# Trimming adaptors from reads

# Trimming adaptors from reads

# SAM and BAM file formats

Sequence Alignment Map, Binary Alignment Map



Head of .sam file

Tail of .sam file

# 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

# 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

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