Suffix Arrays CMSC 423

Suffix Arrays

- Even though Suffix Trees are O(n) space, the constant hidden by the big-Oh notation is somewhat "big": ≈ 20 bytes / character in good implementations.
- If you have a 10Gb genome, 20 bytes / character = 200Gb to store your suffix tree. "Linear" but large.
- Suffix arrays are a more efficient way to store the suffixes that can do most of what suffix trees can do, but just a bit slower.
- Slight space vs. time tradeoff.

Example Suffix Array

- Idea: lexicographically sort all the suffixes.
- Store the starting indices of the suffixes in an array.





suffix of s

sort the suffixes alphabetically

the indices just "come along for the ride"

index of suffix

s = attcatg

Example Suffix Array

- Idea: lexicographically sort all the suffixes.
- Store the starting indices of the suffixes in an array.

sort the suffixes alphabetically

the indices just "come along for the ride"



index of suffix suffix of s

s = attcatg

Another Example Suffix Array

- Idea: lexicographically sort all the suffixes.
- Store the starting indices of the suffixes in an array.



s = cattcat\$



sort the suffixes alphabetically

the indices just "come along for the ride"

index of suffix suffix of s

Another Example Suffix Array

- Idea: lexicographically sort all the suffixes.
- Store the starting indices of the suffixes in an array.

sort the suffixes alphabetically

the indices just "come along for the ride"

s = cattcat

 7
 t\$

 8
 \$

index of suffix suffix of s

cattcat\$

2 attcat\$

3 ttcat\$

tcat\$

cat\$

at\$

4

5

6

Search via Suffix Arrays

s = cattcat



- Does string "at" occur in s?
- Binary search to find "at".
- What about "tt"?

Counting via Suffix Arrays

s = cattcat\$



- How many times does "at" occur in the string?
- All the suffixes that start with "at" will be next to each other in the array.
- Find one suffix that starts with "at" (using binary search).
- Then count the neighboring sequences that start with at.

K-mer counting

Problem: Given a string s, an integer k, output all pairs (b, i) such that b is a length-k substring of s that occurs exactly i times.



I. Build a suffix array.

2.Walk down the suffix array, keeping a CurrentCount count If the current suffix has length < k, skip it</p>

If the current suffix starts with the same length-k string as the previous suffix: increment CurrentCount else output CurrentCount and previous length-k suffix CurrentCount := I

Output CurrentCount & length-k suffix.

Constructing Suffix Arrays

• Easy O(n² log n) algorithm:

sort the n suffixes, which takes $O(n \log n)$ comparisons, where each comparison takes O(n).

- There are several direct O(n) algorithms for constructing suffix arrays that use very little space.
- The Skew Algorithm is one that is based on divide-and-conquer.
- An simple O(n) algorithm: build the suffix tree, and exploit the relationship between suffix trees and suffix arrays (next slide)

Relationship Between Suffix Trees & Suffix Arrays



Red #s = starting position of the suffix ending at that leaf

Leaf labels left to right: 86251743

Edges leaving each node are sorted by label (left-to-right).

Relationship Between Suffix Trees & Suffix Arrays



sorted by label (left-to-right).

Recap

- Suffix arrays can be used to search and count substrings.
- Construction:
 - Easily constructed in O(n² log n)
 - Simple algorithms to construct them in O(n) time using possibly $O(n^2)$ space.
 - More complicated algorithms to construct them in O(n) time using at most O(n) space.
- More space efficient than suffix trees: just storing the original string + a list of integers.