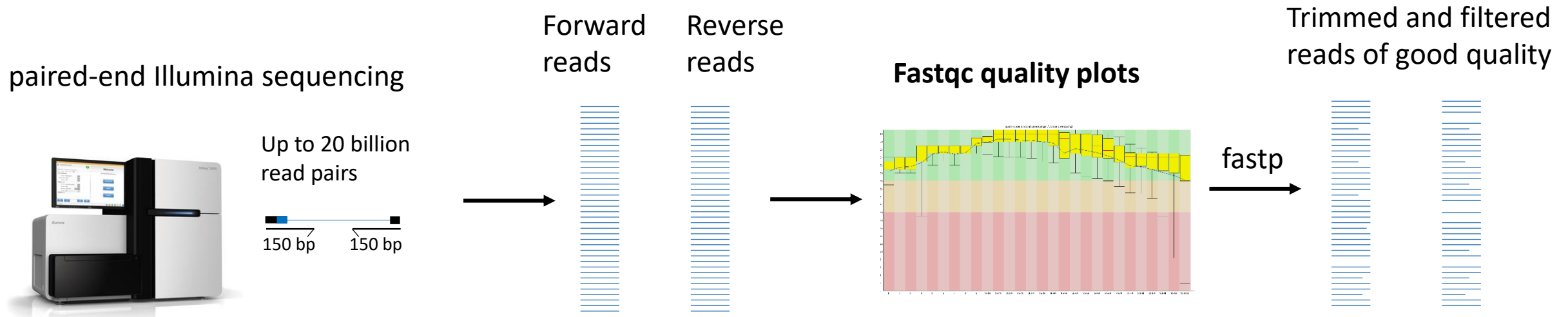


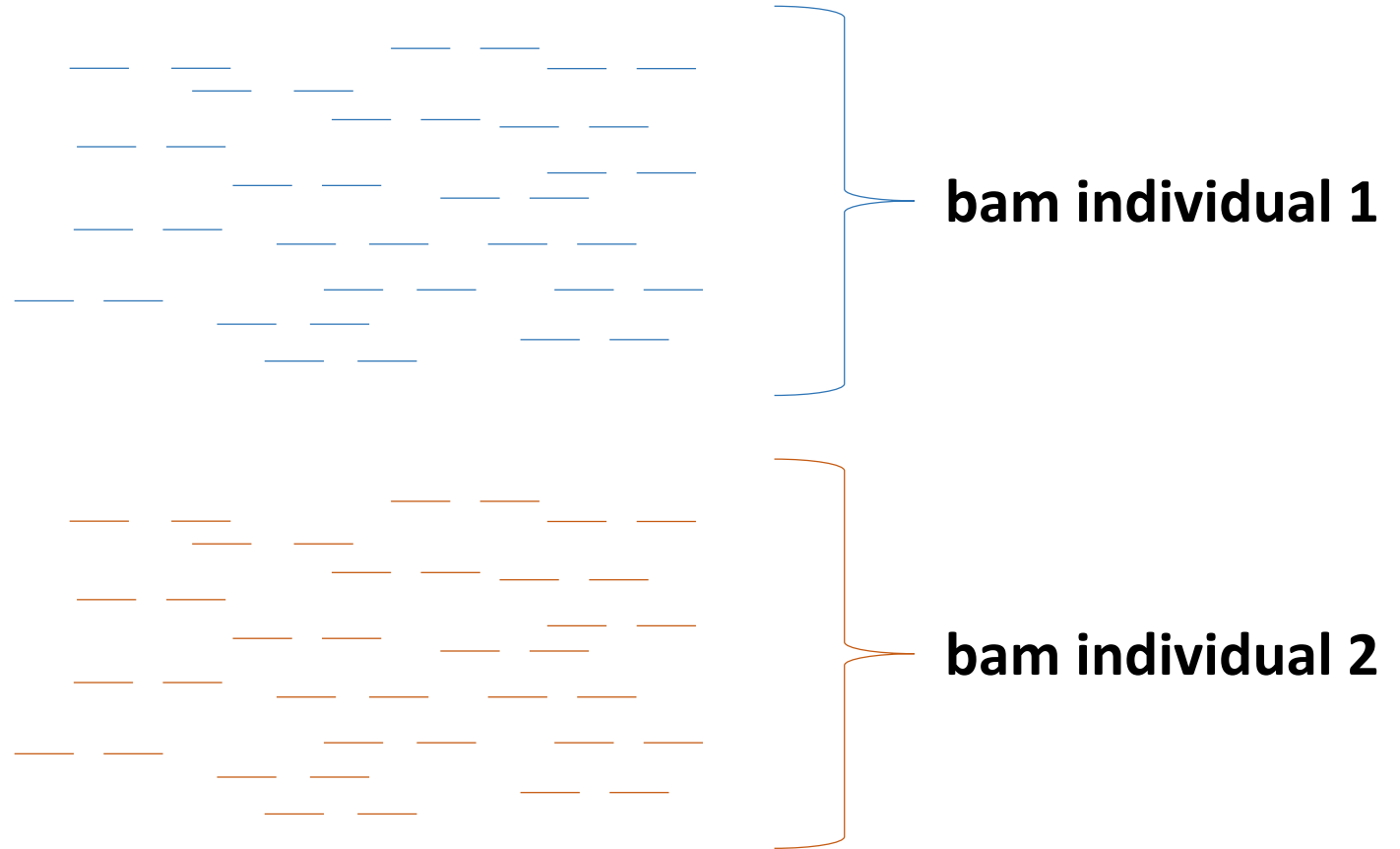
Mapping to a reference
genome

1. Quality check and trimming raw reads



2. Alignment to the reference genome

reference



Tools to align short read data

reference

read

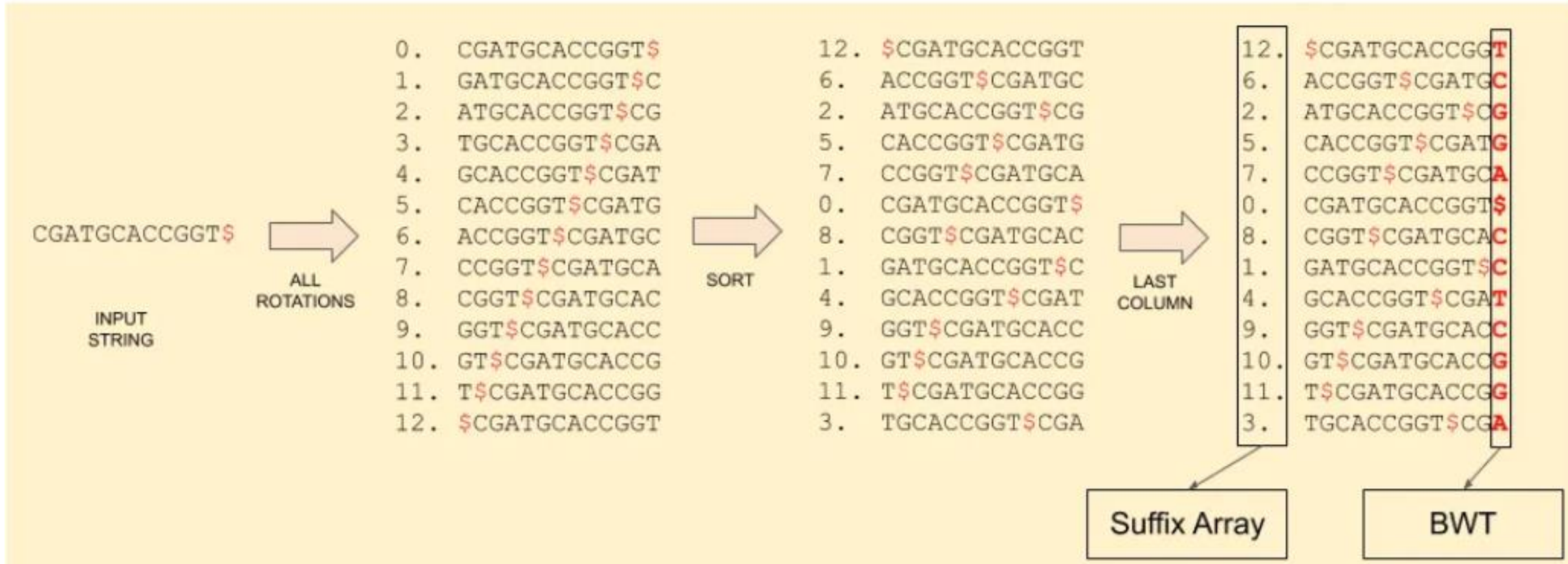


bwa-mem2 (Burroughs-Wheeler alignment – maximum exact matches)

Bowtie2

Burroughs Wheeler Transform

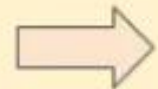
(also used in file compression, e.g. bzip, figures from <https://medium.com/@mr-easy>)



Burroughs Wheeler Alignment

CGATGCACCGGT\$

Read: GCA



LAST
COLUMN

12.	\$	C	G	A	T	G	C	A	C	C	G	G	T
6.	A	C	C	G	G	T	\$	C	G	A	T	G	C
2.	A	T	G	C	A	C	C	G	G	T	\$	C	G
5.	C	A	C	C	G	G	T	\$	C	G	A	T	G
7.	C	C	G	G	T	\$	C	G	A	T	G	C	A
0.	C	G	A	T	G	C	A	C	C	G	G	T	\$
8.	C	G	G	T	\$	C	G	A	T	G	C	A	C
1.	G	A	T	G	C	A	C	C	G	G	T	\$	C
4.	G	C	A	C	C	G	G	T	\$	C	G	A	T
9.	G	G	T	\$	C	G	A	T	G	C	A	C	C
10.	G	T	\$	C	G	A	T	G	C	A	C	C	G
11.	T	\$	C	G	A	T	G	C	A	C	C	G	G
3.	T	G	C	A	C	C	G	G	T	\$	C	G	A

Suffix Array

BWT

Read: GCA

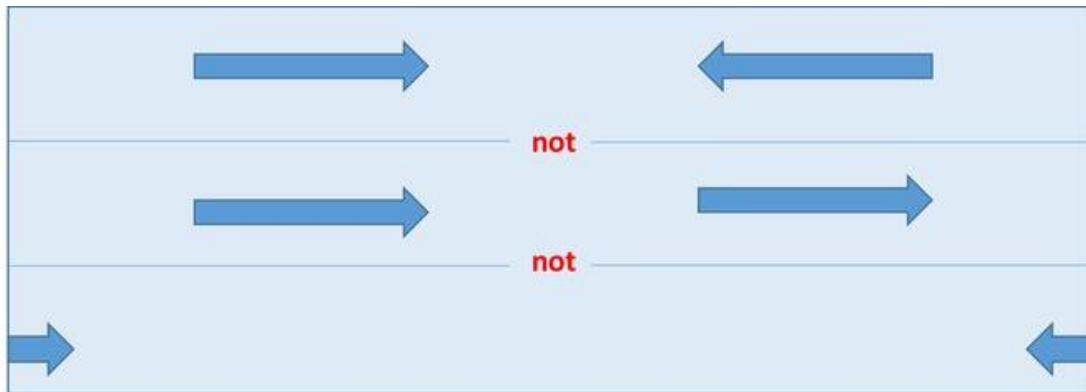
Read: GCA

Read: GCA

12.	0.	\$.	.	.	T	0.	\$.	.	.	T	0.	\$.	.	.	T
6.	1.	A	.	.	.	C	1.	A	.	.	.	C	1.	A	.	.	.	C
2.	2.	A	.	.	.	G	2.	A	.	.	.	G	2.	A	.	.	.	G
5.	3.	C	.	.	.	G	3.	CA	.	.	.	G	3.	CA	.	.	.	G
7.	4.	C	.	.	.	A	4.	C	.	.	.	A	4.	C	.	.	.	A
0.	5.	C	.	.	.	\$	5.	C	.	.	.	\$	5.	C	.	.	.	\$
8.	6.	C	.	.	.	C	6.	C	.	.	.	C	6.	C	.	.	.	C
1.	7.	G	.	.	.	C	7.	G	.	.	.	C	7.	G	.	.	.	C
4.	8.	G	.	.	.	T	8.	G	.	.	.	T	8.	GCA	.	.	.	T
9.	9.	G	.	.	.	C	9.	G	.	.	.	C	9.	G	.	.	.	C
10.	10.	G	.	.	.	G	10.	G	.	.	.	G	10.	G	.	.	.	G
11.	11.	T	.	.	.	G	11.	T	.	.	.	G	11.	T	.	.	.	G
3.	12.	T	.	.	.	A	12.	T	.	.	.	A	12.	T	.	.	.	A

Aligning reads is complicated by:

- Imperfect match to the reference due to mutations or sequencing errors, or errors in the reference genome
- Multiple positions where the read could match (repeated regions)
- Low quality of the read
- With paired-end reads: only one read maps or the other read maps on a different chromosome



Alignment tools such as bwa-mem2 are able to handle all of these complications and give information on the mapping quality.

Regions that are repeated in the genome make it very difficult to map reads

Reference genome



Read to be aligned to the reference



Regions that are repeated in the genome make it very difficult to map reads

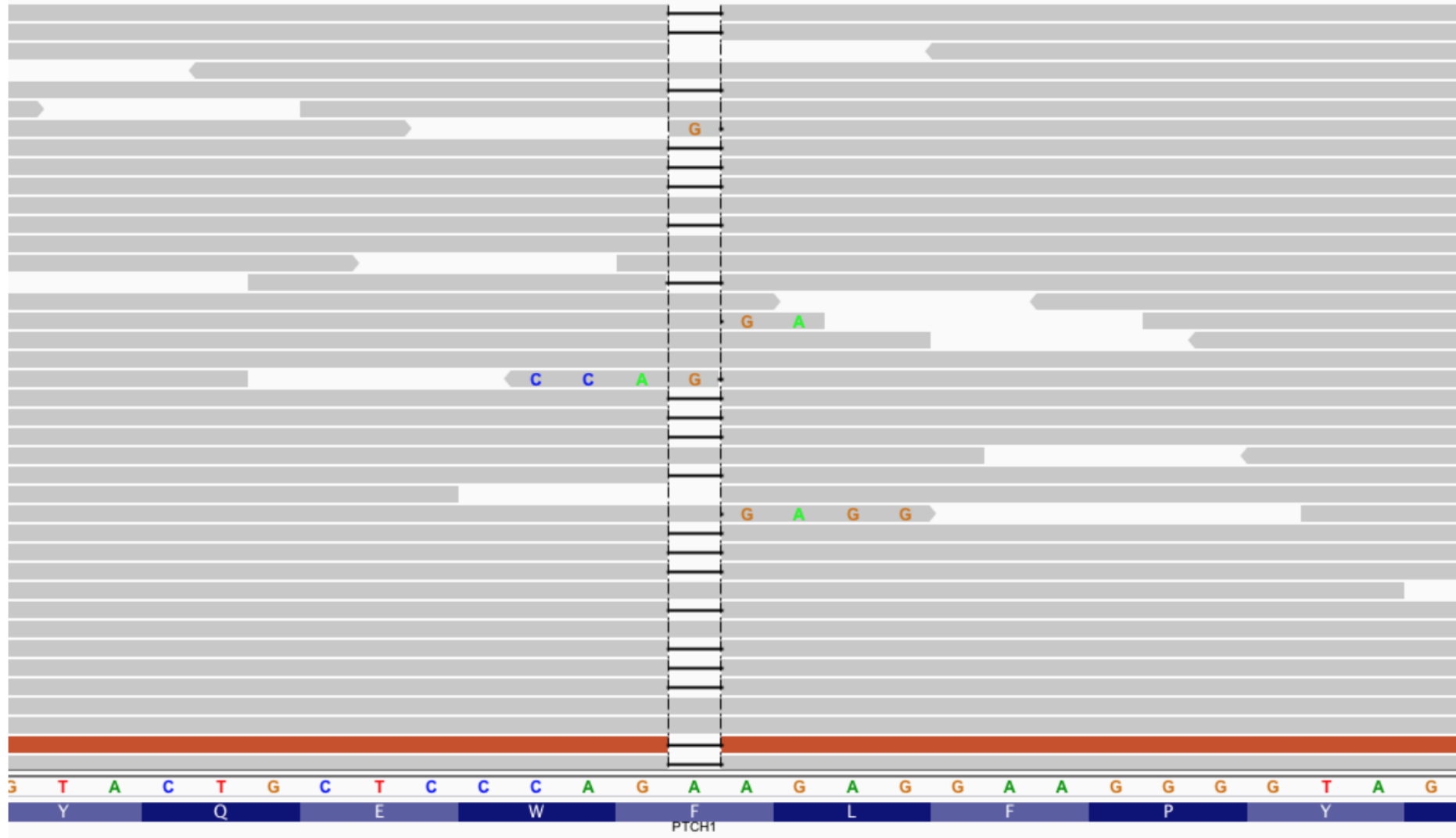
Reference genome



Read to be aligned to the reference

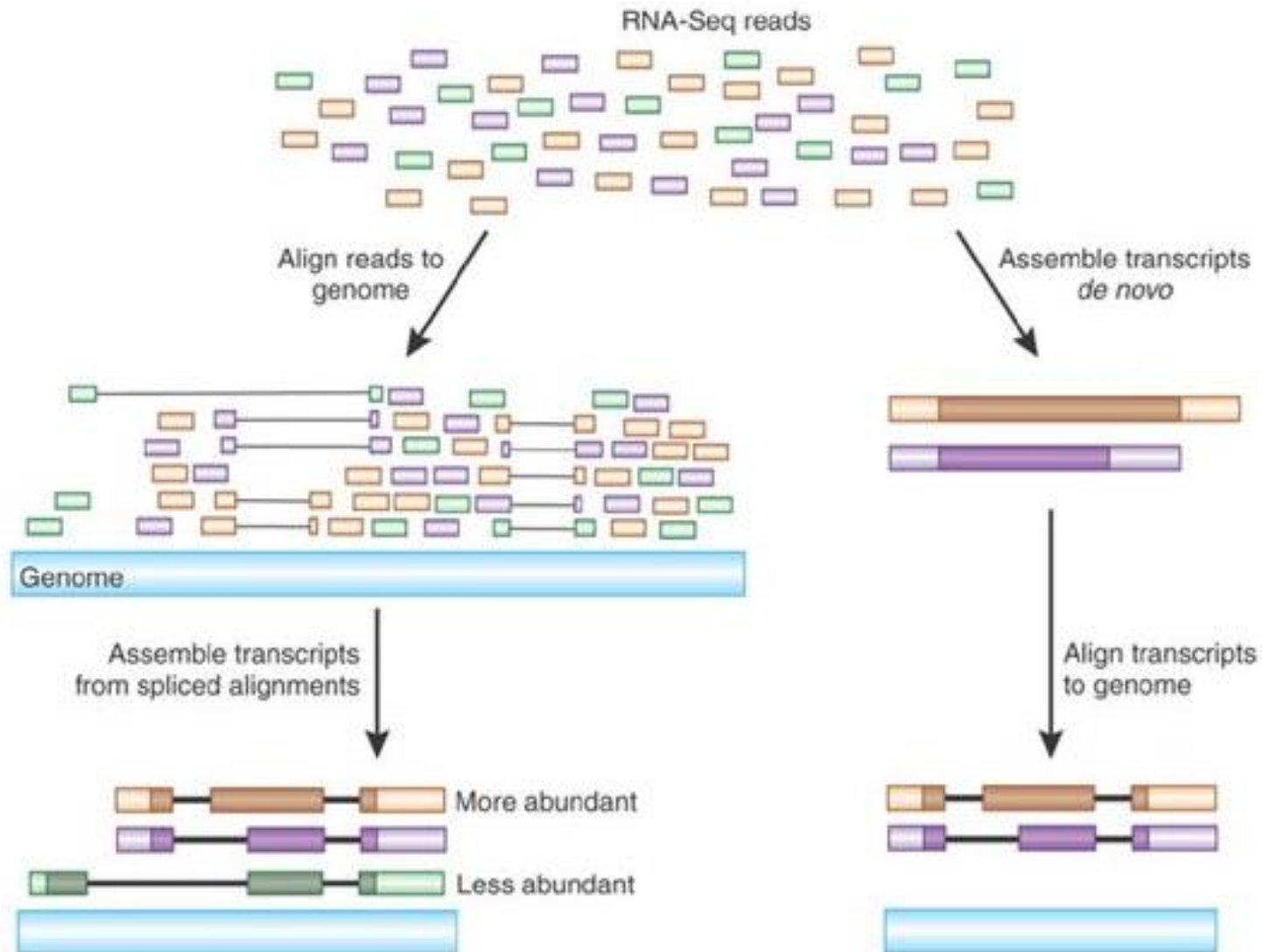


Indels: deletions and insertions



Aligning other type of data: RNA data

Tools:
e.g.
STAR
HISAT2



PacBio mapping

Karin Näsvall



Sequel IIe Instrument

Generate Sequencing Data



Data Transfer/SMRT Link Server



SMRT Link GUI

Access using a web browser

Sequel IIe System output files

```
<your_specified_output_directory>/r64009_20200825_221039/1_A01/  
|-- m64009_200825_222052.baz2bam_1.log  
|-- m64009_200825_222052.ccs.log  
|-- m64009_200825_222052.ccs_reports.json  
|-- m64009_200825_222052.ccs_reports.txt  
|-- m64009_200825_222052.consensusreadset.xml  
|-- m64009_200825_222052.reads.bam  
|-- m64009_200825_222052.reads.bam.pbi  
|-- m64009_200825_222052.sts.xml  
|-- m64009_200825_222052.transferdone  
|-- m64009_200825_222052.zmw_metrics.json.gz
```

Automatic HiFi reads generation (Export Reads)

```
hifi_reads.fastq.gz - FASTQ file containing HiFi Reads  
hifi_reads.fasta.gz - FASTA file containing HiFi Reads  
hifi_reads.bam - BAM file containing HiFi Reads
```

HiFi reads QV > 20

<https://pacbiofileformats.readthedocs.io/en/13.0/>

PacBio

<https://pacbiofileformats.readthedocs.io/en/13.0/BAM.html#hifi-reads>

- `$ samtools view -H
/pacbio/m84093_240426_124306_s1.ccs.bc2033.rmdup.bam
| head -n2`
- `@HD VN:1.6 S0:coordinate pb:5.0.0`
- `@RG ID:f56a67e5/0--0 PL:PACBIO
DS:READTYPE=CCS;Ipd:CodecV1=ip;PulseWidth:CodecV1=pw;B
INDINGKIT=102-739-100;SEQUENCINGKIT=102-118-
800;BASECALLERVERSION=5.0;FRAMERATEHZ=100.000000;Barco
deFile=/lustre/scratch123/tol/resources/barcodes/PacBi
o_ULI_adapter.fasta;BarcodeHash=cf95303a081e62fbaedf29
888b16fdb7;BarcodeCount=1;BarcodeMode=Symmetric;Barcod
eQuality=Score LB:TRAC-2-8589
PU:m84093_240426_124306_s1 SM:Meier Genomes13637824
PM:REVIO BC:AAGCAGTGGTATCAACGCAGAGTACT CM:R/P1-C1/5.0-
25M`

PacBio

<https://pacbiofileformats.readthedocs.io/en/13.0/BAM.html#hifi-reads>

- [illegible]

PacBio

```
cat m84093_240426_124306_s1.ccs.bc2033.stats
```

```
-
```

```
A = 7108650862 (34.9%), C = 3058767823 (15.0%), G = 3061827434 (15.0%),  
T = 7127816130 (35.0%), CpG = 1003156910 (4.9%)  
sum = 20357062249, n = 2218131, mean = 9177.57438537219, largest =  
30301, smallest = 135
```


PacBio mapping – Minimap2

Minimap2 <https://github.com/lh3/minimap2>

Uses **minimizers** – short sequences of length $-k$ [default 15 bases]

1. Extracts minimizers from the reference (target) and index them
2. Match each minimizer in the query sequence against the reference set of minimizers
3. Sorts the position of each minimizer after position
4. Make a chain of the minimizers
5. Repeat for all query sequences

```
Minimap2 -ax map-hifi target.fa query.fa > output.sam
```

```
Align PacBio high-fidelity (HiFi) reads to a reference genome (-k19 -w19 -U50,500 -g10k -A1 -B4 -O6,26 -E2,1 -s200)
```

PacBio

- Preprocessing
 - Filtering adapters (blastn)
- Alignment
 - MINIMAP2
- Alignment post-processing
 - Statistics



<https://pipelines.tol.sanger.ac.uk/readmapping>