## 2. Notation

We use the shorthands $[i, j] = \{i, \ldots, j\}$ and $[i, j) = [i, j-1]$ for ranges of integers and extend to substrings as seen below.

The *input* of a suffix array construction algorithm is a *string* $T = T[0, n) = t_0 t_1 \cdots t_{n-1}$ over the alphabet $[1, n]$, that is, a sequence of $n$ integers from the range $[1, n]$. For convenience, we assume that $t_j = 0$ for $j \geq n$. Sometimes we also assume that $n + 1$ is a multiple of some constant $v$ or a square to avoid a proliferation of trivial case distinctions and $\lceil \cdot \rceil$ operations. An implementation will either spell out the case distinctions or pad (sub)problems with an appropriate number of zero characters. The restriction to the alphabet $[1, n]$ is not a serious one. For a string $T$ over any alphabet, we can first sort the characters of $T$, remove duplicates, assign a rank to each character, and construct a new string $T'$ over the alphabet $[1, n]$ by renaming the characters of $T$ with their ranks. Since the renaming is order preserving, the order of the suffixes does not change.

For $i \in [0, n]$, let $S_i$ denote the *suffix* $T[i, n) = t_i t_{i+1} \cdots t_{n-1}$. We also extend the notation to sets: for $C \subseteq [0, n]$, $S_C = \{S_i \mid i \in C\}$. The goal is to sort the set $S_{[0,n]}$ of suffixes of $T$, where comparison of substrings or tuples assumes the lexicographic order throughout this article. The *output* is the *suffix array SA*$[0, n]$ of $T$, a permutation of $[0, n]$ satisfying $S_{SA[0]} < S_{SA[1]} < \cdots < S_{SA[n]}$.

## 3. Linear-Time Algorithm

We begin with a detailed description of the simple linear-time algorithm, which we call DC3 (for Difference Cover modulo 3, see Section 4). The execution of the algorithm is illustrated with the following example

$$
\begin{array}{c}
\phantom{T[0, n) = }\;\;0\;\;\;1\;\;\;2\;\;\;3\;\;\;4\;\;\;5\;\;\;6\;\;\;7\;\;\;8\;\;\;9\;\;10\;\;11 \\
T[0, n) = \text{y a b b a d a b b a d o},
\end{array}
$$

where we are looking for the suffix array

$$SA = (12, 1, 6, 4, 9, 3, 8, 2, 7, 5, 10, 11, 0).$$

*Step* 0: *Construct a Sample*
For $k = 0, 1, 2$, define

$$B_k = \{i \in [0, n] \mid i \bmod 3 = k\}.$$

Let $C = B_1 \cup B_2$ be the set of *sample positions* and $S_C$ the set of *sample suffixes*.

*Example* 3.1.  $B_1 = \{1, 4, 7, 10\}$, $B_2 = \{2, 5, 8, 11\}$, that is, $C = \{1, 4, 7, 10, 2, 5, 8, 11\}$.

*Step* 1: *Sort Sample Suffixes*

For $k = 1, 2$, construct the strings

$$R_k = [t_k t_{k+1} t_{k+2}][t_{k+3} t_{k+4} t_{k+5}] \ldots [t_{\max B_k} t_{\max B_k+1} t_{\max B_k+2}]$$

whose characters are triples $[t_i t_{i+1} t_{i+2}]$. Note that the last character of $R_k$ is unique because $t_{\max B_k+2} = 0$. Let $R = R_1 \odot R_2$ be the concatenation of $R_1$ and $R_2$. Then, the (nonempty) suffixes of $R$ correspond to the set $S_C$ of sample suffixes: $[t_i t_{i+1} t_{i+2}][t_{i+3} t_{i+4} t_{i+5}] \cdots$ corresponds to $S_i$. The correspondence is order preserving, that is, by sorting the suffixes of $R$ we get the order of the sample suffixes $S_C$.

*Example* 3.2.   $R = [abb][ada][bba][do0][bba][dab][bad][o00]$.

To sort the suffixes of $R$, first radix sort the characters of $R$ and rename them with their ranks to obtain the string $R'$. If all characters are different, the order of characters gives directly the order of suffixes. Otherwise, sort the suffixes of $R'$ using Algorithm DC3.

*Example* 3.3.   $R' = (1, 2, 4, 6, 4, 5, 3, 7)$ and $SA_{R'} = (8, 0, 1, 6, 4, 2, 5, 3, 7)$.

Once the sample suffixes are sorted, assign a rank to each suffix. For $i \in C$, let $rank(S_i)$ denote the rank of $S_i$ in the sample set $S_C$. Additionally, define $rank(S_{n+1}) = rank(S_{n+2}) = 0$. For $i \in B_0$, $rank(S_i)$ is undefined.

*Example* 3.4.

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $rank(S_i)$ | $\perp$ | 1 | 4 | $\perp$ | 2 | 6 | $\perp$ | 5 | 3 | $\perp$ | 7 | 8 | $\perp$ | 0 | 0. |

*Step* 2: *Sort Nonsample Suffixes*

Represent each nonsample suffix $S_i \in S_{B_0}$ with the pair $(t_i, rank(S_{i+1}))$. Note that $rank(S_{i+1})$ is always defined for $i \in B_0$. Clearly, we have, for all $i, j \in B_0$,

$$S_i \leq S_j \iff (t_i, rank(S_{i+1})) \leq (t_j, rank(S_{j+1})).$$

The pairs $(t_i, rank(S_{i+1}))$ are then radix sorted.

*Example* 3.5.   $S_{12} < S_6 < S_9 < S_3 < S_0$ because $(0, 0) < (a, 5) < (a, 7) < (b, 2) < (y, 1)$.

*Step* 3: *Merge*

The two sorted sets of suffixes are merged using a standard comparison-based merging. To compare suffix $S_i \in S_C$ with $S_j \in S_{B_0}$, we distinguish two cases:

$$i \in B_1 : \quad S_i \leq S_j \iff (t_i, rank(S_{i+1})) \leq (t_j, rank(S_{j+1}))$$
$$i \in B_2 : \quad S_i \leq S_j \iff (t_i, t_{i+1}, rank(S_{i+2})) \leq (t_j, t_{j+1}, rank(S_{j+2}))$$

Note that the ranks are defined in all cases.

*Example* 3.6.   $S_1 < S_6$ because $(a, 4) < (a, 5)$ and $S_3 < S_8$ because $(b, a, 6) < (b, a, 7)$.

The time complexity is established by the following theorem.

THEOREM 3.7.   *The time complexity of Algorithm DC3 is $\mathcal{O}(n)$.*

PROOF.   Excluding the recursive call, everything can clearly be done in linear time. The recursion is on a string of length $\lceil 2n/3 \rceil$. Thus, the time is given by the recurrence $T(n) = T(2n/3) + \mathcal{O}(n)$, whose solution is $T(n) = \mathcal{O}(n)$.   $\square$