Vorläufiges Skriptum VL Text-Indexierung und Information Retrieval

Wintersemester 2013/14

Johannes Fischer (TU Dortmund) last update: December 18, 2013

Disclaimer 1

Dieses Skript wird den Studierenden an der TU Dortmund im Voraus zur Verfügung gestellt. Die Inhalte werden im Laufe des Semesters aber noch angepasst, insbesondere in Bezug auf Information Retrieval. Die horizontale Linie kennzeichnet den bisher tatsächlich behandelten Stoff.

Disclaimer 2

Students attending my lectures are often astonished that I present the material in a much livelier form than in this script. The reason for this is the following:

This is a *script*, not a *text book*.

It is meant to *accompany* the lecture, not to *replace* it! We do examples on all concepts, definitions, theorems and algorithms in the lecture, but usually not this script. In this sense, it is **not a good idea** to study the subject soley by reading this script.

1 Recommended Reading

In the order of relevance for this lecture:

- 1. E. Ohlebusch: Bioinformatics Algorithms: Sequence Analysis, Genome Rearrangements, and Phylogenetic Reconstruction. Oldenbusch Verlag, 2013.
- 2. D. Adjeroh, T. Bell, and A. Mukherjee: *The Burrows-Wheeler Transform: Data Compression, Suffix Arrays and Pattern Matching.* Springer, 2008.
- 3. D. Gusfield: Algorithms on Strings, Trees, and Sequences. Cambridge University Press, 1997.

2 Tries

Definition 1. Let $S = \{S_1, \ldots, S_k\}$ be a set of k strings over the alphabet Σ of size $\sigma = |\Sigma|$. A trie on S is a rooted tree S = (V, E) with edge labels from Σ that fulfills the following two constraints:

- 1. $\forall v \in V$: all outgoing edges from v start with a different $a \in \Sigma$.
- 2. For all $S_i \in \mathcal{S}$ there is a leaf ℓ such that S_i is a prefix of the concatenation of the labels on the root-to- ℓ path.
- 3. For all leaves $\ell \in V$ there is a string $S_i \in \mathcal{S}$ such the root-to- ℓ path spells out exactly S_i .

We often deal with *compacted* tries, which can be defined similarly to Def. 1, with the difference that the edge labels are now from Σ^+ , and with an additional constraint:

4. Apart from the root, all nodes have out-degree $\neq 1$.

Tries support existential queries ("Is pattern P one of the strings in S?"), prefix queries ("Which strings in S have P as a prefix?"), and also predecessor queries ("If P is none of the strings in S, which ones are lexicographically closest?"). All of those queries work in a top-down manner, starting at the root and trying to match further characters in P on the way down. The search time of all these operations depends mainly on the way the outgoing edges of a trie node are implemented; this is what we consider next.

Let v be a node in the trie.

- 1. We can simply scan all of v's outgoing edges to find the next character of P. This results in $O(|P| \cdot \sigma)$ search time. The space of the trie is O(n+k) = O(n) for $n = \sum_{i=1}^{k} |s_i|$ being the total size of the strings in S.
- 2. The outgoing edges are implemented as arrays of size σ . This results in optimal O(|P|) search time, but the space shoots up to $O(|P| \cdot \sigma)$.
- 3. We can use either a hash table at every node, or a global hash table using (node, character) pairs as keys. In any case, this results in optimal O(|P|) search time, but only with high probability. Also, predecessor searches are not supported. The space is O(n).
- 4. The outgoing edges are implemented as arrays of size s_v , where s_v denotes the number of v's children. Using binary search over these arrays, this results in total $O(|P|\log \sigma)$ search time. The overall space is O(n) (WHY?). Note that if the trie is dynamic, the arrays can be replaced by balanced binary search trees, yielding the same running times.
- 5. Modifying the previous approach, we can use weight-balanced binary search trees (WB-BST), where each trie node v has a weight w_v equal to the number of leaves below v (hence, the number of strings stored in v's subtree). Then the binary search tree at every trie node v with children v_1, \ldots, v_x is formed as follows (see also Fig. 1). Split the total weights w_{v_1}, \ldots, w_{v_x} exactly in the middle (namely at $\sum w_{v_i}/2$), respecting the lexicographic order of the corresponding characters. This creates the root of the WB-BST (the character touching this middle). The process continues recursively in the left and right children of the root. It is then easy to see that one character comparison in this WB-BST either advances one character in P, or reduces the number of strings to be considered by at least 1/2. Since the latter situation can happen only $\log k$ times, this results in a total search time of $O(|P| + \log k)$, while the space remains linear.
- 6. Here comes the climax! Divide the trie into an upper top tree and several lower bottom trees by declaring all maximally deep nodes with weight at least σ as leaves of the top tree. Then use approach (5) for the nodes in the bottom trees; since their size is now $O(\sigma)$, this results in $O(|P| + \log \sigma)$ time. In the top tree, all branching nodes (meaning thay have at least 2 children) are handled by approach (2) above. Since the number of branching nodes in the top tree are at most $O(n/\sigma)$, this results in O(n) total space for the entire trie. Non-branching nodes of the top tree are simply stored by noting the character of their only outgoing edge. In sum, we get $O(|P| + \log \sigma)$ time, and O(n) space.

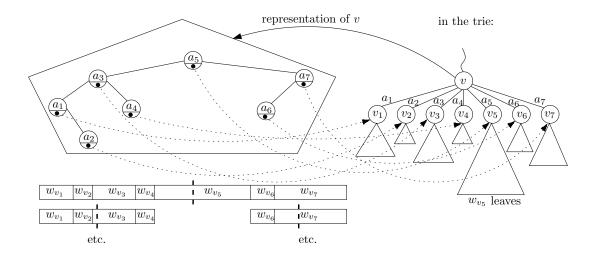


Figure 1: Representation of trie nodes with weight-balanced binary search trees.

3 Suffix Trees and Arrays

In this section we will introduce suffix trees and suffix arrays, which, among many other things, can be used to solve the string matching task: find pattern P of length m in a text T of length n in O(n+m) time. We already know that other methods (Boyer-Moore, e.g.) solve this task in the same time. So why do we need suffix trees?

The advantage of suffix trees and arrays over the other string-matching algorithms (Boyer-Moore, KMP, etc.) is that those structures are an *index* of the text. So, if T is *static* and there are several patterns to be matched against T, the O(n)-task for building the index needs to be done only once, and subsequent matching-tasks can be done in time proportional only to m, and only weakly depends on n ("weakly" meaning, for example, logarithmically). If $m \ll n$, this is a clear advantage over the other algorithms.

Throughout this section, let $T = t_1 t_2 \dots t_n$ be a text over an alphabet Σ of size σ . We use the notation $T_{i...j}$ as an abbreviation of $t_i t_{i+1} \dots t_j$, the substring of T ranging from i to j.

To make this more formal, let P be a pattern of length m. We will be concerned with the two following problems:

Problem 1. Counting: Return the number of matches of P in T. Formally, return the size of $O_P = \{i \in [1, n] : T_{i...i+m-1} = P\}$

Problem 2. Reporting: Return all occurrences of P in T, i. e., return the set O_P .

Definition 2. The i'th suffix of T is the substring $T_{i...n}$ and is denoted by T^i .

3.1 Suffix- and LCP-Arrays

Definition 3. The suffix array A of T is a permutation of $\{1, 2, ..., n\}$ such that A[i] is the i-th smallest suffix in lexicographic order: $T^{A[i-1]} < T^{A[i]}$ for all $1 < i \le n$.

Hence, the suffix array is a compact representation (O(n) space) of the sorted order of all suffixes of a text.

The second array H builds on the suffix array:

Definition 4. The LCP-array H of T is defined such that H[1] = 0, and for all i > 1, H[i] holds the length of the longest common prefix of $T^{A[i]}$ and $T^{A[i-1]}$.

From now on, we assume that T terminates with a \$, and we define \$ to be lexicographically smaller than all other characters in Σ : \$ < a for all $a \in \Sigma$.

3.2 Linear-Time Construction of Suffix Arrays

Now we explain the *induced sorting* algorithm for constructing suffix arrays (called *SAIS* in the literature). Its basic idea is to sort a certain subset of suffixes recursively, and then use this result to *induce* the order of the remaining suffixes.

- Ge Nong, Sen Zhang, Wai Hong Chan: Two Efficient Algorithms for Linear Time Suffix Array Construction. IEEE Trans. Computers **60**(10): 1471–1484 (2011).
- E. Ohlebusch: Bioinformatics Algorithms: Sequence Analysis, Genome Rearrangements, and Phylogenetic Reconstruction. Oldenbusch Verlag, 2013. Chapter 4.1.2.

Definition 5. For $1 \le i < n$, suffix T^i is said to be S-type if $T^i <_{\text{lex}} T^{i+1}$, and L-type otherwise. The last suffix is defined to be S-type. For brevity, we also use the terms S- and L-suffixes for suffixes of the corresponding type.

The type of each suffix can be determined in linear time by a right-to-left scan of T: first, T^n is declared as S-type. Then, for every i from n-1 to 1, T^i is classified by the following rule:

$$T^i$$
 is S-type iff either $t_i < t_{i+1}$, or $t_i = t_{i+1}$ and T^{i+1} is S-type.

We further say that an S-suffix T^i is of type S^* iff T^{i-1} is of type L. (Note that the S-suffixes still include the S^* -suffixes in what follows.)

In A, all suffixes starting with the same character $c \in \Sigma$ form a consecutive interval, called the *c-bucket* henceforth. Observe that in any *c*-bucket, the L-suffixes precede the S-suffixes. Consequently, we can sub-divide buckets into S-type buckets and L-type buckets.

Now the induced sorting algorithm can be explained as follows:

- 1. Sort the S*-suffixes. This step will be explained in more detail below.
- 2. Put the sorted S*-suffixes into their corresponding S-buckets, without changing their order.
- 3. Induce the order of the L-suffixes by scanning A from left to right: for every position i in A, if $T^{A[i]-1}$ is L-type, write A[i]-1 to the current head of the L-type c-bucket ($c=t_{A[i]-1}$), and increase the current head of that bucket by one. Note that this step can only induce "to the right" (the current head of the c-bucket is larger than i).
- 4. Induce the order of the S-suffixes by scanning A from right to left: for every position i in A, if $T^{A[i]-1}$ is S-type, write A[i]-1 to the current end of the S-type c-bucket ($c=t_{A[i]-1}$), and decrease the current end of that bucket by one. Note that this step can only induce "to the left," and might intermingle S-suffixes with S*-suffixes.

It remains to explain how the S*-suffixes are sorted (step 1 above). To this end, we define:

Definition 6. An S*-substring is a substring $T_{i...j}$ with $i \neq j$ of T such that both T^i and T^j are S^* -type, but no suffix in between i and j is also of type S^* .

Let $R_1, R_2, \ldots, R_{n'}$ denote these S*-substrings, and σ' be the number of different S*-substrings. We assign a name $v_i \in [1, \sigma']$ to any such R_i , such that $v_i < v_j$ if $R_i <_{\text{lex}} R_j$ and $v_i = v_j$ if $R_i = R_j$. We then construct a new text $T' = v_1 \ldots v_{n'}$ over the alphabet $[1, \sigma']$, and build the suffix array A' of T' by applying the inducing sorting algorithm recursively to T' if $\sigma' < n'$ (otherwise there is nothing to sort, as then the order of the S*-suffixes is given by the order of the S*-substrings). The crucial property to observe here is that the order of the suffixes in T' is the same as the order of the respective S*-suffixes in T; hence, A' determines the sorting of the S*-suffixes in T. Further, as at most every second suffix in T can be of type S*, the complete algorithm has worst-case running time T(n) = T(n/2) + O(n) = O(n), provided that the naming of the S*-substrings also takes linear time, which is what we explain next.

The naming of the S*-substrings is similar to the inducing of the S-suffixes in the induced sorting algorithm (steps 2–4 above), with the difference that in step 2 we put the *unsorted* S*-suffixes into their corresponding buckets (hence they are only sorted according to their first character). Steps 3 and 4 work exactly as described above. At the end of step 4, we can assign names to the S*-substrings by comparing adjacent S*-suffixes naively until we find a mismatch or reach their end; this takes overall linear time.

Theorem 1. We can construct the suffix array for a text of length n in O(n) time.

3.3 Searching in Suffix Arrays

3.3.1 Exact Searches

- G. Navarro, V. Mäkinen: Compressed Full-Text Indexes. ACM Computing Surveys **39**(1), Article Article No. 2, 2007. Section 3.3.
- E. Ohlebusch: Bioinformatics Algorithms: Sequence Analysis, Genome Rearrangements, and Phylogenetic Reconstruction. Oldenbusch Verlag, 2013. Chapter 5.1.3.

We can use a plain suffix array A to search for a pattern P, using the ideas of binary search, since the suffixes in A are sorted lexicographically and hence the occurrences of P in T form an interval in A. The algorithm below performs two binary searches. The first search locates the starting position s of P's interval in A, and the second search determines the end position r. A counting query returns r-s+1, and a reporting query returns the numbers $A[s], A[s+1], \ldots, A[r]$.

Note that both while-loops in Alg. 2 make sure that either l is increased or r is decreased, so they are both guaranteed to terminate. In fact, in the first while-loop, r always points one position behind the current search interval, and r is decreased in case of equality (when $P = T_{A[q]...\min\{A[q]+m-1,n\}}$). This makes sure that the first while-loop finds the leftmost position of P in A. The second loop works symmetrically. Note further that in the second while-loop it is enough to check for lexicographical equality, as the whole search is done in the interval of A where all suffixes are lexicographically no less than P.

Theorem 2. The suffix array allows to answer counting queries in $O(m \log n)$ time, and reporting queries in $O(m \log n + |O_P|)$ time.

Algorithm 1: function SAsearch $(P_{1...m})$

```
1 l \leftarrow 1; r \leftarrow n + 1;
 2 while l < r do
            \begin{array}{l} q \leftarrow \lfloor \frac{l+r}{2} \rfloor; \\ \text{if } P >_{\text{lex}} T_{A[q]...\min\{A[q]+m-1,n\}} \text{ then} \end{array}
 3
 4
 \mathbf{5}
 6
             else
 7
                  r \leftarrow q;
            end
 8
 9 end
10 s \leftarrow l; l--; r \leftarrow n;
     while l < r do
            q \leftarrow \lceil \tfrac{l+r}{2} \rceil;
12
            if P = \lim_{l \to \infty} T_{A[q]...\min\{A[q]+m-1,n\}} then
13
14
15
                  r \leftarrow q - 1;
16
            end
17
18 end
19 return [s, r];
```

3.3.2 Approximate Searches

• Huynh, T. N., Hon, W. K., Lam, T. W., Sung, W. K.: Approximate string matching using compressed suffix arrays. Theoretical Computer Science, **352**(1), 240-249, 2006.

There are many variants of approximate string matching, both indexed and sequential (i.e., non-indexed). The easiest type of errors are *mismatches*: counting the number of non-equal characters between two strings. This number is usually called the *Hamming distance*.

Smaller distances are usually achieved by the *edit distance*: counting the minimum number of edit operations (insertions, deletions, substitutions) to transform one string into the other. For a given text T of length n and $1 \le i \le n$, we define ed(P, i) as the minimum number of such edit operations on P such that the modified P matches (exactly) at position i in T.

Our algorithmic idea is now to modify the pattern at each possible position and search each of the modified patterns in T, using the suffix array A of T. We focus on the case with one error, but the idea can be generalized to more errors. Since there are exactly $m + m(\sigma - 1) + (m + 1)\sigma$ possible modifications (deletions, substitutions, insertions), this would take $O(m^2\sigma \log n)$ time if we were to use the plain $O(m \log n)$ suffix array search algorithm for each modified pattern.

To speed up the search, we need a small lemma, which allows us to "paste" to suffix arrays intervals together, faster than starting the search from scratch:

Lemma 3. Let $[s_1, e_1]$ be the suffix array interval corresponding to pattern P_1 , and $[s_2, e_2]$ the suffix array interval corresponding to pattern P_2 . Then the suffix array interval for the concatenation P_1P_2 can be found in $O(\log n)$ time.

Proof. We need to find the sub-range [s,e] of $[s_1,e_1]$ such that the suffixes $T^{A[i]}$, $s \le i \le e$, are exactly the suffixes that are prefixed by P_1P_2 . Let $s \le i \le e$. If $T^{A[i]}$ continues with P_2 after the initial $m := |P_1|$ characters, then

$$s_2 \le A^{-1}[A[j] + m_1] \le e_2$$
,

and vice versa. Hence s is the smallest value in $[s_1, e_1]$ such that the above inequality holds, and can be found by a binary search in $O(\log n)$ time. The arguments for e are symmetric.

Using this lemma, the algorithm now works as follows.

```
Algorithm 2: function SAsearchApproximate(P_{1...m})
```

```
1 foreach a \in \Sigma do compute [S_a, E_a], the suffix array interval of a;
                                                                                   // linear scan of A
 2 for 1 \le i \le m do compute [s_i, e_i], the suffix array interval of P_{1...i};
                                                                                       // using Alg. 2
s' \leftarrow 1; e' \leftarrow n;
                                                  // interval of P_{i+1...m} in the following loop
 4 for i = m, ..., 1 do
       // deletion at i > 0:
       from [s', e'] and [s_{i-1}, e_{i-1}], find interval of P_{1...i-1}P_{i+1...m};
                                                                                       // using Lemma 3
5
6
       foreach c \in \Sigma do
           from [s', e'] and [S_c, E_c], find interval [s'', e''] of cP_{i+1...m};
                                                                                       // using Lemma 3
7
           // substitution at i > 0:
           from [s'', e''] and [s_{i-1}, e_{i-1}], find interval of P_{1...i-1}cP_{i+1...m};
                                                                                       // using Lemma 3
8
           // insertion at i \geq 0:
           from [s'', e''] and [s_i, e_i], find interval of P_{1...i}cP_{i+1...m};
                                                                                       // using Lemma 3
9
10
       [s', e'] \leftarrow \text{interval of } P_i P_{i+1...m};
                                                       // using Lemma 3 with [s',e'] and [S_{P_i},E_{P_i}]
12 end
```

The running time is now $O(m\sigma \log n)$; for constant alphabets, this is $O(m \log n)$, which is not too bad. As shown, the algorithm computes all positions i with ed(P, i) = 1 (lines 4, 7, and 8). For Hamming distance, one would just have to consider the case in line 7, and drop all other cases.

3.4 Linear-Time Construction of LCP-Arrays

• E. Ohlebusch: Bioinformatics Algorithms: Sequence Analysis, Genome Rearrangements, and Phylogenetic Reconstruction. Oldenbusch Verlag, 2013. Chapters 4.2.1.

It remains to be shown how the LCP-array H can be constructed in O(n) time. Here, we assume that we are given T and A, the text and the suffix array for T.

We will construct H in text order, which is also the order of the inverse suffix array A^{-1} , the latter defined by $A^{-1}[A[i]] = i$ for all $1 \le i \le n$, which is easily computable from A in linear time. In other words, we aim at filling $H[A^{-1}[i]]$ before $H[A^{-1}[i+1]]$, because in this case we know that H cannot decrease too much, as shown next.

Going from suffix T^i to T^{i+1} , we see that the latter equals the former, but with the first character t_i truncated. Let h = H[i]. Then the suffix T^j , $j = A[A^{-1}[i] - 1]$, has a longest common prefix with T^i of length h. So T^{i+1} has a longest common prefix with T^{j+1} of length h-1. But every suffix T^k that is lexicographically between T^{j+1} and T^{i+1} must have a longest common prefix with

 T^{j+1} that is at least h-1 characters long (for otherwise T^k would not be in lexicographic order). In particular, the suffix right before T^{i+1} in A, which is suffix $T^{A[A^{-1}[i+1]-1]}$, must share a common prefix with S_{i+1} of length at least h-1. Hence, $H[A^{-1}[i+1]] \ge h-1$. We have thus proved the following:

Lemma 4. For all $1 \le i < n$: $H[A^{-1}[i+1]] \ge H[A^{-1}[i]] - 1$.

This gives rise to the following elegant algorithm to construct H:

Algorithm 3: Linear-Time Construction of the LCP-Array

```
1 for i = 1, ..., n do A^{-1}[A[i]] \leftarrow i;

2 h \leftarrow 0, H[1] \leftarrow 0;

3 for i = 1, ..., n do

4 | if A^{-1}[i] \neq 1 then

5 | j \leftarrow A[A^{-1}[i] - 1];

6 | while t_{i+h} = t_{j+h} do h \leftarrow h + 1;

7 | H[A^{-1}[i]] \leftarrow h;

8 | h \leftarrow \max\{0, h - 1\};

9 | end

10 end
```

The linear running time follows because h starts and ends at 0, is always less than n and decreased at most n times in line 8. Hence, the number of times where k is increased in line 6 is bounded by n, so there are at most 2n character comparisons in the whole algorithm. We have proved:

Theorem 5. We can construct the LCP array for a text of length n in O(n) time.

3.4.1 Practical Improvements

Let us do some algorithm engineering on the LCP-array construction algorithm! The problem with this algorithm is its poor locality behavior, resulting in many potential cache misses (4n in total). Our idea is now to rearrange the computations such that in the big for-loop accesses only one array in a random access manner, whereas all other arrays are scanned sequentially. To this end, we first compute a temporary array $\Phi[1,n]$ that at $\Phi[i]$ stores the lexicographic preceding suffix of $T^{[i]}$. (This is exactly the suffix with whom we have to compare $T^{[i]}$ for longest common prefix computation.) Further, in the for-loop we write the computed LCP-values in text order. (This is exactly the order in which they are computed.) The resulting algorithm can be seen in Alg. 4.

Algorithm 4: More Cache-Efficient Linear-Time Construction of the LCP-Array

In total, the algorithm now produces at most 3n cache misses (as opposed to 4n in Alg. 3). The practical running time of Alg. 4 is reported to be 1.5 times faster than Alg. 3.

3.5 Suffix Trees

Definition 7. The suffix tree of T is a compact trie over all suffixes $\{T^1, T^2, \dots, T^n\}$.

The following definitions make it easier to argue about suffix trees and compact tries in general:

Definition 8. Let S = (V, E) be a compact trie.

- For $v \in V$, \overline{v} denotes the concatenation of all path labels from the root of S to v.
- $|\overline{v}|$ is called the string-depth of v and is denoted by d(v).
- S is said to display $\alpha \in \Sigma^*$ iff $\exists v \in V, \beta \in \Sigma^* : \overline{v} = \alpha \beta$.
- If $\overline{v} = \alpha$ for $v \in V, \alpha \in \Sigma^*$, we also write $\overline{\alpha}$ to denote v.
- words(S) denotes all strings in Σ^* that are displayed by S: $words(S) = \{\alpha \in \Sigma^* : S \text{ displays } \alpha\}$

[NB. With these new definitions, an alternative definition of suffix trees would be: "The suffix tree of T is a compact trie that displays exactly the subwords of T."]

It is useful if each suffix ends in a leaf of S. This can again be accomplished by adding a new character $\$ \notin \Sigma$ to the end of T, and build the suffix tree over T\$. This gives a one-to-one correspondence between T's suffixes and the leaves of S, which implies that we can *label the leaves* with a function I by the start index of the suffix they represent: $I(v) = i \iff \overline{v} = T^i$.

The following observations relate the suffix array A and the LCP-array H with the suffix tree S.

Observation 1. If we do a lexicographically-driven depth-first search through S (visit the children in lexicographic order of the first character of their corresponding edge-label), then the leaf-labels seen in this order give the suffix-array A.

To relate the LCP-array H with the suffix tree S, we need to define the concept of lowest common ancestors:

Definition 9. Given a tree S = (V, E) and two nodes $v, w \in V$, the lowest common ancestor of v and w is the deepest node in S that is an ancestor of both v and w. This node is denoted by LCA(v, w).

Observation 2. The string-depth of the lowest common ancestor of the leaves labeled A[i] and A[i-1] is given by the corresponding entry H[i] of the LCP-array, in symbols: $\forall i > 1 : H[i] = d(LCA(\overline{T^{A[i]}}, \overline{T^{A[i-1]}}))$.

An important *implementation detail* is that the edge labels in a suffix tree are represented by a pair (i, j), $1 \le i \le j \le n$, such that $T_{i...j}$ is equal to the corresponding edge label. This ensures that an edge label uses only a constant amount of memory.

From this implementation detail and the fact that S contains exactly n leaves and hence less than n internal nodes, we can formulate the following theorem:

Theorem 6. A suffix tree of a text of length n occupies O(n) space in memory.

3.6 Searching in Suffix Trees

Since the suffix tree is a trie, we can use any of the methods from the section on tries (Sect. 2) for navigation: for example, the *counting* the number of pattern matches can be done in $O(m \log \sigma)$ time (with outgoing-edge representation (4) from the previous chapter on tries): traverse the tree from the root downwards, in each step locating the correct outgoing edge, until P has been scanned completely. More formally, suppose that $P_{1...i-1}$ have already been parsed for some $1 \le i < m$, and our position in the suffix tree S is at node v ($\overline{v} = P_{1...i-1}$). We then find v's outgoing edge e whose label starts with P_i . This takes $O(\log \sigma)$ time. We then compare the label of e character-by-character with $P_{i...m}$, until we have read all of P (i = m), or until we have reached position $j \ge i$ for which $\overline{P_{1...j}}$ is a node v' in S, in which case we continue the procedure at v'. This takes a total of $O(m \log \sigma)$ time. With the more sophisticated representation (6) from Sect. 2 this time can be reduced to $O(m + \log \sigma)$.

Suppose the search procedure has brought us successfully to a node v, or to the incoming edge of node v. We then output the size of S_v , the subtree of S rooted at v. This can be done in constant time, assuming that we have labeled all nodes in S with their subtree sizes. This answers the *counting query*. For *reporting* all positions where the pattern matches, we output the labels of all leaves in S_v (recall that the leaves are labeled with text positions).

Theorem 7. The suffix tree can answer counting queries in $O(m + \log \sigma)$ time, and reporting queries in $O(m + \log \sigma + |O_P|)$ time.

3.7 Linear-Time Construction of Suffix Trees

Assume for now that we are given T, A, and H, and we wish to construct S, the suffix tree of T. We will show in this section how to do this in O(n) time. Later, we will also see how to construct A and H only from T in linear time. In total, this will give us an O(n)-time construction algorithm for suffix trees.

The idea of the algorithm is to insert the suffixes into S in the order of the suffix array: $T^{A[1]}, T^{A[2]}, \ldots, T^{A[n]}$. To this end, let S_i denote the partial suffix tree for $0 \le i \le n$ (S_i is the compact Σ^+ -tree with $words(S_i) = \{T_{A[k]...j} : 1 \le k \le i, A[k] \le j \le n\}$). In the end, we will have $S = S_n$.

We start with S_0 , the tree consisting only of the root (and thus displaying only ϵ). In step i+1, we climb up the rightmost path of S_i (i.e., the path from the leaf labeled A[i] to the root) until we meet the deepest node v with $d(v) \leq H[i+1]$. If d(v) = H[i+1], we simply insert a new leaf x to S_i as a child of v, and label (v,x) by $T^{A[i+1]+H[i+1]}$. Leaf x is labeled by A[i+1]. This gives us S_{i+1} .

Otherwise (i.e., d(v) < H[i+1]), let w be the child of v on S_i 's rightmost path. In order to obtain S_{i+1} , we split up the edge (v, w) as follows.

- 1. Delete (v, w).
- 2. Add a new node y and a new edge (v, y). (v, y) gets labeled by $T_{A[i]+d(v)...A[i]+H[i+1]-1}$.
- 3. Add (y, w) and label it by $T_{A[i]+H[i+1]...A[i]+d(w)-1}$.
- 4. Add a new leaf x (labeled A[i+1]) and an edge (y,x). Label (y,x) by $T^{A[i+1]+H[i+1]}$.

The correctness of this algorithm follows from observations 1 and 2 above. Let us now consider the execution time of this algorithm. Although climbing up the rightmost path could take O(n) time in a single step, a simple amortized argument shows that the running time of this algorithm can be bounded by O(n) in total: each node traversed in step i (apart from the last) is removed from the rightmost path and will not be traversed again for all subsequent steps j > i. Hence, at most 2n nodes are traversed in total.

Theorem 8. We can construct T's suffix tree in linear time from T's suffix- and LCP-array.

4 Repeats

Definition 10. Let $T = t_1 t_2 \dots t_n$ be a text of length n. A triple (i, j, ℓ) with $1 \le i < j \le n - \ell + 1$ is called repeat if $T_{i\dots i+\ell-1} = T_{j\dots j+\ell-1}$.

We look at three different tasks:

- 1. Output all triples (i, j, ℓ) that are a repeat according to Def. 10.
- 2. Output all strings $\alpha \in \Sigma^*$ such that there is a repeat (i, j, ℓ) with $T_{i...i+\ell-1} = \alpha (= T_{j...j+\ell-1})$.
- 3. Like (2), but instead of outputting α just output a pair (l,r) with $\alpha = T_{l...r}$.

Task (1) yields probably many triples (consider $T = a^n$), whereas the returned strings in task (2) may probably be very long. The outuatin task (3) may be seen as a space-efficient representation of (2). Nevertheless, in most cases we are happy with repeats that cannot be extended, as captured in the following section.

4.1 Maximal Repeats

Definition 11. A repeat (i, j, ℓ) is called

• left-maximal, if $t_{i-1} \neq t_{i-1}$.

- right-maximal, if $t_{i+\ell} \neq t_{j+\ell}$.
- maximal, if it is both left- and right-maximal.

To make this definition valid for the text borders, we extend T with $t_0 = \mathcal{L}$ to the left, and with $t_{n+1} = \$$ to the right $(\mathcal{L}, \$ \notin \Sigma)$.

Observation 3. If (i, j, ℓ) is a right-maximal repeat, then there must be an internal node v in T's suffix tree S with $\overline{v} = T_{i...i+\ell-1}$. For otherwise S could not display both $T_{i...i+\ell}$ and $T_{j...j+\ell}$, which must be different due to $t_{i+\ell} \neq t_{j+\ell}$.

Let

$$\mathcal{R}_T = \{ \alpha \in \Sigma^* : \text{ there is a maximal repeat } (i, j, \ell) \text{ with } T_{i...i+\ell-1} = \alpha \}$$

be the set of T's maximal repeat strings (task (2) above). Then the observation above shows that $|\mathcal{R}_T| < n$, as there are only n leaves and hence less than n internal nodes in the suffix tree. Hence, for task (3) we should be able to come up with a O(n)-time algorithm.

It remains to show how left-maximality can be checked efficiently.

Definition 12. Let S be T's suffix tree. A node v in S is called left-diverse if there are at least two leaves b_1 and b_2 below v such that $t_{\ell(b_1)-1} \neq t_{\ell(b_2)-1}$. [Recall that $\ell(\cdot)$ denotes the leaf label (=suffix number)]. Character $t_{\ell(v)-1}$ is called v's left-character.

Lemma 9. A repeat (i, j, ℓ) is maximal iff there is left-diverse node v in T's suffix tree S with $\overline{v} = T_{i...i+\ell-1}$.

Proof. It remains to care about left-maximality.

" \Rightarrow " Let (i, j, ℓ) be maximal. Let v be the node in S with $\overline{v} = T_{i...i+\ell-1}$, which must exist due to right-maximality. Due to left-maximality, we know $t_{i-1} \neq t_{j-1}$. Hence there are two different leaves b_1 and b_2 below with $\ell(b_1) = i$ and $\ell(b_2) = j$. So v is left-diverse.

This yields the following **algorithm** to compute \mathcal{R}_T :

In a depth-first search through S do:

- Let v be the current node.
- If v is a leaf: propagate v's left-character to its parent.
- If v is internal with children v_1, \ldots, v_k :
 - * If one of the v_i 's is left-diverse, or if at least two of the v_i 's have a different left-character: output \overline{v} and propagate "left-diverse" to the parent.
 - * Otherwise, propagate the unique left-character of the v_i 's to the parent.

We formulated the above algorithm to solve task (2) above, but it can easily be adapted to task (1) or (3). Note in particular the *linear* running time for (3). For (1), we also have to propagate lists of positions $L_a(v)$ to the parent, where $L_a(v)$ contains all leaf labels below v that have $a \in \Sigma$ as their left-character. These lists have to be concatenated in linear time (using linked lists with additional pointers to the ends), and in each step we have to output the Cartesian product of $L_a(v)$ and $L_b(v)$ for all $a, b \in \Sigma$, $a \neq b$. The resulting algorithm is still optimal (in an output-sensitive meaning).

4.2 Super-Maximal Repeats

The maximal repeats in Def. 11 can still contain other maximal repeats, as the example T =axybxxyyyaxyb shows (both xy and axyb are maximal repeats, for example). This is prevented by the following definition:

Definition 13. A maximal repeat (i, j, ℓ) is called super-maximal if there is no maximal repeat (i', j', ℓ') such that $T_{i...i+\ell-1}$ is a proper subword of $T_{i'...i'+\ell-1}$.

The algorithmic difference to the previous section is that we only have to consider internal nodes whose children are all *leaves*. Hence, we can also report all k super-maximal repeats in output-optimal time O(n+k).

4.3 Longest Common Substrings

As a last simple example of repeated sequences, consider the following problem: We are given two strings T_1 and T_2 . Our task is to return the longest string $\alpha \in \Sigma^*$ which occurs in both T_1 and T_2 .

Computer-science pioneer Don Knuth conjectured in the late 60's that no linear-time algorithm for this problem can exist. However, he was deeply wrong, as suffix trees make the solution almost trivial: Build a suffix tree S for $T = T_1 \# T_2$. In a DFS through S (where v is the current node), propagate to v's parent from which of the T_i 's the suffixes below v come (either from T_1 , T_2 , or from both). During the DFS, remember the node w of greatest string depth which has suffixes from both T_1 and T_2 below it. In the end, \overline{w} is the solution. Total time is O(n) for $n = |T_1| + |T_2|$.

5 Range Minimum Queries

Range Minimum Queries (RMQs) are a versatile tool for many tasks in exact and approximate pattern matching, as we shall see at various points in this lecture. They ask for the position of the minimum element in a specified sub-array, formally defined as follows.

Definition 14. Given an array H[1,n] of n integers (or any other objects from a totally ordered universe) and two indices $1 \le i \le j \le n$, $\text{RMQ}_H(i,j)$ is defined as the position of the minimum in H's sub-array ranging from i to j, in symbols: $\text{RMQ}_H(i,j) = \operatorname{argmin}_{i < k < j} H[k]$.

We often omit the subscript H if the array under consideration is clear from the context.

Of course, an RMQ can be answered in a trivial manner by scanning H[i, j] (H's sub-array ranging from position i to j) for the minimum each time a query is posed. In the worst case, this takes O(n) query time.

However, if H is static and known in advance, and there are several queries to be answered on-line, it makes sense to preprocess H into an auxiliary data structure (called index or scheme) that allows to answer future queries faster. As a simple example, we could precompute all possible $\binom{n+1}{2}$ RMQs and store them in a table M of size $O(n^2)$ — this allows to answer future RMQs in O(1) time by a single lookup at the appropriate place in M.

We will show in this section that this naive approach can be dramatically improved, as the following proposition anticipates:

Proposition 10. An array of length n can be preprocessed in time O(n) such that subsequent range minimum queries can be answered in optimal O(1) time.

5.1 Linear Equivalence of RMQs and LCAs

Recall the definition of range minimum queries (RMQs): $RMQ_D(\ell, r) = argmin_{\ell \leq k \leq r} D[k]$ for an array D[1, n] and two indices $1 \leq \ell \leq r \leq n$. We show in this section that a seemingly unrelated problem, namely that of computing *lowest common ancestors* (LCAs) in static rooted trees, can be reduced quite naturally to RMQs.

Definition 15. Given a rooted tree T with n nodes, $LCA_T(v, w)$ for two nodes v and w denotes the unique node ℓ with the following properties:

- 1. Node ℓ is an ancestor of both v and w.
- 2. No descendant of ℓ has property (1).

Node ℓ is called the lowest common ancestor of v and w.

The reduction of an LCA-instance to an RMQ-instance works as follows:

- Let r be the root of T with children u_1, \ldots, u_k .
- Define T's inorder tree walk array I = I(T) recursively as follows:
 - If k = 0, then I = [r].
 - If k = 1, then $I = I(T_{u_1}) \circ [r]$.
 - Otherwise, $I = I(T_{u_1}) \circ [r] \circ I(T_{u_2}) \circ [r] \circ \cdots \circ [r] \circ I(T_{u_k})$, where "o" denotes array concatenation. Recall that T_v denotes T's subtree rooted at v.
- Define T's depth array D = D(T) (of the same length as I) such that D[i] equals the tree-depth of node I[i].
- Augment each node v in T with a "pointer" p_v to an arbitrary occurrence of v in I ($p_v = j$ only if I[j] = v).

Lemma 11. The length of I (and of D) is between n (inclusively) and 2n (exclusively).

Proof. By induction on n.

n=1: The tree T consists of a single leaf v, so I=[v] and |I|=1<2n.

 $\leq n \to n+1$: Let r be the root of T with children u_1, \ldots, u_k . Let n_i denote the number of nodes in T_{u_i} . Recall $I = I(T_{u_1}) \circ [r] \circ \cdots \circ [r] \circ I(T_{u_k})$. Hence,

$$|I| = \max(k - 1, 1) + \sum_{1 \le i \le k} |I(T_{u_i})|$$

$$\le \max(k - 1, 1) + \sum_{1 \le i \le k} (2n_i - 1) \qquad \text{(by the induction hypothesis)}$$

$$= \max(k - 1, 1) - k + 2 \sum_{1 \le i \le k} n_i$$

$$\le 1 + 2 \sum_{1 \le i \le k} n_i$$

$$= 1 + 2(n - 1)$$

$$< 2n . \square$$

Here comes the desired connection between LCA and RMQ:

Lemma 12. For any pair of nodes v and w in T, $LCA_T(v, w) = I[RMQ_D(p_v, p_w)]$.

Proof. Consider the inorder tree walk I = I(T) of T. Assume $p_v \leq p_w$ (otherwise swap). Let ℓ denote the LCA of v and w, and let u_1, \ldots, u_k be ℓ 's children. Look at

$$I(T_{\ell}) = I(T_{u_1}) \circ \cdots \circ I(T_{u_x}) \circ [\ell] \circ \cdots \circ [\ell] \circ I(T_{u_y}) \circ \cdots \circ I(T_{u_k})$$

such that $v \in T_{u_x}$ and $w \in T_{u_y}$ ($v = \ell$ or $w = \ell$ can be proved in a similar manner).

Note that $I(T_{\ell})$ appears in I exactly the same order, say from a to b: $I[a,b] = I(T_{\ell})$. Now let d be the tree depth of ℓ . Because ℓ 's children u_i have a greater tree depth than d, we see that D attains its minima in the range [a,b] only at positions i where the corresponding entry I[i] equals ℓ . Because $p_v, p_w \in [a,b]$, and because the inorder tree walk visits ℓ between u_x and u_y , we get the result.

To summarize, if we can solve RMQs in O(1) time using O(n) space, we also have a solution for the LCA-problem within the same time- and space-bounds.

Interestingly, this reduction also works the other way around: a linear-space data structure for O(1) LCAs implies a linear-space data structure for O(1) RMQs. To this end, we need the concept of Cartesian Trees:

Definition 16. Let A[1, n] be an array of size n. The Cartesian Tree C(A) of A is a labelled binary tree, recursively defined as follows:

- Create a root node r and label it with $p = \operatorname{argmin}_{1 < i < n} A[i]$.
- The left and right children of r are the roots of the Cartesian Trees C(A[1, p-1]) and C(A[p+1, n]), respectively (if existent).

Constructing the Cartesian Tree according to this definition requires $O(n^2)$ time (scanning for the minimum in each recursive step), or maybe $O(n \log n)$ time after an initial sorting of A. However, there is also a linear time **algorithm** for constructing $\mathcal{C}(A)$, which we describe next.

Let C_i denote the Cartesian Tree for A[1,i]. Tree C_1 just consists of a single node r labelled with 1. We now show how to obtain C_{i+1} from C_i . Let the *rightmost path* of C_i be the path v_1, \ldots, v_k in C_i , where v_1 is the root, and v_k is the node labelled i. Let l_i be the label of node v_i for $1 \le i \le k$.

To get C_{i+1} , climb up the rightmost path (from v_k towards the root v_1) until finding the first node v_y where the corresponding entry in A is not larger than A[i+1]:

$$A[l_y] \leq A[i+1], \text{and } A[l_z] > A[i+1] \text{ for all } y < z \leq k$$
 .

Then insert a new node w as the right child of v_y (or as the root, if v_y does not exist), and label w with i + 1. Node v_{y+1} becomes the left child of w. This gives us C_{i+1} .

The linear running time of this algorithm can be seen by the following amortized argument: each node is inserted onto the rightmost path exactly once. All nodes on the rightmost path (except the last, v_y) traversed in step i are removed from the rightmost path, and will never be traversed again in steps j > i. So the running time is proportional to the total number of removed nodes from the rightmost path, which is O(n), because we cannot remove more nodes than we insert.

How is the Cartesian Tree related to RMQs?

Lemma 13. Let A and B be two arrays with equal Cartesian Trees. Then $RMQ_A(\ell, r) = RMQ_B(\ell, r)$ for all $1 \le \ell \le r \le n$.

Proof. By induction on n.

n=1: $\mathcal{C}(A)=\mathcal{C}(B)$ consists of a single node labelled 1, and RMQ(1,1) = 1 in both arrays.

 $\leq n \to n+1$: Let v be the root of $\mathcal{C}(A) = \mathcal{C}(B)$ with label μ . By the definition of the Cartesian Tree,

$$\underset{1 \le k \le n}{\operatorname{argmin}} A[k] = \mu = \underset{1 \le k \le n}{\operatorname{argmin}} B[k] . \tag{1}$$

Because the left (and right) children of $\mathcal{C}(A)$ and $\mathcal{C}(B)$ are roots of the same tree, this implies that the Cartesian Trees $\mathcal{C}(A[1,\mu-1])$ and $\mathcal{C}(B[1,\mu-1])$ (and $\mathcal{C}(A[\mu+1,n])$ and $\mathcal{C}(B[\mu+1,n])$) are equal. Hence, by the induction hypothesis,

$$RMQ_A(\ell, r) = RMQ_B(\ell, r) \forall 1 \le \ell \le r < \mu, \text{ and } RMQ_A(\ell, r) = RMQ_B(\ell, r) \forall \mu < \ell \le r \le n. \quad (2)$$

In total, we see that $RMQ_A(\ell, r) = RMQ_B(\ell, r)$ for all $1 \le \ell \le r \le n$, because a query must either contain position μ (in which case, by (1), μ is the answer to both queries), or it must be completely to the left/right of μ (in which case (2) gives what we want).

5.2 O(1)-RMQs with $O(n \log n)$ Space

We already saw that with $O(n^2)$ space, O(1)-RMQs are easy to realize by simply storing the answers to all possible RMQs in a two-dimensional table of size $n \times n$. We show in this section a little trick that lowers the space to $O(n \log n)$.

The basic idea is that it suffices to precompute the answers only for query lengths that are a power of 2. This is because an arbitrary query $RMQ_D(l,r)$ can be decomposed into two overlapping sub-queries of equal length 2^h with $h = \lfloor \log_2(r-l+1) \rfloor$:

$$m_1 = \text{RMQ}_D(l, l + 2^h - 1)$$
 and $m_2 = \text{RMQ}_D(r - 2^h + 1, r)$

The final answer is then given by $\text{RMQ}_D(l,r) = \operatorname{argmin}_{\mu \in \{m_1,m_2\}} D[\mu]$. This means that the precomputed queries can be stored in a two-dimensional table $M[1,n][1,|\log_2 n|]$, such that

$$M[x][h] = \text{RMQ}_D(x, x + 2^h - 1)$$

whenever $x + 2^h - 1 \le n$. Thus, the size of M is $O(n \log n)$. With the identity

$$\begin{split} M[x][h] &= \operatorname{RMQ}_D(x, x + 2^h - 1) \\ &= \operatorname{argmin}\{D[i] : i \in \{x, \dots, x + 2^h - 1\}\} \\ &= \operatorname{argmin}\{D[i] : i \in \{\operatorname{RMQ}_D(x, x + 2^{h-1} - 1), \operatorname{RMQ}_D(x + 2^{h-1}, x + 2^h - 1)\}\} \\ &= \operatorname{argmin}\{D[i] : i \in \{M[x][h - 1], M[x + 2^{h-1}][h - 1]\}\} \ , \end{split}$$

we can use dynamic programming to fill M in optimal $O(n \log n)$ time.

5.3 O(1)-RMQs with O(n) Space

We divide the input array D into blocks B_1, \ldots, B_m of size $s := \frac{\log_2 n}{4}$ (where $m = \lceil \frac{n}{s} \rceil$ denotes the number of blocks): $B_1 = D[1, s]$, $B_2 = D[s + 1, 2s]$, and so on. The reason for this is that any query $RMQ_D(l, r)$ can be decomposed into at most three non-overlapping sub-queries:

- At most one query spanning exactly over several blocks.
- At most two queries completely inside of a block.

We formalize this as follows: Let $i = \lceil \frac{l}{s} \rceil$ and $j = \lceil \frac{r}{s} \rceil$ be the block numbers where l and r occur, respectively. If i = j, then we only need to answer one in-block-query to obtain the final result. Otherwise, $\text{RMQ}_D(l,r)$ is answered by $\text{RMQ}_D(l,r) = \operatorname{argmin}_{\mu \in \{m_1,m_2,m_3\}} D[\mu]$, where the m_i 's are obtained as follows:

- $m_1 = \text{RMQ}_D(l, is)$
- $m_2 = \text{RMQ}_D(is+1,(j-1)s)$ (only necessary if j > i+1)
- $m_3 = \text{RMQ}_D((j-1)s+1,r)$

We first show how to answer queries spanning exactly over several blocks (i.e., finding m_2).

5.3.1 Queries Spanning Exactly over Blocks

Define a new array D'[1, m], such that D'[i] holds the minimum inside of block B_i : $D'[i] = \min_{(i-1)s < j \le is} D[j]$. We then prepare D' for constant-time RMQs with the algorithm from Sect. 5.2, using

$$O(m\log m) = O(\frac{n}{s}\log(\frac{n}{s})) = O(\frac{n}{\log n}\log\frac{n}{\log n}) = O(n)$$

space.

We also define a new array W[1, m], such that W[i] holds the position where D'[i] occurs in D: $W[i] = \operatorname{argmin}_{(i-1)s < j \le is} D[j]$. A query of the form $\operatorname{RMQ}_D(is+1, (j-1)s)$ is then answered by $W[\operatorname{RMQ}_{D'}(i+1, j-1)]$.

5.3.2 Queries Completely Inside of Blocks

We are left with answering "small" queries that lie completely inside of blocks of size s. These are actually more complicated to handle than the "long" queries from Sect. 5.3.1.

The consequence of this is that we only have to precompute in-block RMQs for blocks with different Cartesian Trees, say in a table called P. But how do we know in O(1) time where to look up the results for block B_i ? We need to store a "number" for each block in an array T[1, m], such that T[i] gives the corresponding row in the lookup-table P.

Lemma 14. A binary tree T with s nodes can be represented uniquely in 2s + 1 bits.

Proof. We first label each node in T with a '1' (these are not the same labels as for the Cartesian Tree!). In a subsequent traversal of T, we add "missing children" (labelled '0') to every node labelled '1', such that in the resulting tree T' all leaves are labelled '0'. We then list the 0/1-labels of T' level-wise (i.e., first for the root, then for the nodes at depth 1, then for depth 2,

etc.). This uses 2s + 1 bits, because in a binary tree without nodes of out-degree 1, the number of leaves equals the number of internal nodes plus one.

It is easy to see how to reconstruct T from this sequence. Hence, the encoding is unique. \Box So we perform the following steps:

- 1. For every block B_i , we compute the bit-encoding of $C(B_i)$ and store it in T[i]. Because $s = \frac{\log n}{4}$, every bit-encoding can be stored in a single computer word.
- 2. For every *possible* bit-vector t of length 2s + 1 that describes a binary tree on s nodes, we store the answers to all RMQs in the range [1, s] in a table:

 $P[t][l][r] = \text{RMQ}_B(l,r)$ for some array B of size s whose Cartesian Tree has bit-encoding t

Finally, to answer a query $\text{RMQ}_D(l,r)$ which is completely contained within a block $i = \lceil \frac{l}{s} \rceil = \lceil \frac{r}{s} \rceil$, we simply look up the result in P[T[i]][l - (i-1)s][r - (i-1)s].

To analyze the space, we see that T occupies $m = n/\log n = O(n)$ words. It is perhaps more surprising that also P occupies only a linear number of words, namely order of

$$2^{2s} \cdot s \cdot s = \sqrt{n} \cdot \log^2 n = O(n) .$$

Construction time of the data structures is O(ms) = O(n) for T, and $O(2^{2s} \cdot s \cdot s \cdot s) = O(\sqrt{n} \cdot \log^3 n) = O(n)$ for P (the additional factor s accounts for finding the minimum in each precomputed query interval).

This finishes the description of the algorithm.

6 Lempel-Ziv Compression

6.1 Longest Common Prefixes and Suffixes

An indispensable tool in pattern matching are efficient implementations of functions that compute $longest\ common\ prefixes$ and $longest\ common\ suffixes$ of two strings. We will be particularly interested in longest common prefixes of suffixes from the same string T:

Definition 17. For a text T of length n and two indices $1 \le i, j \le n$, $LCP_T(i, j)$ denotes the length of the longest common prefix of the suffixes starting at position i and j in T, in symbols: $LCP_T(i, j) = \max\{\ell \ge 0 : T_{i...i+\ell-1} = T_{j...j+\ell-1}\}.$

Note that LCP(·) only gives the *length* of the matching prefix; if one is actually interested in the *prefix* itself, this can be obtained by $T_{i...i+\text{LCP}(i,j)-1}$.

Note also that the LCP-array H from Sect. 3.1 holds the lengths of longest common prefixes of lexicographically consecutive suffixes: H[i] = LCP(A[i], A[i-1]). Here and in the remainder of this chapter, A is again the suffix array of text T.

But how do we get the lcp-values of suffixes that are *not* in lexicographic neighborhood? The key to this is to employ RMQs over the LCP-array, as shown in the next lemma (recall that A^{-1} denotes the *inverse suffix array* of T).

Lemma 15. Let $i \neq j$ be two indices in T with $A^{-1}[i] < A^{-1}[j]$ (otherwise swap i and j). Then $LCP(i,j) = H[RMQ_H(A^{-1}[i] + 1, A^{-1}[j])].$

Proof. First note that any common prefix ω of T^i and T^j must be a common prefix of $T^{A[k]}$ for all $A^{-1}[i] \leq k \leq A^{-1}[j]$, because these suffixes are lexicographically between T^i and T^j and must hence start with ω . Let $m = \text{RMQ}_H(A^{-1}[i] + 1, A^{-1}[j])$ and $\ell = H[m]$. By the definition of H, $T_{i...i+\ell-1}$ is a common prefix of all suffixes $T^{A[k]}$ for $A^{-1}[i] \leq k \leq A^{-1}[j]$. Hence, $T_{i...i+\ell-1}$ is a common prefix of T^i and T^j .

Now assume that $T_{i...i+\ell}$ is also a common prefix of T^i and T^j . Then, by the lexicographic order of A, $T_{i...i+\ell}$ is also a common prefix of $T^{A[m-1]}$ and $T^{A[m]}$. But $|T_{i...i+\ell}| = \ell + 1$, contradicting the fact that $H[m] = \ell$ tells us that $T^{A[m-1]}$ and $T^{A[m]}$ share no common prefix of length more than

The above lemma implies that with the inverse suffix array A^{-1} , the LCP-array H, and constant-time RMQs on H, we can answer lcp-queries for arbitrary suffixes in O(1) time.

Now consider the "reverse" problem, that of finding longest common suffixes of prefixes.

Definition 18. For a text T of length n and two indices $1 \le i, j \le n$, $LCS_T(i, j)$ denotes the length of the longest common suffix of the prefixes ending at position i and j in T, in symbols: $LCS_T(i, j) = \max\{k \ge 0 : T_{i-k+1...i} = T_{j-k+1...j}\}.$

For this, it suffices to build the *reverse* string \tilde{T} , and prepare it for lcp-queries as shown before. Then $LCS_T(i,j) = LCP_{\tilde{T}}(n-i+1,n-j+1)$.

6.2 Longest Previous Substring

We now show how to compute an array L of longest previous substrings, where L[i] holds the length of the longest prefix of T^i that has another occurrence in T starting strictly before i.

Definition 19. The longest-previous-substring-array L[1, n] is defined such that $L[i] = \max\{\ell \geq 0 : \exists k < i \text{ with } T_{i...i+\ell-1} = T_{k...k+\ell-1}\}.$

Note that for a character $a \in \Sigma$ which has its first occurrence in T at position i, the above definition correctly yields L[i] = 0, as in this case any position k < i satisfies $T_{i...i-1} = \epsilon = T_{k...k-1}$.

If we are also interested in the *position* of the longest previous substring, we need another array:

Definition 20. The array O[1, n] of previous occurrences is defined by:

$$O[i] = \begin{cases} k & if \ T_{i...i+L[i]-1} = T_{k...k+L[i]-1} \neq \epsilon \\ \bot & otherwise \end{cases}$$

A first approach for computing L is given by the following lemma, which follows directly from the definition of L and LCP:

Lemma 16. For all
$$2 \le i \le n$$
: $L[i] = \max\{LCP(i, j) : 1 \le j < i\}$.

For convenience, from now on we assume that both A and H are padded with 0's at their beginning and end: A[0] = H[0] = A[n+1] = H[n+1] = 0. We further define T^0 to be the empty string ϵ .

Definition 21. Given the suffix array A and an index $1 \le i \le n$ in A, the previous smaller value function $PSV_A(\cdot)$ returns the nearest preceding position where A is strictly smaller, in symbols: $PSV_A(i) = \max\{k < i : A[k] < A[i]\}$. The next smaller value function $NSV(\cdot)$ is defined similarly for nearest succeeding positions: $NSV_A(i) = \min\{k > i : A[k] < A[i]\}$.

The straightforward solution that stores the answers to all PSV-/NSV-queries in two arrays P[1,n] and N[1,n] is sufficient for our purposes. Both arrays can be computed from left to right, setting P[i] to i-1 if A[i-1] < A[i]. Otherwise, continue as follows: if A[P[i-1]] < A[i], set P[i] to P[i-1]. And so on $(P[P[i-1]], P[P[P[i-1]]], \ldots)$, until reaching the beginning of the array (set $P[0] = -\infty$ for handling the border case). By a similar argument we used for constructing Cartesian trees, this algorithms takes O(n) time.

The next lemma shows how PSVs/NSVs can be used to compute L efficiently:

Lemma 17. For all
$$1 \le i \le n$$
, $L[A[i]] = \max(LCP(A[PSV_A(i)], A[i]), LCP(A[i], A[NSV_A(i)]))$.

Proof. Rewriting the claim of Lemma 16 in terms of the suffix array, we get

$$L[A[i]] = \max\{LCP(A[i], A[j]) : A[j] < A[i]\}$$

for all $1 \le i \le n$. This can be split up as

$$L[A[i]] = \max(\max\{LCP(A[i], A[j]) : 0 \le j < i \text{ and } A[j] < A[i]\},$$

$$\max\{LCP(A[i], A[j]) : i < j \le n \text{ and } A[j] < A[i]\}).$$

To complete the proof, we show that $LCP(A[PSV(i)], A[i]) = \max\{LCP(A[i], A[j]) : 0 \le j < i \text{ and } A[j] < A[i]\}$ (the equation for NSV follows similarly). To this end, first consider an index j < PSV(i). Because of the lexicographic order of A, any common prefix of $T^{A[j]}$ and $T^{A[i]}$ is also a prefix of $T^{A[PSV(i)]}$. Hence, the indices j < PSV(i) need not be considered for the maximum. For the indices j with PSV(i) < j < i, we have $A[j] \ge A[i]$ by the definition of PSV. Hence, the maximum is given by LCP(A[PSV(i)], A[i]).

To summarize, we build the array L of longest common substrings in O(n) time as follows:

- Build the suffix array A and the LCP-array H.
- Calculate two arrays P and N such that $PSV_A(i) = P[i]$ and $NSV_A(i) = N[i]$.
- Prepare H for O(1)-RMQs, as $LCP(A[PSV(i)], A[i]) = H[RMQ_H(P[i] + 1, i)]$ by Lemma 15.
- Build L by applying Lemma 17 to all positions i.

The array O of previous occurrences can be filled along with L, by writing to O[A[i]] the value A[P[i]] if $LCP(A[P[i]], A[i]) \ge LCP(A[N[i]], A[i])$, and the value A[N[i]] otherwise.

6.3 Lempel-Ziv Factorization

Although the Lempel-Ziv factorization is usually introduced for data compression purposes (gzip, WinZip, etc. are all based on it), it also turns out to be useful for efficiently finding repetitive structures in texts, due to the fact that it "groups" repetitions in some useful way.

Definition 22. Given a text T of length n, its LZ-decomposition is defined as a sequence of k strings s_1, \ldots, s_k , $s_i \in \Sigma^+$ for all i, such that $T = s_1 s_2 \ldots s_k$, and s_i is either a single letter not occurring in $s_1 \ldots s_{i-1}$, or the longest factor occurring at least twice in $s_1 s_2 \ldots s_i$.

Note that the "overlap" in the definition above exists on purpose, and is not a typo!

We describe the LZ-factorization by a list of k pairs $(b_1, e_1), \ldots, (b_k, e_k)$ such that $s_i = T_{b_i \ldots e_i}$. We now observe that given our array L of longest previous substrings from the previous section, we can obtain the LZ-factorization quite easily in linear time:

Algorithm 5: O(n)-computation of the LZ-factorization

```
1 i \leftarrow 1, e_0 \leftarrow 0;

2 while e_{i-1} < n do

3 | b_i \leftarrow e_{i-1} + 1;

4 | e_i \leftarrow b_i + \max(0, L[b_i] - 1);

5 | ++i;

6 end
```

6.4 A More Space Efficient Algorithm

The idea for a more space efficient algorithm is to factor the string T by the usual pattern matching algorithm. Suppose that a prefix $T_{1,...,i-1}$ is already factored into $s_1s_2...s_{k-1}$. To find the next factor s_k , we start matching t_i in the text T itself, with the help of the suffix array A. Suppose the t_i -interval in A is $[\ell, r]$. Then t_i occurs before position i iff there is a value in $A[\ell, r]$ that is less than i, in particular if the minimum in $A[\ell, r]$ is less than i. This can be efficiently checked by range minimum queries over A. We then continue and find the $t_i t_{i+1}$ -interval in A, and so on, until the minimum in the suffix array range equals i. Since a single search step in the suffix array takes $O(\log n)$ time, the whole algorithm takes $O(n \log n)$ time.

7 Burrows Wheeler Transformation

The Burrows-Wheeler Transformation was originally invented for text *compression*. Nonetheless, it was noted soon that it is also a very useful tool in text *indexing*.

7.1 The Transformation

Definition 23. Let $T = t_1 t_2 \dots t_n$ be a text of length n, where $t_n = \$$ is a unique character lexicographically smaller than all other characters in Σ . Then the i-th cyclic shift of T is $T_{i...n}T_{1...i-1}$. We denote it by $T^{(i)}$.

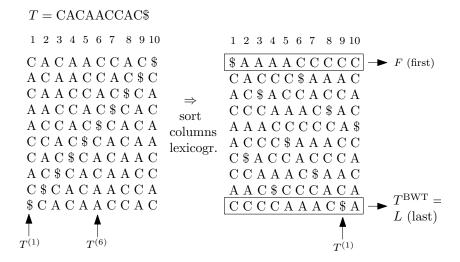
Example 1.

$$T = {\rm CACAACCAC\$}$$
 $T^{(6)} = {\rm CCAC\$CAC\$A}$

The Burrows-Wheeler-Transformation (BWT) is obtained by the following steps:

- 1. Write all cyclic shifts $T^{(i)}$, $1 \le i \le n$, column-wise next to each other.
- 2. Sort the columns lexicographically.
- 3. Output the last row. This is T^{BWT} .

Example 2.



The text T^{BWT} in the last row is also denoted by L (last), and the text in the first row by F (first). Note:

- Every row in the BWT-matrix is a permutation of the characters in T.
- Row F is a sorted list of all characters in T.
- In row $L=T^{\text{BWT}}$, similar characters are grouped together. This is why T^{BWT} can be compressed more easily than T.

7.2 Construction of the BWT

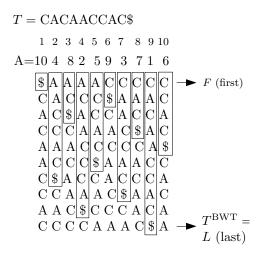
The BWT-matrix needs not to be constructed explicitly in order to obtain T^{BWT} . Since T is terminated with the special character \$, which is lexicographically smaller than any $a \in \Sigma$, the shifts $T^{(i)}$ are sorted exactly like T's suffixes. Because the last row consists of the characters preceding the corresponding suffixes, we have

$$T^{\text{BWT}}[i] = t_{A[i]-1} (= T^{(A[i])}[n]) ,$$

where A denotes again T's suffix array, and t_0 is defined to be t_n (read T cyclically!). Because the suffix array can be constructed in linear time (Thm. 5), we get:

Theorem 18. The BWT of a text length-n text over an integer alphabet can be constructed in O(n) time.

Example 3.



7.3 The Reverse Transformation

The amazing property of the BWT is that it is not a random permutation of T's letters, but that it can be $transformed\ back$ to the original text T. For this, we need the following definition:

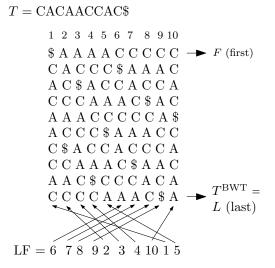
Definition 24. Let F and L be the strings resulting from the BWT. Then the last-to-front mapping LF is a function LF: $[1, n] \rightarrow [1, n]$, defined by

$$\operatorname{LF}(i) = j \iff T^{(A[j])} = (T^{(A[i])})^{(n)} (\iff A[j] = A[i] + 1) \ .$$

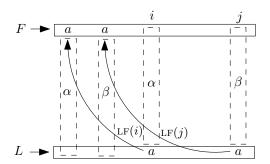
(Remember that $T^{(A[i])}$ is the i'th column in the BWT-matrix, and $(T^{(A[i])})^{(n)}$ is that column rotated by one character downwards.)

Thus, $\mathrm{LF}(i)$ tells us the position in F where L[i] occurs.

Example 4.



Observation 4. Equal characters preserve the same order in F and L. That is, if L[i] = L[j] and i < j, then LF(i) < LF(j). To see why this is so, recall that the BWT-matrix is sorted lexicographically. Because both the LF(i)'th and the LF(j)'th column start with the same character a = L[i] = L[j], they must be sorted according to what follows this character a, say a and b. But since a < b, we know a < b, hence a < b.



This observation allows us to compute the LF-mapping without knowing the suffix array of T.

Definition 25. Let T be a text of length n over an alphabet Σ , and let $L = T^{\text{BWT}}$ be its BWT.

- Define $C: \Sigma \to [1, n]$ such that C(a) is the number of occurrences in T of characters that are lexicographically smaller than $a \in \Sigma$.
- Define OCC: $\Sigma \times [1, n] \to [1, n]$ such that OCC(a, i) is the number of occurrences of a in L's length-i-prefix L[1, i].

Lemma 19. With the definitions above,

$$LF(i) = C(L[i]) + OCC(L[i], i) .$$

Proof: Follows immediately from the observation above.

This gives rise to the following algorithm to recover T from $L = T^{\text{BWT}}$.

- 1. Scan $L = T^{\text{BWT}}$ and compute array $C[1, \sigma]$.
- 2. Compute the first row F from C; as F consists of all characters in L sorted lexicographically, this step is trivial.
- 3. Compute OCC(L[i], i) for all $1 \le i \le n$.
- 4. Recover $T = t_1 t_2 \dots t_n$ from right to left: we know that $t_n = \$$, and the corresponding cyclic shift $T^{(n)}$ appears in column 1 in BWT. Hence, $t_{n-1} = L[1]$. Shift $T^{(n-1)}$ appears in column LF(1), and thus $t_{n-2} = L[\text{LF}(1)]$. This continues until the whole text has been recovered:

$$t_{n-i} = L[\underbrace{\mathrm{LF}(\mathrm{LF}(\ldots(\mathrm{LF}(1))\ldots))}_{i-1 \text{ applications of LF}}]$$

Example 5.

$$C = \stackrel{\$}{0} \stackrel{\text{A C}}{1} \stackrel{\text{C}}{5}$$
 $F = \$ \text{ A A A A C C C C C }$
 $L = \text{C C C C A A A C \$ A}$
 $\text{OCC}(L[i], i) = 1 \ 2 \ 3 \ 4 \ 1 \ 2 \ 3 \ 5 \ 1 \ 4$

$$T_{n} = \$, k = 1$$

$$L[1] = C \Rightarrow T_{n-1} = C, k = LF(1) = 6$$

$$L[6] = C \Rightarrow T_{n-2} = A, k = LF(6) = 3$$

$$L[3] = C \Rightarrow T_{n-3} = C, k = LF(3) = 8$$

$$L[8] = C \Rightarrow T_{n-4} = C, k = LF(8) = 10$$

$$L[10] \text{ etc.}$$

7.4 Compression

Storing T^{BWT} plainly needs the same space as storing the original text T. However, because equal characters are grouped together in T^{BWT} , we can compress T^{BWT} in a second stage.

We can directly exploit that T^{BWT} consists of many equal-letter runs. Each such run a^{ℓ} can be encoded as a pair (a, ℓ) with $a \in \Sigma, \ell \in [1, n]$. This is known as run-length encoding.

Example 6.

$$T^{\text{\tiny BWT}} = \overset{1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\ 10}{\text{\tiny CCCCAAAC$$A$}}$$

$$\Rightarrow \text{RLE}(T^{\text{\tiny BWT}}) = (\text{C},4), (\text{A},4), (\text{C},1), (\$,1), (\text{A},1)$$

A different possibility for compression is to proceed in two steps: first, we perform a move-to-front encoding of the BWT. Then, we review different methods for compressing the output of the move-to-front algorithm. Both steps are explained in the following sections.

7.4.1 Move-to-front (MTF)

- Initialize a list Y containing each character in Σ in alphabetic order.
- In a left-to-right scan of T^{BWT} , (i = 1, ..., n), compute a new array R[1, n]:
 - Write the position of character $T_i^{\mbox{\tiny BWT}}$ in Y to R[i].
 - Move character $T_i^{\mbox{\tiny BWT}}$ to the front of Y.

MTF is easy to reverse.

Observation 5. MTF produces "many small" numbers for equal characters that are "close together" in T^{BWT} . These can be compressed using an order-0 compressor, as explained next.

7.4.2 0-Order Compression

We looked at unary, Elias- γ and Elias- δ code.

8 Backwards Search and FM-Indices

We are now going to explore how the BW-transformed text is helpful for (indexed) pattern matching. Indices building on the BWT are called FM-indices, most likely in honor of their inventors P. Ferragina and G. Manzini. From now on, we shall always assume that the alphabet Σ is goodnatured: $\sigma = o(n/\log \sigma)$.

8.1 Model of Computation and Space Measurement

For the rest of this lecture, we work with the word-RAM model of computation. This means that we have a processor with registers of width w (usually w = 32 or w = 64), where usual arithmetic operations (additions, shifts, comparisons, etc.) on w-bit wide words can be computed in constant time. Note that this matches all current computer architectures. We further assume that n, the input size, satisfies $n \leq 2^w$, for otherwise we could not even address the whole input.

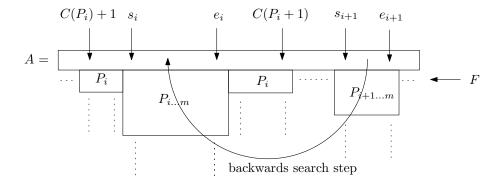
From now on, we measure the space of all data structures in *bits* instead of words, in order to be able to differentiate between the various text indexes. For example, an array of n numbers from the range [1, n] occupies $n\lceil \log n \rceil$ bits, as each array cell stores a binary number consisting of $\lceil \log n \rceil$ bits. As another example, a length-n text over an alphabet of size σ occupies $n\lceil \log \sigma \rceil$ bits. In this light, all text indexes we have seen so far (suffix trees, suffix arrays, suffix trays) occupy $O(n \log n + n \log \sigma)$ bits. Note that the difference between $\log n$ and $\log \sigma$ can be quite large, e.g., for the human genome with $\sigma = 4$ and $n = 3.4 \times 10^9$ we have $\log \sigma = 2$, whereas $\log n \approx 32$. So the suffix array occupies about 16 times more memory than the genome itself!

8.2 Backward Search

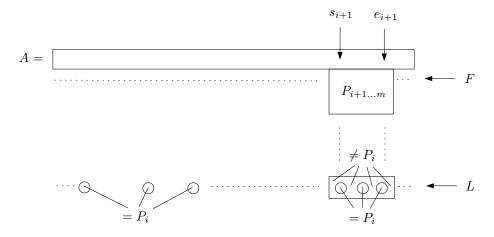
We first focus our attention on the *counting problem* (p. 3); i.e., on finding the number of occurrences of a pattern $P_{1...m}$ in $T_{1...n}$. Recall from Chapter 7 that

- A denotes T's suffix array.
- L/F denotes the first/last row of the BWT-matrix.
- $LF(\cdot)$ denotes the last-to-front mapping.
- C(a) denotes the number of occurrences in T of characters lexicographically smaller than $a \in \Sigma$.
- OCC(a, i) denotes the number of occurrences of a in L[1, i].

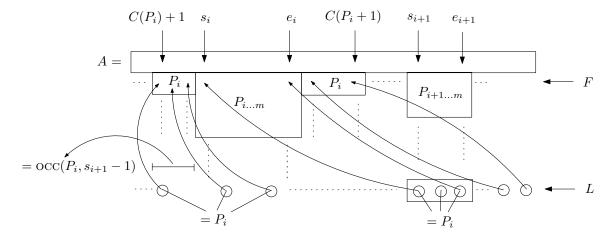
Our aim is identify the interval of P in A by searching P from right to left (= backwards). To this end, suppose we have already matched $P_{i+1...m}$, and know that the suffixes starting with $P_{i+1...m}$ form the interval $[s_{i+1}, e_{i+1}]$ in A. In a backwards search step, we wish to calculate the interval $[s_i, e_i]$ of $P_{i...m}$. First note that $[s_i, e_i]$ must be a sub-interval of $[C(P_i) + 1, C(P_i + 1)]$, where $(P_i + 1)$ denotes the character that follows P_i in Σ .



So we need to identify, from those suffixes starting with P_i , those which continue with $P_{i+1...m}$. Looking at row L in the range from s_{i+1} to e_{i+1} , we see that there are exactly $e_i - s_i + 1$ many positions $j \in [s_{i+1}, e_{i+1}]$ where $L[j] = P_i$.



From the BWT decompression algorithm, we know that characters preserve the same order in F and L. Hence, if there are x occurrences of P_i before s_{i+1} in L, then s_i will start x positions behind $C(P_i) + 1$. This x is given by $OCC(P_i, s_{i+1} - 1)$. Likewise, if there are y occurrences of P_i within $L[s_{i+1}, e_{i+1}]$, then $e_i = s_i + y - 1$. Again, y can be computed from the OCC-function.



Algorithm 6: function backwards-search($P_{1...m}$)

```
1 s \leftarrow 1; e \leftarrow n;

2 for i = m \dots 1 do

3 \begin{vmatrix} s \leftarrow C(P_i) + \text{OCC}(P_i, s - 1) + 1; \\ e \leftarrow C(P_i) + \text{OCC}(P_i, e); \\ \text{if } s > e \text{ then} \\ | \text{return "no match"}; \\ \text{7} & | \text{end} \\ \text{8 end} \\ \text{9 return } [s, e]; \\ \end{vmatrix}
```

This gives rise to the following, elegant algorithm for backwards search:

The reader should compare this to the "normal" binary search algorithm in suffix arrays. Apart from matching backwards, there are two other notable deviations:

- 1. The suffix array A is not accessed during the search.
- 2. There is no need to access the input text T.

Hence, T and A can be deleted once T^{BWT} has been computed. It remains to show how array C and occ are implemented. Array C is actually very small and can be stored plainly using $\sigma \log n$ bits. Because $\sigma = o(n/\log n)$, |C| = o(n) bits. For occ, we have several options that are explored in the rest of this chapter. This is where the different FM-Indices deviate from each other. In fact, we will see that there is a natural trade-off between time and space: using more space leads to a faster computation of the occ-values, while using less space implies a higher query time.

Theorem 20. With backwards search, we can solve the counting problem in $O(m \cdot t_{OCC})$ time, where t_{OCC} denotes the time to answer an $OCC(\cdot)$ -query.

8.3 First Ideas for Implementing Occ

For answering OCC(c, i), there are two simple possibilities:

- 1. Scan L every time an $OCC(\cdot)$ -query has to be answered. This occupies no space, but needs O(n) time for answering a single $OCC(\cdot)$ -query, leading to a total query time of O(mn) for backwards search.
- 2. Store all answers to OCC(c, i) in a two-dimensional table. This table occupies $O(n\sigma \log n)$ bits of space, but allows constant-time $OCC(\cdot)$ -queries. Total time for backwards search is *optimal* O(m).

For more more practical implementation between these two extremes, let us define the following:

Definition 26. Given a bit-vector B[1, n], $rank_1(B, i)$ counts the number of 1's in B's prefix B[1, i]. Operation $rank_0(B, i)$ is defined similarly for 0-bits.

¹More precisely, we should say $\sigma[\log n]$ bits, but we will usually omit floors and ceilings from now on.

In the lecture "Advanced Data Structures" (every winter semester) it is shown that a bit-vector B, together with additional information for constant-time rank-operations, can be stored in n+o(n) bits. This can be used as follows for implementing OCC: For each character $c \in \Sigma$, store an indicator bit vector $B_c[1, n]$ such that $B_c[i] = 1$ iff L[i] = c. Then

$$OCC(c, i) = rank_1(B_c, i)$$
.

The total space for all σ indicator bit vectors is thus $\sigma n + o(\sigma n)$ bits. Note that for reporting queries, we still need the suffix array to output the values in A[s,e] after the backwards search.

Theorem 21. With backwards search and constant-time rank operations on bit-vectors, we can answer counting queries in optimal O(m) time. The space (in bits) is $\sigma n + o(\sigma n) + \sigma \log n$.

Example 7.

 $L = { \begin{array}{cccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ CCCCAAAC\$A & & & & & \\ \end{array} }$

 $B_{\$} = 0000000010$ $B_A = 0000111001$ $B_C = 1111000100$

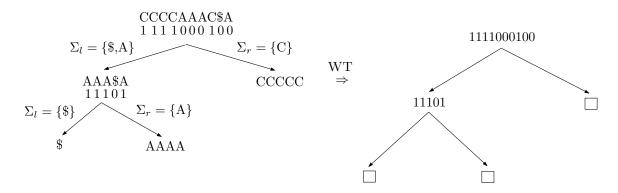
8.4 Wavelet Trees

Armed with constant-time *rank*-queries, we now develop a more space-efficient implementation of the occ-function, sacrificing the optimal query time. The idea is to use a *wavelet tree* on the BW-transformed text.

The wavelet tree of a sequence L[1,n] over an alphabet $\Sigma[1,\sigma]$ is a balanced binary search tree of height $O(\log \sigma)$. It is obtained as follows. We create a root node v, where we divide Σ into two halves $\Sigma_l = \Sigma[1, \lceil \frac{\sigma}{2} \rceil]$ and $\Sigma_r = \Sigma[\lceil \frac{\sigma}{2} \rceil + 1, \sigma]$ of roughly equal size. Hence, Σ_l holds the lexicographically first half of characters of Σ , and Σ_r contains the other characters. At v we store a bit-vector B_v of length n (together with data structures for O(1) rank-queries), where a '0' of position i indicates that character L[i] belongs to Σ_l , and a '1' indicates the it belongs to Σ_r . This defines two (virtual) sequences L_v and R_v , where L_v is obtained from L by concatenating all characters L[i] where $B_v[i] = 0$, in the order as they appear in L. Sequence R_v is obtained in a similar manner for positions i with $B_v[i] = 1$. The left child l_v is recursively defined to be the root of the wavelet tree for L_v , and the right child r_v to be the root of the wavelet tree for R_v . This process continues until a sequence consists of only one symbol, in which case we create a leaf.

Example 8.

L=CCCCAAAC\$A $\Sigma = \{\$,A,C\}$



Note that the sequences themselves are not stored explicitly; node v only stores a bit-vector B_v and structures for O(1) rank-queries.

Theorem 22. The wavelet tree for a sequence of length n over an alphabet of size σ can be stored in $n \log \sigma \times (1 + o(1))$ bits.

Proof: We concatenate all bit-vectors at the same depth d into a single bit-vector B_d of length n, and prepare it for O(1)-rank-queries. Hence, at any level, the space needed is n + o(n) bits. Because the depth of the tree is $\lceil \log \sigma \rceil$ the claim on the space follows. In order to "know" the sub-interval of a particular node v in the concatenated bit-vector B_d at level d, we can store two indices α_v and β_v such that $B_d[\alpha_v, \beta_v]$ is the bit-vector B_v associated to node v. This accounts for additional $O(\sigma \log n)$ bits. Then a rank-query is answered as follows $(b \in \{0, 1\})$:

$$rank_b(B_v, i) = rank_b(B_d, \alpha_v + i - 1) - rank_b(B_d, \alpha_v - 1) ,$$

where it is assumed that $i \leq \beta_v - \alpha_v + 1$, for otherwise the result is not defined.

How does the wavelet tree help for implementing the OCC-function? Suppose we want to compute OCC(c,i), i.e., the number of occurrences of $c \in \Sigma$ in L[1,i]. We start at the root r of the wavelet tree, and check if c belongs to the first or to the second half of the alphabet. In the first case, we know that the c's are "stored" in the left child of the root, namely L_r . Hence, the number of c's in L[1,i] corresponds to the number of c's in $L_r[1, rank_0(B_r, i)]$. If, on the hand, c belongs to the second half of the alphabet, we know that the c's are "stored" in the subsequence R_r that corresponds to the right child of r, and hence compute the number of occurrences of c in $R_r[1, rank_1(B_r, i)]$ as the number of c's in L[1, i]. This leads to the following recursive procedure for computing OCC(c, i), to be invoked with WT-OCC $(c, i, 1, \sigma, r)$, where r is the root of the wavelet tree. (Recall that we assume that the characters in Σ can be accessed as $\Sigma[1], \ldots, \Sigma[\sigma]$.)

Due to the depth of the wavelet tree, the time for $\operatorname{WT-occ}(\cdot)$ is $O(\log \sigma)$. This leads to the following theorem.

Theorem 23. With backward-search and a wavelet-tree on T^{BWT} , we can answer counting queries in $O(m \log \sigma)$ time. The space (in bits) is

$$\underbrace{O(\sigma \log n)}_{|C| + \text{ space for } \alpha_v\text{'s}} + \underbrace{n \log \sigma}_{\text{wavelet tree}} + \underbrace{o(n \log \sigma)}_{\text{rank data structure}}$$

Algorithm 7: function $\forall \mathsf{T-occ}(c,i,\sigma_l,\sigma_r,v)$

```
1 if \sigma_l = \sigma_r then
2 | return i;
3 end
4 \sigma_m = \lfloor \frac{\sigma_l + \sigma_r}{2} \rfloor;
5 if c \leq \Sigma[\sigma_m] then
6 | return WT-occ(c, rank_0(B_v, i), \sigma_l, \sigma_m, l_v);
7 else
8 | return WT-occ(c, rank_1(B_v, i), \sigma_m + 1, \sigma_r, r_v);
9 end
```

8.5 Sampling the Suffix Array

If we also want to solve the reporting problem (outputting all starting positions of P in T, see p. 3), we do need the actual suffix array values. A simple way to solve this is to sample regular text positions in A, and use the LF-function to recover unsampled values. More precisely, we choose a sampling parameter s, and in an array A' we write the values $1, s, 2s, 3s, \ldots$ in the order as they appear in the full suffix array A. Array A' takes $O(n/s\log n)$ bits. In a bit-vector S of length n, we mark the sampled suffix array values with a '1', and augment S with constant-time rank information. Now let i be a position for which we want to find the value of A[i]. We first check if S[i] = 1, and if so, return the value $A'[rank_1(S,i)]$. If not (S[i] = 0), we go to position LF(i) in time t_{LF} , making use of the fact that if A[i] = j, then A[LF(i)] = j - 1. This processes continues until we hit a sampled position d, which takes at most s steps. We then add the number of times we followed LF to the sampled value of A'[d]; the result is A[i]. The overall time for this process is $O(s \cdot t_{OCC})$ for a single suffix array value. Choosing $s = \log_{\sigma} n$ and wavelet trees for implementing the OCC-function, we get an index of $O(n\log \sigma)$ space, $O(m\log \sigma)$ counting time, and $O(k\log n)$ reporting time for k occurrences to be reported.

9 Simulation of Suffix Trees

So far, we have seen compressed text indices that have only one functionality: locating all occurrences of a search pattern P in a text T. In some cases, however, more functionality is required. From other courses you might know that many sequence-related problems are solved efficiently with suffix trees (e.g., computing tandem repeats, MUMs, ...). However, the space requirement of a suffix tree is huge: it is at least 20–40 times higher then the space of the text itself, using very proprietary implementations that support only a very small number of all conceivable suffix tree operations. In this chapter, we present a generic approach that allows for the simulation of all suffix tree operations, by using only compressed data structures. More specifically, we will build on the compressed suffix array from Chapter 8, and show how all suffix tree operations can be simulated by computations on suffix array intervals (the same intervals that we used for suffix trays). Space-efficient data structures that facilitate these computations will be handled in subsequent chapters.

9.1 Basic Concepts

The reader is encouraged to recall the definitions from Sect. 3.5, in particular Def. 8. From now on, we regard the suffix tree as an abstract data type that supports the following operations.

Definition 27. A suffix tree S supports the following operations.

- ROOT(): returns the root of the suffix tree.
- ISLEAF(v): true iff v is a leaf.
- Leaflabel(v): returns l(v) if v is a leaf, and null otherwise.
- ISANCESTOR(v, w): true iff v is an ancestor of w.
- SDEPTH(v): returns d(v), the string-depth of v.
- Count(v): the number of leaves in S_v .
- Parent (v): the parent node of v.
- FIRSTCHILD(v): the alphabetically first child of v.
- NextSibling (v): the alphabetically next sibling of v.
- LCA(v): the lowest common ancestor of v and w.
- CHILD(v, a): node w such that the edge-label of (v, w) starts with $a \in \Sigma$.
- EDGELABEL(v, i) the i'th letter on the edge (PARENT(v), v).

We recall from from previous chapters that A denotes the suffix array, H the LCP-array, and RMQ a range minimum query. Because we will later be using *compressed* data structures (which not necessarily have constant access times), we use variables $t_{\rm SA}$, $t_{\rm LCP}$ and $t_{\rm RMQ}$ for the access time to the corresponding array/function. E. g., with uncompressed (plain) arrays, we have $t_{\rm SA} = t_{\rm LCP} = t_{\rm RMQ} = O(1)$, while with the sampled suffix array from Sect. 8.5 we have $t_{\rm SA} = O(\log n)$.

We represent a suffix tree node v by the interval $[v_{\ell}, v_r]$ such that $A[v_{\ell}], \ldots, A[v_r]$ are exactly the labels of the leaves below v. For such a representation we have the following basic lemma (from now on we assume H[1] = H[n+1] = -1 for an easy handling of border cases):

Lemma 24. Let $[v_{\ell}, v_r]$ be the interval of an internal node v. Then

- (1) For all $k \in [v_{\ell} + 1, v_r] : H[k] \ge d(v)$.
- (2) $H[v_{\ell}] < d(v)$ and $H[v_r + 1] < d(v)$.
- (3) There is a $k \in [v_{\ell} + 1, v_r]$ with H[k] = d(v).

Proof: Condition (1) follows because all suffixes $T^{A[k]}$, $k \in [v_{\ell}, v_r]$, have \overline{v} as their prefix, and hence $H[k] = \text{LCP}(T^{A[k]}, T^{A[k-1]}) \geq |\overline{v}| = d(v)$ for all $k \in [v_{\ell} + 1, v_r]$. Property (2) follows because otherwise suffix $T^{A[v_{\ell}]}$ or $T^{A[v_r+1]}$ would start with \overline{v} , and hence leaves labeled $A[v_{\ell}]$ or $A[v_r+1]$ would also be below v. For proving property (3), for the sake of contradiction assume H[k] > d(v) for all $k \in [v_{\ell} + 1, v_r]$. Then all suffixes $T^{A[k]}$, $k \in [v_{\ell}, v_r]$, would start with $\overline{v}a$ for some $a \in \Sigma$.

Hence, v would only have one outgoing edge (whose label starts with a), contradicting the fact that the suffix tree is compact (has no unary nodes).

As a side remark, this is actually an "if and only if" statement, as every interval satisfying the three conditions from Lemma 24 corresponds to an internal node.

Definition 28. Let $[v_{\ell}, v_r]$ be the interval of an internal node v. Any position $k \in [v_{\ell} + 1, v_r]$ satisfying point (3) in Lemma 24 is called a d(v)-index of v.

Our aim is to simulate all suffix tree operations by computations on suffix intervals: given the interval $[v_{\ell}, v_r]$ corresponding to node v, compute the interval of w = f(v) from the values v_{ℓ} and v_r alone, where f can be any function from Def. 27; e.g., f = PARENT. We will see that most suffix tree operations follow a generic approach: first locate a d(w)-index p of w, and then search for the (yet unknown) delimiting points w_{ℓ} and w_r of w's suffix interval. For this latter task (computation of w_{ℓ} and w_r from p), we also need the previous- and next-smaller-value functions as already defined in Def. 21 in Sect. 6.2. However, this time we define them to work on the LCP-array:

Definition 29. Given the LCP-array H and an index $1 \le i \le n$, the previous smaller value function $PSV_H(i) = \max\{k < i : H[k] < H[i]\}$. The next smaller value function $NSV_H(i)$ is defined similarly for succeeding positions: $NSV_H(i) = \min\{k > i : H[k] < H[i]\}$.

We use t_{PNSV} to denote the time to compute a value $NSV_H(i)$ or $PSV_H(i)$. In what follows, we often use simply PSV and NSV instead of PSV_H and NSV_H , implicitly assuming that array H is the underlying array. The following lemma shows how these two functions can be used to compute the delimiting points w_ℓ and w_r of w's suffix interval:

Lemma 25. Let p be a d(w)-index of an internal node w. Then $w_{\ell} = PSV(p)$, and $w_r = NSV(p) - 1$.

Proof: Let l = PSV(p), and r = NSV(p). We must show that all three conditions in Lemma 27 are satisfied by [l, r-1]. Because H[l] < H[p] by the definition of PSV, and likewise H[r] < H[p], point (1) is clear. Further, because l and r are the closest positions where H attains a smaller value, condition (2) is also satisfied. Point (3) follows from the assumption that p is a d(w)-index. We thus conclude that $w_l = l$ and $w_r = r - 1$.

9.2 Suffix Tree Operations

We now step through the operations from Def. 27 and show how they can be simulated by computations on the suffix array intervals. Let $[v_{\ell}, v_r]$ denote the interval of an arbitrary node v. The most easy operations are:

- ROOT(): returns the interval [1, n].
- IsLeaf(v): true iff $v_{\ell} = v_r$.
- Count(v): returns $v_r v_\ell + 1$.
- ISANCESTOR(v, w): true iff $v_{\ell} \leq w_r \leq v_r$.

Time is O(1) for all four operations.

• Leaflabel(v): If $v_{\ell} \neq v_r$, return null. Otherwise, return $A[v_{\ell}]$ in $O(t_{SA})$ time.

- SDEPTH(v): If $v_{\ell} = v_r$, return $n A[v_{\ell}] + 1$ in time $O(t_{\text{SA}})$, as this is the length of the $A[v_{\ell}]$ 'th suffix. Otherwise from Lemma 24 we know that d(v) is the minimum LCP-value in $H[v_{\ell} + 1, v_r]$. We hence return $H[\text{RMQ}_H(v_{\ell} + 1, v_r)]$ in time $O(t_{\text{RMQ}} + t_{\text{LCP}})$.
- PARENT(v): Because S is a compact tree, either $H[v_{\ell}]$ or $H[v_r+1]$ equals the string-depth of the parent-node, whichever is greater. Hence, we first set $p = \operatorname{argmax}\{H[k] : k \in \{v_{\ell}, v_r+1\}\}$, and then, by Lemma 25, return [PSV(p), NSV(p) 1]. Time is $O(t_{LCP} + t_{PNSV})$.
- FIRSTCHILD(v): If v is a leaf, return NULL. Otherwise, locate the first d(v)-value in $H[v_{\ell}, v_r]$ by $p = \text{RMQ}_H(v_{\ell} + 1, v_r)$. Here, we assume that RMQ returns the position of the *leftmost* minimum, if it is not unique. The final result is $[v_{\ell}, p-1]$, and the total time is $O(t_{\text{RMQ}})$.
- NEXTSIBLING(v): First, compute v's parent as w = Parent(v). Now, if $v_r = w_r$, return NULL, since v does not have a next sibling in this case. If $w_r = v_r + 1$, then v's next sibling is a leaf, so we return $[w_r, w_r]$. Otherwise, try to locate the first d(w)-value after $v_r + 1$ by $p = \text{RMQ}_H(v_r + 2, w_r)$. If H[p] = d(w), we return $[v_r + 1, p 1]$ as the final result. Otherwise (H[p] > d(w)), the final result is $[v_r + 1, w_r]$. Time is $O(t_{\text{LCP}} + t_{\text{PNSV}} + t_{\text{RMQ}})$.
- LCA(v, w): First check if one of v or w is an ancestor of the other, and return that node in this case. Otherwise, assume $v_r < w_\ell$ (otherwise swap v and w). Let u denote the (yet unknown) LCA of v and w, so that our task is to compute u_ℓ and u_r . First note that all suffixes $T^{A[k]}$, $k \in [v_\ell, v_r] \cup [w_\ell, w_r]$, must be prefixed by \overline{u} , and that u is the deepest node with this property. Further, because none of v and w is an ancestor of the other, v and w must be contained in subtrees rooted at two different children \hat{u} and \hat{u} of u, say v is in \hat{u} 's subtree and w in the one of \hat{u} . Because $v_r \leq w_\ell$, we have $\hat{u}_r \leq \hat{u}_\ell$, and hence there must be a d(u)-index in H between \hat{u}_r and \hat{u}_ℓ , which can be found by $p = \text{RMQ}_H(v_r + 1, w_\ell)$. The endpoints of u's interval are again located by $u_\ell = PSV(p)$ and $u_r = NSV(p) 1$. Time is $O(t_{\text{RMQ}} + t_{\text{PNSV}})$.
- EDGELABEL(v, i): First, compute the string-depth of v by $d_1 = \text{SDEPTH}(v)$, and that of u = PARENT(v) by $d_2 = \text{SDEPTH}(u)$, in total time $O(t_{\text{RMQ}} + t_{\text{LCP}} + t_{\text{PNSV}})$. Now if $i > d_1 d_2$, return NULL, because i exceeds the length of the label of (u, v) in this case. Otherwise, the result is given by $t_{A[v_\ell]+d_2+i-1}$, since the edge-label of (u, v) is $T_{A[k]+d_2...A[k]+d_1-1}$ for an arbitrary $k \in [v_\ell, v_r]$. Total time is thus $O(t_{\text{SA}} + t_{\text{RMQ}} + t_{\text{LCP}} + t_{\text{PNSV}})$.

A final remark is that we can also simulate many other operations in suffix trees not listed here, e.g. suffix links, Weiner links, level ancestor queries, and many more.

9.3 Compressed LCP-Arrays

We now show how to reduce the space for the LCP-array H from $n \log n$ to O(n) bits. To this end, we first note that the LCP-value can decrease by at most 1 when moving from suffix A[i] - 1 to A[i] in H (i. e., when enumerating the LCP-values in text order):

Lemma 26. For all
$$1 < i \le n$$
, $H[i] \ge H[A^{-1}[A[i] - 1] - 1$.

Proof: If H[i] = 0, the claim is trivial. Hence, suppose H[i] > 0, and look at the two suffixes starting at positions A[i] and A[i-1], which must start with the same character. Suppose $T^{A[i]} = a\alpha$ and $T^{A[i-1]} = a\beta$ for $a \in \Sigma$, $\alpha, \beta \in \Sigma^*$.

Because the suffixes are sorted lexicographically in A, and $a\alpha >_{\text{lex}} a\beta$, we know $\alpha >_{\text{lex}} \beta$, and that α and β share a common prefix of length H[i]-1, call it γ . Now note that all suffixes between β and α in A must also start with γ , as otherwise the suffixes would not be in lexicographic order. In particular, suffix $T^{A[A^{-1}[A[i]+1]-1]}$ must be prefixed by γ , and hence $H[A^{-1}[A[i]+1]=1]$ LCP $(T^{A[i]+1}, T^{A[A^{-1}[A[i]+1]-1]}) = \text{LCP}(\alpha, T^{A[A^{-1}[A[i]+1]-1]}) \geq |\gamma| = H[i]-1$.

From the above lemma, we can conclude that $I[1,n] = [H[A^{-1}[1]]+1, H[A^{-1}[2]]+2, H[A^{-1}[3]]+3, \ldots, H[A^{-1}[n]]+n]$ is an array of *increasing* integers. Further, because no LCP-value can exceed the length of corresponding suffixes, we see that $H[A^{-1}[i]] \leq n-i+1$. Hence, sequence I must be in range [1,n]. We encode I differentially: writing $\Delta[i] = I[i] - I[i-1]$ for the difference between entry i and i-1, and defining I[0] = 0 for handling the border case, we encode $\Delta[i]$ in unary as $0^{\Delta[i]}1$. Let the resulting sequence be S.

Note that the number of 1's in S is exactly n, and that the number of 0's is at most n, as the $\Delta[i]$'s sum up to at most n. Hence, the length of S is at most 2n bits. We further prepare S for constant-time $rank_0$ - and $select_1$ -queries, using additional o(n) bits. Then H[i] can be retrieved by

$$H[i] = rank_0(S, select_1(S, A[i])) - A[i]$$
.

This is because the *select*-statement points to the position of the terminating '1' of $0^{\Delta[A[i]]}1$ in S, and the *rank*-statement counts the sum of Δ -values before that position, which is I[A[i]]. From this, in order to get H[i], we need to subtract A[i], which has bin "artificially" added when deriving I from H.

By noting that there are exactly A[i] 1's up to position $select_1(S, A[i])$ in S (and therefore $select_1(S, A[i]) - A[i]$ 0's), the calculation can be further simplified to

$$H[i] = select_1(S, A[i]) - 2A[i]$$
.

We have proved:

Theorem 27. The LCP-array H can be stored in 2n + o(n) bits such that retrieving an arbitrary entry H[i] takes $t_{\text{LCP}} = O(t_{\text{SA}})$ time.

Note that with the sampled suffix array from Sect. 8.5, this means that we no more have constant-time access to H, as $t_{SA} = O(\log n)$ in this case.

10 Succinct Data Structures for RMQs and PSV/NSV Queries

This chapter shows that O(n) bits are sufficient to answer RMQs and PSV/NSV-queries in constant time. For our compressed suffix tree, we assume that all three queries are executed on the LCP-array H, although the data structures presented in this chapter are applicable to any array of ordered objects.

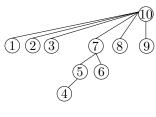
10.1 2-Dimensional Min-Heaps

We first define a tree that will be the basis for answering RMQs and NSV-queries. The solution for PSV-queries is symmetric. The following definition assumes that H[n+1] is always the smallest value in H, what can be enforced by introducing a "dummy" element $H[n+1] = -\infty$.

Definition 30. Let H[1, n+1] be an array of totally ordered objects, with the property that H[n+1] < H[i] for all $1 \le i \le n$. The 2-dimensional Min-Heap \mathcal{M}_H of H is a tree an n nodes $1, \ldots, n$, defined such that NSV(i) is the parent-node of i for $1 \le i \le n$.

Note that \mathcal{M}_H is a well-defined tree whose root is n+1.

Example 9.

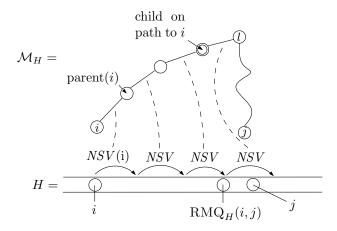


$$H = -1 \ 0 \ 0 \ 3 \ 1 \ 2 \ 0 \ 1 \ 1 \ -\infty$$

From the definition of \mathcal{M}_H , it is immediately clear that the value NSV(i) is given by the parent node of i $(1 \le i \le n)$. The next lemma shows that \mathcal{M}_H is also useful for answering RMQs on H.

Lemma 28. For $1 \le i < j \le n$, let $l = LCA_{\mathcal{M}_H}(i,j)$. Then if l = j, $RMQ_H(i,j) = j$. Otherwise, $RMQ_H(i,j)$ is given by the child of l that is on the path from l to i.

Proof: "graphical proof":



Example 10. Continuing the example above, let i = 4 and j = 6. We have $LCA_{\mathcal{M}_H}(4,6) = 7$, and 5 is the child of 7 on the path to 4. Hence, $RMQ_H(4,6) = 5$.

10.2 Balanced Parentheses Representation of Trees

Any ordered tree T on n nodes can be represented by a sequence B of 2n parentheses as follows: in a depth-first traversal of T, write an opening parenthesis '(' when visiting a node v for the first time, and a closing parenthesis ')' when visiting v for the last time (i. e., when all nodes in T_v have been traversed).

Example 11. Building on the 2d-Min-Heap from the Example 9, we have B = (()()()((())())()).

In a computer, a '(' could be represented by a '1'-bit, and a ')' by a '0'-bit, so the space for B is 2n bits. In the lecture "Advanced Data Structures" it is shown that this representation allows us to answer queries like $rank_{\ell}(B,i)$ and $select_{\ell}(B,i)$, by using only o(n) additional space.

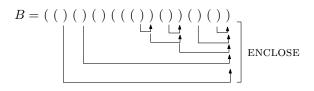
Note that the sequence B is balanced, in the sense that in each prefix the number of closing parentheses is no more than the number of opening parenthesis, and that there are n opening and closing parentheses each in total. Hence, this representation of trees is called balanced parentheses sequence (BPS).

We also need the following operation.

Definition 31. Given a sequence B[1,2n] of balanced parentheses and a position i with B[i] = ')', enclose(B,i) returns the position of the closing parenthesis of the nearest enclosing '()'-pair.

In other words, if v is a node with closing parenthesis at position i < 2n in B, and w is the parent of v with closing parenthesis at position j in B, then enclose(B, i) = j. Note that enclose(i) > i for all i, because of the order in which nodes are visited in a depth first traversal.

Example 12.



We state the following theorem that is also shown in the lecture "Advanced Data Structures."

Theorem 29. There is a data structure of size $O\left(\frac{n \log \log n}{\log n}\right) = o(n)$ bits that allows for constant-time enclose-queries.

(The techniques are roughly similar to the techniques for rank- and select-queries.)

Now look at an arbitrary position i in B, $1 \le i \le 2n$. We define the excess-value E[i] at position i as the number of opening parenthesis in B[1,i] minus the number of closing parenthesis in B[1,i]. Note that the excess-values do not have to be stored explicitly, as

$$\begin{aligned} |E[i]| &= rank_{(}(B,i) - rank_{)}(B,i) \\ &= i - rank_{)}(B,i) - rank_{)}(B,i) \\ &= i - 2rank_{)}(B,i) \ . \end{aligned}$$

Example 13.

$$B = (\ (\)\ (\)\ (\)\ (\ (\ (\)\)\ (\)\$$

Note:

- 1. E[i] > 0 for all $1 \le i < 2n$
- 2. E[2n] = 0
- 3. If i is the position of the closing parenthesis of node v, then E[i] is the depth of v. (Counting starts at 0, so the root has depth 0.)

We also state the following theorem without proof.

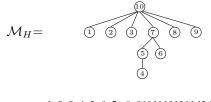
Theorem 30. Given a sequence B of balanced parentheses, there is a data structure of size $O\left(\frac{n \log \log n}{\log n}\right) = o(n)$ bits that allows to answer RMQs on the associated excess-sequence E in constant time.

(The techniques are again similar to rank and select: blocking and table-lookups. Note in particular that $\frac{\log n}{2}$ excess-values $E[x], E[x+1], \ldots, E\left[x+\frac{\log n}{2}-1\right]$ are encoded in a single computerword $B\left[x, x+\frac{\log n}{2}-1\right]$, and hence it is again possible to apply the Four-Russians-Trick!)

10.3 Answering Queries

We represent \mathcal{M}_H by its BPS B, and identify each node i in \mathcal{M}_H by the position of its *closing* parenthesis in B.

Example 14.



Note that the (closing parenthesis of) nodes appear in B in sorted order - this is simply because in \mathcal{M}_H node i hast post-order number i, and the closing parenthesis appear in post-order by the definition of the BPS. This fact allows us to jump back and forth between indices in H and positions of closing parentheses ')' in B, by using rank- and select-queries in the appropriate sequences.

Answering NSV-queries is now simple. Suppose we wish to answer $NSV_H(i)$. We then move to the position of the i'th ')' by

$$x \leftarrow select_{)}(B, i)$$
,

and then call

$$y \leftarrow enclose(B, x)$$

in order to move to the position y of the closing parenthesis of the parent j of i in \mathcal{M}_H . The (yet unknown) value j is computed by

$$j \leftarrow rank_{j}(B, y)$$
.

Example 15. We want to compute NSV(7). First compute $x \leftarrow select_j(B,7) = 15$, and then $y \leftarrow enclose(15) = 20$. The final result is $j \leftarrow rank_j(B,20) = 10$.

Answering RMQs is only slightly more complicated. Suppose we wish to answer RMQ $_H(i, j)$ for $1 \le i < j \le n$. As before, we go to the appropriate positions in B by

$$x \leftarrow select_{1}(B, i)$$
 and

$$y \leftarrow select_1(B, j)$$
.

We then compute the position of the minimum excess-value in the range [x, y] by

$$z \leftarrow \text{RMQ}_E(x, y)$$
,

and map it back to a position in H by

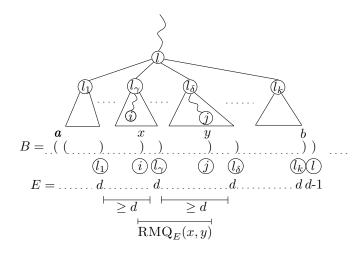
$$m \leftarrow rank_1(B, z)$$
.

This is the final answer.

Example 16. We want to compute $RMQ_H(4,9)$. First, compute $x \leftarrow select_{\uparrow}(B,4) = 11$ and $y \leftarrow enclose(B,9) = 19$. The range minimum query yields $z \leftarrow RMQ_E(11,19) = 15$. Finally, $m \leftarrow rank_{\uparrow}(B,15) = 7$ is the result.

We now justify the correctness of this approach. First assume that $\ell = \text{LCA}_{\mathcal{M}_H}(i,j)$ is different from j. Let ℓ_1, \ldots, ℓ_k be the children of ℓ , and assume $i \in T_{\ell_{\gamma}}$ and $j \in T_{\ell_{\delta}}$ for some $1 \leq \gamma < \delta \leq k$. By Lemma 28, we thus need to show that the position of the closing parenthesis of ℓ_{γ} is the position where E attains the minimum in E[x, y].

Example 17.



Let d-1 be the tree-depth of ℓ , and let B[a,b] denote the part of B that "spells out" T_{ℓ} (i.e., B[a,b] is the BPS of the sub-tree of T rooted at ℓ). Note that a < x < y < b, as i and j are both below ℓ in T.

Because B[a] is the opening parenthesis of node ℓ , we have E[a] = d. Further, because B is balanced, we have $E[c] \geq d$ for all a < c < b. But E assumes the values d at the positions of the closing parenthesis of nodes ℓ_{β} ($1 \leq \beta \leq k$), in particular for ℓ_{γ} . Hence, the leftmost minimum in E[x, y] is attained at the position z of the closing parenthesis of node ℓ_{γ} , which is computed by an RMQ in E. The case where $\ell = j$ is similar (and even simpler to prove). Thus, we get:

Theorem 31. With a data structure of size 2n + o(n) bits, we can answer RMQs and NSV-queries on an array of n ordered objects on O(1) time.

The drawback of the 2d-Min-Heap, however, is that it is inherently asymmetric (as the parentrelationship is defined by the minimum to the right), and cannot be used for answering PSV-queries as well. For this, we could build another 2d-Min-Heap \mathcal{M}_H^R on the reversed sequence H^R , using another 2n + o(n) bits. (Note that an interesting side-effect of this \mathcal{M}_H^R is that it would allow to compute the rightmost minimum in any query range, instead of the leftmost, which could have interesting applications in compressed suffix trees.)

In the lecture we also discussed the possibility to just add another bit-vector of length n bits — however, this seems only to work if we represent the 2d-Min-Heap by DFUDS (instead of BPS). If we plug all these structures into the compressed suffix tree from Chapter 9 (which was indeed the reason for developing the solutions for RMQs and PNSVs), we get:

Theorem 32. A suffix tree on a text of length n over an alphabet of size σ can be stored in |SA| + 3n + o(n) bits of space (where |SA| denotes the space for the suffix array), such that operations ROOT, ISLEAF, COUNT, ISANCESTOR, FIRSTCHILD, and LCA take O(1) time, and operations LEAFLABEL, SDEPTH, PARENT, NEXTSIBLING and EDGELABEL take $O(t_{SA})$ time (where t_{SA} denotes the time to retrieve an element from the suffix array).

11 Inside Google*

11.1 The Task

You are given a collection $S = \{S_1, \ldots, S_m\}$ of sequences $S_i \in \Sigma^*$ (web pages, protein or DNA-sequences, or the like). Your task is to build an index on S such that the following type of on-line queries can be answered *efficiently*:

given: a pattern $P \in \Sigma^*$.

return: all $j \in [1, m]$ such that S_j contains P.

Exercise: What has this to do with Google?

11.2 The Straight-Forward Solution

Define a string

$$T = S_1 \# S_2 \# \dots \# S_m \#$$

of length $n := \sum_{1 \le i \le m} (|S_i| + 1) = m + \sum_{1 \le i \le m} |S_i|$. Build the suffix array A on T. In an array D[1, n] remember from which string in S the corresponding suffix comes from:

$$D[i] = j \text{ iff } \sum_{k=1}^{j-1} (|S_k| + 1) < A[i] \le \sum_{k=1}^{j} (|S_k| + 1) .$$

When a query pattern P arrives, first locate the interval $[\ell, r]$ of P in A. Then output all numbers in $D[\ell, r]$, removing the duplicates (how?).

11.3 The Problem

Even if we can efficiently remove the duplicates, the above query algorithm is *not* output sensitive. To see why, consider the situation where P occurs many (say x) times in S_1 , but never in S_j for j > 1. Then the query takes O(|P| + x) time, just to output *one* sequence identifier (namely nr. 1). Note that x can be as large as $\Theta(n)$, e.g., if $|S_1| \ge \frac{n}{2}$.

11.4 An Optimal Solution

The following algorithm solves the queries in optimal O(|P|+d) time, where d denotes the number of sequences in S where P occurs.

We set up a new array E[1, n] such that E[i] points to the nearest previous occurrence of D[i] in D:

$$E[i] = \left\{ \begin{array}{ll} j & \text{if there is a } j < i \text{ with } D[j] = D[i], \text{ and } D[k] \neq D[i] \text{ for all } j < k < i \text{ }, \\ -1 & \text{if no such } j \text{ exists.} \end{array} \right.$$

It is easy to compute E with a single left-to-right scan of D. We further process E for constant-time RMQs.

When a query pattern P arrives, we first locate P's interval $[\ell, r]$ in A in O(|P|) time (as before). We then call $report(\ell, r)$, which is a procedure defined as follows.

Algorithm 8: Document Reporting

```
\begin{array}{lll} \textbf{1} \  \, & \textbf{procedure report} \ (i,j); \\ \textbf{2} & m \leftarrow \text{RMQ}_E(i,j); \\ \textbf{3} & \textbf{if} \ E[m] \leq \ell \ \textbf{then} \\ \textbf{4} & \text{output} \ D[m]; \\ \textbf{5} & \textbf{if} \ m-1 \geq i \ \textbf{then report}(i,m-1); \\ \textbf{6} & \textbf{if} \ m+1 \leq j \ \textbf{then report}(m+1,j); \\ \textbf{7} \ \ \textbf{end} \end{array}
```

The claimed O(d) running time of the call to $\operatorname{report}(\ell,r)$ relies on the following observation. Consider the range $[\ell,r]$. Note that P is a prefix of $T^{A[i]}$ for all $\ell \leq i \leq r$. The idea is that the algorithm visits and outputs only those suffixes $T^{A[i]}$ with $i \in [\ell,r]$ such that the corresponding suffix σ_i of $S_{D[i]}$ ($\sigma_i = T^{A[i]...e}$, where $e = \sum_{1 \leq j \leq D[i]} (|S_j| + 1)$ is the end position of $S_{D[j]}$ in T) is the lexicographically smallest among those suffixes of $S_{D[i]}$ that are prefixed by P. Because the suffix array orders the suffixes lexicographically, we must have $E[i] \leq \ell$ for such suffixes σ_i . Further, there is at most one such position i in $[\ell,r]$ for each string S_j . Because the recursion searches the whole range $[\ell,r]$ for such positions i, no string $S_j \in \mathcal{S}$ is missed by the procedure.

Finally, when the recursion stops (i.e., $E[m] > \ell$), because E[m] is the minimum in E[i,j], we must have that the identifiers of the strings $S_{D[k]}$ for all $k \in [i,j]$ have already been output in a previous call to $\mathtt{report}(i',j')$ for some $\ell \leq i' \leq j' < i$. Hence, we can safely stop the recursion at this point.