

Mapping

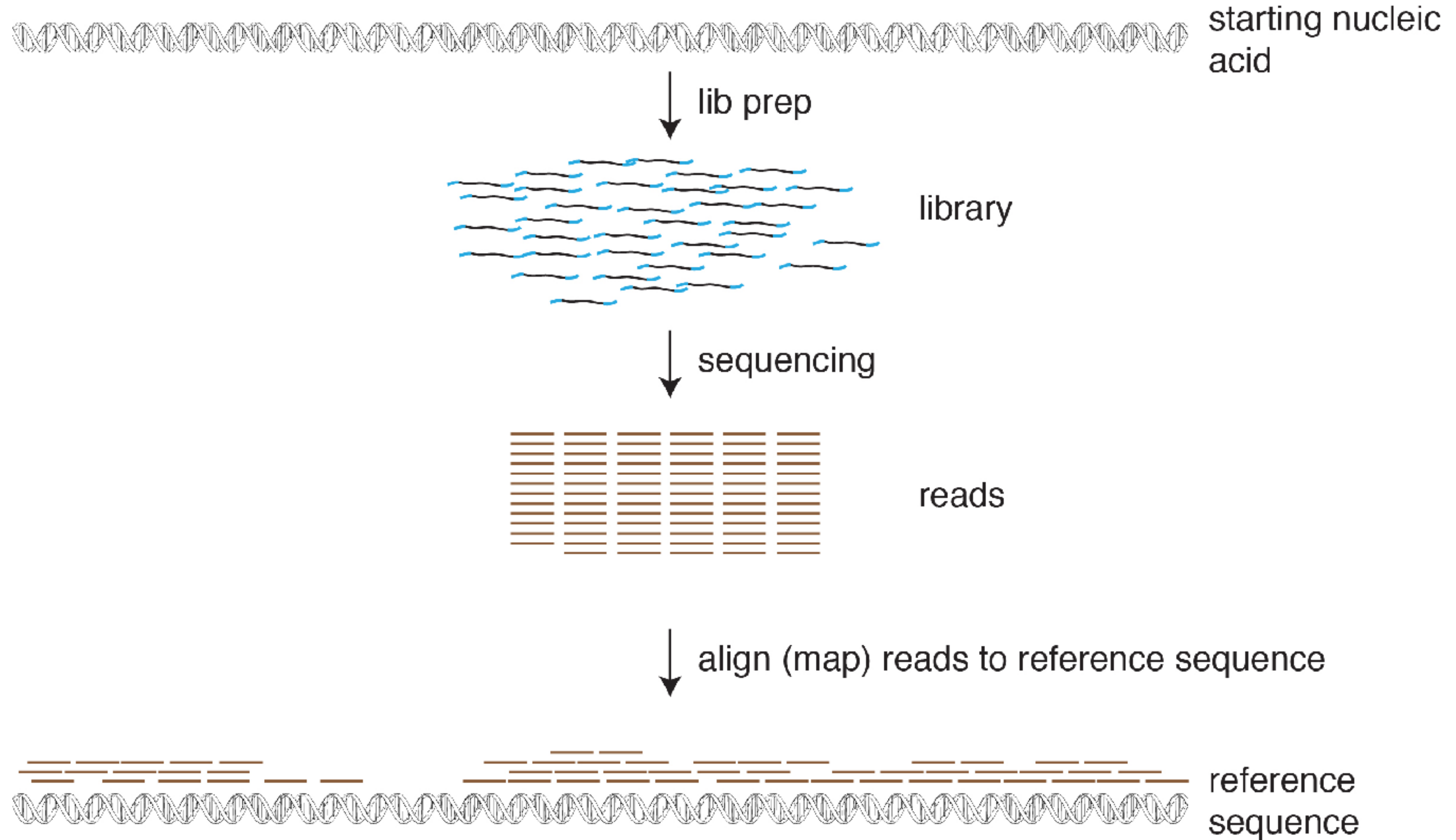
Mark Stenglein



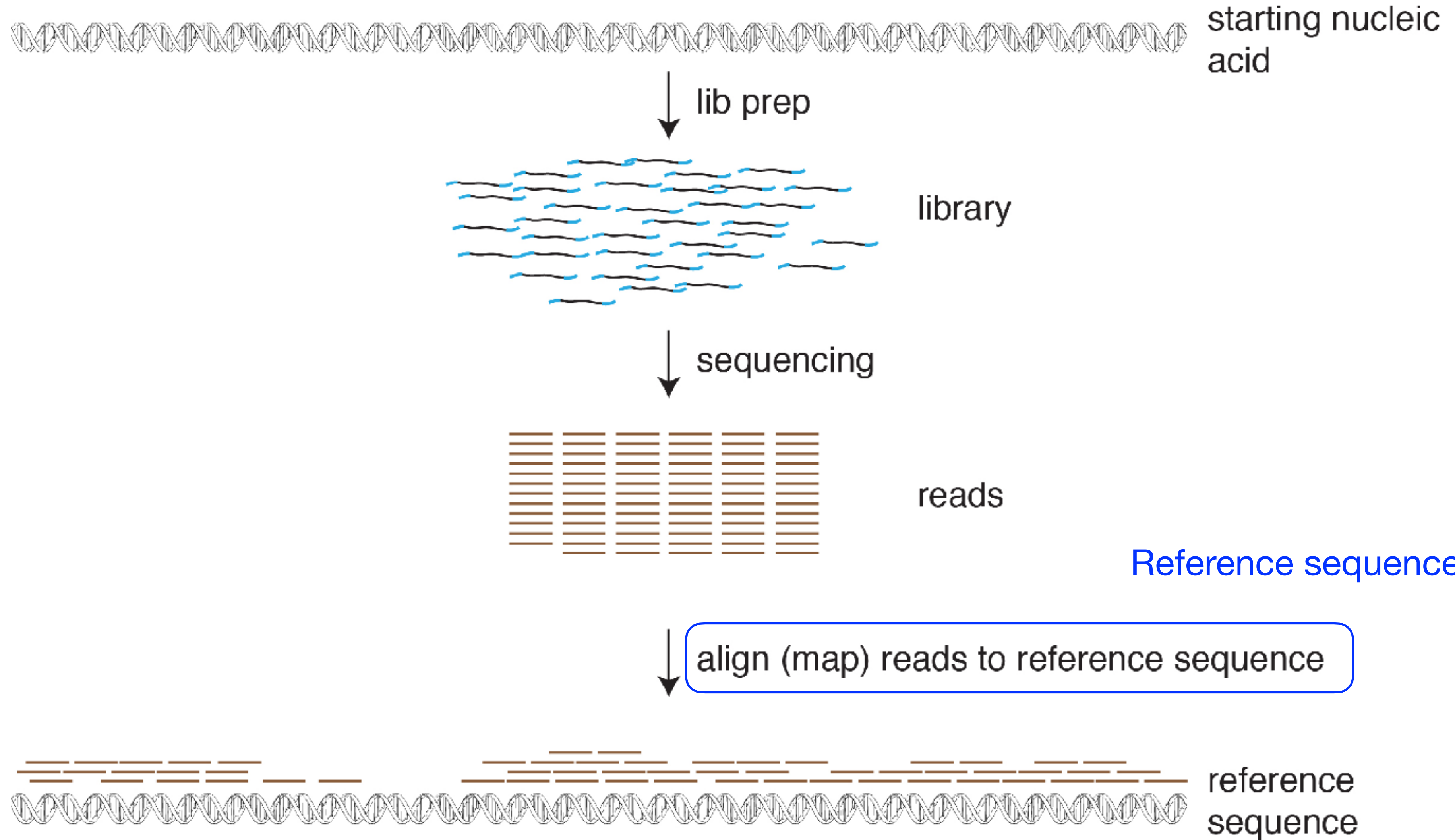
Computational Biology
Workshop

Todos Santos Center
May 9-12, 2022

Mapping is the process by which sequencing reads are aligned to the region of a genome from which they derive.



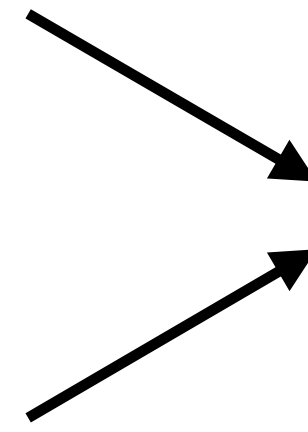
You need to have an existing reference sequence to map to



Mapping inputs

Sequence Reads

Reference sequence(s)



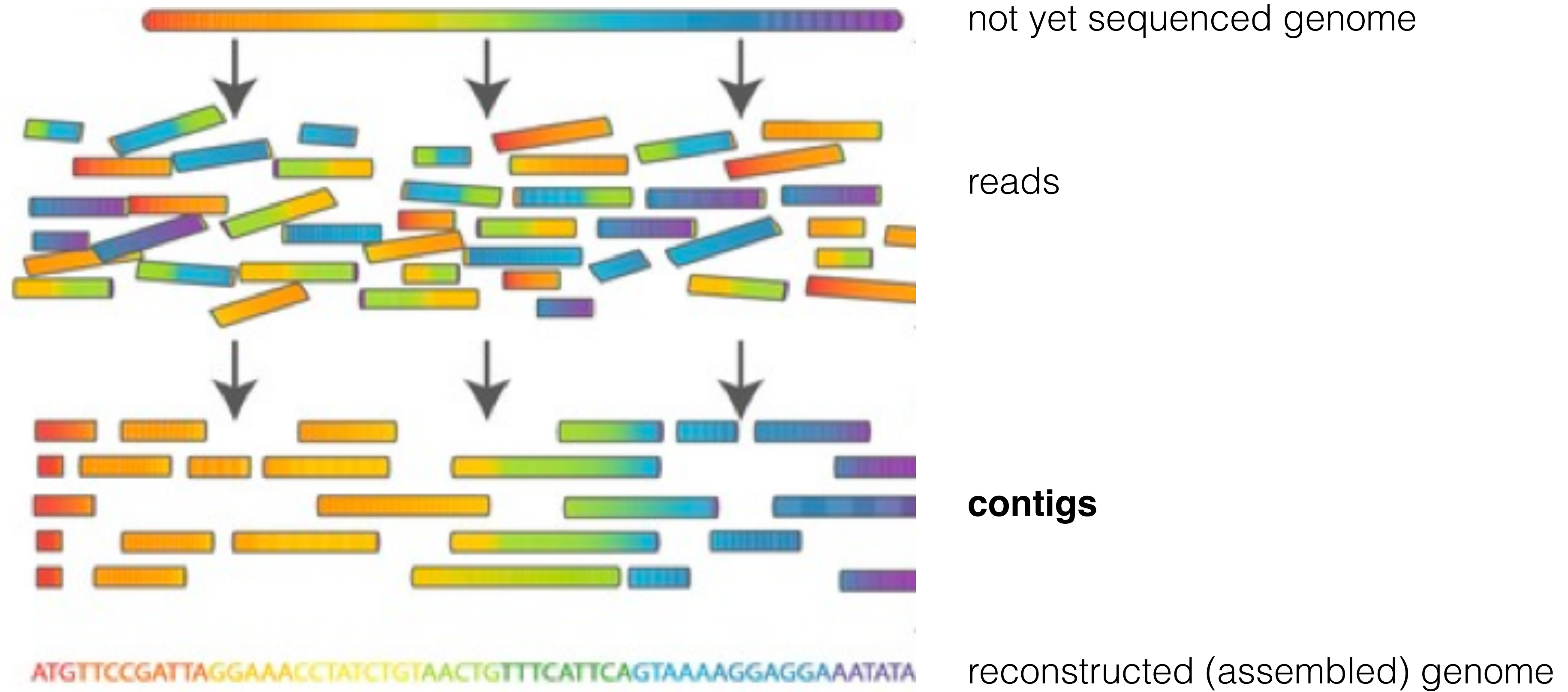
Mapping output

Does each read map?

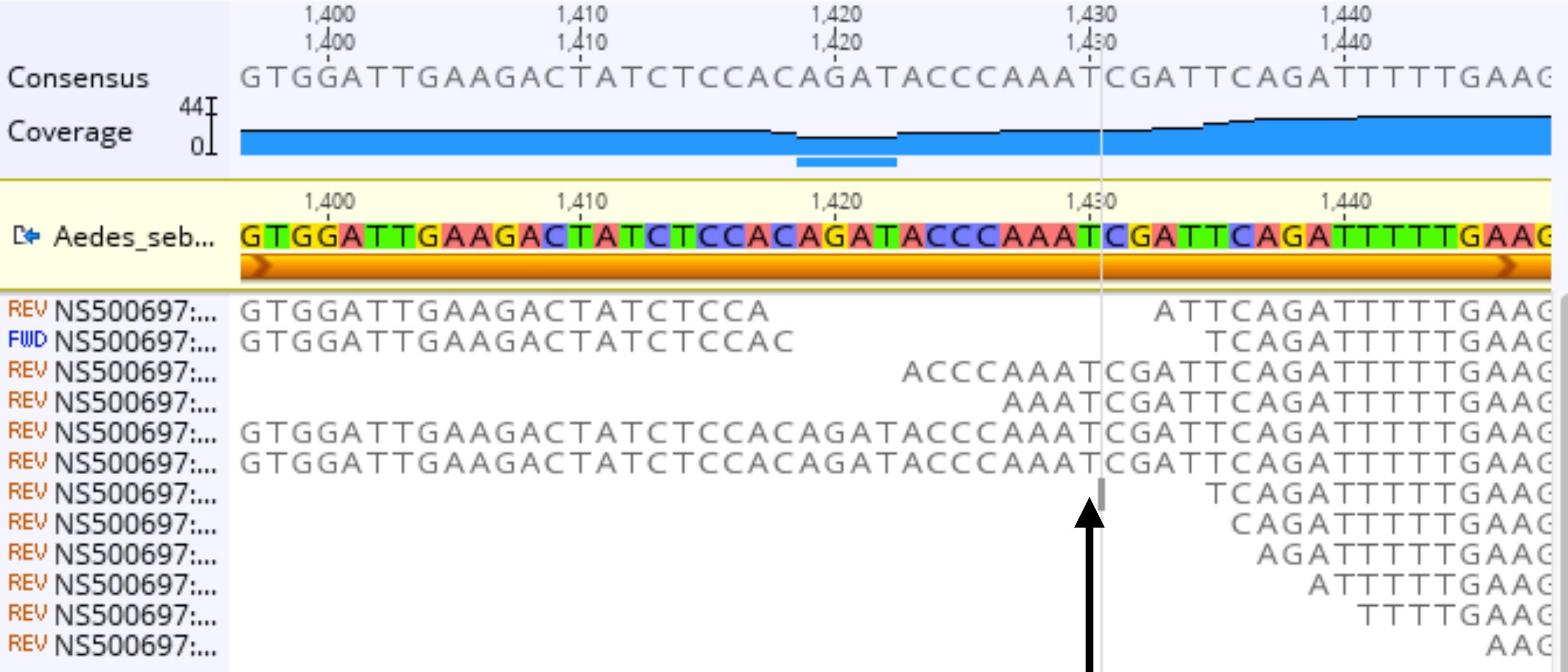
Where on the ref. seq. does it map?

How well does it map?

Genome assembly is the process of trying to reconstruct a genome sequence from reads (making a new reference sequence)



Coverage is the number of individual mapped reads that support a particular nucleotide in a reference sequence

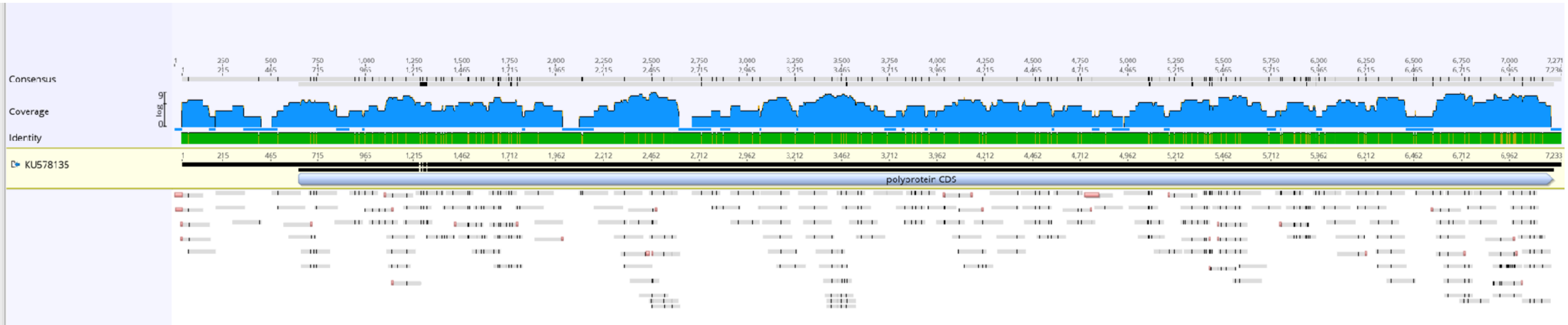


coverage is often referred to as 'depth' or 'depth of coverage'

This T at position 1430 in the reference sequence has 4x coverage

Coverage is also used to describe the fraction of a genome with >0 coverage depth

reads from human oral swab RNA aligned to a coxsackie virus genome



96% genome coverage (96% of bases have >0x coverage)

3.4x average coverage depth (range 0-9x)



hand foot and mouth disease



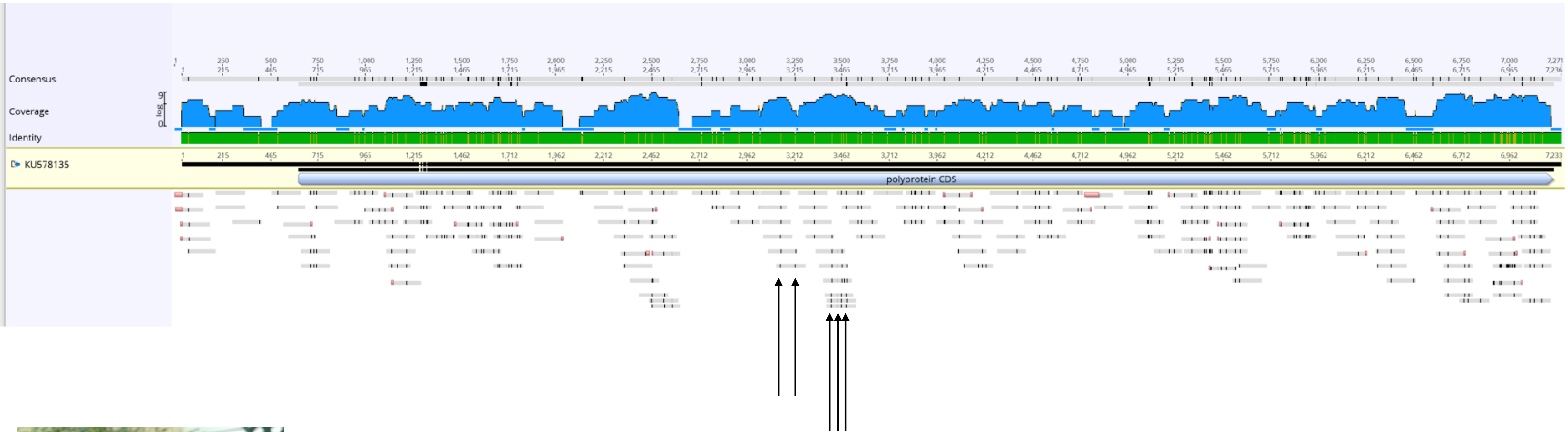
(Mayo clinic)

Applications of mapping

- **Quantification:** using reads for counting: sequence itself not important per se
 - RNA-seq: reads mapping to a particular transcript proportional to its abundance in the sample
 - ChIP-seq: coverage levels proportional to binding of proteins to DNA sequences
- **Variant identification**
 - Single nucleotide variants (SNVs aka SNPs)
 - Structural variants
 - Consensus-changing or sub-consensus
- **Remove sequences** of specific origins
 - Contaminating organisms
 - Plasmid
 - Organellar

There are variants in the reads relative to the co reference sequence:
these differences are the basis for 'variant calling'

reads from human oral swab RNA aligned to a coxsackie virus genome

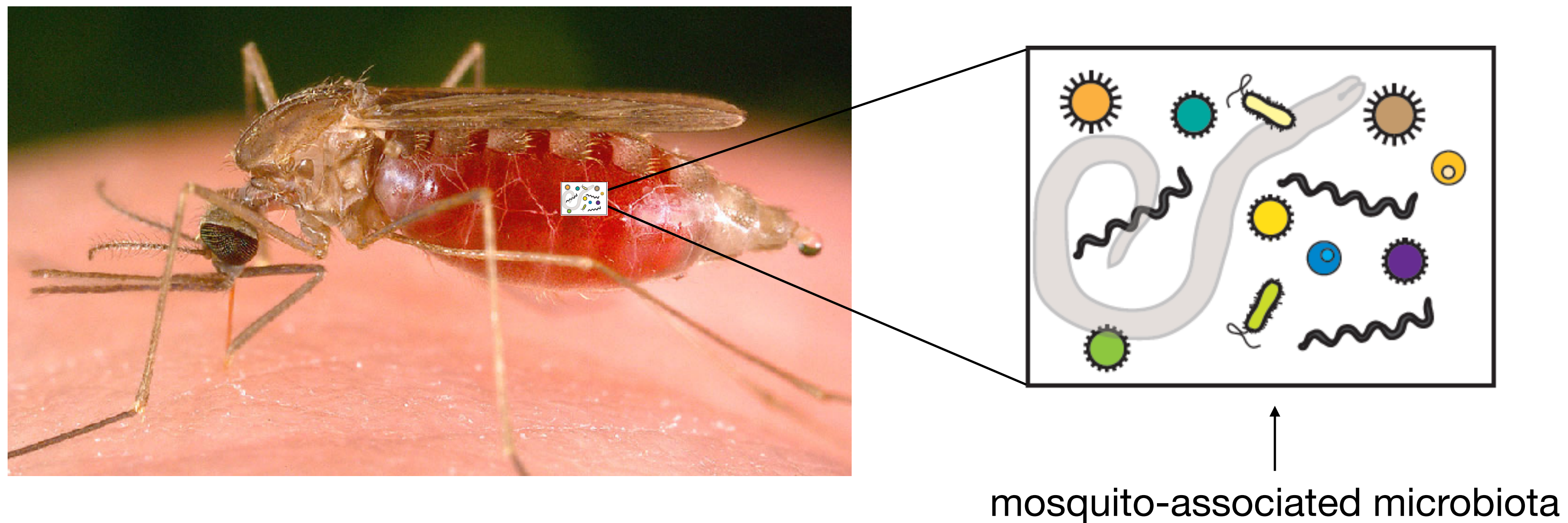


mismatches between reads and the reference sequence

hand foot and mouth disease



Mapping can be used to remove reads that derive from an organism you don't care about



You can sequence everything and use a mosquito reference genome to remove all the mosquito reads leaving reads from all the other non-mosquito organisms

Mapping Exercise

Work in pairs to map reads to the provided 'genome'



Questions to consider while doing this mapping exercise

Coverage

What was the (approximate) average coverage depth?

What was the maximum coverage depth?

What was the minimum coverage depth?

Was coverage across the 'genome' even?

What percent of the genome was covered by at least one read?

Mapping

Were all of your reads mappable?

Where did the unmappable reads come from?

In a real sequencing dataset, why might there be unmappable reads?

What fraction (approximately) of reads mapped unambiguously (uniquely)?

Did you identify any sequencing errors?

Did you identify any variants (SNPs)?

Speed

What was your mapping speed (how many reads per minute did you map)?

How do you think that speed compares to the speed of mapping software like bowtie or bwa?

Could you have mapped faster with more workers in your group?

We choose to go to the moon. We choose to go to the moon in this decade and do the other
 o the moon choose_to_o_tosthe_m n_in_this_ nd_do_the_
 go_to_the_ the_moon_ decade_an otzur_thi
 hoosx_to_g moon_we_ch _to_jo_to_ is_decade_ o_the_othg
 oose_to_go e_moon_we_ o_the_moon_in_this_d do_the_oth
 choose_to_ _moon_we_c e_to_go_to_ _in_this_dnc _and_do_th
 hoose_to_g o_the_moon_n_this_dnc _and_do_th
 o_tosthe_m _mo_go_to_ his_decade_ nd_do_the_
 go_to_the_ hoose_to_g
 hoose_to_g
 se_to_go_t

Unmapped reads:
 coronaviru
 us_vaccine

Uniquely mapped
 Ambiguously mapped

59 reads:
 • 57 mapped (97%):
 • 32 mapped uniquely (54%)
 • 2 unmapped (3%)

Coverage:
 • 57 mapped reads x 10 'bases' / read = 570 bases of data
 • 570 base / 150 base genome = 3.8x avg coverage

things, not because they are easy, but because they are hard
 ings_not_b y_are_easy ey_are_har
 yr_things_ ecause_the _easy_but_ cause_they
 r_things_n ey_are_eas use_they_a
 ngs_not_be they_are_e _but_becau they_are_h
 _things_no y_are_easy _bacaue_t
 easy_but_b
 y_are_easy ut_because
 re_easy_bu
 _easy_but_
 _easy_but_
 ↑
 10x coverage of this e
 0x coverage of this d

Mapping tools like bowtie2 map fast!

Bowtie2 mapping **1M** 50nt reads to the human genome (3B bp)
1 CPU (1 thread/1 core):

```
mdstengl@cctsi-104:~/analyses/test_human_mapping$ time ./run_bowtie ERR3252925_1_1M.fastq GCA_000001405.15_GRCh38_no_alt_an|
alysis_set.fna.bowtie_index
bowtie2 -x GCA_000001405.15_GRCh38_no_alt_analysis_set.fna.bowtie_index -q -U ERR3252925_1_1M.fastq --local --score-min C,1
00,1 --no-unal --threads 1 -S ERR3252925_1_1M.fastq.GCA_000001405.15_GRCh38_no_alt_analysis_set.fna.bowtie_index.sam

real    1m16.427s
user    1m13.836s
sys     0m12.844s
```

1 minute 16 seconds: 13,000 reads per second

Bowtie2 mapping **1M** 50nt reads to the human genome (3B bp)
24 CPUs (24 thread/24 core):

```
mdstengl@cctsi-104:~/analyses/test_human_mapping$ time ./run_bowtie_multiple_threads ERR3252925_1_1M.fastq GCA_000001405.15
_GRCh38_no_alt_analysis_set.fna.bowtie_index
bowtie2 -x GCA_000001405.15_GRCh38_no_alt_analysis_set.fna.bowtie_index -q -U ERR3252925_1_1M.fastq --local --score-min C,1
00,1 --no-unal --threads 24 -S ERR3252925_1_1M.fastq.GCA_000001405.15_GRCh38_no_alt_analysis_set.fna.bowtie_index.sam

real    0m9.641s
user    1m38.696s
sys     0m33.124s
```

9.6 seconds: ~100,000 reads per second

Like airport security, computers can run tasks in parallel to make jobs go faster



Jobs awaiting processing

Done being processed 😊

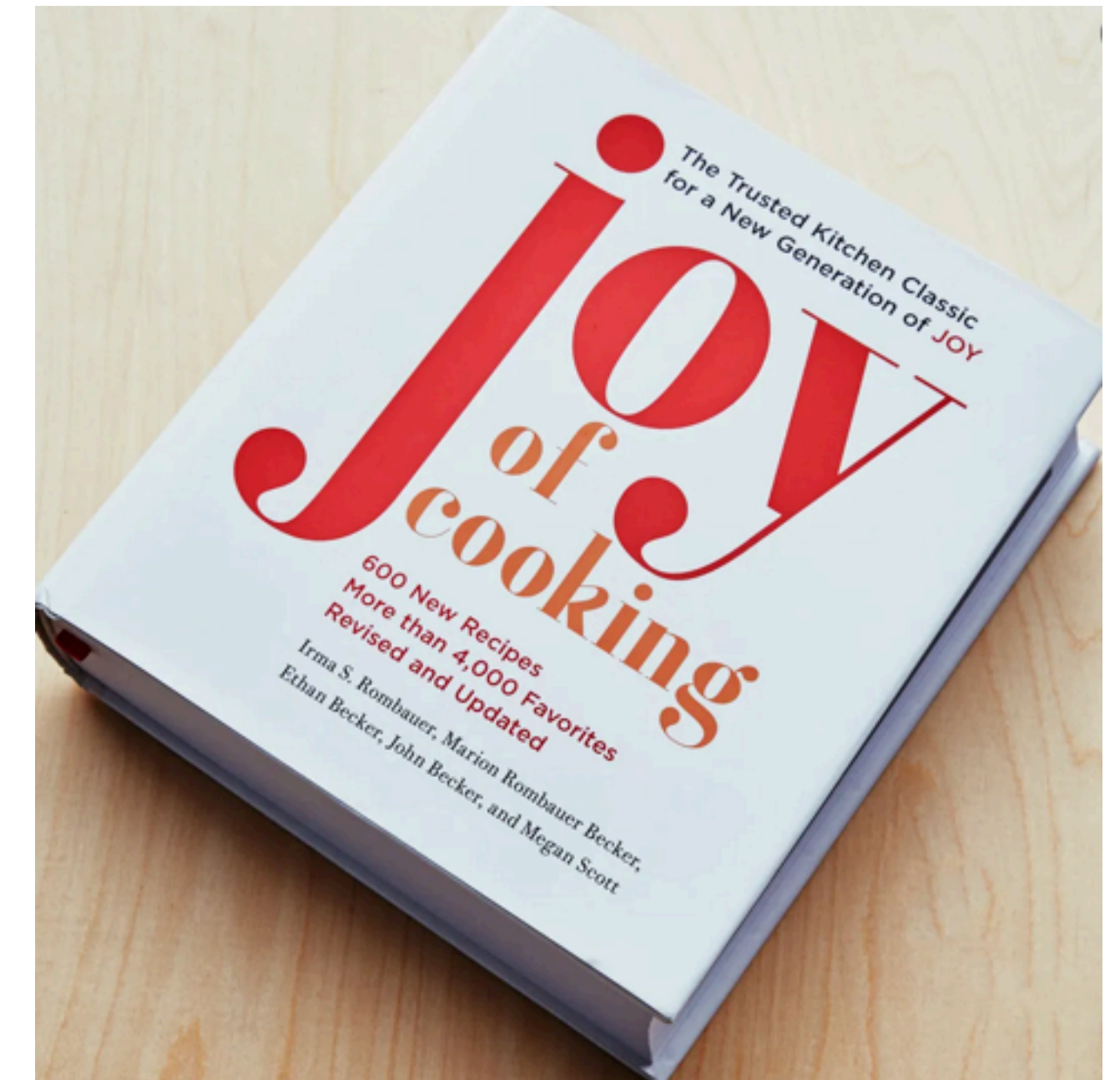
Mapping tools map fast because they **pre-index** reference sequences

cookies

- almond macaroons, 771
- almond pretzels (mandelplättchen), 774
- angel slices, 764–765
- baking, 760–761
- bar. see bars
- bar and square, about, 762
- biscotti, 774
- book club brownies, 762
- brandy snaps, 781
- brownies Cockaigne, 762
- brown sugar sand tarts, 773
- butterscotch brownies or blondies, 762
- butterscotch icebox, 776
- butterscotch nut, 771
- cheesecake brownies, 763
- chocolate chip, 766–767
- chocolate chip icebox, 776
- chocolate coconut macaroons, 771
- chocolate icebox, 776
- chocolate shortbread, 775–776
- Christmas, 762
- cinnamon stars, 774–775
- cocoa meringue kisses, 771

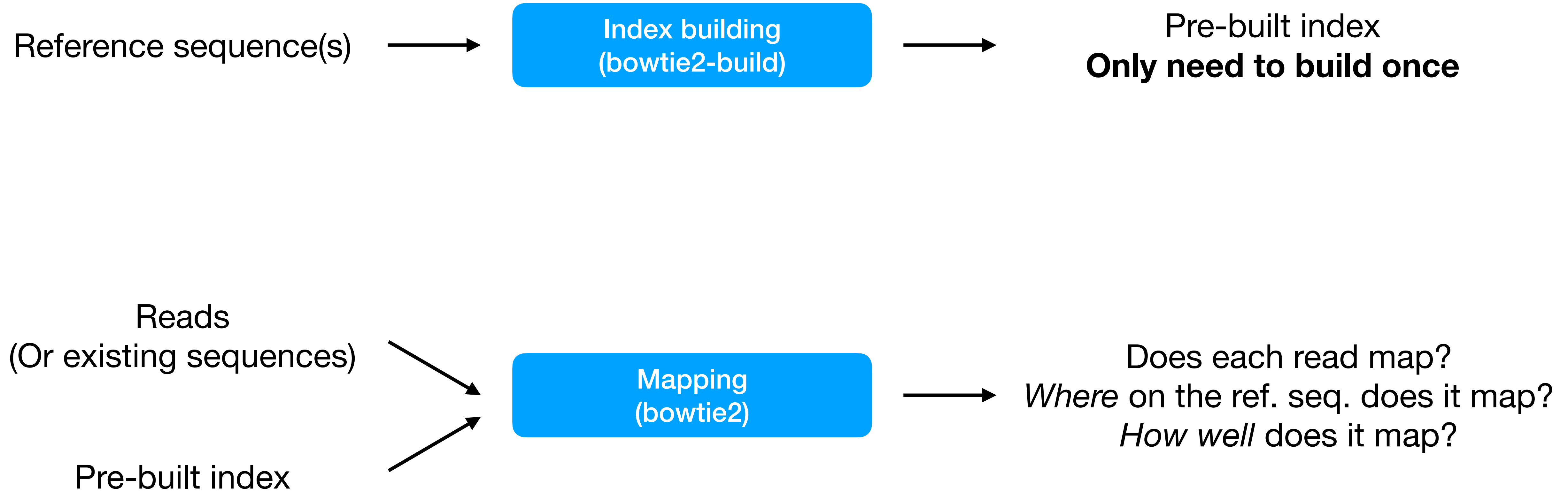
Indexes help you find things faster

chocolate chip
cookies on page
766



BLAST databases are another example of pre-built indexes

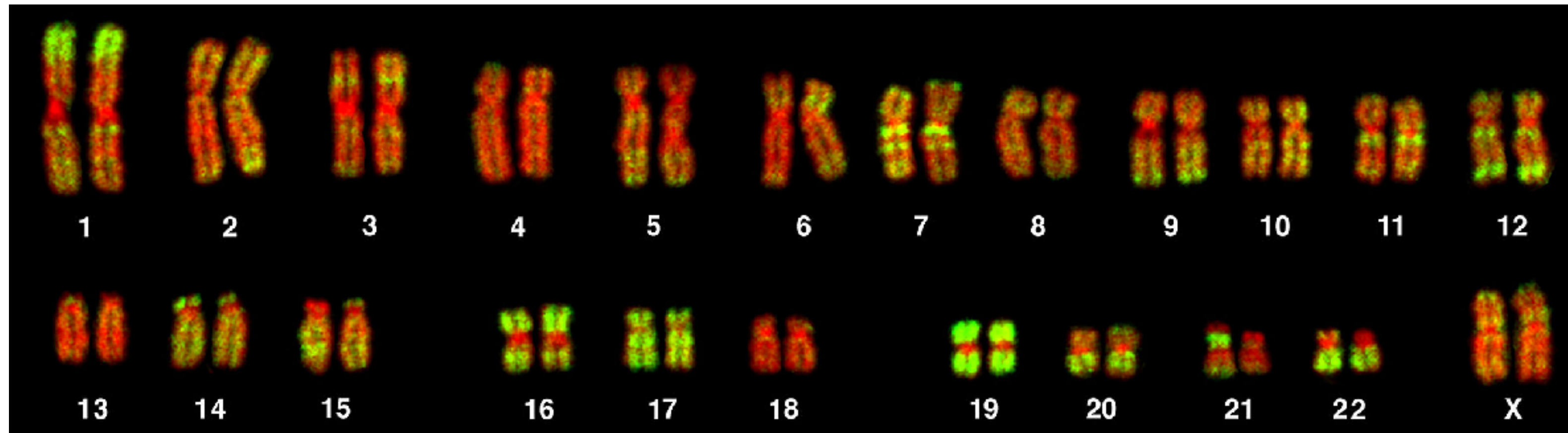
Mapping software includes tools to build the index



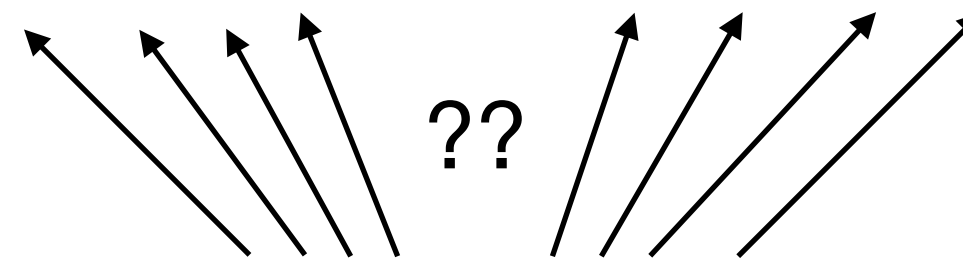
Reads might map to more than one location in a genome

This is especially a problem for reads that derive from repeat elements

Alu sequences in the human genome
1 million copies, ~10% of the mass



Bolzer et al (2005) PLoS Biol



If you had a read derived from an Alu sequence,
which of these million copies should you map it to?

Mapping tools like bowtie2 have options that define how they will deal with ambiguously mapping reads and provide information about whether a read mapped uniquely or not

Mapping quality measures whether a read maps uniquely or not

Mapping quality: higher = more unique

The aligner cannot always assign a read to its point of origin with high confidence. For instance, a read that originated inside a repeat element might align equally well to many occurrences of the element throughout the genome, leaving the aligner with no basis for preferring one over the others.

Aligners characterize their degree of confidence in the point of origin by reporting a mapping quality: a non-negative integer $Q = -10 \log_{10} p$, where p is an estimate of the probability that the alignment does not correspond to the read's true point of origin. Mapping quality is sometimes abbreviated MAPQ, and is recorded in the `SAM` MAPQ field.

Mapping quality is related to "uniqueness." We say an alignment is unique if it has a much higher alignment score than all the other possible alignments. The bigger the gap between the best alignment's score and the second-best alignment's score, the more unique the best alignment, and the higher its mapping quality should be.

Accurate mapping qualities are useful for downstream tools like variant callers. For instance, a variant caller might choose to ignore evidence from alignments with mapping quality less than, say, 10. A mapping quality of 10 or less indicates that there is at least a 1 in 10 chance that the read truly originated elsewhere.

From: <http://bowtie-bio.sourceforge.net/bowtie2/manual.shtml>

Mapping quality scores are like basecall quality scores in FASTQ files

$$\text{Quality score} = -10 \log_{10} (p)$$

Basecall Q score = $-10 \log_{10}$ (probability basecall is incorrect)

Mapping Q score = $-10 \log_{10}$ (probability that the read is not mapped to its true location)

Q score	P
10	0.1 = 1/10
20	0.01 = 1/100
30	0.001 = 1/1,000
40	0.0001 = 1/10,000

bowtie2 mapped reads from *Drosophila melanogaster* to the *D. melanogaster* reference genome

```
bowtie2 -x /home/databases/fly/fly_genome -q -U dros_pool_R1_fu.fastq --local --score-min C,120,1 --no-unal --time --al dros_pool_R1_fu.
astq.fly_genome.hits.fastq --threads 12 -S dros_pool_R1_fu.fastq.fly_genome.sam
Time loading reference: 00:00:00
Time loading forward index: 00:00:01
Time loading mirror index: 00:00:00
Multiseed full-index search: 00:00:08
186708 reads; of these:
  186708 (100.00%) were unpaired; of these:
    20301 (10.87%) aligned 0 times
    87911 (47.08%) aligned exactly 1 time
    78496 (42.04%) aligned >1 times
89.13% overall alignment rate
Time searching: 00:00:09
Overall time: 00:00:09
```

42% of reads mapped non-uniquely

47% of reads mapped uniquely

10% of reads didn't map, what are these?

SAM files are used as input to downstream tools that use mapping data

