

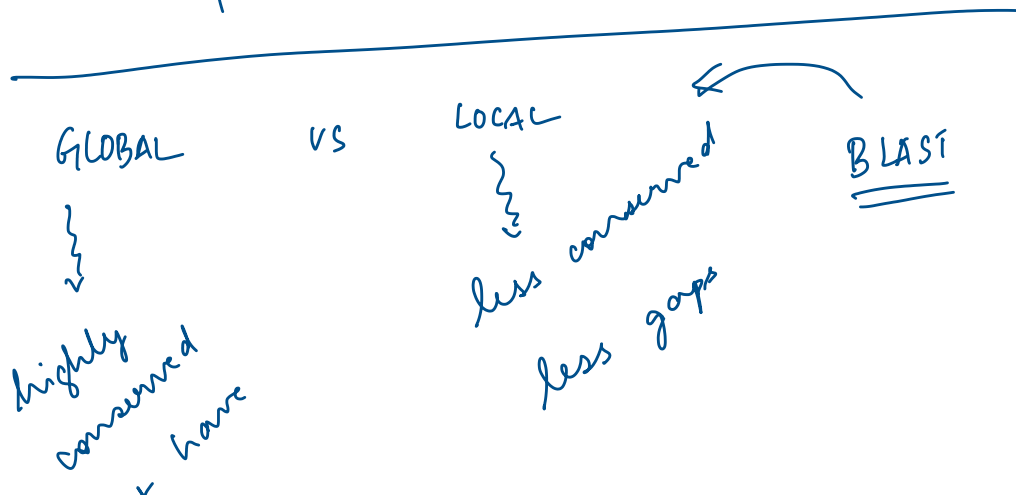
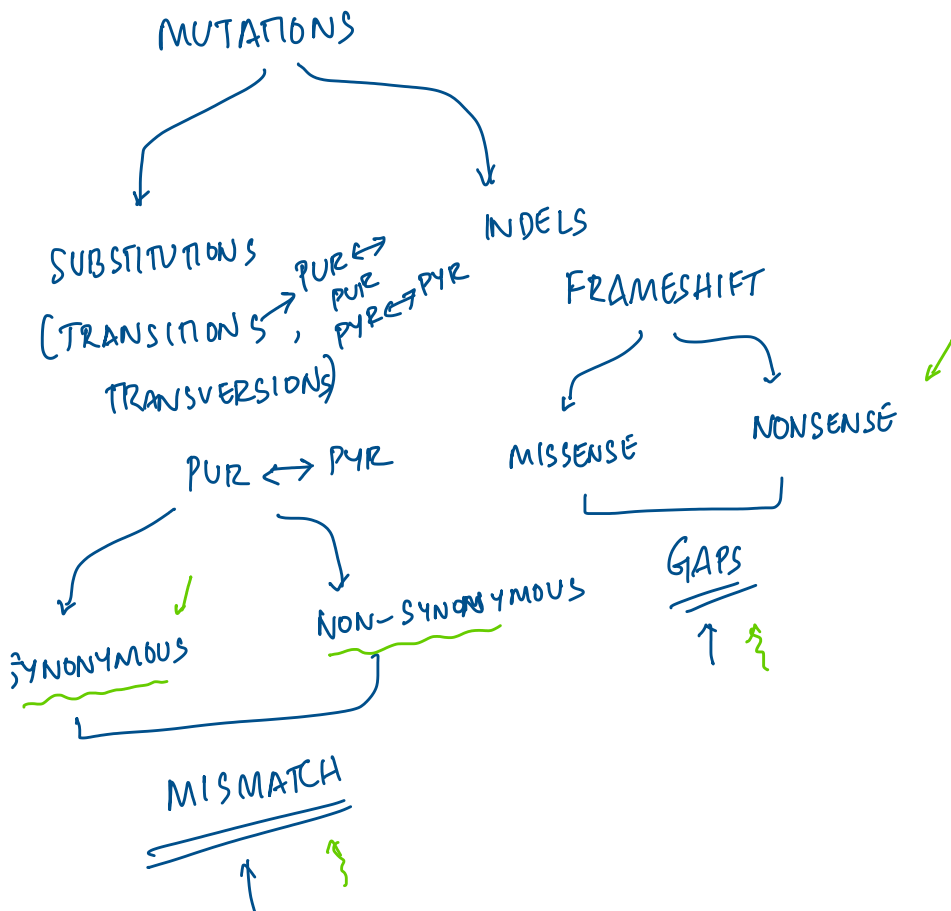
Week 3,4 Guided Assembly

Friday, 4 February 2022 8:08 AM

BOWTIE, BWA-MEM

⇒ efficiently align reads to a reference

ALIGNMENT → Homology → MATCH



high
(gap)

NEEDLEMAN & WUNSCH

seq. 1 - m

seq 2 - n

STEP 1: create a matrix $(m+1) \times (n+1)$ ✓

STEP 2: fill up gap penalties

STEP 3: use a scoring scheme:

$$S_{i,j} = \max \begin{cases} F(i-1, j-1) + s_{i,j} \\ F(i-1, j) - \text{gap} \\ F(i, j-1) - \text{gap} \end{cases} \quad \begin{matrix} \text{(match)} \\ \text{DIAGONAL} \end{matrix}$$

eg. MATCH = +1, MISMATCH = -1, GAP = -2

seq 1 = AGC (m)

seq 2 = AAAC (n)

		A	G	C
1	0	-2	-4	-6
A	-2	-1	-4	-3
A	-4	-1	0	-2
A	-6	-3	-2	-1

A	A	A	C
	X		
A	G	—	C
		⋮	

C | -8 | -5 | -4 | -1

S < W \rightarrow if you a -ve score, you push it to 0

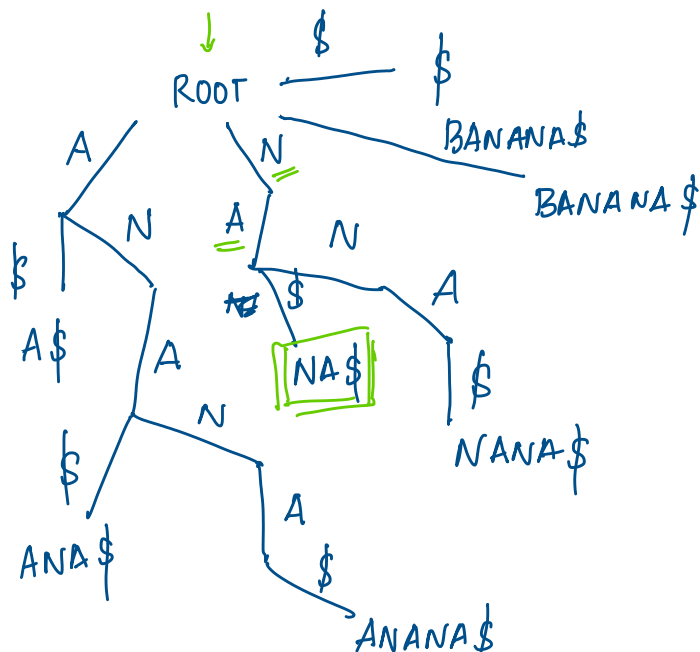
PROBLEM: align 'billions' of reads to a reference genome

ERGO
we are going to build a dictionary or database

SUFFIX ARRAY / TREE

GENOME: BANANA\$

all possible substrings \rightarrow
BANANA\$, ANANA\$, NANA\$, ANA\$, NA\$, A\$, \$]



NA\$
n = 1

BURROW - WHEELER TRANSFORM

0 BANANA\$

1 \$BANANA

2 A\$BANAN

3 NA\$BANA

4 ANA\$BAN

5 NANA\$BA

6 ANANA\$B

STEP 1

STEP 2: arrange the words alphabetically

1 \$BANANA

2 A\$BANAN

4 ANA\$BAN

6 ANANA\$B

0 BANANA\$

3 NA\$BANA

5 NANA\$BA

SA = I column
= \$AABNN

BWT = LAST COLUMN
= ANNBSAA

STEP 3

STEP 1

STEP 2

SUFFIX
ARRAY

BWT

2 1 5

POSITION (INDEX)	1	2	4	6	0	3	5
BWT	A	N	N	B	\$	A	A
→ SA	\$	A	A	A	B	N	N

eg. NA (read)

ALGORITHM 1 - BISECTION

⇒ 2 matches → Position ③ & ⑤

BA^③NA^⑤\$

Backward First Algorithm

- ① Look for A in SA [2, 4, 6]
- ② Look for all N prefixes in BWT [2, 4]
- ③ Whichth occurrence of N's in BWT? [1st, 2nd]
- ④ Look for those occurrences of N in SA → positions 3, ⑤ (backwards) STOP

REF GENOME :

ATGCATGC\$

READ : ATG

- ⑤ ATGC ATGC \$
- ① \$ ATGCATGC
- ② C\$ ATGCATG
- ③ GC\$ ATGCAT
- ④ TGC \$ ATGCA
- ⑤ ATGC \$ ATGC
- ⑥ CATGC \$ ATG
- ⑦ GCATGC \$ AT
- ⑧ TGCATGC \$ A

A
↓
Z

- ① \$ ATGCATGC
- ⑤ ATGC \$ ATGC
- ⑥ ATGCATGC \$
- ② C \$ ATGCATG
- ⑥ CATGC \$ ATG
- ③ GC\$ ATGCAT
- ⑦ GCATGC \$ AT
- ④ TGC \$ ATGCA
- ⑧ TGCATGC \$ A

POSITION	1	⑤	⑥	2	6	3	7	4	8
BWT	C	C	\$	G	G	<u>T_{1st}</u>	<u>T_{2nd}</u>	A [↓]	A [←]
SA	\$	A	A	C	C	G	G	T	T

↙ ↘ ↙
ATG
==

⇒ INDEX 0 and 5 backwards
i.e. 5th and 9th positions
backwards

- STEP 1: Look for G in SA → 3, 7 ;
- STEP 2: which of those are prefixed by T? → both ;
whichth T occurrences? 1st 2nd
- STEP 3: 1st and 2nd occurrences of T in SA → 4, 8
- STEP 4: which of those are prefixed by A? → both
...th A occurrences? 1st, 2nd

STEP 5 : which A occurs

STEP 6 : 1st, 2nd A's in SA \rightarrow 5, 0 backwards (STOP)

2 types of coins in bag \rightarrow BIASED OR UNBIASED

$$P(\text{BIASED}) = P(\text{UNBIASED}) = \frac{1}{2}$$

OUTCOMES

BIASED

$$P(H | \text{BIASED}) = 1$$

$$P(T | \text{BIASED}) = 0$$

UNBIASED

$$P(H | \text{UNBIASED}) = \frac{1}{2}$$

$$P(T | \text{UNBIASED}) = \frac{1}{2}$$

$$P(\text{BIASED} | H) = ?$$

BAYES' THEOREM

$$P(A | B)$$

POSTERIOR PROB.

$$= \frac{\overbrace{P(B|A) P(A)}^{\text{LIKELIHOOD}}}{P(B) \leftarrow \text{PRIOR PROBS.}}$$