

Suffix Array

The suffix array of a text T is a lexicographically ordered array of the set $T_{[0..n]}$ of all suffixes of T . More precisely, the suffix array is an array $SA[0..n]$ of integers containing a permutation of the set $[0..n]$ such that $T_{SA[0]} < T_{SA[1]} < \dots < T_{SA[n]}$.

A related array is the **inverse suffix array** SA^{-1} which is the inverse permutation, i.e., $SA^{-1}[SA[i]] = i$ for all $i \in [0..n]$. The value $SA^{-1}[j]$ is the **lexicographical rank** of the suffix T_j .

As with suffix trees, it is common to add the end symbol $T[n] = \$$. It has no effect on the suffix array assuming $\$$ is smaller than any other symbol.

Example 4.7: The suffix array and the inverse suffix array of the text $T = \text{banana}\$$.

i	$SA[i]$	$T_{SA[i]}$	j	$SA^{-1}[j]$	T_j
0	6	\$	0	4	banana\$
1	5	a\$	1	3	anana\$
2	3	ana\$	2	6	nana\$
3	1	anana\$	3	2	ana\$
4	0	banana\$	4	5	na\$
5	4	na\$	5	1	a\$
6	2	nana\$	6	0	\$

177

Suffix array is much simpler data structure than suffix tree. In particular, the type and the size of the alphabet are usually not a concern.

- The size of the suffix array is $\mathcal{O}(n)$ on any alphabet.
- We will later see that the suffix array can be constructed in the same asymptotic time it takes to **sort the characters** of the text.

Suffix array construction algorithms are quite fast in practice too. Probably the fastest way to construct a suffix tree is to construct a suffix array first and then use it to construct the suffix tree. (We will see how in a moment.)

Suffix arrays are rarely used alone but are augmented with other arrays and data structures depending on the application. We will see some of them in the next slides.

178

Exact String Matching

As with suffix trees, **exact string matching** in T can be performed by a **prefix search** on the suffix array. The answer can be conveniently given as a **contiguous interval** $SA[b..e]$ that contains the suffixes with the given prefix. The interval can be found using string binary search.

- If we have the additional arrays $LLCP$ and $RLCP$, the result interval can be computed in $\mathcal{O}(|P| + \log n)$ time.
- Without the additional arrays, we have the same time complexity on average but the worst case time complexity is $\mathcal{O}(|P| \log n)$.
- We can then count the number of occurrences in $\mathcal{O}(1)$ time, list all occ occurrences in $\mathcal{O}(occ)$ time, or list a sample of k occurrences in $\mathcal{O}(k)$ time.

An alternative algorithm for computing the interval $SA[b..e]$ is called **backward search**. It is commonly used with compressed representations of suffix arrays and will be covered in the course **Data Compression Techniques**.

179

LCP Array

Efficient string binary search uses the arrays $LLCP$ and $RLCP$. However, for many applications, the suffix array is augmented with the lcp array of Definition 1.11 (Lecture 2). For all $i \in [1..n]$, we store

$$LCP[i] = \text{lcp}(T_{SA[i]}, T_{SA[i-1]})$$

Example 4.8: The LCP array for $T = \text{banana}\$$.

i	$SA[i]$	$LCP[i]$	$T_{SA[i]}$
0	6		\$
1	5	0	a\$
2	3	1	ana\$
3	1	3	anana\$
4	0	0	banana\$
5	4	0	na\$
6	2	2	nana\$

180

Using the solution of Exercise 2.4 (construction of compact trie from sorted array and LCP array), the suffix tree can be constructed from the suffix and LCP arrays in linear time.

However, many suffix tree applications can be solved using the suffix and LCP arrays directly. For example:

- The **longest repeating factor** is marked by the maximum value in the LCP array.
- The **number of distinct factors** can be computed by the formula

$$\frac{n(n+1)}{2} + 1 - \sum_{i=1}^n LCP[i]$$

since it equals the number of nodes in the uncompact suffix trie, for which we can use Theorem 1.17.

- **Matching statistics** of S with respect to T can be computed in linear time using the generalized suffix array of S and T (i.e., the suffix array of $S\$T\$$) and its LCP array (exercise).

181

LCP Array Construction

The LCP array is easy to compute in linear time using the suffix array SA and its inverse SA^{-1} . The idea is to compute the lcp values by comparing the suffixes, but skip a prefix based on a known lower bound for the lcp value obtained using the following result.

Lemma 4.9: For any $i \in [0..n]$, $LCP[SA^{-1}[i]] \geq LCP[SA^{-1}[i-1]] - 1$

Proof. For each $j \in [0..n]$, let $\Phi(j) = SA[SA^{-1}[j] - 1]$. Then $T_{\Phi(j)}$ is the immediate lexicographical predecessor of T_j and $LCP[SA^{-1}[j]] = \text{lcp}(T_j, T_{\Phi(j)})$.

- Let $\ell = LCP[SA^{-1}[i-1]]$ and $\ell' = LCP[SA^{-1}[i]]$. We want to show that $\ell' \geq \ell - 1$. If $\ell = 0$, the claim is trivially true.
- If $\ell > 0$, then for some symbol c , $T_{i-1} = cT_i$ and $T_{\Phi(i-1)} = cT_{\Phi(i-1)+1}$. Thus $T_{\Phi(i-1)+1} < T_i$ and $\text{lcp}(T_i, T_{\Phi(i-1)+1}) = \text{lcp}(T_{i-1}, T_{\Phi(i-1)}) - 1 = \ell - 1$.
- If $\Phi(i) = \Phi(i-1) + 1$, then $\ell' = \text{lcp}(T_i, T_{\Phi(i)}) = \text{lcp}(T_i, T_{\Phi(i-1)+1}) = \ell - 1$.
- If $\Phi(i) \neq \Phi(i-1) + 1$, then $T_{\Phi(i-1)+1} < T_{\Phi(i)} < T_i$ and $\ell' = \text{lcp}(T_i, T_{\Phi(i)}) \geq \text{lcp}(T_i, T_{\Phi(i-1)+1}) = \ell - 1$. \square

182

The algorithm computes the lcp values in the order that makes it easy to use the above lower bound.

Algorithm 4.10: LCP array construction

Input: text $T[0..n]$, suffix array $SA[0..n]$, inverse suffix array $SA^{-1}[0..n]$

Output: LCP array $LCP[1..n]$

```

(1)  $\ell \leftarrow 0$ 
(2) for  $i \leftarrow 0$  to  $n-1$  do
(3)    $k \leftarrow SA^{-1}[i]$ 
(4)    $j \leftarrow SA[k-1]$  //  $j = \Phi(i)$ 
(5)   while  $T[i+\ell] = T[j+\ell]$  do  $\ell \leftarrow \ell + 1$ 
(6)    $LCP[k] \leftarrow \ell$ 
(7)   if  $\ell > 0$  then  $\ell \leftarrow \ell - 1$ 
(8) return  $LCP$ 

```

The time complexity is $\mathcal{O}(n)$:

- Everything except the while loop on line (5) takes clearly linear time.
- Each round in the loop increments ℓ . Since ℓ is decremented at most n times on line (7) and cannot grow larger than n , the loop is executed $\mathcal{O}(n)$ times in total.

183

RMQ Preprocessing

The **range minimum query** (RMQ) asks for the smallest value in a given range in an array. Any array can be preprocessed in linear time so that RMQ for any range can be answered in constant time.

In the LCP array, RMQ can be used for computing the lcp of any two suffixes.

Lemma 4.11: The length of the **longest common prefix** of two suffixes $T_i < T_j$ is $\text{lcp}(T_i, T_j) = \min\{LCP[k] \mid k \in [SA^{-1}[i] + 1..SA^{-1}[j]]\}$.

The lemma can be seen as a generalization of Lemma 1.31 (Lecture 3) and holds for any sorted array of strings. The proof is left as an exercise.

- The RMQ preprocessing of the LCP array supports the same kind of applications as the LCA preprocessing of the suffix tree, but RMQ preprocessing is simpler than LCA preprocessing.
- The RMQ preprocessed LCP array can also replace the LLCP and RLCP arrays in binary searching.

184

We will next describe the RMQ data structure for an arbitrary array $L[1..n]$ of integers.

- We precompute and store the minimum values for the following collection of ranges:
 - Divide $L[1..n]$ into blocks of size $\log n$.
 - For all $0 \leq \ell \leq \log(n/\log n)$, include all ranges that consist of 2^ℓ blocks. There are $\mathcal{O}(\log n \cdot \frac{n}{\log n}) = \mathcal{O}(n)$ such ranges.
 - Include all prefixes and suffixes of blocks. There are a total of $\mathcal{O}(n)$ of them.
- Now any range $L[i..j]$ that overlaps or touches a block boundary can be exactly covered by at most **four ranges** in the collection.



The minimum value in $L[i..j]$ is the minimum of the minimums of the covering ranges and can be computed in constant time.

Enhanced Suffix Array

The enhanced suffix array adds two more arrays to the suffix and LCP arrays to make the data structure fully equivalent to suffix tree.

- The idea is to represent a suffix tree node v representing a factor S_v by the suffix array interval of the suffixes that begin with S_v . That interval contains exactly the suffixes that are in the subtree rooted at v .
- The additional arrays support navigation in the suffix tree using this representation: one array along the regular edges, the other along suffix links.

With all the additional arrays the suffix array is not very space efficient data structure any more. Nowadays suffix arrays and trees are often replaced with **compressed text indexes** that provide the same functionality in much smaller space. These will be covered in the course **Data Compression Techniques**.

Here are some of the key properties of the BWT.

- The BWT is easy to compute using the suffix array:

$$L[i] = \begin{cases} \$ & \text{if } SA[i] = 0 \\ T[SA[i] - 1] & \text{otherwise} \end{cases}$$
- The BWT is **invertible**, i.e., T can be reconstructed from the BWT L alone. The inverse BWT can be computed in the same time it takes to sort the characters.
- The BWT L is typically **easier to compress** than the text T . Many text compression algorithms are based on compressing the BWT.
- The BWT supports **backward searching**, a different technique for indexed exact string matching. This is used in many **compressed text indexes**.

BWT will be covered in more detail in the course **Data Compression Techniques**.

Prefix Doubling

Our first specialized suffix array construction algorithm is a conceptually simple algorithm achieving $\mathcal{O}(n \log n)$ time.

Let T_i^ℓ denote the text factor $T[i.. \min\{i + \ell, n + 1\}]$ and call it an ℓ -factor. In other words:

- T_i^ℓ is the factor starting at i and of length ℓ except when the factor is cut short by the end of the text.
- T_i^ℓ is the **prefix** of the suffix T_i of length ℓ , or T_i when $|T_i| < \ell$.

The idea is to sort the sets $T_{[0..n]}^\ell$ for ever increasing values of ℓ .

- First sort $T_{[0..n]}^1$, which is equivalent to sorting individual characters. This can be done in $\mathcal{O}(n \log n)$ time.
- Then, for $\ell = 1, 2, 4, 8, \dots$, use the sorted set $T_{[0..n]}^\ell$ to sort the set $T_{[0..n]}^{2\ell}$ in $\mathcal{O}(n)$ time.
- After $\mathcal{O}(\log n)$ rounds, $\ell > n$ and $T_{[0..n]}^\ell = T_{[0..n]}$, so we have sorted the set of all suffixes.

Ranges $L[i..j]$ that are completely inside one block are handled differently.

- Let $NSV(i) = \min\{k > i \mid L[k] < L[i]\}$ (NSV=Next Smaller Value). Then the position of the minimum value in the range $L[i..j]$ is the last position in the sequence $i, NSV(i), NSV(NSV(i)), \dots$ that is in the range. We call these the NSV positions for i .
- For each i , store the NSV positions for i up to the end of the block containing i as a bit vector $B(i)$. Each bit corresponds to a position within the block and is one if it is an NSV position. The size of $B(i)$ is $\log n$ bits and we can assume that it fits in a single machine word. Thus we need $\mathcal{O}(n)$ words to store $B(i)$ for all i .
- The position of the minimum in $L[i..j]$ is found as follows:
 - Turn all bits in $B(i)$ after position j into zeros. This can be done in constant time using bitwise shift -operations.
 - The right-most 1-bit indicates the position of the minimum. It can be found in constant time using a lookup table of size $\mathcal{O}(n)$.

All the data structures can be constructed in $\mathcal{O}(n)$ time (exercise).

Burrows–Wheeler Transform

The Burrows–Wheeler transform (BWT) is an important technique for **text compression**, **text indexing**, and their combination **compressed text indexing**.

Let $T[0..n]$ be the text with $T[n] = \$$. For any $i \in [0..n]$, $T[i..n]T[0..i]$ is a **rotation** of T . Let \mathcal{M} be the matrix, where the rows are all the rotations of T in lexicographical order. All columns of \mathcal{M} are **permutations** of T . In particular:

- The first column F contains the text characters in order.
- The last column L is the BWT of T .

Example 4.12: The BWT of $T = \text{banana}\$$ is $L = \text{anb}\$aa$.

F		L
$\$$	b a n a n a	a
a	$\$$ b a n a n	n
a	n a $\$$ b a n	n
a	n a n a n a $\$$	b
b	a n a n a n a $\$$	b
n	a $\$$ b a n a	a
n	a n a $\$$ b a	a

Suffix Array Construction

Suffix array construction means simply sorting the set of all suffixes.

- Using standard sorting or string sorting the time complexity is $\Omega(\Sigma LCP(T_{[0..n]}))$.
- Another possibility is to first construct the suffix tree and then traverse it from left to right to collect the suffixes in lexicographical order. The time complexity is $\mathcal{O}(n)$ on a constant alphabet.

Specialized suffix array construction algorithms are a better option, though.

We still need to specify, how to use the order for the set $T_{[0..n]}^\ell$ to sort the set $T_{[0..n]}^{2\ell}$. The key idea is assigning **order preserving names** (lexicographical names) for the factors in $T_{[0..n]}^\ell$. For $i \in [0..n]$, let N_i^ℓ be an integer in the range $[0..n]$ such that, for all $i, j \in [0..n]$:

$$N_i^\ell \leq N_j^\ell \text{ if and only if } T_i^\ell \leq T_j^\ell.$$

Then, for $\ell > n$, $N_i^\ell = SA^{-1}[i]$.

For smaller values of ℓ , there can be many ways of satisfying the conditions and any one of them will do. A simple choice is

$$N_i^\ell = \{j \in [0, n] \mid T_j^\ell < T_i^\ell\}.$$

Example 4.13: Prefix doubling for $T = \text{banana}\$$.

N^1	N^2	N^4	$N^8 = SA^{-1}$
4 b	4 ba	4 bana	4 banana\$
1 a	2 an	3 anan	3 anana\$
5 n	5 na	6 nana	6 nana\$
1 a	2 an	2 ana\$	2 ana\$
5 n	5 na	5 na\$	5 na\$
1 a	1 a\$	1 a\$	1 a\$
0 \$	0 \$	0 \$	0 \$

Now, given N^ℓ , for the purpose of sorting, we can use

- N_i^ℓ to represent T_i^ℓ
- the pair $(N_i^\ell, N_{i+\ell}^\ell)$ to represent $T_i^{2\ell} = T_i^\ell T_{i+\ell}^\ell$.

Thus we can sort $T_{[0..n]}^{2\ell}$ by sorting pairs of integers, which can be done in $\mathcal{O}(n)$ time using LSD radix sort.

Theorem 4.14: The suffix array of a string $T[0..n]$ can be constructed in $\mathcal{O}(n \log n)$ time using prefix doubling.

- The technique of assigning order preserving names to factors whose lengths are powers of two is called the [Karp–Miller–Rosenberg naming technique](#). It was developed for other purposes in the early seventies when suffix arrays did not exist yet.
- The best practical variant is the [Larsson–Sadakane algorithm](#), which uses ternary quicksort instead of LSD radix sort for sorting the pairs, but still achieves $\mathcal{O}(n \log n)$ total time.

193

Let us return to the first phase of the prefix doubling algorithm: assigning names N_i^1 to individual characters. This is done by sorting the characters, which is easily within the time bound $\mathcal{O}(n \log n)$, but sometimes we can do it faster:

- On an ordered alphabet, we can use ternary quicksort for time complexity $\mathcal{O}(n \log \sigma_T)$ where σ_T is the number of distinct symbols in T .
- On an integer alphabet of size n^c for any constant c , we can use LSD radix sort with radix n for time complexity $\mathcal{O}(n)$.

After this, we can replace each character $T[i]$ with N_i^1 to obtain a new string T' :

- The characters of T' are integers in the range $[0..n]$.
- The character $T'[n] = 0$ is the unique, smallest symbol, i.e., \$.
- The suffix arrays of T and T' are **exactly the same**.

Thus we can construct the suffix array using T' as the text instead of T .

As we will see next, the suffix array of T' can be constructed in linear time. Then [sorting the characters](#) of T to obtain T' is the asymptotically **most expensive operation** in the suffix array construction of T for any alphabet.

194