

Another approach to string problems is to build a more general data structure that represents the strings in a way that allows for certain queries about the string to be answered quickly. There are a lot of such indexes (e.g. wavelet trees, FM-index, etc.) that make different tradeoffs and support different queries. Today, we're going to see two of the most common string index data structures: suffix trees and suffix arrays.

1 Suffix Trees

Consider a string T of length t (long). Our goal is to preprocess T and construct a data structure that will allow various kinds of queries on T to be answered efficiently. The most basic example is this: given a pattern P of length p , find all occurrences of P in the text T . What is the performance we aiming for?

- The time to find all occurrences of pattern P in T should be $O(p + k)$ where k is the number of occurrences of P in T .
- Moreover, ideally we would require $O(t)$ time to do the preprocessing, and $O(t)$ space to store the data structure.

Suffix trees are a solution to this problem, with all these ideal properties.¹ They can be used to solve many other problems as well. In this lecture, we'll consider the alphabet size to be $|\Sigma| = O(1)$.

1.1 Tries

The first piece of the puzzle is a *trie*, a data structure for storing a set of strings. This is a tree, where each edge of the tree is labeled with a character of the alphabet. Each node then implicitly represents a certain string of characters. Specifically, a node v represents the string of letters on the edges we follow to get from the root to v . (The root represents the empty string.) Each node has a bit in it that indicates whether the path from the root to this node is a member of the set—if the bit is set, we say the node is marked.

Since our alphabet is small, we can use an array of pointers at each node to point at the subtrees of it. So to determine if a pattern P occurs in our set we simply traverse down from the root of the tree one character at a time until we either (1) walk off the bottom of the tree, in which case P does not occur, or (2) we stop at some node v . We now know that P is a prefix of some string in our set. And if v is marked, then P is in our set, otherwise it is not.

This search process takes $O(p)$ time because each step simply looks up the next character of P in an array of child pointers from the current node. (We used that $|\Sigma| = O(1)$ here.)

1.2 Tries → Suffix Trees

Our first attempt to build a data structure that solves this problem is to build a trie which stores all the strings that are suffixes of the given text T . It's going to be useful to avoid having one suffix match the beginning of another suffix. So in order to avoid this we will affix a special

¹Suffix trees were invented by Peter Weiner.

character denoted “\$” to the end of the text T , which occurs nowhere else in T . (This character is lexicographically less than any other character.)

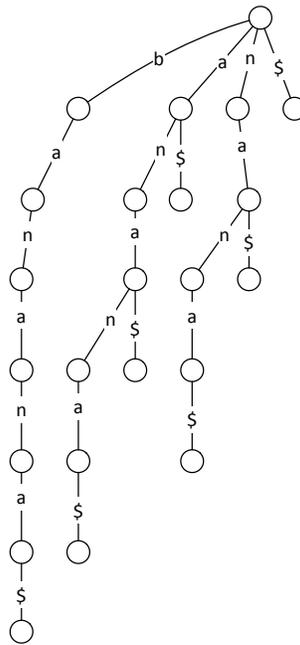
For example if the text were $T = \text{banana}\$,$ the suffixes of T are then

```

banana$
anana$
nana$
ana$
na$
a$
$

```

(Since we appended the \$ sign to the end, we’re not including the empty suffix here.) And the trie of these suffixes would be:



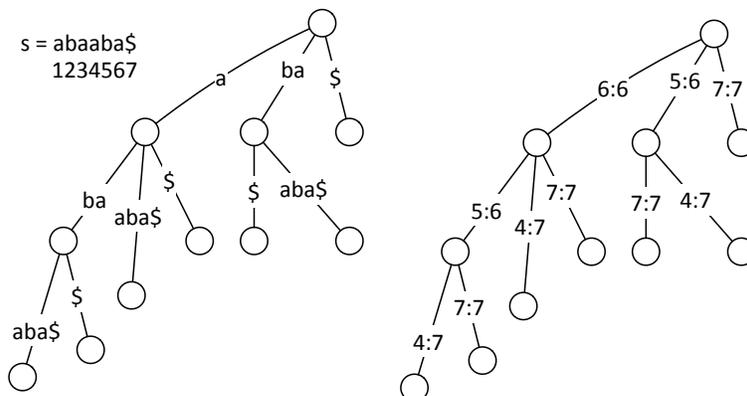
Suppose we have constructed the trie containing all these suffixes of T and we store at each node the count of the number of leaves in the subtree rooted at that node. Now given a pattern P , we can count the number of occurrences of P in T in $O(|P|)$ time: we just walk down the trie and when we run out of P we look at the count of the node v we’re sitting on. It’s our answer.

But there are a number of problems with this solution. First of all, the space to store this trie could be as large as $\Theta(t^2)$. Also, it’s unsatisfactory in that it does not tell us where in s these patterns occur. Finally, it will also take too long to build it.

Shrinking the tree to $O(t)$ space. The first two issues are easy to handle: we can create a “compressed” tree of size only $O(t)$. Since no string occurs as a prefix of any other, we can divide the nodes of our trie into *internal* and *leaf* nodes.

The leaf nodes represent a suffix of T . We can have each leaf node point to the place in T where the given suffix begins.

If there is a long path in the trie with branching factor 1 at each node on that path. We can compress that path into a single edge that represents the entire path. Moreover, the string of characters corresponding to that path must occur in T , so we can represent it implicitly by a pair of pointers into the string T . So an edge is now labeled with a pair of indices into T instead of just a single character (or a string). Here's an example (with the substrings labeling the edges on the left, and the start-end pairs labeling them on the right):



This representation uses $O(t)$ space. (We count pointers as $O(1)$ space.) Why? Each internal node now has degree at least 2, hence the total number of nodes in the tree is at most twice the number of leaves. But each leaf corresponds to some suffix of T , and there are t suffixes.

Building the tree — preview. What about the time to build the data structure? Let's first look at the naïve construction, by adding suffixes into it one at a time. To add a new suffix, we walk down the current tree until we come to a place where the path leads off of the current tree. (This must occur because the suffix is not already in the tree.) This could happen in the middle of an edge, or at an already existing node. In the former case, we split the edge in two and add a new node with a branching factor of 2 in the middle of it. In the latter case we simply add a new edge from an already existing node. In either case the process terminates with a tree containing $O(t)$ nodes, and the running time of this naïve construction algorithm is $O(t^2)$.

In fact, it is possible to build a suffix tree on a string of length t in time $O(t)$. We may not have time to go over this construction algorithm in detail.

2 Applications of Suffix Trees

We've seen how suffix trees can do exact search in time proportional to query string, once the tree is built. There are many other applications of suffix trees to practical problems on strings. Gusfield discusses many of these in his book. We'll just mention just a few here.

2.1 Simple Queries

Suffix trees make it easy to answer common (and less common) kinds of queries about strings. For example:

- Check whether P is a suffix of T : follow the path for q starting from the root and check whether you end at a leaf.
- Count the number of occurrences of P in T : follow the path for q ; the number of leaves under the node you end up at is the number of occurrences of P . If you are going to answer this kind of query a lot, you can store the number of leaves under each node in the nodes.
- Find the lexicographically (alphabetically) first suffix: start at the root, repeatedly follow the edge labeled with the lexicographically (alphabetically) smallest letter.
- Find the longest repeat in T . That is, find the longest string r such r occurs at least twice in T : Find the deepest node that has ≥ 2 leaves under it.

2.2 Longest Common Substring of Two Strings

Given two strings S and T , what is the longest substring that occurs in both of them? For example if $S = \text{boogie}$ and $T = \text{ogre}$ then the answer is og . How can one compute this efficiently? The answer is to use suffix trees. Here's how.

Construct a new string $U = S\%T$. That is, concatenate S and T together with an intervening special character that occurs nowhere else (indicated here by “%”). Let n be the sum of the lengths of the two strings. Now construct the suffix tree for U . Every leaf of the suffix tree represents a suffix that begins in S or in T . Mark every internal node with two bits: one that indicates if this subtree contains a substring of S , and another for T . These bits can be computed by depth first search in linear time. Now take the deepest node in the suffix tree (in the sense of the longest string in the suffix tree) that has both marks. This tells you the the longest common substring.

This is a very elegant solution to a natural problem. Before suffix trees, an algorithm of Karp, Miller, and Rosenberg gave an $O(n \log n)$ time solution, and Knuth had even conjectured a lower bound of $\Omega(n \log n)$.

2.3 Searching for Matching Strings in a Database

The above idea can be extended to more than 2 strings. This is called a *generalized suffix tree*. To represent a set of strings $D = \{S_1, \dots, S_m\}$, you concatenate the strings together with a unique character and label the leaves with the index of the string in which that suffix begins. To find which strings in D contain a query string q you follow the path for q and report the indices that occur in the subtree under the node at which you stopped.

3 Suffix Arrays

Imagine that you write down all the suffixes of a string T of length t . The i^{th} suffix is the one that begins at position i . Now imagine that you sort all these suffixes. And you write down the indices of them in an array in their sorted order. This is the suffix array. For example, suppose $T = \text{banana\$}$ and we've sorted the suffixes:

0	1	2	3	4	5	6	
	b	a	n	a	n	a	\$
6:	\$						
5:	a\$						

```

3: ana$
1: anana$
0: banana$
4: na$
2: nana$

```

The numbers to the left are the indices of these suffixes. So the *suffix array* is:

6 5 3 1 0 4 2

This array can be computed in a straightforward way by sorting the suffixes directly. This takes $O(t^2 \log t)$ time because each comparison of two strings takes $O(t)$ time. In fact, the suffix array can be constructed in $O(t)$ time. We will see faster construction algorithms in a minute.

It's often handy to have an auxiliary array called the “*LCP*” array around. Each successive suffix in this order matches the previous one in some number of letters. (Maybe zero letters.) This is recorded in the common prefix lengths array, or the LCP array. In this case we have:

```

suffix array is:           6 5 3 1 0 4 2
common prefix lengths array 0 1 3 0 0 2

```

3.1 Searching T using the suffix array

Consider the standard string search problem: we have the suffix array A for a string T and we want to find where pattern P occurs in A .

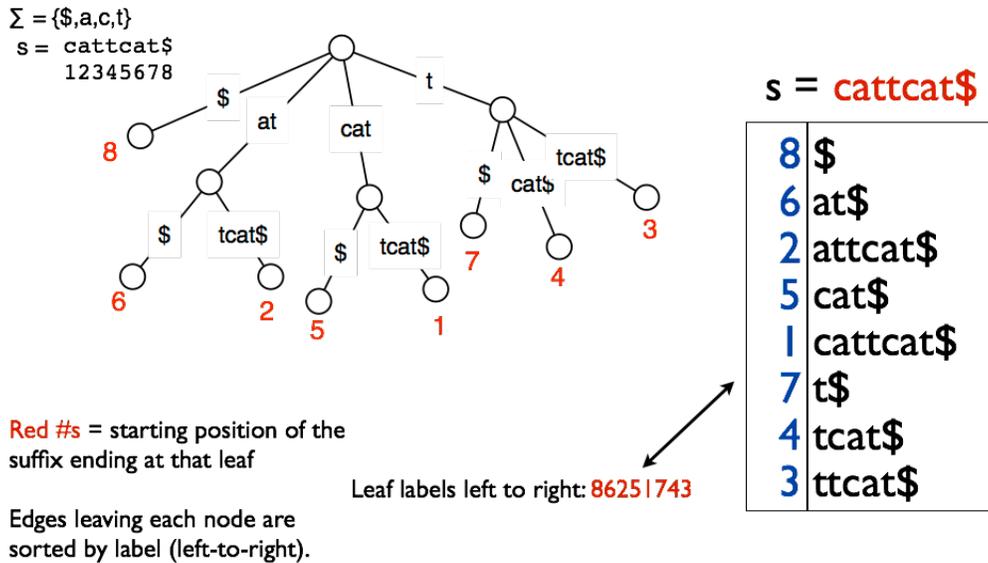
We can do this in $O(|P| \log |T|)$ time using binary search using the suffix array: Maintain a range $[U, D]$ of candidate positions in the array; initially the range is the entire suffix array. Let $M(U, D)$ be the midpoint of that range. Repeatedly check whether P comes before or after the suffix at position $M(U, D)$ of the array, and update U and D accordingly. This takes $O(|P| \log |T|)$ since we're doing a binary search over $|T|$ elements, and each comparison of P against the string at suffix $M(U, D)$ takes $O(P)$ time.

How do we find *all* the occurrences of P in T ? The thing to note is that these will all be adjacent in the suffix array since all the suffixes that start with P obviously start with the same sequence. We can find the range that starts with P in two ways: The first way: we could do 2 binary searches: one that takes the “left” range when there is a tie — which will find the start of the range — and one which takes the “right” range when there is a tie — which will find the end of the range. These searches will give you the range which contains suffixes starting with P . The second way: using the LCP array, we can walk left and right from a suffix that starts with P , continuing as long as the LCP is $\geq |P|$.

Note that $O(|P| \log |T|)$ time to search is slower than we got with suffix trees — we can do better, however. In fact, there is an $O(|P|)$ algorithm to search, matching the time for suffix trees. If there's time, we will see an $O(|P| + \log |T|)$ time algorithm that is almost as good.

3.2 Suffix Array \leftrightarrow Suffix Tree

Suffix Tree \rightarrow Suffix Array. The suffix array can be computed from the suffix tree by doing an in-order traversal of the tree, where in-order means that we visit children in lexicographic (alphabetical) order.



This takes $O(|T|)$ time.

Suffix Array \rightarrow Suffix Tree. We add the suffixes one at a time into a partially built suffix tree in the order that they appear in the suffix array. At any point in time, we keep track of the sequence of nodes on the path from the most recently added leaf to the root. To add the next suffix, we find where this suffix's path deviates from the current path we're keeping track of. To do this, we just use the common prefix length value. We walk up the path until we pass this prefix length. This tells us where to add the new node.

A potential argument can be used to see that this process runs in linear time. Imagine a token on each of the edges on the path from the current leaf to the root. We use these tokens to pay for walking up the tree until we find the branch point where a new child is added. The tokens on the path pay for the steps we take up the tree. We'll need a new token for the edge that connects to the new leaf. We may also need another token in case we have to split an edge. So in all, at most two new tokens are needed to pay for the work. This proves that the running time is linear.

So how do we compute the suffix array and the common prefix lengths array? There are linear time algorithms for this, but here I will describe a probabilistic method that is $O(n \log^2 n)$.

It's based on Karp-Rabin fingerprinting. If we could compare two suffixes in $O(1)$ time we could then just sort them in $O(n \log n)$ time. Instead we'll use a method for comparing two suffixes that works in $O(\log n)$ time.

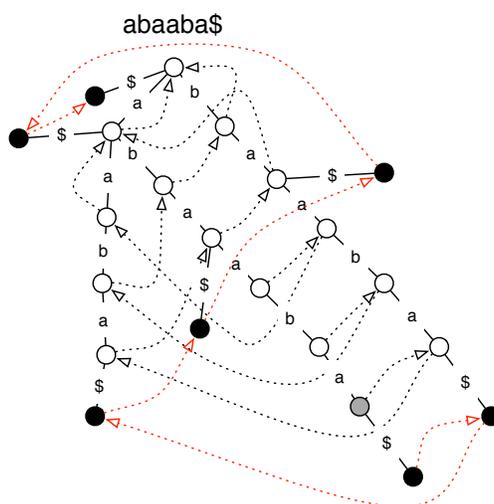
Using Karp-Rabin fingerprinting we can in $O(1)$ time (see the previous lecture) compare two substrings for equality. To compare two suffixes for lexicographic order, we use binary search to find the shortest length R such that the first R characters of each of the suffixes differ, but the first $R-1$ characters of them are the same. Then the lexicographic order is determined by the R th character of them. Furthermore this also tells us the common prefix length between the two strings.

4 Another Application: Longest Common Substring

Suppose we can preprocess our large database text T in anticipation of many future (unknown) queries q where we want to find the longest common substring between T and q . To solve this problem, we need to introduce the concept of a *suffix link*:

Definition 1 A suffix link is an extra pointer leaving each node in a suffix tree. If a node u represents string $x\alpha$ (where x is a character and α is a string) then u 's suffix link connects from u to the node representing the string α .

Every node has a suffix link. Why? Every node represents the prefix of some suffix. Suppose node u represents the prefix of suffix i of length m . Then u 's suffix link should point to the node on the path representing the prefix of length $m - 1$ of suffix $i + 1$. A trie with all the suffix links shown as dashed arrows is given below:



If node representing string $x\alpha$ exists in our tree (even after compression), then a node representing α must exist. The reason for this is since the node $x\alpha$ exists, it has two children y and z . That means both the strings αy and αz must exist in the string. So the node representing by α must have 2 children too.

Using suffix links we can get an $O(|q|)$ -time algorithm to find the longest common substring between T and q :

1. Walk down the tree, following q
2. If you hit a dead-end at node u ,
 - (a) save the depth of node u if it is deeper than the deepest dead end you found so far.
 - (b) follow the suffix link out of u .
 - (c) Go to step 1 (starting from your current place in q)
3. When you exhaust q , return the string represented by the deepest node you visited.

The first idea here is the the longest common substring starts at some suffix — we just don't know which suffix, so we try them all starting with suffix 0. The second idea here is that following a suffix link chops off the first character of q taking you to the place in the tree that you would have been at if the query had started with the 2nd character of q .

5 Summary

We saw different approaches to exact string matching problems: given a pattern P and a text T , find all occurrences of P within T . We saw:

- The randomized Karp-Rabin fingerprinting scheme which has a linear running time. It is a versatile idea and extends to different settings (like 2-dimensional pattern matching).
- the Knuth-Morris-Pratt algorithm which runs in time $O(t + p)$. Here if you have a pattern you want to find in many texts, you can preprocess the pattern in $O(p)$ time and space, and then search over multiple texts, the search in text T_i taking the time $O(|T_i|)$.
- The suffix-tree construction. Here you preprocess the text in $O(t)$ time and space, and then you can perform many different operations (including searching for different patterns) in time $O(p + n_{p,t})$, where $n_{p,t}$ is the number of occurrences of pattern p in text t .

Each one its advantages.