

Algorithms in Bioinformatics

More suffix trees

Inexact Matching



Recap

- Office hours today start at 4
(plus I'll leave a little time
at the end of today
for questions)

- Next assignment: essay

Officially, due Thursday
(after break), but
feel free to submit
the next Tuesday if
you need more time

- No class in 1 week

Last time(s): Exact matching

We focused on:

- hashing : import something
- suffix trees

There are many variants,
but key tradeoff:

fixed P vs
fixed T

More suffix trees:

- Generalized suffix trees:
a tree that recognizes suffixes
for > 1 string

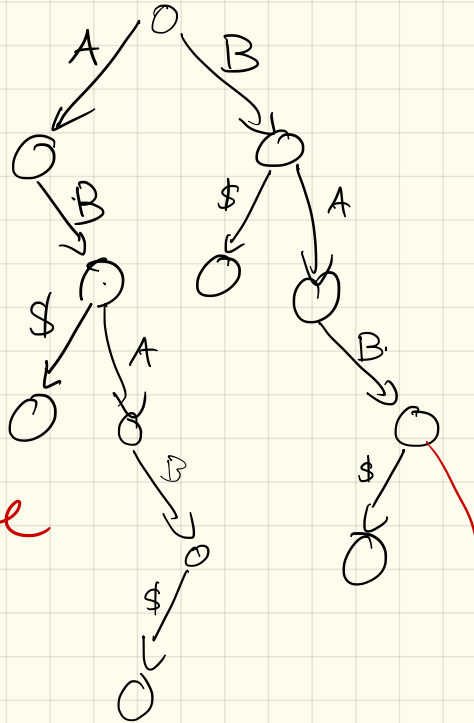
Idea:

- Build suffix tree for S_1

Eg: $S_1 = \underline{A} \underline{B} A B \$$

Then parse
for S_2 .

$S_2 = \underline{B} \underline{A} B A \#$

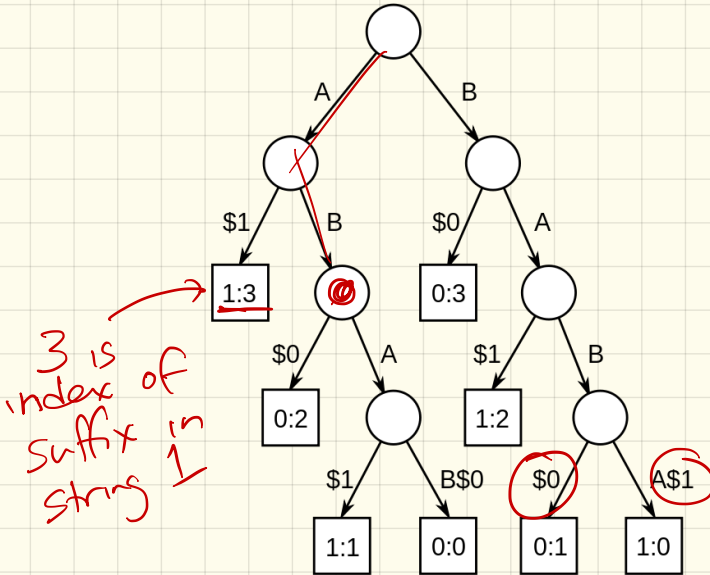


$O(|S_2|)$ time
to add to
tree

Final tree:

S₁ = BABA (\$0)

S₂ = ABAB (\$1)



(This can be done efficiently, but we won't unpack those details yet...)

Can actually do this for
even more strings:

- Make each terminate with
a different end
character. (\$1, \$2, \$3)
- Annotate tree as you
go

↳ encodes all suffixes
of all ~~trees~~

strings in a single
tree.

More applications:

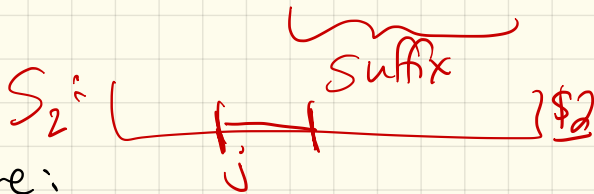
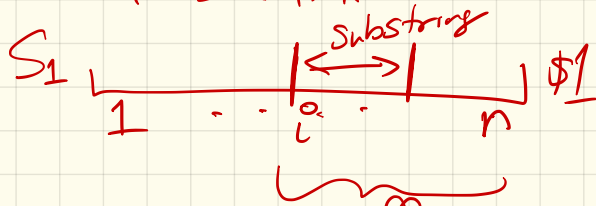
- Longest common substring
(not sequence) \hookrightarrow contiguous set of characters

How?

Use generalized suffix tree.

Mark internal nodes:

a substring in S_1
is a prefix of
a suffix.



Runtime:

Linear time

(to build gen. suffix tree)

• DNA contamination problem:
(Fun dinosaur story)

Problem: Given string S_1 ,
& a known string S_2 ,
find all substrings of S_2
that occur in S_1 and
are longer than some length l .

(These are candidates for
unwanted pieces of S_2
that have contaminated S_1 .)

How? (Similar to last one!)

Build gen. suffix tree.

$O(n)$

Same Example:

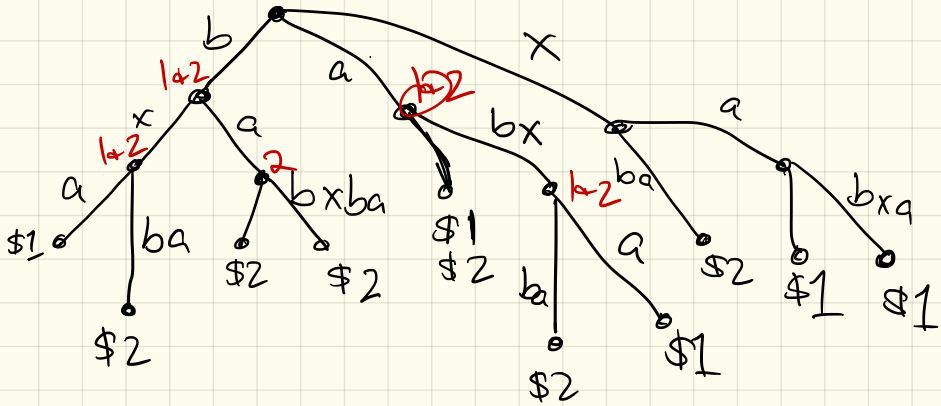
$S_1 = xabxc$ (\$1)

unknown DNA

$S_2 = babxba$ (\$2)

my DNA

Tree:



Set $l=2$:

Look for any nodes of
depth ≥ 2 .

Longest Common Extensions:

Input: Two strings S_1 + S_2

$$|S_1|, |S_2| = n$$

and ^{long} sequence of index
↑ pairs (i, j)

Output:

For each (i, j) , return length
of longest common prefix
of $S_1(i..n)$ matching
a prefix of $S_2(j..n)$

How?

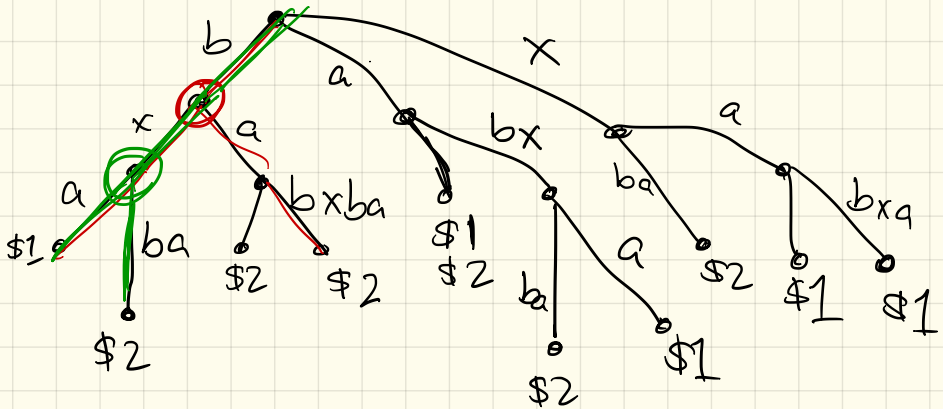
- Use lowest common ancestor in gen. suffix tree:

Same Example:

$$S_1 = x a b x c \text{ (\$1)}$$

$$S_2 = b a b x b a \text{ (\$2)}$$

Tree:



$$\text{Let } i = 3 \text{ \& } j = \underline{1}$$

$$1 : b$$

$$\text{or } \underline{i = 3} \text{ \& } \underline{j = 3}$$

$$2 : bx$$

Inexact Matching:

the k -mismatch problem

Idea: For some simpler variations on matching, don't need full-on dynamic programming.

Why bother?

Not going to allow insertions or deletions

$\sim O(n)$ instead of $O(mn)$

Problem:

K -mismatch Problem:

Approximate Pattern Matching Problem:

Find all approximate occurrences of a pattern in a text.

Input: A pattern $\mathbf{p} = p_1p_2 \dots p_n$, text $\mathbf{t} = t_1t_2 \dots t_m$, and parameter k , the maximum number of mismatches.

Output: All positions $1 \leq i \leq m - n + 1$ such that $t_it_{i+1} \dots t_{i+n-1}$ and $p_1p_2 \dots p_n$ have at most k mismatches (i.e., $d_H(\mathbf{t}_i, \mathbf{p}) \leq k$).

Ex: P = bend

T = abent banana end

k = 2

Then 3 2-mismatches:

• at position 2,
bent has 1 mismatch

at pos 6,
bana w/ 2 mismatches

at pos 11,
aend 1 mismatch
↓
bend

First approach: brute force:

APPROXIMATEPATTERNMATCHING(p, t, k)

```
1  $n \leftarrow$  length of pattern  $p$ 
2  $m \leftarrow$  length of text  $t$ 
3 for  $i \leftarrow 1$  to  $m - n + 1$ 
4    $dist \leftarrow 0$ 
5   for  $j \leftarrow 1$  to  $n$ 
6     if  $t_{i+j-1} \neq p_j$ 
7        $dist \leftarrow dist + 1$ 
8   if  $dist \leq k$ 
9     output  $i$ 
```

Runtime: $O(mn)$
 $O((m-n+1)n)$

Goal: $O(km)$

where $|P| = n$
 $|T| = m$ ($m > n$)

\rightarrow instead of $O(mn)$

Method: use longest common extension idea

For any position i in T
(the longer string),
do a longest common
extension query.

Say it has length l .
if $l = |P|$, done.

otherwise, position $i+l$
is a mismatch:

Ex:

	1	2	...	l	$l+1$	
P:	A	G	...	T	G	...
T =	...	A	G	...	T	A
		i	$i+1$...	$i+l-1$	$i+l$...

Then use a mismatch
(so $k-1$ left), reset to

$T[i+l+1]$ & $P[l+2]$

& continue, looking for $k-1$
match

Runtime: $O(km)$