

15/01/2024

Next-Generation Sequencing.

- Roche 454 (2006-2013)
- Illumina → most popular (90%)
- PacBio SMRT → single-molecule real-time
- Ion Torrent
- Oxford Nanopore → latest technology
- Data Analysis
- Genomics
- Transcriptomics
- Assembly Algorithms

Sanger sequencing method:

di-deoxy termination method.



↓
capillary electrophoresis

A schematic diagram showing a series of vertical lanes, each containing a series of short horizontal dashes. This represents the output of a capillary electrophoresis instrument, where different sized DNA fragments migrate at different speeds through a gel matrix.

- ✓ one DNA fragment at a time → Drawbacks
- ✓ we get around 1 kb (10^3 bp) → Drawbacks
- ✓ took 13 years to sequence the entire human genome!

Human Genome Project.

Shotgun sequencing

NGS Technologies:

- ✓ direct sequencing
- ✓ high-throughput, accurate → ~96 fragments together
- ✓ cost-effective ~\$1000 & can go down

Atlas Project

expression of genes
changes during
diseases

dataset to correlate
genetics & drugs

→ The Cancer Genome Atlas
(TCGA).

Precision
Medicine

useful in patient-
specific therapy.

Reads:

- ✓ output sequences from sequencing platform
- ✓ of different lengths

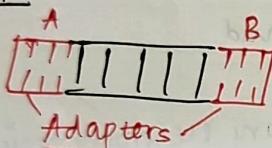
Roche 454 Sequencing:

Pyrosequencing

||||| DNA

↓ fragmentation

|||



synthetically-designed DNA sequences of known sequences

primer-binding sites

as we do not know the seq. of fragment yet

Index sequences
DNA barcode

to identify each sample differently

Multiplexing

mixing multiple samples & sequencing them together

Pyrosequencing:

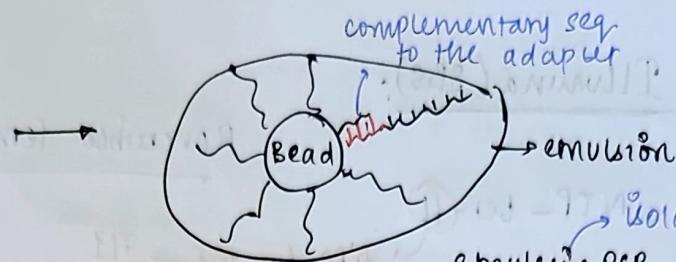
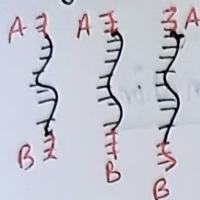
relies on synthesizing the complementary strand

→ A

every time there is an addition of a base, a light signal is emitted

if two same bases simultaneously → double the intensity

fragment

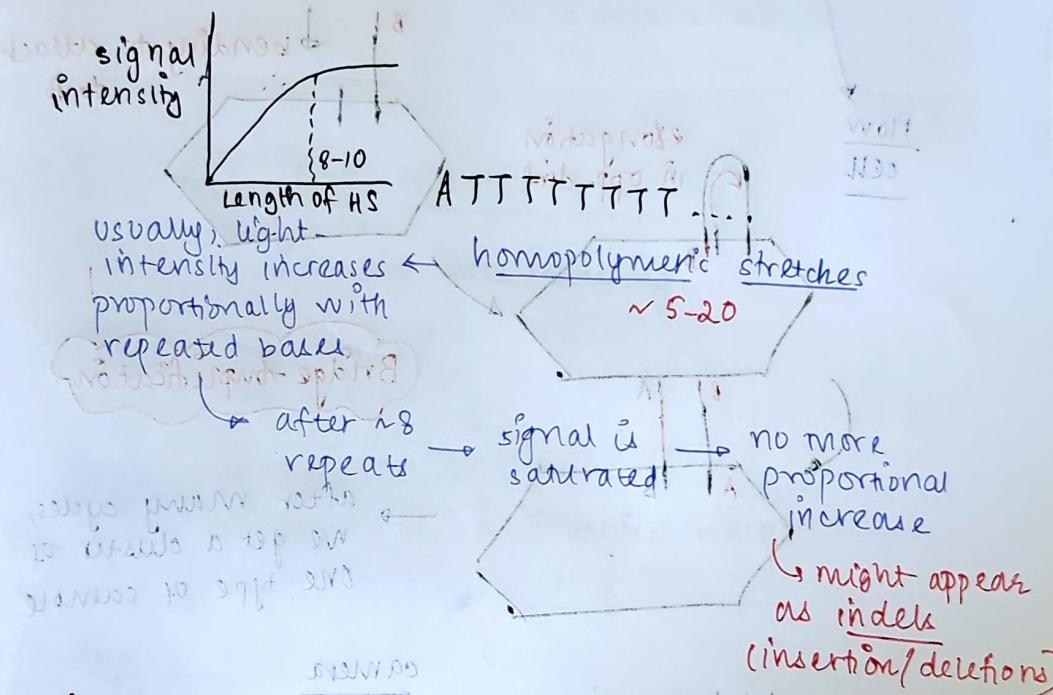


picotitre plates → 100,000 fragments

put in sequencer

Drawbacks

- ✓ many enzymes & substrates required
- ✓ expensive compared to Illumina, etc.



Illumina - Sequencing by Synthesis (SBS)

Reversible termination

(Sanger has irreversible termination)

modified dNTP

dNTP-ter-fl

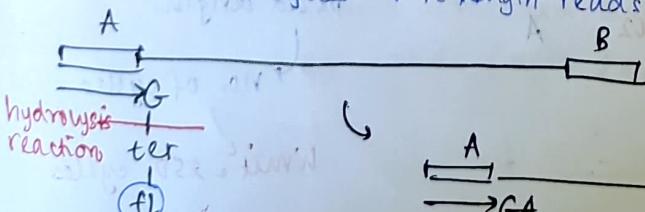
fluorophore

terminator

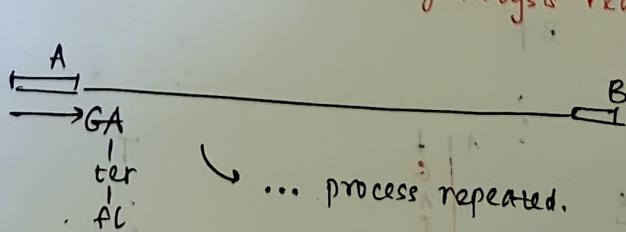
can be removed by a hydrolysis reaction

one base at a time

n cycles → n length reads



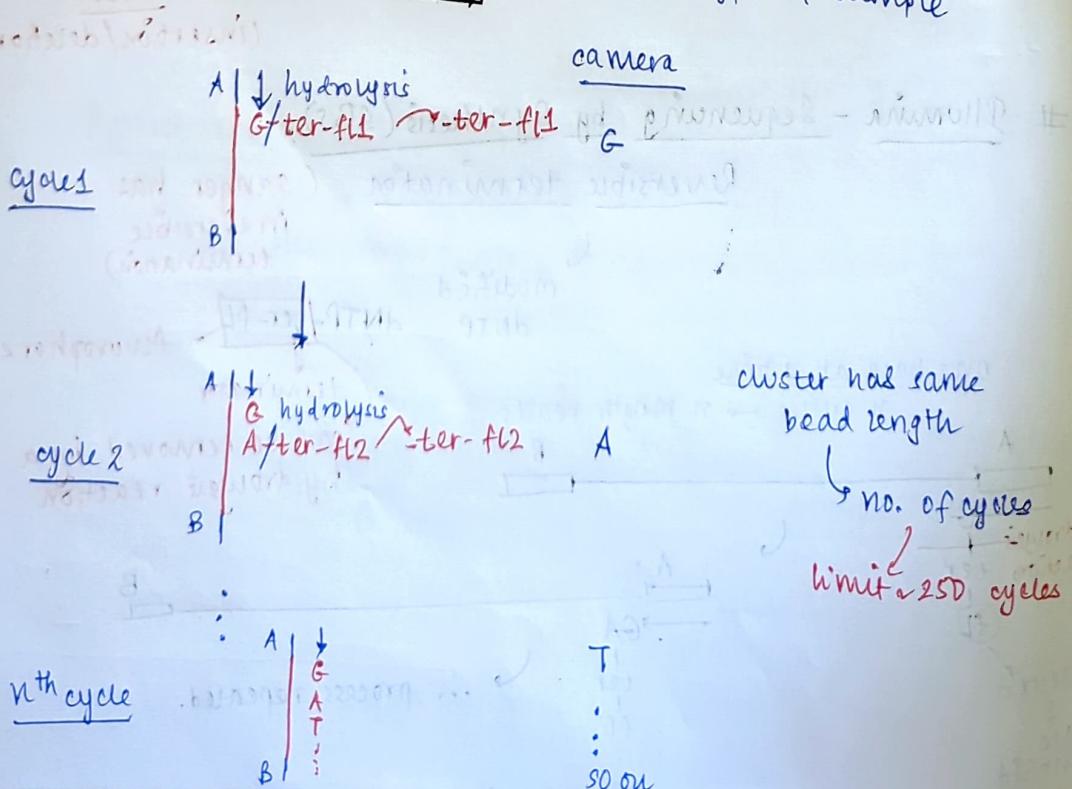
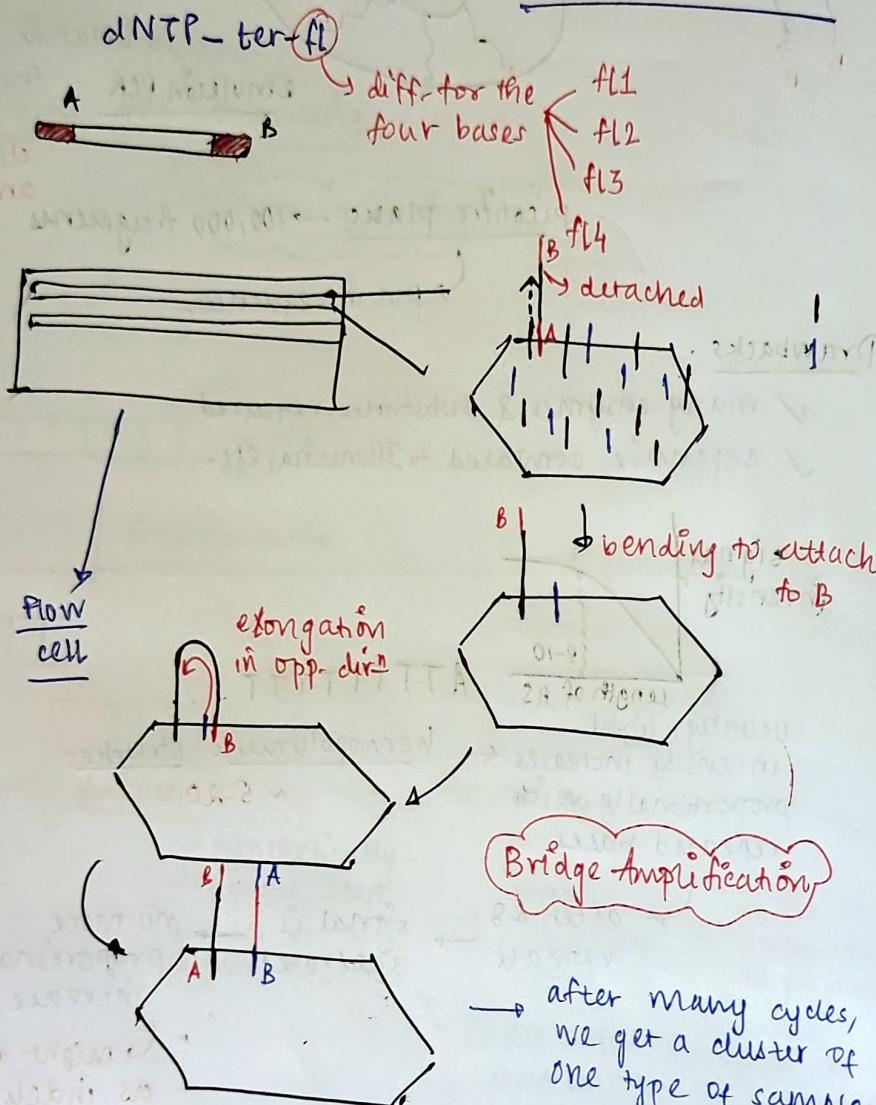
detected
base identified



29/1/24

Illumina (SBS):

Reversible Termination



in a cluster → n100 molecules
that are identical → detector reads average signal from a cluster
so after 250 signals/base,
these avg. signals can be confusing & hard to decipher

in the adaptor seq.

we add index sequences

eg: ATTGCC

helps to identify sequences in a mixed sample

Multiplexing

single-end (SE) and paired-end (PE) sequencing:

Illumina Data Processing.

↑ Throughput: high reads/
lot of data

150 bp: SE seq.

2x150 bp: PE seq. (use in Illumina video).

Advantages:

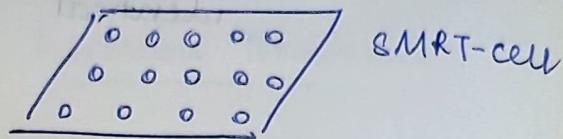
- cost-effective (not many enzymes are required).
- polymeric stretch → seq. one by one, not all at a time like Roche 454.
- highly accurate.

Drawbacks:

- expensive compared to the newer NGS even a little costs in thousands.
- short reads (upto 2x250 bp)

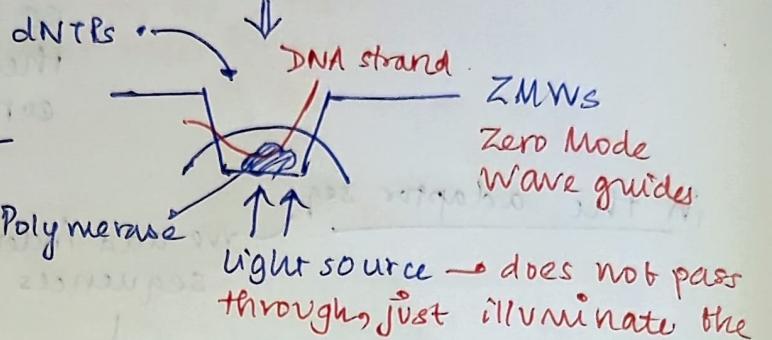
Single Molecule Real Time Sequencing (SMRT): → [SBS]

sequencing happens inside ZMWS.



everything else is similar to Illumina ↪

immobilized at the bottom



→ no amplification step is required.

→ no cluster formation

→ no averaging of signals → that's why 'single molecule'

more sensitive detected

signal detected only at illumination

so incoming dNTPs do not produce signal

Advantages:

- long reads
- avg. read length of 2-3 kb, some 10-20 kb
- single-molecule resolution

Drawbacks:

- high error rate (~14%)

Circular Consensus Sequence (CCS) reads

→ 99% accuracy

1. ATGCCCTA TT

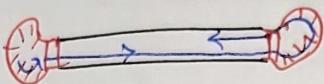
2. ATCCCTA TT

3. ATGCCAATT

4. ATGCATATT

take the most freq. at a place

ATG CCTATT } Consensus sequence



multiple repeats

↓
single base

circular synthesis

polymerase → always illuminated by light → leads to adverse effect on fidelity

↳ high energy state

short reads but

→ high accuracy

use of Illumina & SMRT together

Hybrid Sequencing

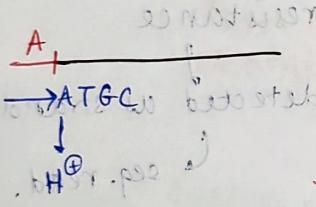
↳ large reads but

low accuracy.

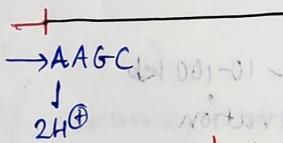
5/2/24

Ion Torrent Sequencing : → [SBS],

whenever there is a base addition, a hydrogen ion is released.

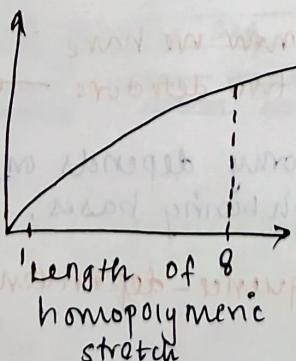


sensor



- add one base at a time. ↗ just like 454
- add a dNTP → check change in H⁺ ion concentration.
- same drawback as 454.

this drawback is related to the sensor



eg: A → T → G → C

↓
if an A is present in the template, a change in pH is observed at this stage.
washed ↗ repeated.

Advantages:

- ✓ read length 400-700 bp
- ✓ accuracy ~ 99%
- ✓ no expensive optical detection unit required
- ✓ no modified dNTP required.

Drawbacks:

- ✓ lower throughput compared to Illumina

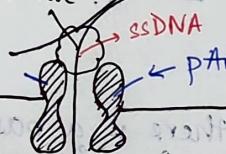
Oxford Nanopore Technology (ONT):

↳ Direct DNA Sequencing (not SBS).

- no dNTPs
- no polymerase

} → very economical.

hectase



each base has different resistance

detected as strand passes

seq. read.

Advantages:

- ✓ long read seq. ~ 10-100 kb

- ✓ real time observation

Drawbacks:

- ✓ ↑ error rate (~20%)

base detection from the current signal

now we have

two detectors → improved accuracy.

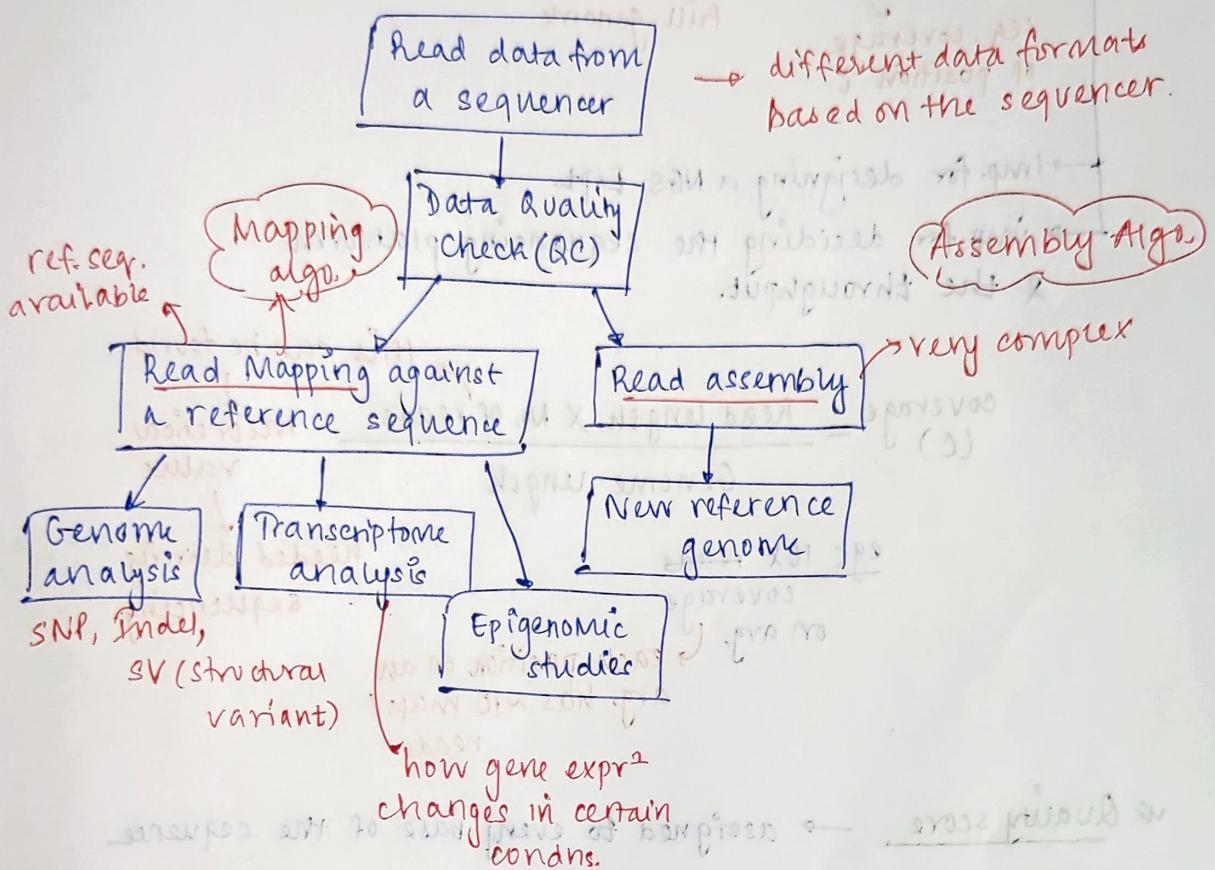
- ✓ current change not only depends on individual bases but also on the neighboring bases.

individual bases

sequence-dependent pattern

Data formats & Data quality:

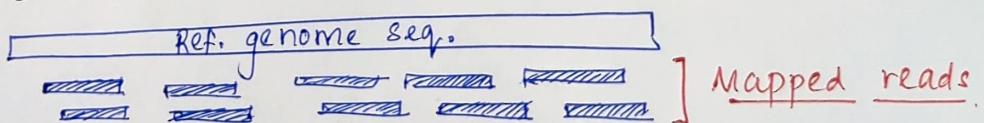
NGS Data analysis:



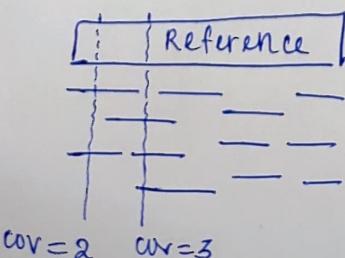
The Assembly Problem:

↳ if no ref genome, how do we know where the genome starts from?

The Mapping Problem:



Sequencing Coverage / Sequencing Depth:



no. of reads mapped to that particular position

\times Mean coverage → mean of coverage for all positions across the whole genome

$$= \frac{\sum_{i=1}^L cov_i}{L}$$

seq. coverage of position i

length of the full genome

→ imp. for designing a NGS Expt.

→ imp. for deciding the sequencing platform & the throughput.

$$\text{coverage} = \frac{\text{Read length} \times \text{No. of reads}}{\text{Genome length}}$$

this can be found

theoretical value

eg: 10X reads
coverage on avg. ↗ each position on an arg. has ~10 maps/ reads
needed during sequencing

\times Quality score → assigned to every base of the sequence

12/2/24

FASTA
FASTQ → Illumina

SMRT
Don Torrent
Nanopore

Read Mapping → the mapping problem is "string search"
Burrows-Wheeler Transform ← Algorithm

PacBio SMRT seq. data format

HDF5 - Hierarchical Data Format version 5

BAM - a binary format → machine-readable → not human-readable
↳ can be converted to FASTQ
bam2fastx
as most tools are based on this

Don Torrent data format

raw data → DAT files - windows format
↳ pH change to store files

seq. data → VCF →

→ HDF5;

FAST5 data format

Fast5
↓
raw data
↓
picurrent changes
analysis data
↓
sequence

Read Mapping :

Challenges:

- large genome → human genome: 3×10^9 bp
- huge no. of short reads → length of read $\sim 10^3$ bp
(millions or billions)
- time requirement → a long time for so many reads
- mutations & seq. error
- repetitive regions in the genome.

→ a read can map to multiple regions if seq. is similar

Mapping Algorithms:

Data Structures	① Hash-table based mapping algorithm	very fast but take up a lot of memory.
	② Suffix tree or suffix array based	
stores the ref genome multiple times	③ Burrows-Wheeler transform based. ref. genome → very small only once → amt. of memory	cannot be run on mere Desktop.

Burrows-Wheeler Transform (BWT):

Reference : G A C G T A C G T C A A \$

✓ \$ G A C G T A C G T C A A

✓ A. \$ G A C G T A C G T C A

✓ A A \$ G A C G T A C G T C

✓ C A A \$ G A C G T A C G T

✓ T C A A \$ G A C G T A C G

✓ G T C A A \$ G A C G T A C

✓ C G T C A A \$ G A C G T A

✓ A C G T C A A \$ G A C G T

✓ T A C G T C A A \$ G A C G

✓ G T A C G T C A A \$ G A C

✓ C G T A C G T C A A \$ G A

✓ A C G T A C G T C A A \$ G

✓ T A C G T C A A \$ G A C G

STEP:1

Alphabetical sort: [\\$' comes before A]

\$	A
A	A
A	C
A	G
A	T
C	T
C	A
C	A
G	\$
G	C
G	C
T	G
T	G

Burrows-Wheeler String

$\text{BWT}(S) = \boxed{\text{AACGTTAA\$ CCGG}}$

STEP:2

last column $\xrightarrow{\text{sort alphabetically}}$ first column

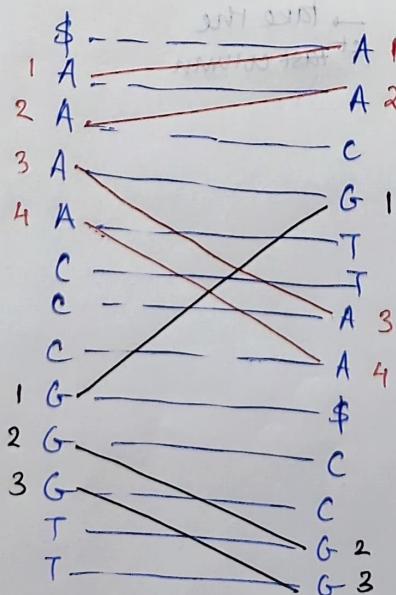
position info. also stored.

(say) search read AGGT

→ how do we take into account mismatches? AGGT is closest to ACGT but with one mismatch

Last-first (LF) mapping:

the i^{th} occurrence of character c in the last column & the i^{th} occurrence of character c in the first column correspond to the same character in the ref. seq.



if we no. the As in first column & no the As in last column, the 1-1, 2-2, 3-3, --, n-n correspond & occupy the same position in the reference genome

FM Index:

Full-Text Minute-space (Ferragina, Manzini)

first column of the BWT matrix can simply be stored as numbers of letters of each type.

A	C	G	T
4	3	3	2

this is enough to generate the 1st column

Searching a read position:

First column	<u>LF map:</u>	BWT	<u>Position</u>
\$		A ₁	12
A ₁		A ₂	11
A ₂		C ₁	10
A ₃	X	G ₁	1
A ₄		T ₁	5
C ₁	X	T ₂	9
C ₂		A ₃	2
C ₃		A ₄	6
G ₁		\$	0
G ₂		C ₂	3
G ₃		C ₃	7
T ₁		G ₂	4
T ₂		G ₃	8

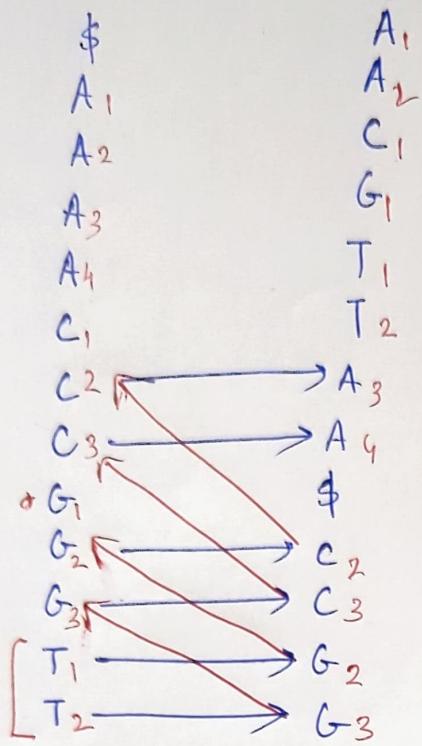
Read: TAC

→ LF mapping
 → strings that come before C in seq.

Read: T A C

start from the last we have an & coming!

take the 1st test column



ACG[T]

both ways/tracks
give us ACGT

repeat sequences.