

# Suffix Arrays

1 2 3 4 5 6 7 8  
A B A B A A B A

( 8 5 6 3 1 7 4 2 )

we can get  $O(n)$  time by  
using suffix trees!

$O(n)$  time direct methods exist.

q: BA

naive querying

$$O(|q| \log n)$$



$$O(|q| + \log n)$$

n Suffixe

$$O(\underline{n \log n}) \times O(\underline{n})$$

Cost of  
Comparison

naive

$$O(n^2 \log n)$$

1 2 3 4 5 6 7 8  
A B A B A A B A

1 2 3 4 5 6 7 8  
 8 5 6 3 1 7 4 2  
 ↑  
 L<sub>1</sub>  
 ↑  
 M<sub>1</sub>  
 L<sub>1</sub>  
 ↑  
 M<sub>2</sub>  
 R

q = BA

l = LCP(q, L) = 0  
 Longer common prefix length

r = LCP(q, R) = 2

1) maintain mlr

start any comparison at <sup>Position</sup> mlr + 1 mlr = 0

mlr = min(l, r)

we find BA.

2) start at max of l, r

L → 10  
 15  
 2  
 M → 1  
 L → 3  
 M → 3  
 5  
 R → 20  
 ABA  
 ABAC  
 ABBC  
 ABAC  
 ABC  
 ABCA  
 ABCD

q = ABC

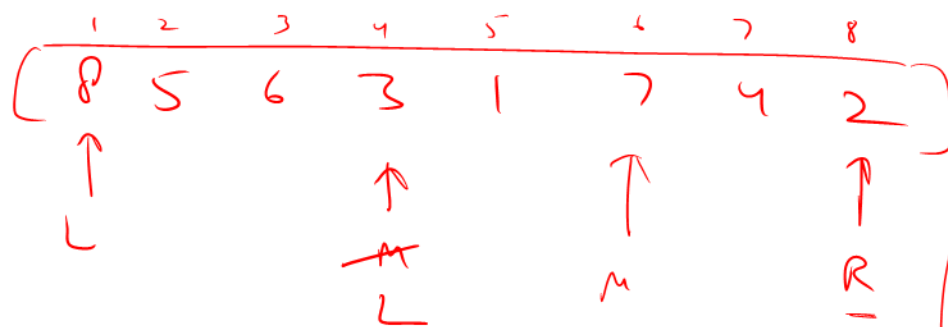
l = 2 = LCP(q, L)

r = 3 = LCP(q, R)

mlr + 1 = 2 + 1 = 3

O(|q| · l, n)

1	2	3	4	5	6	7	8
A	B	A	B	A	A	B	A


$$q = BA$$
$$LCP(q, L)$$
$$l = 0$$
$$LCP(q, r)$$
$$r = 2$$
$$LCP(L, M) = 0$$
$$L^p(M, \mathbb{R}) = 2$$

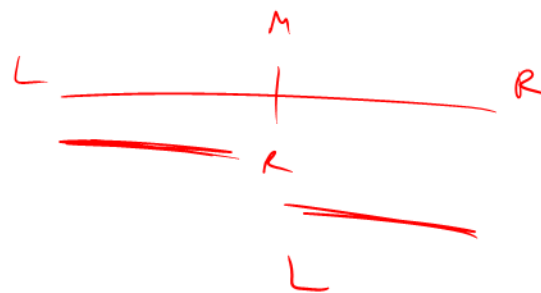
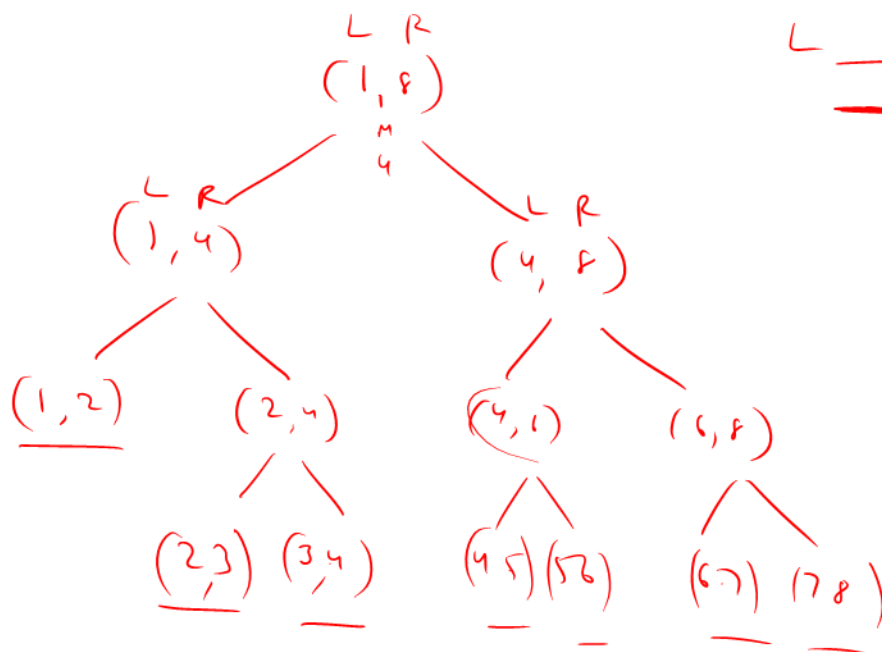
Pre computed  
Value

9 matches only 7

$$\partial(|q| + \underline{b_{s,n}})$$

look only  
"once"

## binary search



$n-1$  internal nodes

7 leaves  
n leave

$$O(n + n - 1) = O(n) \times O(1) = O(n)$$

# Exact Matching via Hashing

$l \leftarrow$  motif length

$d \leftarrow$  # of mismatches allowed

LSH : locality sensitive hashing

$h_1$ : choose any  $k$  position at random  
 $\rightarrow$  hash key

DB: H

A	C	A	G	T	A	C	A	C	T	T	G	C
1	2	3	4	5	6	7	8	9	10	11	12	13

$q$ :  
DB[4]

A	T	A	C	A	T	C	A
1	2	3	4	5	6	7	8

$h_1: (3, 4, 5)$  e.g.  $h(1, 2, 3)$

$l = 8$

$d = 4$

$k = 3$

HT

ACA $\rightarrow 4, 6$

Prob of a match between  
 $q$  & DB for a  
single position

Answer: DB[4]  
DB[6]

$$\left(1 - \frac{d}{l}\right)$$

Prob of a match over  $k$  position is

$$\left(1 - \frac{d}{l}\right)^k$$

1) false positive

→ do something, <sup>DB[i]</sup> hashes to a cell, but  
the true distance  $> d$

✓

solution: just compare the g with DB[i]

2) false negative:

there exists a match, but they don't hash  
to the same cell?

solution: we cannot eliminate this

however, we can minimize this!

do  $m$  hash keys

$$\begin{array}{l} h_1(1, 2, 3) \\ h_2(3, 4, 5) \\ h_3(\dots) \\ \vdots \\ h_m(\dots) \end{array} \longrightarrow P(\text{of mismatch}) = \left(1 - \left(1 - \frac{d}{x}\right)^k\right)^m$$

Prob of a false negative even after  $m$  hash functions  
mismatch

$$\left(1 - \left(1 - \frac{d}{x}\right)^k\right)^{(m)} < \epsilon$$

$\Sigma$  is the desired prob of a mismatch.

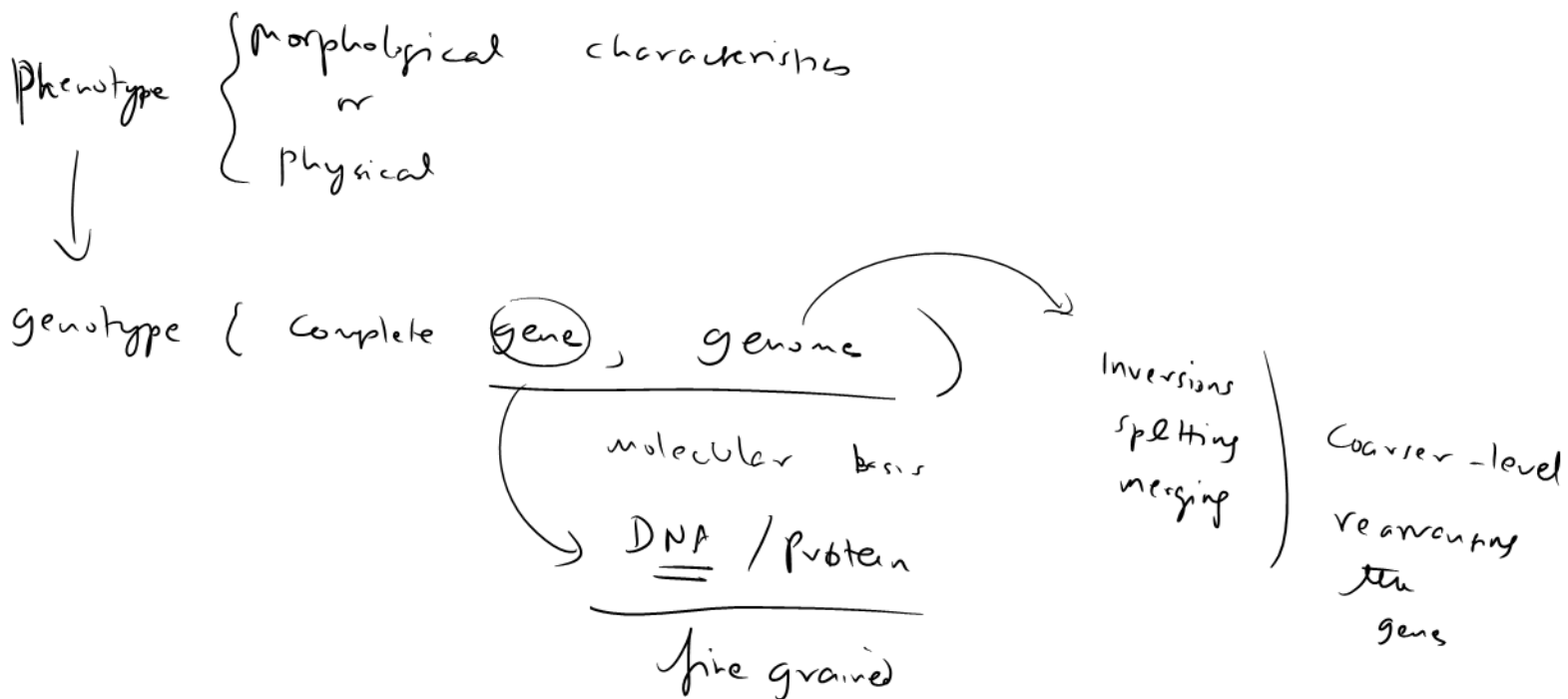
e.g.  $\Sigma = \underline{\underline{0.001}}$

solve for  $m$ !

---

## The Tree of life : Evolutionary Tree

---



---

### 1) Parsimony - based methods

most economical tree from # of changes

assumption: all positions are independent

⇒ add up the cost across all positions  
# of changes

$s_1: \begin{array}{|c|c|c|} \hline A & B & A \\ \hline \end{array}$   
 $s_2: \begin{array}{|c|c|c|} \hline B & A & A \\ \hline \end{array}$   
 $s_3: \begin{array}{|c|c|c|} \hline A & A & B \\ \hline \end{array}$   
 $s_4: \begin{array}{|c|c|c|} \hline B & B & A \\ \hline \end{array}$

Given: search over binary trees

1) search over different topologies

$(2n-3)!!$  rooted trees

$$\equiv (2n-3)(2n-5)(2n-7)\dots(1)$$

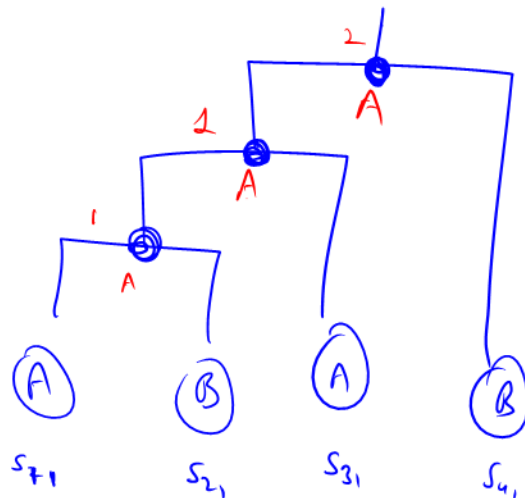
Infeasible to enumerate all trees

→ sampling

e.g. MCMC

Markov Chain Monte Carlo

2) how to compute the cost of a given tree



Upper-bound (3)

bottom up:

node: If  $L \cap R = \emptyset$

then LOR is kept  
else

keep LOR

top-down phase:

Starting at the root, for any node

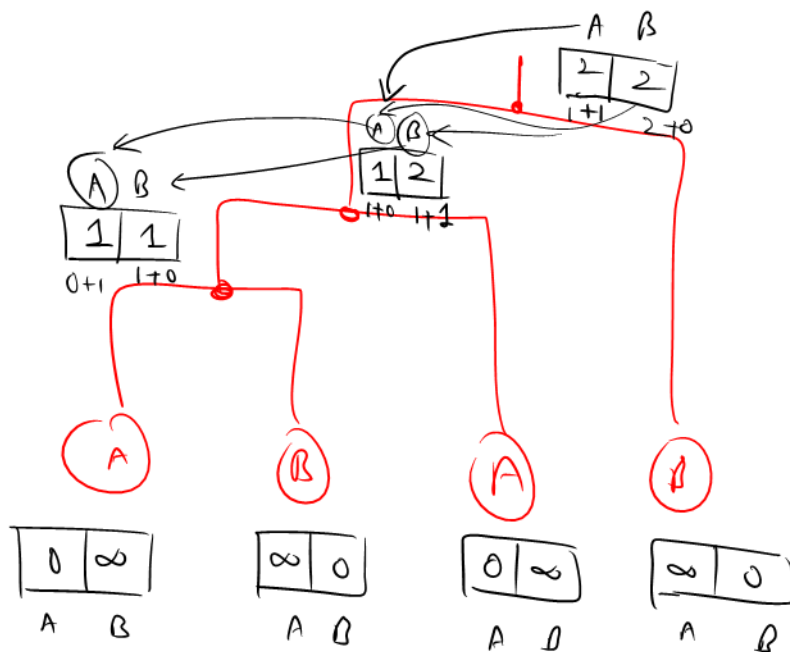
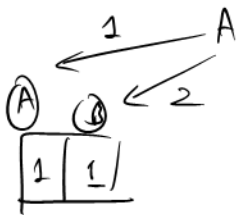
If more than one choice

use parent's symbol

	A	B
A	0	1
B	1	0

Dynamic programming

find the least cost tree.



2 2  
A B