# Breaking the $\mathcal{O}(n)$-Barrier in the Construction of Compressed Suffix Arrays and Suffix Trees

Dominik Kempa[*]
Stony Brook University
kempa@cs.stonybrook.edu

Tomasz Kociumaka[†]
Max Planck Institute for Informatics
tomasz.kociumaka@mpi-inf.mpg.de

## Abstract

The suffix array, describing the lexicographical order of suffixes of a given text, and the suffix tree, a path-compressed trie of all suffixes, are the two most fundamental data structures for string processing, with plethora of applications in data compression, bioinformatics, and information retrieval. For a length-$n$ text, however, they use $\Theta(n \log n)$ bits of space, which is often too costly. To address this, Grossi and Vitter [STOC 2000] and, independently, Ferragina and Manzini [FOCS 2000] introduced space-efficient versions of the suffix array, known as the *compressed suffix array* (CSA) and the *FM-index*. Sadakane [SODA 2002] then showed how to augment them to obtain the *compressed suffix tree* (CST). For a length-$n$ text over an alphabet of size $\sigma$, these structures use only $\mathcal{O}(n \log \sigma)$ bits. Nowadays, these structures are part of the standard toolbox: modern textbooks spend dozens of pages describing their applications, and they almost completely replaced suffix arrays and suffix trees in space-critical applications. The biggest remaining open question is how efficiently they can be constructed. After two decades, the fastest algorithms still run in $\mathcal{O}(n)$ time [Hon et al., FOCS 2003], which is $\Theta(\log_{\sigma} n)$ factor away from the lower bound of $\Omega(n/\log_{\sigma} n)$ (following from the necessity to read the input).

In this paper, we make the first in 20 years improvement in $n$ for this problem by proposing a new compressed suffix array and a new compressed suffix tree which admit $o(n)$-time construction algorithms while matching the space bounds and the query times of the original CSA/CST and the FM-index. More precisely, our structures take $\mathcal{O}(n \log \sigma)$ bits, support SA queries and full suffix tree functionality in $\mathcal{O}(\log^{\epsilon} n)$ time per operation, and can be constructed in $\mathcal{O}(n \min(1, \log \sigma/\sqrt{\log n}))$ time using $\mathcal{O}(n \log \sigma)$ bits of working space. (For example, if $\sigma = 2$, the construction time is $\mathcal{O}(n/\sqrt{\log n}) = o(n)$.) We derive this result as a corollary from a much more general reduction: We prove that all parameters of a compressed suffix array/tree (query time, space, construction time, and construction working space) can essentially be reduced to those of a data structure answering new query types that we call *prefix rank* and *prefix selection*. Using the novel techniques, we also develop a new index for pattern matching.

# 1 Introduction

Let $T$ be a text of length $n$. A *suffix tree* [81] of $T$ is a trie of all suffixes of $T$, in which every unary path has been replaced with a single edge labeled by a text substring. The resulting tree has less than $2n$ nodes and thus can be encoded in $\mathcal{O}(n \log n)$ bits. Related to suffix trees are *suffix arrays* [60]. The suffix array $SA[1 . . n]$ of $T$ stores the permutation of $\{1, \ldots, n\}$ such that $SA[i]$ is the starting position of the $i$th lexicographically smallest suffix of $T$. Consider now the following problem: Construct a data structure that, given any length-$m$ pattern $P$, counts the number of occurrences of $P$ in $T$. To solve it using a suffix tree, it suffices to descend the tree in $\mathcal{O}(m)$ time and report the precomputed number of leaves below the reached node. Using a suffix array, it suffices to perform an $\mathcal{O}(m \log n)$-time binary search resulting in the range $SA[b . . e]$ of suffixes of $T$ having $P$ as a prefix. Then, $e - b$ is the number of occurrences of $P$ in $T$ (and $SA[b . . e]$ contains their starting positions). The advantage of suffix array is that it is more space efficient: it only needs $n \lceil \log n \rceil$ bits. The queries, however, are usually slightly slower.

The above is a canonical application of suffix arrays/trees. It is, however, only the tip of the iceberg. Suffix trees and suffix arrays are widely considered to be the two most fundamental data structures for string processing. As written by Gusfield in his classical textbook [44]: *"Suffix trees can be used to solve the exact matching problem in linear time (...), but their real virtue comes from their use in linear-time solutions to many string problems more complex than exact matching"*. This includes well-studied problems like Maximal Repeats, Longest Repeated Factor, Minimal Absent Word, Longest Common Substring, Matching Statistics, Maximal Unique Matches, LZ77 Factorization, BWT Compression, and many more (see, e.g., [1, 44, 58, 65, 72]).

With the increasing size of datasets that need processing, plain suffix arrays and suffix trees, however, have become expensive to use, particularly in applications where the input text is over a small alphabet $[0 . . \sigma)$. Such text requires $n \lceil \log \sigma \rceil$ bits, whereas the suffix array/tree uses at least $n \lceil \log n \rceil$ bits of space. In some applications, the gap $\frac{\log n}{\log \sigma}$ can be quite large, e.g., in computational biology, where we usually have $\sigma = 4$, the gap is typically between 16 and 32. This shortcoming was addressed by Grossi and Vitter and, independently, Ferragina and Manzini at the turn of the millennium. They introduced space-efficient versions of the suffix array, known as the *compressed suffix array (CSA)* [42, 43] and the *FM-index* [27, 28]. For a length-$n$ text over an alphabet of size $\sigma$, these data structures use only $\mathcal{O}(n \log \sigma)$ bits, and they can answer SA queries (asking for $SA[i]$ given $i \in [1 . . n]$) in $\mathcal{O}(\log^\epsilon n)$ time, where $\epsilon > 0$ is an arbitrary predefined constant. With such data structure, one can execute any algorithm that uses the suffix array, but consuming less space and only incurring a factor of $\mathcal{O}(\log^\epsilon n)$ penalty in the runtime.[1] Shortly after these discoveries, Sadakane [79] extended CSA/FM-index into a *compressed suffix tree (CST)*, supporting all suffix tree operations in $\mathcal{O}(\log^\epsilon n)$ time (while still using $\mathcal{O}(n \log \sigma)$ bits of space). This powerful structure can be plugged into an even larger set of algorithms [37].

Nowadays, CSAs and CSTs are widely used in practice. Modern string algorithms textbooks focus on the use and applications of CSAs/CSTs and related data structures [1, 58], or even entirely on the emerging notion of *compressed data structures* [65]. The FM-index occupies the central role in some of the most commonly used bioinformatics tools, like `Bowtie` [55], `BWA` [56], and `Soap2` [57], and mature and highly engineered implementations of CSAs and CSTs are available through the `sdsl` library[2] of Gog et al. [37, 38]. Despite these developments in functionality and practical adoption of CSAs/CSTs, the time complexity of their construction remains an open problem. The original paper of Grossi and Vitter [42], describes a method that, given a length-$n$ text over alphabet $\Sigma = [0 . . \sigma)$, constructs the CSA in $\mathcal{O}(n \log \sigma)$ time and using $\mathcal{O}(n \log n)$ bits of working space. In 2003, a celebrated result of Hon et al. [46] lowered the time

---

[1]This is often acceptable: a slower algorithm remains usable, but insufficient memory can thwart it entirely.
[2]The library is available at `https://github.com/simongog/sdsl-lite`.

complexity to $\mathcal{O}(n \log \log \sigma)$ and the space to the optimal $\mathcal{O}(n \log \sigma)$ bits. Note, however, that, e.g., for $\sigma = 2$, this algorithm still runs in $\Theta(n)$ time, which is slower by a $\Theta(\log n)$ factor than the lower bound of $\Omega(n/\log n)$, following simply from the necessity to read the entire input. Recently, Belazzougui [5] improved the time complexity of the CSA/CST construction to randomized $\mathcal{O}(n)$ (while using the optimal space of $\mathcal{O}(n \log \sigma)$ bits), making it independent of the alphabet size $\sigma$. Shortly after, Munro, Navarro, and Nekrich [61] proposed a deterministic solution. Despite these advances, 20 years after the result of Hon et al. [46], the bound of $\Omega(n)$ still stands on the construction of CSAs/CSTs. Given the fundamental role of CSAs and CSTs, we thus ask:

> *Given a text over alphabet $\Sigma = [0 \mathinner{.\,.} \sigma)$ represented using $\mathcal{O}(n \log \sigma)$ bits,*
> *can we construct a compressed suffix array/tree of $T$ in $o(n)$ time?*

**Our Results**   We answer the above question affirmatively by describing a new data structure that takes $\mathcal{O}(n \log \sigma)$ bits, supports all operations of CSA and CST in $\mathcal{O}(\log^\epsilon n)$ time, and can be constructed in $\mathcal{O}(n \min(1, \log \sigma/\sqrt{\log n}))$ time using $\mathcal{O}(n \log \sigma)$ bits of space (Theorems 5.31 and 7.81). Thus, our solution matches the size and the query time of [27, 42, 79] (as well as more recent CSTs [14, 16, 31, 34, 75, 77]) but, unlike those, admits a sublinear-time construction for small $\sigma$. In particular, we achieve $\mathcal{O}(n/\sqrt{\log n}) = o(n)$ time for $\sigma = 2$, constituting the first improvement in $n$ since 2003 [46].

In addition to a new CSA/CST, we also present a new pattern matching index. We show (in Theorem 6.29) how, given a length-$n$ text $T$ stored using $\mathcal{O}(n \log \sigma)$ bits, to construct in $\mathcal{O}(n \min(1, \log \sigma/\sqrt{\log n}))$ time an index of size $\mathcal{O}(n \log \sigma)$ bits that, given the packed representation (i.e., using $\mathcal{O}(m \log \sigma)$ bits) of any pattern $P[1 \mathinner{.\,.} m]$, counts the occurrences of $P$ in $T$ in $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$ time (where $\epsilon > 0$ is an arbitrary predefined constant). The best previous solutions using compact space (i.e., $\mathcal{O}(n \log \sigma)$ bits) achieve $\mathcal{O}(n \log \sigma/\sqrt{\log n})$-time[3] construction and $\mathcal{O}(m/\log_\sigma n + \log n \cdot \log_\sigma n)$-time queries [63], or $\mathcal{O}(n)$-time construction and $\mathcal{O}(m/\log_\sigma n + \log_\sigma^\epsilon n)$-time queries [62]. Thus, for the most difficult case of $\sigma = 2^{\mathcal{O}(\sqrt{\log n})}$, our construction subsumes both these indexes in both aspects.[4] Since our pattern-matching index not only returns the number of occurrences but also the range in SA containing all suffixes prefixed with $P$, combining the above result with our CSA yields the structure that can additionally *report* all occurrences of $P$ in $\mathcal{O}(\log^\epsilon n)$ time per occurrence (Theorem 6.30).

The query times of all our data structures (i.e., both the CSA/CST and the pattern matching index) are worst-case, and all our algorithms are deterministic.

Our data structures differ significantly from the CSA of Grossi and Vitter [42], the FM-index of Ferragina and Manzini [27], and the CST of Sadakane [79], which are based on the so-called $\Psi$ function [42] or the Burrows–Wheeler transform [15]. We instead rely on the combination of the recently developed notion of *string synchronizing sets (SSS)* [50] and the new type of queries we call *prefix rank* and *prefix selection* queries. Although the prior work on SSS [50, 52] laid out its basic properties, it cannot be turned into an efficient CSA, because it heavily relies on *orthogonal range counting* queries [18], which are *provably incapable* of supporting SA queries as fast as the CSA or FM-index: Pătraşcu [76] showed a lower bound $\Omega(\frac{\log n}{\log \log n})$ on the query time of any structure using near-linear space. On the other hand, the $\mathcal{O}(n \log \sigma/\sqrt{\log n})$-time BWT construction from [50] is not sufficient to obtain an implementation of CSA since the classical BWT-based CSA, in addition to BWT, requires SA samples, i.e., a set containing all pairs $(\mathrm{SA}^{-1}[j], j)$ such that $j$ is a multiple of $\log n$, and it is not known how to obtain such a sequence using prior techniques. The key difficulty is computing the (global) rank $\mathrm{SA}^{-1}[j]$ of

---

[3]Although CSA lets us implement pattern counting queries, an index implementing pattern counting queries does not let us implement SA queries; thus, although built in $o(n)$ time, [63] cannot be used to answer SA queries.

[4]Note that if $\log \sigma = \mathcal{O}(\sqrt{\log n})$, then $\log_\sigma^\epsilon n = \Theta(\log^{\epsilon'} n)$ holds for $\epsilon' = \frac{\epsilon}{2}$.

each sampled suffix; an easy application of sparse suffix sorting gives, in $\mathcal{O}(n/\log_\sigma n)$ time, the lexicographic order of the sampled suffixes, but this is insufficient for placing each sampled suffix among the $n$ suffixes of the original string.

We sidestep these obstacles and demonstrate that *general* orthogonal range counting queries [19, 18] are in fact not needed at all, and each of their uses can either be: (1) eliminated completely (see the proof in Section 5.3.6), (2) replaced with prefix rank/selection queries (see Section 4.3), or (3) improved, utilizing the fact that the instances arising in our construction have properties that permit a fast custom solution (see Section 4.4). More details are provided in the Technical Overview (Section 3). As a result, we obtain a general set of reductions for the construction of CSA/CST and pattern-matching indexes, stated in Theorems 5.32, 6.31, and 7.82. In a single theorem, we can summarize it as follows; note that our reduction achieves *near-perfect efficiency*, i.e., it incurs no overhead (compared to the optimal solution) in space, preprocessing time, and preprocessing space, and only has an extra $\mathcal{O}(\log\log n)$ term in the query time. Everything else depends entirely on prefix rank and selection queries.

**Theorem 1.1** (Main result of this paper)**.** *Consider a data structure answering prefix rank and selection queries (Section 2.1) that, for any string of length $m$ over alphabet $[0\mathinner{.\,.}\sigma)^\ell$ (or equivalently, a sequence of $m$ length-$\ell$ strings over alphabet $[0\mathinner{.\,.}\sigma)$), achieves:*

1. *Space usage $S(m,\ell,\sigma)$ (measured in $\Theta(\log m)$-bit machine words),*
2. *Preprocessing time $P_t(m,\ell,\sigma)$,*
3. *Preprocessing space $P_s(m,\ell,\sigma)$,*
4. *Query time $Q(m,\ell,\sigma)$.*

*For every $T \in [0\mathinner{.\,.}\sigma)^n$ with $2 \le \sigma < n^{1/7}$, there exist $m = \mathcal{O}(n/\log_\sigma n)$ and $\ell = \mathcal{O}(\log_\sigma n)$ such that, given the packed representation of $T$, we can in $\mathcal{O}(n/\log_\sigma n + P_t(m,\ell,\sigma))$ time and $\mathcal{O}(n/\log_\sigma n + P_s(m,\ell,\sigma))$ working space build a structure of size $\mathcal{O}(n/\log_\sigma n + S(m,\ell,\sigma))$ that:*

- *Supports* SA *and inverse* SA *queries in $\mathcal{O}(\log\log n + Q(m,\ell,\sigma))$ time;*
- *Supports all suffix tree operations (Table 1) in $\mathcal{O}(\log\log n + Q(m,\ell,\sigma))$ time;*
- *Given the packed representation of any pattern $P \in [0\mathinner{.\,.}\sigma)^p$, returns:*

  - *The range* SA$[b\mathinner{.\,.}e)$ *of suffixes of $T$ with prefix $P$ in $\mathcal{O}(p/\log_\sigma n + \log\log n + Q(m,\ell,\sigma))$ time;*
  - *All occ starting positions of $P$ in $T$ in $\mathcal{O}(p/\log_\sigma n + (occ+1)(\log\log n + Q(m,\ell,\sigma)))$ time.*

Using this general reduction, we obtain the specific tradeoffs for CSA/CST and pattern matching queries we announced earlier by plugging in the data structure for prefix rank/selection queries from Theorem 2.2.

**Related Work** In parallel to efforts to improve the complexity of CSA/CST construction, were the efforts to make it more practical [38, 39, 40, 48, 73, 75]. This resulted in libraries of compressed data structures such as `sdsl` [38], `sux`, and `libcds`. More recently, some of these data structures have been extended to the dynamic setting, e.g., in the `DYNAMIC` [74] library.

In addition to CSA/CST and indexes using the optimal space of $\mathcal{O}(n\log\sigma)$ bits, previous work addressed the problem of designing structures using $\omega(n\log\sigma)$ but still $o(n\log n)$ bits [28, 35, 43]. We formulated our main result (Theorem 1.1) as a general reduction so that techniques from these and similar future studies could be easily combined with ours, potentially yielding new tradeoffs for pattern matching and CSA/CST queries.

In recent years, there has also been progress in the query time of $\mathcal{O}(n\log n)$-bit pattern matching indexes. The suffix trees support $\mathcal{O}(m)$-time pattern search after $\mathcal{O}(n)$-time randomized or $\mathcal{O}(n\log\log\sigma)$-time deterministic construction [26, 78]. Fischer and Gawrychowski [30] achieved $\mathcal{O}(m+\log\log\sigma)$-time queries after $\mathcal{O}(n)$-time deterministic construction, improving upon [22, 60].

If the pattern is given using $\mathcal{O}(m \log \sigma)$ bits, Bille et al. [12] achieved $\mathcal{O}(m/\log_\sigma n + \log m + \log\log \sigma)$ time, which Navarro and Nekrich [69] improved to $\mathcal{O}(m/\log_\sigma n + 1)$.

Surprisingly, the size of some CSAs, CSTs, and compact indexes can be reduced *below* $n\lceil\log\sigma\rceil$ bits for statistically compressible texts. For example, already the original FM-index [27] takes only $\mathcal{O}(nH_k(T)) + o(n\log\sigma)$ bits, where $H_k(T)$ denotes the *empirical kth-order entropy* of the text $T$ [23]. Currently, the smallest indexes reach $nH_k(T) + o(n(H_k(T) + 1))$ bits [4, 8]. Navarro and Mäkinen [68], and Belazzougui and Navarro [7] survey the achievable tradeoffs for such *fully compressed* indexes. Chan et al. [17], and Mäkinen and Navarro [59] describe *dynamic* compressed pattern-matching indexes maintaining a collection of texts supporting insertions/deletions.

Compressed indexes based on LZ77 [82] and run-length BWT [15] rapidly gain popularity. The early indexes [6, 11, 13, 32, 33] support only pattern search and random-access operations. Subsequent works generalized them to other dictionary compressors [20, 53, 70] and added dynamism [36, 71]. Support for SA queries is a recent addition of Gagie et al. [34]. Navarro surveys these indexes [67] and the intricate network of the underlying compressibility measures [66]. Interestingly, some of these pattern matching indexes can be constructed in compressed time. For example, the index of [36] can be constructed in $\mathcal{O}(z\log^3 n)$ time from the LZ77 representation of $T$ (with $z$ phrases), and then it locates pattern occurrences in $\mathcal{O}(m + occ\log n)$ time. On the other hand, the only compressed index supporting SA queries [34] is only constructible in $\Omega(n)$ time, but it can be built in compressed space $\mathcal{O}(r\log(n/r))$ given the run-length BWT of $T$ (with $r$ runs).

**Organization of the Paper**   After introducing the basic notation and tools in Section 2, we give a technical overview of the paper in Section 3. In Section 4, we then introduce some auxiliary tools utilized in our data structures. Section 5 describes our data structure answering SA and $\text{SA}^{-1}$ queries. In Section 6, we present our index for counting and reporting occurrences of patterns given using packed representation. Finally, in Section 7, we extend the functionality of our CSA into that of a CST.

## 2   Preliminaries

A *string* is a finite sequence of characters from a given *alphabet*. The length of a string $S$ is denoted $|S|$. For $i \in [1 \mathinner{..} |S|]$,[5] the $i$th character of $S$ is denoted $S[i]$. A *substring* of $S$ is a string of the form $S[i \mathinner{..} j] = S[i]S[i+1]\cdots S[j-1]$ for some $1 \le i \le j \le |S| + 1$. *Prefixes* and *suffixes* are substrings of the form $S[1 \mathinner{..} j]$ and $S[i \mathinner{..} |S|]$, respectively. We use $\overline{S}$ to denote the *reverse* of $S$, i.e., $S[|S|]\cdots S[2]S[1]$. We denote the *concatenation* of two strings $U$ and $V$, that is, $U[1]\cdots U[|U|]V[1]\cdots V[|V|]$, by $UV$ or $U \cdot V$. Furthermore, $S^k = \bigodot_{i=1}^k S$ is the concatenation of $k \in \mathbb{Z}_{\ge 0}$ copies of $S$; note that $S^0 = \varepsilon$ is the *empty string*. For a non-empty string $S \in \Sigma^+$, we define the special infinite string $S^\infty$ such that $S^\infty[i] = S[1 + (i-1) \bmod |S|]$ holds for every $i \in \mathbb{Z}$; in particular, $S^\infty[1 \mathinner{..} |S|] = S[1 \mathinner{..} |S|]$. An integer $p \in [1 \mathinner{..} |S|]$ is a *period* of $S$ if $S[i] = S[i+p]$ holds for every $i \in [1 \mathinner{..} |S|-p]$. We denote the shortest period of $S$ as $\mathrm{per}(S)$.

Throughout the paper, we consider a string (called the *text*) $T$ of length $n \ge 2$ over an integer alphabet $\Sigma = [0 \mathinner{..} \sigma)$,

| $i$ | SA$[i]$ | $T[\text{SA}[i] \mathinner{..} n]$ |
|---|---|---|
| 1 | 18 | \$ |
| 2 | 17 | a\$ |
| 3 | 12 | aababa\$ |
| 4 | 3 | aababababaababa\$ |
| 5 | 15 | aba\$ |
| 6 | 10 | abaababa\$ |
| 7 | 1 | abaababababaababa\$ |
| 8 | 13 | ababa\$ |
| 9 | 8 | ababaababa\$ |
| 10 | 6 | abababaababa\$ |
| 11 | 4 | ababababaababa\$ |
| 12 | 16 | ba\$ |
| 13 | 11 | baababa\$ |
| 14 | 2 | baababababaababa\$ |
| 15 | 14 | baba\$ |
| 16 | 9 | babaababa\$ |
| 17 | 7 | bababaababa\$ |
| 18 | 5 | babababaababa\$ |

**Figure 1:** A list of all sorted suffixes of $T = $ abaababababaababa\$ along with the suffix array of $T$.

---

[5]For $i, j \in \mathbb{Z}$, denote $[i \mathinner{..} j] = \{k \in \mathbb{Z} : i \le k \le j\}$, $[i \mathinner{..} j) = \{k \in \mathbb{Z} : i \le k < j\}$, and $(i \mathinner{..} j] = \{k \in \mathbb{Z} : i < k \le j\}$.

where $\sigma = n^{\mathcal{O}(1)}$. We assume $T[n] = 0$, and that $0$ (also denoted with \$) does not appear elsewhere in $T$. We use $\preceq$ to denote the order on $\Sigma$, extended to the *lexicographic* order on $\Sigma^*$ (the set of strings over $\Sigma$) so that $U, V \in \Sigma^*$ satisfy $U \preceq V$ if and only if either $U$ is a prefix of $V$, or $U[1 .. i] = V[1 .. i]$ and $U[i] \prec V[i]$ holds for some $i \in [1 .. \min(|U|, |V|)]$. The *suffix array* $\mathrm{SA}[1 .. n]$ of $T$ is a permutation of $[1 .. n]$ such that $T[\mathrm{SA}[1] .. n] \prec T[\mathrm{SA}[2] .. n] \prec \cdots \prec T[\mathrm{SA}[n] .. n]$, i.e., $\mathrm{SA}[i]$ is the starting position of the lexicographically $i$th suffix of $T$; see Fig. 1 for an example. The *inverse suffix array* $\mathrm{ISA}[1 .. n]$ (also denoted $\mathrm{SA}^{-1}[1 .. n]$) is the inverse permutation of SA, i.e., $\mathrm{ISA}[j] = i$ holds if and only if $\mathrm{SA}[i] = j$. Intuitively, $\mathrm{ISA}[j]$ stores the lexicographic *rank* of a suffix $T[j .. n]$ among the suffixes of $T$. By $\mathrm{lcp}(U, V)$ we denote the length of the longest common prefix of $U$ and $V$. For $j_1, j_2 \in [1 .. n]$, we let $\mathrm{LCE}(j_1, j_2) = \mathrm{lcp}(T[j_1 ..], T[j_2 ..])$. For any $P, S \in \Sigma^*$, we let

$$\mathrm{Occ}(P, S) = \{j \in [1 .. |S|] : j + |P| \le |S| + 1 \text{ and } S[j .. j+|P|) = P\},$$
$$\mathrm{RangeBeg}(P, S) = |\{i \in [1 .. |S|] : S[i .. |S|] \prec P\}|,$$
$$\mathrm{RangeEnd}(P, S) = \mathrm{RangeBeg}(P, S) + |\mathrm{Occ}(P, S)|.$$

Observe that the following equality holds for every $P \in \Sigma^*$:

$$\mathrm{Occ}(P, T) = \{\mathrm{SA}[i] : i \in (\mathrm{RangeBeg}(P, T) .. \mathrm{RangeEnd}(P, T)]\}.$$

We use the word RAM model of computation [45] with $w$-bit *machine words*, where $w \ge \log n$. In this model, strings are typically represented as arrays, with each character occupying one memory cell. A single character, however, only needs $\lceil \log \sigma \rceil$ bits, which might be much less than $w$. We can therefore store (the *packed representation* of) a text $T \in [0 .. \sigma)^n$ using $\mathcal{O}(\lceil \frac{n \log \sigma}{w} \rceil)$ memory cells.

## 2.1 (Prefix) Rank and Selection Queries

Let us recall the (ordinary) rank and selection queries on a string $S \in \Sigma^n$:

**Rank query** $\mathrm{rank}_{S,a}(j)$: Given $a \in \Sigma$ and $j \in [0 .. n]$, compute $|\{i \in [1 .. j] : S[i] = a\}|$.

**Selection query** $\mathrm{select}_{S,a}(r)$: Given $a \in \Sigma$ and $r \in [1 .. \mathrm{rank}_{S,a}(n)]$, find the $r$th smallest element of $\{i \in [1 .. n] : S[i] = a\}$.

**Theorem 2.1** (Rank and selection queries in bitvectors [3, 21, 47, 64]). *For every string $S \in \{0, 1\}^*$, there exists a data structure of $\mathcal{O}(|S|)$ bits answering rank and selection queries in $\mathcal{O}(1)$ time. Moreover, given the packed representations of $m$ binary strings of total length $n$, the data structures for all these strings can be constructed in $\mathcal{O}(m + n/\log n)$ time.*

Next, we provide a generalization of rank and selection queries specific to sequences of strings (strings whose characters are strings themselves). Let $W \in (\Sigma^*)^m$ be a sequence of $m$ strings.

**Prefix rank query** $\mathrm{rank}_{W,X}(j)$: Given $X \in \Sigma^*$ and $j \in [0 .. m]$, compute $|\{i \in [1 .. j] : X \text{ is a prefix of } W[i]\}|$.

**Prefix selection query** $\mathrm{select}_{W,X}(r)$: Given $X \in \Sigma^*$ and $r \in [1 .. \mathrm{rank}_{W,X}(m)]$, find the $r$th smallest element of $\{i \in [1 .. m] : X \text{ is a prefix of } W[i]\}$.

The following result, proved in Section 4.3 by building on the results of Belazzougui and Puglisi [9], provides an efficient implementation of prefix rank and selection queries. Note that we require $W$ to consist of same-length strings over an integer alphabet.

**Theorem 2.2.** *For all integers $m, \ell, \sigma \in \mathbb{Z}_{\ge 1}$ satisfying $m \ge \sigma^\ell \ge 2$, every constant $\epsilon > 0$, and every string $W \in ([0 .. \sigma)^\ell)^{\le m}$, there exists a data structure of size $\mathcal{O}(m)$ answering prefix rank*

queries in $\mathcal{O}(\ell^{\epsilon/2} \log \log m) = \mathcal{O}(\log^\epsilon m)$ *time and prefix selection queries in* $\mathcal{O}(\ell^{\epsilon/2}) = \mathcal{O}(\log^\epsilon m)$ *time. Moreover, it can be constructed in* $\mathcal{O}(m \min(\ell, \sqrt{\log m}))$ *time using* $\mathcal{O}(m)$ *working space given the packed representation of* $W$ *and the constant parameter* $\epsilon > 0$.

## 2.2 Range Counting and Selection

Let $A[1 \mathinner{\ldotp\ldotp} m]$ be an array of nonnegative integers. We define the following queries on $A$:

**Range counting query** $\mathsf{rcount}_A(v, j)$**:** Given an integer $v \geq 0$ and a position $j \in [0 \mathinner{\ldotp\ldotp} m]$, compute $|\{i \in [1 \mathinner{\ldotp\ldotp} j] : A[i] \geq v\}|$.

**Range selection query** $\mathsf{rselect}_A(v, r)$**:** Given integers $v \geq 0$ and $r \in [1 \mathinner{\ldotp\ldotp} \mathsf{rcount}_A(v, m)]$, find the $r$th smallest element of $\{i \in [1 \mathinner{\ldotp\ldotp} m] : A[i] \geq v\}$.

The currently fastest general-purpose data structure for range counting/selection queries is described in [18, Theorems 2.3 and 3.3]. The instances in our construction, however, satisfy an additional property, namely, that the sum $\sum_{i=1}^m A[i]$ is bounded. This lets us obtain a solution with faster queries and smaller construction time; see Section 4.4.

**Proposition 2.3.** *An array* $A[1 \mathinner{\ldotp\ldotp} m']$ *of* $m' \in [2 \mathinner{\ldotp\ldotp} m]$ *nonnegative integers satisfying* $\sum_{i=1}^{m'} A[i] = \mathcal{O}(m \log m)$ *can be preprocessed in* $\mathcal{O}(m)$ *time so that range counting and selection queries can be answered in* $\mathcal{O}(\log \log m)$ *time and* $\mathcal{O}(1)$ *time, respectively.*

## 2.3 String Synchronizing Sets

**Definition 2.4** ($\tau$-synchronizing set [50])**.** Let $T \in \Sigma^n$ be a string and let $\tau \in [1 \mathinner{\ldotp\ldotp} \lfloor \frac{n}{2} \rfloor]$ be a parameter. A set $\mathsf{S} \subseteq [1 \mathinner{\ldotp\ldotp} n - 2\tau + 1]$ is called a $\tau$-*synchronizing set* of $T$ if it satisfies the following *consistency* and *density* conditions:

1. If $T[i \mathinner{\ldotp\ldotp} i+2\tau] = T[j \mathinner{\ldotp\ldotp} j+2\tau]$, then $i \in \mathsf{S}$ holds if and only if $j \in \mathsf{S}$ (for $i, j \in [1 \mathinner{\ldotp\ldotp} n-2\tau+1]$),
2. $\mathsf{S} \cap [i \mathinner{\ldotp\ldotp} i + \tau] = \emptyset$ if and only if $i \in \mathsf{R}(\tau, T)$ (for $i \in [1 \mathinner{\ldotp\ldotp} n - 3\tau + 2]$), where

$$\mathsf{R}(\tau, T) := \{i \in [1 \mathinner{\ldotp\ldotp} |T| - 3\tau + 2] : \mathrm{per}(T[i \mathinner{\ldotp\ldotp} i + 3\tau - 2]) \leq \tfrac{1}{3}\tau\}.$$

In most applications, we want to minimize $|\mathsf{S}|$. Note, however, that the density condition imposes a lower bound $|\mathsf{S}| = \Omega(\frac{n}{\tau})$ for strings of length $n \geq 3\tau - 1$ that do not contain substrings of length $3\tau - 1$ which are periodic with period $\leq \frac{1}{3}\tau$. Thus, we cannot hope to achieve an upper bound improving in the worst case upon the following one.

**Theorem 2.5** ([50, Proposition 8.10])**.** *For any string* $T$ *of length* $n$ *and parameter* $\tau \in [1 \mathinner{\ldotp\ldotp} \lfloor \frac{n}{2} \rfloor]$, *there exists a* $\tau$-*synchronizing set* $\mathsf{S}$ *of size* $|\mathsf{S}| = \mathcal{O}\left(\frac{n}{\tau}\right)$. *Moreover, if* $T \in [0 \mathinner{\ldotp\ldotp} \sigma)^n$, *where* $\sigma = n^{\mathcal{O}(1)}$, *such* $\mathsf{S}$ *can be deterministically constructed in* $\mathcal{O}(n)$ *time.*

Note that when $\tau = \omega(1) \cap \mathcal{O}(\log_\sigma n)$ and $T \in [0 \mathinner{\ldotp\ldotp} \sigma)^n$ is given in the packed representation, the first part of Theorem 2.5 opens the possibility of an algorithm running in $\mathcal{O}(\frac{n}{\tau}) = o(n)$ time. In [50], it was shown that this lower bound is achievable (the upper bound $\tau = \mathcal{O}(\log_\sigma n)$ follows from the fact that every algorithm needs to at least read the input, which takes $\Theta(n / \log_\sigma n)$ time; thus, for larger $\tau$, the algorithm cannot run in $\mathcal{O}(\frac{n}{\tau})$ time).

**Theorem 2.6** ([50, Theorem 8.11])**.** *For every constant* $\mu < \frac{1}{5}$, *given the packed representation of a text* $T \in [0 \mathinner{\ldotp\ldotp} \sigma)^n$ *and a positive integer* $\tau \leq \mu \log_\sigma n$, *one can deterministically construct in* $\mathcal{O}(\frac{n}{\tau})$ *time a* $\tau$-*synchronizing set of size* $\mathcal{O}(\frac{n}{\tau})$.

## 3 Technical Overview

In this section, we give an overview of our data structures to answer SA and ISA queries (Section 3.1), pattern matching queries (Section 3.2), and suffix tree queries (Section 3.3). Each subsection contains a summary of the key new techniques.

### 3.1 SA and ISA Queries

Let $\epsilon \in (0, 1)$ and $T \in [0 \mathinner{..} \sigma)^n$, where $2 \leq \sigma < n^{1/7}$. In this section, we give an overview of our data structure to compute the value of $\mathrm{SA}[i]$ (resp. $\mathrm{ISA}[j]$) given any $i \in [1 \mathinner{..} n]$ (resp. $j \in [1 \mathinner{..} n]$) in $\mathcal{O}(\log^{\epsilon} n)$ time. The data structure uses $\mathcal{O}(n / \log_{\sigma} n)$ space. We assume $\sigma < n^{1/7}$ since for larger $\sigma$ the plain representations of SA and ISA use $\mathcal{O}(n \log n) = \mathcal{O}(n \log \sigma)$ bits and can be constructed in $\mathcal{O}(n)$ time [49].

Let $\tau = \lfloor \mu \log_{\sigma} n \rfloor$, where $\mu < \frac{1}{6}$ is a positive constant chosen so that $\tau \geq 1$ (such $\mu$ exists by $\sigma < n^{1/7}$). We use $\mathsf{R}$ as a shorthand for $\mathsf{R}(\tau, T)$ (see Definition 2.4). Our data structure to compute $\mathrm{SA}[i]$ (resp. $\mathrm{ISA}[j]$) works differently depending on whether $\mathrm{SA}[i] \in \mathsf{R}$ (resp. $j \in \mathsf{R}$). To check if $\mathrm{SA}[i] \in \mathsf{R}$, we store a bitvector $B_{3\tau-1}$ marking boundaries between the blocks of suffixes in SA sharing the length-$(3\tau-1)$ prefix. We also store the sequence $A_{\mathrm{short}}$ of those prefixes (by $\mu < \frac{1}{6}$, it needs $\mathcal{O}(\sigma^{3\tau-1}) = \mathcal{O}(n^{3\mu}) = o(n / \log_{\sigma} n)$ space). Given any $i \in [1 \mathinner{..} n]$, we can then check if $\mathrm{SA}[i] \in \mathsf{R}$ by first computing the block $k$ containing position $i$ using a rank query on $B_{3\tau-1}$, and then checking (using a lookup table) if $X = A_{\mathrm{short}}[k]$ satisfies $\mathrm{per}(X) \leq \frac{1}{3}\tau$. As for an ISA query, checking if $j \in \mathsf{R}$ only needs the lookup table (since we store $T$).

**The Nonperiodic Positions** We first focus on computing $\mathrm{ISA}[j]$. Let $j \notin \mathsf{R}$ and let $\mathsf{S}$ be a $\tau$-synchronizing set of $T$ of size $n' := |\mathsf{S}| = \mathcal{O}(n/\tau)$ (such $\mathsf{S}$ exists and can be quickly constructed using Theorem 2.6). The query algorithm relies on the following two observations:

*Observation 1: $\mathsf{S}$ induces a partitioning of SA into blocks.* The density condition of $\mathsf{S}$ implies $\mathsf{S} \cap [j \mathinner{..} j + \tau) \neq \emptyset$, i.e., the successor of $j$ in $\mathsf{S}$, denoted $s := \mathrm{succ}_{\mathsf{S}}(j)$, satisfies $s < j + \tau$. Hence, the string $X := T[j \mathinner{..} s + 2\tau)$, called the *distinguishing prefix* of $T[j \mathinner{..} n]$, is of length $|X| \leq 3\tau - 1$. By the local consistency of $\mathsf{S}$, the set $\mathcal{D}$ of distinguishing prefixes of all suffixes $T[j \mathinner{..} n]$ with $j \notin \mathsf{R}$ is prefix-free (i.e., no string in $\mathcal{D}$ is a prefix of another). All positions in $[1 \mathinner{..} n] \setminus \mathsf{R}$ in the SA of $T$ can thus be partitioned into disjoint blocks according to distinguishing prefixes. Since $\mathcal{D} \subseteq [0 \mathinner{..} \sigma)^{\leq 3\tau-1}$, the number of blocks is $\mathcal{O}(\sigma^{3\tau-1})$, so we can store their boundaries in a lookup table of size $\mathcal{O}(\sigma^{3\tau-1}) = \mathcal{O}(n^{3\mu}) = o(n / \log_{\sigma} n)$. To efficiently determine $\mathrm{succ}_{\mathsf{S}}(j)$, we store a bitvector marking positions in $\mathsf{S}$, augmented with $\mathcal{O}(1)$-time rank and select queries. Once the block $\mathrm{SA}(b \mathinner{..} e]$ for $X = T[j \mathinner{..} \mathrm{succ}_{\mathsf{S}}(j) + 2\tau)$ is found, it remains to locate $j$ within that block.

*Observation 2: The order in each block is consistent with $\mathsf{S}$.* Assume $\mathrm{SA}(b \mathinner{..} e]$ represents all suffixes of $T$ having $X = T[j \mathinner{..} \mathrm{succ}_{\mathsf{S}}(j) + 2\tau)$ as a prefix. By the consistency condition, letting $\delta_{\mathrm{text}} = |X| - 2\tau$, for every $i \in (b \mathinner{..} e]$, we have $\mathrm{succ}_{\mathsf{S}}(\mathrm{SA}[i]) = \mathrm{SA}[i] + \delta_{\mathrm{text}}$. Thus, letting $(s_i^{\mathrm{lex}})_{i \in [1 \mathinner{..} n']}$ contain $\mathsf{S}$ sorted by the corresponding suffixes $T[s_i^{\mathrm{lex}} \mathinner{..} n]$, positions in $\mathrm{SA}(b \mathinner{..} e]$ increased by $\delta_{\mathrm{text}}$ occur in $(s_i^{\mathrm{lex}})_{i \in [1 \mathinner{..} n']}$ in the same relative order. Hence, if we define $W[i] = \overline{X_i}$, where $X_i = T^{\infty}[s_i^{\mathrm{lex}} - \tau \mathinner{..} s_i^{\mathrm{lex}} + 2\tau)$, and select $W[y]$ as the $k$th string in $W$ having $\overline{X}$ as a prefix, then $s_y^{\mathrm{lex}} - \delta_{\mathrm{text}} = \mathrm{SA}[b + k]$ is the $k$th position in $\mathrm{SA}(b \mathinner{..} e]$. Thus, to obtain $\mathrm{ISA}[j]$, it suffices to find the index $y$ such that $s_y^{\mathrm{lex}} = \mathrm{succ}_{\mathsf{S}}(j)$. Then, the offset of $j$ in the block $\mathrm{SA}(b \mathinner{..} e]$ is $\mathrm{rank}_{W,\overline{X}}(y)$. To efficiently determine $y$, we store the permutation that maps elements of $\mathsf{S}$ sorted left-to-right to elements of $(s_i^{\mathrm{lex}})_{i \in [1 \mathinner{..} n']}$. This lets us determine $y$ in $\mathcal{O}(1)$ time. We then compute $\mathrm{rank}_{W,\overline{X}}(y)$ in $\mathcal{O}(\log^{\epsilon} n)$ time using Theorem 2.2; see Proposition 5.6.

Let us now turn to the computation of $\mathrm{SA}[i]$ when $\mathrm{SA}[i] \notin \mathsf{R}$. Using a rank query on $B_{3\tau-1}$ and an access to $A_{\mathrm{short}}$, we first determine the length-$(3\tau-1)$ prefix of $T[\mathrm{SA}[i] \mathinner{.\,.} n]$. A lookup table lets us retrieve the prefix $X \in \mathcal{D}$ of $T[\mathrm{SA}[i] \mathinner{.\,.} n]$ and the boundaries of the corresponding block $\mathrm{SA}(b \mathinner{.\,.} e)$. By Observation 2 above, it remains to determine the index $y$ of the $(i-b)$th leftmost string in $W$ having $\overline{X}$ as a prefix, which we compute in $\mathcal{O}(\log^\epsilon n)$ time as $y = \mathsf{select}_{W,\overline{X}}(i-b)$ using Theorem 2.2. We then have $\mathrm{SA}[i] = s_y^{\mathrm{lex}} - \delta_{\mathrm{text}}$, where $\delta_{\mathrm{text}} = |X| - 2\tau$. See Proposition 5.8 for details.

**The Periodic Positions** Let $j \in \mathsf{R}$. We again first focus on computing $\mathrm{ISA}[j]$. As before, we aim to find the location of $j$ in the block $\mathrm{SA}(b \mathinner{.\,.} e)$ containing all suffixes of $T$ prefixed with $X = T[j \mathinner{.\,.} j + 3\tau - 1)$. Note that all positions in $\mathrm{SA}(b \mathinner{.\,.} e)$ are in $\mathsf{R}$. The challenge is thus to devise a way to compare suffixes starting in $\mathsf{R}$. The problem with applying a similar approach as before is that the size of $\mathsf{R}$ can reach $\Theta(n)$. There exists, however, a subset of $\mathsf{R}$ that, when combined with a bitvector representing remaining positions in $\mathsf{R}$, can be applied here. We derive it as follows:

*Structure of $\mathsf{R}$ in the left-to-right (text) order:* The gap between $|X| = 3\tau - 1$ and $\mathrm{per}(X) \le \frac{1}{3}\tau$ ensures that every maximal block of positions in $\mathsf{R}$ corresponds to a $\tau$-*run*, i.e., a maximal substring of $T$ of length $\ge 3\tau - 1$ whose shortest period is $\le \frac{1}{3}\tau$ (Lemma 5.10). Since any two $\tau$-runs overlap by at most $\frac{2}{3}\tau$ positions (see the proof of Lemma 5.13), their number is $\mathcal{O}(n/\tau)$. We can thus succinctly encode $\mathsf{R}$ by storing the set $\mathsf{R}'$ of $\tau$-run starting positions.

*Structure of $\mathsf{R}$ in the lexicographic order:* For $x \in \mathsf{R}$, let $e(x)$ denote the position following the $\tau$-run containing $x$. Observe that, for every $x \in \mathsf{R}$, we can uniquely write $T[x \mathinner{.\,.} e(x)) = H'H^kH''$, where $H$ is the lexicographically smallest rotation of $T[x \mathinner{.\,.} x + p)$, $p = \mathrm{per}(T[x \mathinner{.\,.} e(x)))$, and $H'$ (resp. $H''$) is a proper suffix (resp. prefix) of $H$ (Section 5.3.1). Denote $\text{L-root}(x) = H$, $\text{L-head}(x) = |H'|$, $\text{L-exp}(x) = k$, and $\text{L-tail}(x) = |H''|$. Let also $\mathrm{type}(x) = -1$ if $T[e(x)] \prec T[e(x) - |H|]$ and $\mathrm{type}(x) = +1$ otherwise. Then, in $\mathrm{SA}(b \mathinner{.\,.} e)$, all positions $x$ with $\mathrm{type}(x) = -1$ precede all $x$ with $\mathrm{type}(x) = +1$. Moreover, the value of $e(x) - x$ is non-decreasing (resp. non-increasing) among the positions $x$ with $\mathrm{type}(x) = -1$ (resp. $\mathrm{type}(x) = +1$); see Lemma 5.11.

Denote $\mathsf{R}^- = \{x \in \mathsf{R} : \mathrm{type}(x) = -1\}$, $\mathsf{R}'^- = \mathsf{R}' \cap \mathsf{R}^-$, $\mathsf{R}_H = \{x \in \mathsf{R} : \text{L-root}(x) = H\}$, $\mathsf{R}'_H = \mathsf{R}'^- \cap \mathsf{R}_H$, $\mathsf{R}_{s,H} = \{x \in \mathsf{R}_H : \text{L-head}(x) = s\}$, and $\mathsf{R}^-_{s,H} = \mathsf{R}^- \cap \mathsf{R}_{s,H}$. Assume $\mathrm{type}(j) = -1$ (the case of $\mathrm{type}(j) = +1$ is symmetric), $\text{L-head}(j) = s$, and $\text{L-root}(j) = H$. Given the above structural insights, we can phrase locating $j$ in $\mathrm{SA}(b \mathinner{.\,.} e)$ as counting the positions $x \in \mathsf{R}^-_{s,H}$ satisfying $T[x \mathinner{.\,.} n] \preceq T[j \mathinner{.\,.} n]$. By the analysis above, all such positions satisfy $e(x) - x \le e(j) - j$ and hence $\text{L-exp}(x) \le \text{L-exp}(j)$. We first compute the size of $\mathrm{Pos}^{\mathsf{a}}(j) := \{x \in \mathsf{R}^-_{s,H} : \text{L-exp}(x) \le \text{L-exp}(j)\}$ and then subtract the size of $\mathrm{Pos}^{\mathsf{s}}(j) := \{x \in \mathsf{R}^-_{s,H} : \text{L-exp}(x) = \text{L-exp}(j) \text{ and } T[x \mathinner{.\,.} n] \succ T[j \mathinner{.\,.} n]\}$ as follows:[6]

- Since $|\mathrm{Pos}^{\mathsf{a}}(j)|$ only depends on $\text{L-exp}(j)$, it suffices to store a bitvector $B_{\mathrm{exp}}$ marking the boundaries in $\mathrm{SA}$ between blocks of positions with subsequent values of $\text{L-exp}$. Computing $|\mathrm{Pos}^{\mathsf{a}}(j)|$ then reduces to $\mathcal{O}(1)$-time rank and selection queries on $B_{\mathrm{exp}}$ (Proposition 5.19).
- As for $|\mathrm{Pos}^{\mathsf{s}}(j)|$, we observe that $\mathrm{Pos}^{\mathsf{s}}(j)$ contains at most one position in each $\tau$-run (Lemma 5.20). We store all $x \in \mathsf{R}'^-_H$ sorted by the suffix starting right after the last occurrence of $H$, i.e., $T[e^{\mathrm{full}}(x) \mathinner{.\,.} n]$, where $e^{\mathrm{full}}(x) := e(x) - \text{L-tail}(x)$. In a separate array, we also record $e^{\mathrm{full}}(x) - x$ at the corresponding position. This lets us compute $|\mathrm{Pos}^{\mathsf{s}}(j)|$ by first locating the block of positions $x \in \mathsf{R}'^-_H$ for which $T[e^{\mathrm{full}}(x) \mathinner{.\,.} n] \succ T[e^{\mathrm{full}}(j) \mathinner{.\,.} n]$, and then counting the ones with $e^{\mathrm{full}}(x) - x \ge e^{\mathrm{full}}(j) - j$ (Proposition 5.21). Since the sum of

---

[6]Note that here we first overestimate the number of smaller suffixes in $\mathrm{SA}(b \mathinner{.\,.} e)$ and then subtract the larger suffixes. We explain the reason for this counterintuitive approach in Remark 5.22.

$e^{\text{full}}(x) - x$ over all $x \in \mathsf{R}'$ is $\mathcal{O}(n)$ (Section 5.3.2), we use a specialized structure for range counting (Proposition 2.3), bypassing the general $\Omega(\frac{\log n}{\log \log n})$-time lower bound [76].

Let us now turn to an $\mathrm{SA}[i]$ query with $\mathrm{SA}[i] \in \mathsf{R}^-$. First, we determine $T[\mathrm{SA}[i] \mathinner{..} \mathrm{SA}[i]+3\tau-1]$ using $B_{3\tau-1}$ and $A_{\text{short}}$, as well as $s = \text{L-head}(\mathrm{SA}[i])$ and $H = \text{L-root}(\mathrm{SA}[i])$ using a lookup table (Proposition 5.16). Then, rank and selection queries on $B_{\text{exp}}$ let us easily determine $\text{L-exp}(\mathrm{SA}[i])$ and $|\text{Pos}^{\mathsf{s}}(\mathrm{SA}[i])|$ (Proposition 5.24). As explained above, to compute $\mathrm{SA}[i]$, it remains to first select the $k$th (where $k = |\text{Pos}^{\mathsf{s}}(\mathrm{SA}[i])| + 1$) largest element $j \in \mathsf{R}_H'^-$ according to the string $T[e^{\text{full}}(j) \mathinner{..} n]$, among positions $j'' \in \mathsf{R}_H'^-$ satisfying $e^{\text{full}}(j'') - j'' \geq \text{L-head}(\mathrm{SA}[i]) + \text{L-exp}(\mathrm{SA}[i]) \cdot |H|$. The position $j' \in [j \mathinner{..} e(j) - 3\tau + 2)$ with $\text{L-exp}(j') = \text{L-exp}(\mathrm{SA}[i])$ and $\text{L-head}(j') = \text{L-head}(\mathrm{SA}[i])$ must then satisfy $\mathrm{SA}[i] = j'$. To compute $j$, we use our specialized data structure for range queries (Proposition 2.3). Position $j'$ is then obtained by subtracting $\text{L-head}(\mathrm{SA}[i]) + \text{L-exp}(\mathrm{SA}[i]) \cdot |H|$ from $e^{\text{full}}(j)$ (Proposition 5.25).

In total, the query time for periodic positions is $\mathcal{O}(\log \log n)$ (Propositions 5.23 and 5.26).

**Summary of New Techniques**  The key distinctive feature of our technique is the use of local consistency without general orthogonal range queries, present in prior approaches [20, 50, 51, 52]. This lets us sidestep Pătraşcu's $\Omega(\frac{\log n}{\log \log n})$ lower bound [76], which is achieved in three steps:

- We replace range counting/selection in the nonperiodic case with *prefix rank* and *prefix selection*, for which we propose a new tradeoff by plugging in the technique of Belazzougui and Puglisi [9]. This reveals the power of our reduction: we achieve a non-trivial tradeoff for complex queries by solving a simple bit-permuting problem (note that the tradeoff behind Theorem 2.2, e.g., occupies only 1.5 pages in Section 4.3; the bulk of our paper is the reduction).
- We replace range counting/selection in the periodic case by observing that the sum of coordinates is small (Section 5.3.2). This lets us use a specialized solution (Section 4.4).
- The above two cases occur at query time. The third case concerns the construction of the structure. More precisely, we completely eliminate range queries naturally occurring during the construction of components for periodic positions [50, 51] by a complex bit-optimal algorithm for the construction of the bitvector $B_{\text{exp}}$ (see Section 5.3.6).

As a result, we obtain a very general reduction stated in Theorem 5.32.

## 3.2 Pattern Matching Queries

Let $\epsilon \in (0, 1)$ and $T \in [0 \mathinner{..} \sigma)^n$ be as in Section 3.1. We now give an overview of our data structure that, given a packed representation of any pattern $P \in [0 \mathinner{..} \sigma)^m$, returns the pair of indexes $(b, e) = (\text{RangeBeg}(P, T), \text{RangeEnd}(P, T))$, i.e., the boundaries of the SA block containing all suffixes having $P$ as a prefix, in $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$ time. Note that having this range immediately gives us $|\text{Occ}(P, T)| = e - b$, i.e., it implements *pattern counting*. Moreover, combined with the result from Section 3.1, we obtain *pattern reporting*, i.e., we can enumerate $\text{Occ}(P, T)$ in $\mathcal{O}(m/\log_\sigma n + (|\text{Occ}(P, T)| + 1) \log^\epsilon n)$ time.

Let $\tau$ be as in Section 3.1. Our structure to compute $(\text{RangeBeg}(P, T), \text{RangeEnd}(P, T))$ works differently depending on whether $m \geq 3\tau - 1$ and whether $\text{per}(P[1 \mathinner{..} 3\tau - 1]) \leq \frac{1}{3}\tau$ (such $P$ is called *periodic*) or not. Checking if $P$ is periodic is easily implemented via a lookup table.

**The Nonperiodic Patterns**  Let us assume $m \geq 3\tau - 1$ (shorter patterns are handled using a precomputed array) and let $\mathsf{S}$ be as in Section 3.1. The basic idea of the pattern matching query is to decompose $P = XY$, where $X \in \mathcal{D}$, and then utilize the following observation about $\mathsf{S}$ (generalizing the second observation in Section 3.1):

*Observation: The order in the suffix array range corresponding to* $\mathrm{Occ}(P,T)$ *is consistent with* $\mathsf{S}$. Let $(b,e) = (\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$. By the consistency condition, for every $i \in (b\mathinner{.\,.}e]$, we have $\mathrm{succ}_{\mathsf{S}}(\mathrm{SA}[i]) = \mathrm{SA}[i] + \delta_{\text{text}}$, where $\delta_{\text{text}} = |X| - 2\tau$. Thus, letting $(s_i^{\text{lex}})_{i \in [1\mathinner{.\,.}n']}$ be defined as in Section 3.1, the positions in $\mathrm{SA}(b\mathinner{.\,.}e]$ increased by $\delta_{\text{text}}$ occur in $(s_i^{\text{lex}})_{i \in [1\mathinner{.\,.}n']}$ in the same relative order. Therefore, to compute $|\mathrm{Occ}(P,T)|$, it suffices to first locate a range of $(s_i^{\text{lex}})_{i \in [1\mathinner{.\,.}n']}$ consisting of positions followed by $P(\delta_{\text{text}}\mathinner{.\,.}m]$ in $T$, and then count those which are additionally preceded with $X[1\mathinner{.\,.}\delta_{\text{text}}]$. The first goal is implemented in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time via a compact trie over the set of strings $\{T[s_i^{\text{lex}}\mathinner{.\,.}n]\}_{i \in [1\mathinner{.\,.}n']}$, reinterpreted as strings over alphabet of size $n^{\Theta(1)}$. The second step reduces to a prefix rank query over $W[1\mathinner{.\,.}n']$ (Section 3.1). This approach easily generalizes to return $(b,e)$ instead of $|\mathrm{Occ}(P,T)|$; see Lemma 6.5.

**The Periodic Patterns**  Let us now assume that $m \geq 3\tau - 1$ and $\mathrm{per}(P[1\mathinner{.\,.}3\tau-1]) \leq \frac{1}{3}\tau$. We first generalize the notion of L-root$(x)$, $e(x)$, and all other functions from positions to strings (Section 6.3.1). The main idea is to decompose $P$ into the periodic prefix $P[1\mathinner{.\,.}e(P))$ and the remaining suffix $P[e(P)\mathinner{.\,.}|P|]$. Let us consider the harder case when $e(P) = |P| + 1$ (see Lemma 6.14). We define $\mathrm{Occ}^{\mathsf{a}}(P,T) = \{j \in \mathsf{R}_{s,H} \cap \mathrm{Occ}(P,T) : \text{L-exp}(j) > \text{L-exp}(P)\}$ and $\mathrm{Occ}^{\mathsf{s}}(P,T) = \{j \in \mathsf{R}_{s,H} \cap \mathrm{Occ}(P,T) : \text{L-exp}(j) = \text{L-exp}(P)\}$, where $H = \text{L-root}(P)$ and $s = \text{L-head}(P)$. The value $|\mathrm{Occ}(P,T)|$ is determined in two steps:

- First, we compute the size of $\mathrm{Occ}^{\mathsf{a}-}(P,T) := \mathrm{Occ}^{\mathsf{a}}(P,T) \cap \mathsf{R}^-$ (the size of $\mathrm{Occ}^{\mathsf{a}+}(P,T)$ is computed symmetrically) as in Section 3.1 by utilizing rank and selection queries on $B_{\exp}$ (Proposition 6.15). This only requires knowing $\text{L-exp}(P)$, which can be retrieved in $\mathcal{O}(1 + m/\log_\sigma n)$ time (Proposition 6.12).
- Next, we compute the size of $\mathrm{Occ}^{\mathsf{s}-}(P,T) := \mathrm{Occ}^{\mathsf{s}}(P,T) \cap \mathsf{R}^-$. We first show that $\mathrm{Occ}^{\mathsf{s}-}(P,T)$ contains at most one position in every $\tau$-run (Lemma 6.16), i.e., an analogue of Lemma 5.20. The computation is similar as in Section 3.1, except that we use a trie over meta-symbols to find the range of positions $x \in \mathsf{R}'^-_H$ with $T[e^{\text{full}}(x)\mathinner{.\,.}n]$ prefixed by $P[e^{\text{full}}(P)\mathinner{.\,.}m]$. We then perform an $\mathcal{O}(\log\log n)$-time range query (Proposition 6.17). A small complication is to separate positions $x \in \mathsf{R}'^-$ with different L-root$(x)$ in the trie. For this, we insert into the trie suffixes starting slightly earlier than $e^{\text{full}}(x)$; see Section 6.3.2.

The above algorithm generalizes to the computation of $(b,e)$, rather than $|\mathrm{Occ}(P,T)|$, except that handling "fully periodic" patterns (with $e(P) = |P|+1$) requires some care (see Remark 6.25).

**Summary of New Techniques**  Our key technical contributions are as follows:

- We show how to directly apply string synchronizing sets [50] to the problem of pattern matching (not by simply using SA/ISA queries) and consequently obtain a very efficient reduction from pattern matching to prefix rank and prefix selection queries.
- To achieve this, we prove several new combinatorial results for periodic patterns (Section 6.3.1), and then show for efficiently apply them (Sections 6.3.2 to 6.3.4).
- As a result, we obtain the first optimal-size pattern matching index that is constructible in $o(n)$ time and supports:
  - pattern occurrence counting in $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$ time, and
  - pattern occurrence reporting in $\mathcal{O}(m/\log_\sigma n + (|\mathrm{Occ}(P,T)| + 1)\log^\epsilon n)$ time.

This improves over [62, 63] in either construction or query time and, perhaps more importantly, provides a very general reduction that enables achieving further time-space tradeoffs much easier (Theorem 6.31).

## 3.3 Suffix Tree Queries

Let $\epsilon \in (0,1)$ and $T \in [0 \mathinner{.\,.} \sigma)^n$ be as in Section 3.1. In this section, we outline how to extend the techniques presented in Sections 3.1 and 3.2 to obtain the full suffix tree functionality in optimal $\mathcal{O}(n/\log_\sigma n)$ space. All operations (see Table 1) are supported in $\mathcal{O}(\log^\epsilon n)$ time, and the data structure can be constructed in $\mathcal{O}(n \min(1, \log \sigma/\sqrt{\log n}))$ time and $\mathcal{O}(n/\log_\sigma n)$ working space. Thus, e.g., for $\sigma = 2$, our construction takes $\mathcal{O}(n/\sqrt{\log n}) = o(n)$ time.

Each node $v$ of the suffix tree of $T$, denoted $\mathcal{T}_{\mathrm{st}}$, is encoded either as $(j, \ell)$ such that $T[j \mathinner{.\,.} j + \ell] = \mathrm{str}(v)$ or as $(b, e) = (\mathrm{RangeBeg}(\mathrm{str}(v), T), \mathrm{RangeEnd}(\mathrm{str}(v), T))$, where $\mathrm{str}(v)$ is the string represented by $v$. Since the latter representation is more common (e.g. [31]), we adopt it as the default interface and denote $\mathrm{repr}(v)$.

In this overview, we focus on the child operation, which illustrates some of our key techniques. We remark, however, that other operations require different combinatorial insights.

Let $\tau$ be as in Section 3.1. Our structure works differently depending on whether the operation is performed on a node $v$ such that $\mathrm{str}(v)$ is periodic or not (see Section 3.2).

**The Nonperiodic Nodes**  Let $v$ be a node of $\mathcal{T}_{\mathrm{st}}$ such that $\mathrm{str}(v)$ is nonperiodic and let $c \in \Sigma$. Our goal is to compute $\mathrm{repr}(\mathrm{child}(v, c))$ given $\mathrm{repr}(v)$. Let $\mathsf{S}$ again be as in Section 3.1. The basic idea is to reduce the computation concerning the SA-interval for string $\mathrm{str}(v)$ to the computation involving only positions in $\mathsf{S}$. For this purpose, we store the compact trie $\mathcal{T}_{\mathsf{S}}$ for $\{T[s_i^{\mathrm{lex}} \mathinner{.\,.} n]\}_{i \in [1 \mathinner{.\,.} n']}$. Typically, each operation on $\mathcal{T}_{\mathrm{st}}$ then involves the following steps:

1. Map the input node $v$ of $\mathcal{T}_{\mathrm{st}}$ (given as $\mathrm{repr}(v)$) to some node $u$ of $\mathcal{T}_{\mathsf{S}}$,
2. Perform some operation in $\mathcal{T}_{\mathsf{S}}$ resulting in a node $u'$ (in our case, $u' = \mathrm{child}(u, c)$),
3. Map $u'$ back to some node $v'$ of $\mathcal{T}_{\mathrm{st}}$ (producing $\mathrm{repr}(v')$ as output).

*Mapping from $\mathcal{T}_{\mathrm{st}}$ to $\mathcal{T}_{\mathsf{S}}$:* Let $(b, e) = \mathrm{repr}(v)$ and let $X \in \mathcal{D}$ be the distinguishing prefix of $\mathrm{str}(v)$. By the observation for nonperiodic patterns in Section 3.2, the left-to-right order of leaves of $\mathcal{T}_{\mathsf{S}}$ corresponding to suffixes in $\mathrm{SA}(b \mathinner{.\,.} e]$ shifted by $\delta_{\mathrm{text}} = |X| - 2\tau$ is consistent with their order in $\mathrm{SA}(b \mathinner{.\,.} e]$. Moreover, since we know how to compute the position in $(s_i^{\mathrm{lex}})_{i \in [1 \mathinner{.\,.} n']}$ corresponding to suffix $T[\mathrm{SA}[i] \mathinner{.\,.} n]$ for any $i \in [1 \mathinner{.\,.} n]$ such that $\mathrm{SA}[i] \in [1 \mathinner{.\,.} n] \setminus \mathsf{R}$ (see Section 3.1), we can compute the pointer to the leaf of $\mathcal{T}_{\mathsf{S}}$ corresponding to any suffix in $\mathrm{SA}(b \mathinner{.\,.} e]$. This implies that: (1) there exists a node $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v) := u$ in $\mathcal{T}_{\mathsf{S}}$ such that $\mathrm{str}(v) = X[1 \mathinner{.\,.} \delta_{\mathrm{text}}] \cdot \mathrm{str}(u)$, and (2) a pointer to $u$ can be computed via a lowest common ancestor (LCA) query from the leaves of $\mathcal{T}_{\mathsf{S}}$ corresponding to first and last suffix in $\mathrm{SA}(b \mathinner{.\,.} e]$ (see Section 7.2.2).

*Mapping from $\mathcal{T}_{\mathsf{S}}$ to $\mathcal{T}_{\mathrm{st}}$:* After computing $u' = \mathrm{child}(u, c)$, the next step is to go back to $\mathcal{T}_{\mathrm{st}}$. Our approach exploits a similar principle as when computing $\mathrm{ISA}[j]$: knowing $X \in \mathcal{D}$ and a position in $(s_i^{\mathrm{lex}})_{i \in [1 \mathinner{.\,.} n']}$ lets us determine the corresponding position in SA. More generally, given $X \in \mathcal{D}$ and an interval of $(s_i^{\mathrm{lex}})_{i \in [1 \mathinner{.\,.} n']}$, we can retrieve the corresponding interval in SA. The former is precomputed and stored with each node of $\mathcal{T}_{\mathsf{S}}$. Note, however, the following complication: $\mathcal{T}_{\mathsf{S}}$ may have extra nodes between $u = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v)$ and $\widehat{u} = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(\mathrm{child}(v, c))$, and then $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(\mathrm{child}(v, c)) \neq \mathrm{child}(\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v), c)$ (see also Remark 7.23). We thus need to first prove that applying the inverse mapping to *any* node between $u$ and $\widehat{u}$ (in particular, to $u'$) yields $\mathrm{repr}(\mathrm{child}(v, c))$ (Lemma 7.22). This exploits properties of $\mathsf{S}$ specific to the child operation; we omit the details here but remark that this step differs among core operations (see, e.g., Lemmas 7.20 and 7.27).

**The Periodic Nodes**  Let us now assume that $\mathrm{str}(v)$ is periodic. The basic idea is similar as above: we keep a compact trie (denoted $\mathcal{T}_{\mathsf{Z}}$) letting us search the suffixes in the set $\{T[e^{\mathrm{full}}(j) \mathinner{.\,.} n]\}_{j \in \mathsf{R}'^-}$. Although the implementation of mapping and combinatorial proofs are

more technical, establishing these higher-level navigation primitives results in simpler and more concise implementation of queries (see, e.g., Lemma 7.52).

**Summary of New Techniques**   Our key technical contributions are as follows:

- We show how to reduce all operations of a suffix tree to prefix rank and selection queries, resulting in the first $o(n)$-time construction of optimal-size compressed suffix tree, with all operations simultaneously matching the state-of-the-art [14, 16, 31, 34, 75, 77, 80].
- To achieve this, we first define and efficiently implement mappings between the nodes of $\mathcal{T}_{\mathrm{st}}$ and of two auxiliary tries $\mathcal{T}_{\mathsf{S}}$ and $\mathcal{T}_{\mathsf{Z}}$. We then prove new combinatorial results (see, e.g., Lemmas 7.20, 7.22, 7.25, 7.27, 7.43, 7.47, 7.50, 7.52, and 7.67) showing that these high-level navigation primitives correctly handle all suffix tree operations.

## 4   Auxiliary Tools

### 4.1   Weighted Ancestors

Consider a rooted tree $\mathcal{T}$. Let $\mathrm{root}(\mathcal{T})$ denote the node at depth 0 and let $\mathrm{parent}(v)$ denote the immediate ancestor of each node $v \neq \mathrm{root}(\mathcal{T})$. We let $\mathrm{parent}(\mathrm{root}(\mathcal{T})) = \bot$. Assume that each node $v$ has an associated weight $w(v)$ such that $w(\mathrm{root}(\mathcal{T})) = 0$ and, for every $v \neq \mathrm{root}(\mathcal{T})$, it holds $w(\mathrm{parent}(v)) < w(v)$. We then say that the weight function $w$ is *monotone*. Given any node $v$ of $\mathcal{T}$ and an integer $0 \leq d \leq w(v)$, we define $v' = \mathrm{WA}(v, d)$ (the *weighted ancestor* [25]) as the (unique) ancestor of $v$ in $\mathcal{T}$ that satisfies $w(v') \geq d$ and for which $w(v')$ is minimized.

**Theorem 4.1** ([2, Section 6.2.1]). *Let $\mathcal{T}$ be a rooted tree with $n \leq N$ nodes and a monotone weight function $w$ mapping nodes to $[0 . . N)$. There exists a data structure of size $\mathcal{O}(n)$ that answers weighted ancestor queries in $\mathcal{T}$ in $\mathcal{O}(\log \log N)$ time after $\mathcal{O}(n \log_n N)$-time preprocessing.*

*Proof.* A solution for $N = n$ was presented in [2, Section 6.2.1]. To generalize it to the case of $N \geq n$, it suffices to map all node weights to their ranks. For this, we sort all the node weights in $\mathcal{O}(n \log_n N)$ time (radix sort) and build a deterministic predecessor structure [30, Proposition 2] in $\mathcal{O}(n)$ time so that the rank of a query threshold can be computed in $\mathcal{O}(\log \log N)$ time.   $\square$

### 4.2   Tries and Compact Tries

A set of strings $\mathcal{S} \subseteq \Sigma^+$ is *prefix-free* if there are no $S, S' \in \mathcal{S}$ such that $S$ is a proper prefix of $S'$. For any prefix-free set of string $\mathcal{S} \subseteq \Sigma^+$, its *trie* is a minimal rooted tree $\mathcal{T}$, with each edge labelled by some $c \in \Sigma$, such that: (1) no two edges outgoing from the same node have the same label, (2) for each $S \in \mathcal{S}$ there exists a path from $\mathrm{root}(\mathcal{T})$ to some node such that the concatenation of edge-labels on that path is equal to $S$, and (3) children of every node are ordered according to the lexicographical rank of the connecting edge. A *compact trie* of $\mathcal{S}$ is a trie of $\mathcal{S}$ in which all maximal unary paths have been replaced with edges labelled by substrings of elements of $\mathcal{S}$. The nodes of the trie omitted in the compact trie are referred to as *implicit*. All other nodes are *explicit*. Unless explicitly stated otherwise, by *node* we always mean an explicit node.

For any node $v$ of a (compact) trie $\mathcal{T}$, by $\mathrm{str}(v)$ we denote the *label* of $v$, i.e., the string obtained by concatenating the labels of all edges on the path from $\mathrm{root}(\mathcal{T})$ to $v$. We denote $\mathrm{sdepth}(v) = |\mathrm{str}(v)|$. The parent of $v$ in denoted $\mathrm{parent}(v)$. For any $c \in \Sigma$, we define $\mathrm{child}(v, c)$ as a child $v'$ of $v$ such that $\mathrm{str}(v')[|\mathrm{str}(v)| + 1] = c$, or $\bot$ if no such node exists. For any $c \in \Sigma$, we also define $\mathrm{pred}(v, c)$ as follows:

- If there exists $c' < c$ such that $\mathrm{child}(v, c') \neq \bot$, then we let $\mathrm{pred}(v, c) = \mathrm{child}(v, c_{\max})$, where $c_{\max} = \max\{c' \in [0 \mathinner{\ldotp\ldotp} c) : \mathrm{child}(v, c') \neq \bot\}$.
- Otherwise, we let $\mathrm{pred}(v, c) = \bot$.

We define $(\mathrm{lrank}(v), \mathrm{rrank}(v))$ as a pair of integers satisfying $\mathrm{lrank}(v) = |\{S \in \mathcal{S} : S \prec \mathrm{str}(v)\}|$ and $\mathrm{rrank}(v) - \mathrm{lrank}(v) = |\{S \in \mathcal{S} : \mathrm{str}(v) \text{ is a prefix of } S\}|$. Observe that then collecting every $i$th leftmost leaf of $\mathcal{T}$, where $i \in (\mathrm{lrank}(v) \mathinner{\ldotp\ldotp} \mathrm{rrank}(v)]$ results in precisely the set of leaves in the subtree rooted in $v$. Given any node $v$ of $\mathcal{T}$ and an integer $0 \leq d \leq |\mathrm{str}(v)|$, we let $v' = \mathrm{WA}(v, d)$ to be the weighted ancestor of $v$ assuming the weight of each node is defined as $w(v) = \mathrm{sdepth}(v)$. Thus, $v'$ is the (unique) ancestor of $v$ in $\mathcal{T}$ that satisfies $\mathrm{sdepth}(v') \geq d$ and for which $\mathrm{sdepth}(v')$ is minimized. For any two nodes $v_1$ and $v_2$, the node $v = \mathrm{LCA}(v_1, v_2)$ (the *lowest common ancestor*) is defined as the (unique) ancestor of both $v_1$ and $v_2$ with the maximal depth.

**Observation 4.2.** *If $v_1$ and $v_2$ are nodes of a (compact) trie, then letting $v = \mathrm{LCA}(v_1, v_2)$ and $\ell = \mathrm{lcp}(\mathrm{str}(v_1), \mathrm{str}(v_2))$, it holds $\mathrm{sdepth}(v) = \ell$ and $\mathrm{str}(v) = \mathrm{str}(v_1)[1 \mathinner{\ldotp\ldotp} \ell] = \mathrm{str}(v_2)[1 \mathinner{\ldotp\ldotp} \ell]$.*

### 4.2.1 Small Alphabet

**Proposition 4.3.** *Given a packed representation of $T \in [0 \mathinner{\ldotp\ldotp} \sigma)^n$ with $2 \leq \sigma \leq n$ and an array $A[1 \mathinner{\ldotp\ldotp} q]$ of $q$ positions in $T$ such that, for any $1 \leq i < j \leq q$, it holds $T[A[i] \mathinner{\ldotp\ldotp} n] \prec T[A[j] \mathinner{\ldotp\ldotp} n]$, we can in $\mathcal{O}(q + n/\log_\sigma n)$ time construct a representation of the compact trie $\mathcal{T}$ of the set $\{T[A[i] \mathinner{\ldotp\ldotp} n] : i \in [1 \mathinner{\ldotp\ldotp} q]\}$, augmented with auxiliary data structures to support the following operations on $\mathcal{T}$ in $\mathcal{O}(1)$ time:*

- *Given $i \in [1 \mathinner{\ldotp\ldotp} q]$ return the $i$th leftmost leaf of $\mathcal{T}$,*
- *Given pointers to nodes $v_1$ and $v_2$ return a pointer to $\mathrm{LCA}(v_1, v_2)$,*
- *Given a pointer to node $v$, return $(\mathrm{lrank}(v), \mathrm{rrank}(v))$ and $\mathrm{sdepth}(v)$.*

*It also supports the following operations in $\mathcal{O}(\log\log n)$ time:*

- *Given a pointer to node $v$ and $d$ such that $0 \leq d \leq |\mathrm{str}(v)|$, return the pointer to $\mathrm{WA}(v, d)$,*
- *Given a pointer to node $v$ and $c \in \Sigma$, check if $\mathrm{pred}(v, c) \neq \bot$ (resp. $\mathrm{child}(v, c) \neq \bot$), and if so, return the pointer to $\mathrm{pred}(v, c)$ (resp. $\mathrm{child}(v, c)$).*

*Proof.* The data structure consists of five components:

1. The packed representation of $T$ using $\mathcal{O}(n/\log_\sigma n)$ space.
2. The compact trie $\mathcal{T}$. Since we assumed that $T[n]$ is unique in $T$ (see Section 2), the set $\{T[A[i] \mathinner{\ldotp\ldotp} n]\}_{i \in [1 \mathinner{\ldotp\ldotp} q]}$ is prefix-free, and hence $\mathcal{T}$ has exactly $q$ leaves. Each node $v$ of $\mathcal{T}$ stores the precomputed values $\mathrm{sdepth}(v)$, $(\mathrm{lrank}(v), \mathrm{rrank}(v))$, and a predecessor data structure that, given any $c \in [0 \mathinner{\ldotp\ldotp} \sigma)$, returns a pointer to $\mathrm{pred}(v, c)$. Using the structure from [30, Proposition 2], we achieve linear space and $\mathcal{O}(\log\log n)$ query time.
3. The array of pointers $L[1 \mathinner{\ldotp\ldotp} q]$ such that $L[i]$ is the pointer to the $i$th leftmost leaf of $\mathcal{T}$.
4. The data structure of Bender and Farach-Colton [10] that augments $\mathcal{T}$ with support for LCA queries. The structure needs $\mathcal{O}(q)$ space and answers queries in $\mathcal{O}(1)$ time.
5. The data structure from Theorem 4.1 for $\mathcal{T}$ with the weight function $\mathrm{sdepth}(v) \in [0 \mathinner{\ldotp\ldotp} n]$. The structure needs $\mathcal{O}(q)$ space and answers queries in $\mathcal{O}(\log\log n)$ time.

In total, the data structure takes $\mathcal{O}(q + n/\log_\sigma n)$ space.

Using the above structures, the implementation of all queries in the claim follows immediately (note that $\mathrm{child}(v, c)$ can be determined using $\mathrm{pred}(v, c + 1)$ and the packed representation of $T$).

*Construction algorithm* The data structure is constructed as follows. We start by building a data structure that supports LCE queries for suffixes of $T$. Using [50, Theorem 5.4], the construction takes $\mathcal{O}(n/\log_\sigma n)$ time, and the resulting data structure answers queries in $\mathcal{O}(1)$

time. We then construct $\mathcal{T}$ by inserting elements of $\{T[A[i] \mathinner{\ldotp\ldotp} n]\}_{i \in [1 \mathinner{\ldotp\ldotp} q]}$ in the order given by $A$. We maintain a stack containing the internal nodes on the rightmost path, with the deepest node on top. When inserting each string, we first determine the depth at which that string branches from the rightmost path using LCE queries on $T$. We then update the rightmost path of the trie. Adding each string first removes some elements from the stack, and then adds at most two new elements. Since the total number of elements pushed on stack is $\mathcal{O}(q)$, the construction of $\mathcal{T}$ takes $\mathcal{O}(q)$ time. During the construction, we record the value sdepth$(v)$ in each node and collect pointers to consecutive leaves in the $L[1 \mathinner{\ldotp\ldotp} q]$ array. With the single traversal of $\mathcal{T}$, we then precompute in $\mathcal{O}(q)$ time the values lrank$(v)$ and rrank$(v)$ for every node $v$ (note that at this point, pointers to all children of each node are stored simply using a list, since this traversal does not require fast lookups or predecessor queries). Next, we construct the predecessor data structure for every node. Since the keys in every node are sorted, using [30, Proposition 2], over all nodes of $\mathcal{T}$, the construction takes $\mathcal{O}(q)$ time. Finally, we construct the data structures supporting LCA and WA queries on $\mathcal{T}$. Using [10] and Theorem 4.1, this takes $\mathcal{O}(q)$ and $\mathcal{O}(q \log_q n) = \mathcal{O}((q + \sqrt{n}) \log_{q+\sqrt{n}} n) = \mathcal{O}(q + \sqrt{n}) = \mathcal{O}(q + n/\log_\sigma n)$ time, respectively. $\qquad\square$

### 4.2.2 Large Alphabet

**Proposition 4.4.** *Given a packed representation of $T \in [0 \mathinner{\ldotp\ldotp} \sigma)^n$ with $2 \le \sigma < n^{1/7}$ and an array $A[1 \mathinner{\ldotp\ldotp} q]$ of $q$ positions in $T$ such that, for any $1 \le i < j \le q$, it holds $T[A[i] \mathinner{\ldotp\ldotp} n] \prec T[A[j] \mathinner{\ldotp\ldotp} n]$, we can in $\mathcal{O}(q + n/\log_\sigma n)$ time construct a data structure that, given the packed representation of any $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$, returns in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time a pair of integers $(b_{\mathrm{pre}}, e_{\mathrm{pre}})$ satisfying:*

- $b_{\mathrm{pre}} = |\{i \in [1 \mathinner{\ldotp\ldotp} q] : T[A[i] \mathinner{\ldotp\ldotp} n] \prec P\}|$, *and*
- $(b_{\mathrm{pre}} \mathinner{\ldotp\ldotp} e_{\mathrm{pre}}] = \{i \in [1 \mathinner{\ldotp\ldotp} q] : P \text{ is a prefix of } T[A[i] \mathinner{\ldotp\ldotp} n]\}$.

*Proof.* The basic idea is to construct the compact trie of strings in $\{T[A[i] \mathinner{\ldotp\ldotp} n]\}_{i \in [1 \mathinner{\ldotp\ldotp} q]}$ converted into strings over the alphabet of metasymbols (of $\Theta(\log_\sigma n)$ original symbols each). Our mapping of symbols to metasymbols does not, however, simply group symbols into blocks. We instead introduce a special mapping that will allow us to deduce the output range $(b_{\mathrm{pre}}, e_{\mathrm{pre}})$ using two predecessor queries in the image of the set $\{T[A[i] \mathinner{\ldotp\ldotp} n]\}_{i \in [1 \mathinner{\ldotp\ldotp} q]}$ and some carefully crafted patterns. This will allow us to use the augmentation of tries of Fischer and Gawrychowski [30, Theorem 1] in a black-box manner.

*Definitions* First, we introduce the mapping of strings over $[0 \mathinner{\ldotp\ldotp} \sigma)$ into strings over metasymbols. Let $\tau = \lfloor \frac{1}{7} \log_\sigma n \rfloor$ and $\kappa = 3\tau - 1$. For any $X \in [0 \mathinner{\ldotp\ldotp} \sigma)^{\le 3\tau - 1}$, let int$(X)$ denote an integer constructed by appending $6\tau - 2|X|$ zeros and $|X|$ cs (where $c = \sigma - 1$) to $X$, and then interpreting the resulting string as a base-$\sigma$ representation of a number in $[0 \mathinner{\ldotp\ldotp} \sigma^{6\tau})$. Note that $X \ne X'$ implies int$(X) \ne$ int$(X')$. Let also int$'(X)$ denote an integer constructed by appending $6\tau - 2|X|$ cs (where $c = \sigma - 1$) and $|X|$ zeros to $X$, and then interpreting the resulting string as a base-$\sigma$ representation of a number in $[0 \mathinner{\ldotp\ldotp} \sigma^{6\tau})$. For any string $S \in [0 \mathinner{\ldotp\ldotp} \sigma)^*$ of length $\ell \ge 0$, we define mstr$(S)$ as a string of length $\ell' = \lceil \frac{\ell+1}{\kappa} \rceil > 0$ over alphabet $[0 \mathinner{\ldotp\ldotp} \sigma^{6\tau})$ such that, for any $i \in [1 \mathinner{\ldotp\ldotp} \ell']$, it holds mstr$(S)[i] =$ int$(S((i-1) \cdot \kappa \mathinner{\ldotp\ldotp} \min(\ell, i \cdot \kappa)])$. For $S \in [0 \mathinner{\ldotp\ldotp} \sigma)^*$ of length $\ell \ge 0$, we define mstr$'(S)$ as a string of length $\ell' = \lceil \frac{\ell+1}{\kappa} \rceil > 0$ over alphabet $[0 \mathinner{\ldotp\ldotp} \sigma^{6\tau})$ such that mstr$'(S)[1 \mathinner{\ldotp\ldotp} \ell'] =$ mstr$(S)[1 \mathinner{\ldotp\ldotp} \ell']$ and mstr$'(S)[\ell'] =$ int$'(S((\ell'-1) \cdot \kappa \mathinner{\ldotp\ldotp} \min(\ell, \ell' \cdot \kappa)])$. Note that if $\ell$ is a multiple of $\kappa$, then the last symbol of mstr$(S)$ is int$(\varepsilon) = 0$, whereas the last symbol of mstr$'(S)$ is int$'(\varepsilon) = \sigma^{6\tau} - 1$. Observe that

- For any set of strings $\mathcal{S} \subseteq [0 \mathinner{\ldotp\ldotp} \sigma)^*$, the set $\{\text{mstr}(X) : X \in \mathcal{S}\}$ is prefix-free.
- For any strings $X, Y \in [0 \mathinner{\ldotp\ldotp} \sigma)^*$, $X \prec Y$ holds if and only if mstr$(X) \prec$ mstr$(Y)$.

14

- A string $P \in [0 \mathinner{.\,.} \sigma)^*$ is a prefix of $X \in [0 \mathinner{.\,.} \sigma)^*$ if and only if $\mathrm{mstr}(P) \preceq \mathrm{mstr}(X) \prec \mathrm{mstr}'(P)$.

For any set of strings $\mathcal{S}$ and any string $Y$, denote $\mathrm{rank}_{\mathcal{S}}(Y) := |\{X \in \mathcal{S} : X \prec Y\}|$. Observe that by the above properties, letting $P_1 = \mathrm{mstr}(P)$, $P_2 = \mathrm{mstr}'(P)$, and $\mathcal{A} = \{\mathrm{mstr}(T[A[i] \mathinner{.\,.} n])\}_{i \in [1 \mathinner{.\,.} q]}$, we have $(b_{\mathrm{pre}}, e_{\mathrm{pre}}) = (\mathrm{rank}_{\mathcal{A}}(P_1), \mathrm{rank}_{\mathcal{A}}(P_2))$. Let $\mathcal{T}$ denote the compact trie of the set $\mathcal{A}$.

*Components* The data structure consists of two components:

1. The packed representation of $T$ using $\mathcal{O}(n/\log_{\sigma} n)$ space.
2. The trie $\mathcal{T}$ augmented using [30, Theorem 1]. Note that this result requires that the alphabet of strings in $\mathcal{A}$ is of size $|\mathcal{A}|^{\mathcal{O}(1)}$, which may be violated for $q = n^{o(1)}$. Thus, we actually define the alphabet to be $[0 \mathinner{.\,.} \sigma')$, where $\sigma' = \sigma^{6\tau} + \lceil\sqrt{n}\rceil$, and insert to $\mathcal{A}$ additional $\lceil\sqrt{n}\rceil$ dummy length-1 strings corresponding to the $\lceil\sqrt{n}\rceil$ largest characters in $[0 \mathinner{.\,.} \sigma')$. As a result, we must have $\sigma' = \mathcal{O}(n) = \mathcal{O}(|\mathcal{A}|^2)$. At the same time, the dummy strings do not change $\mathrm{rank}_{\mathcal{A}}(Q)$ for any $Q \in [0 \mathinner{.\,.} \sigma^{6\tau})^*$. By [30, Theorem 1], such augmented $\mathcal{T}$ needs $\mathcal{O}(q + \sqrt{n})$ space.

In total, the data structure takes $\mathcal{O}(q + n/\log_{\sigma} n)$ space.

*Implementation of queries* Using $T$ and $\mathcal{T}$, given the packed representation of $P \in [0 \mathinner{.\,.} \sigma)^m$, we compute the output pair $(b_{\mathrm{pre}}, e_{\mathrm{pre}})$ as follows. First, note that, given the packed representation of $P$ and $T$, we can in $\mathcal{O}(1)$ time access any symbol of $\mathrm{mstr}(P)$, $\mathrm{mstr}'(P)$, and $\mathrm{mstr}(T[i \mathinner{.\,.} n])$ for any $i \in [1 \mathinner{.\,.} n]$. We start by computing $P_1$ and $P_2$ from $P$ in $\mathcal{O}(m/\log_{\sigma} n + 1)$ time. Then, using [30, Theorem 1], we compute $\mathrm{rank}_{\mathcal{A}}(P_1)$ and $\mathrm{rank}_{\mathcal{A}}(P_2)$ in $\mathcal{O}(|P_1| + \log\log\sigma') = \mathcal{O}(m/\log_{\sigma} n + \log\log n)$ and $\mathcal{O}(|P_2| + \log\log\sigma') = \mathcal{O}(m/\log_{\sigma} n + \log\log n)$ time, respectively. By the above discussion, this gives us the output pair $(b_{\mathrm{pre}}, e_{\mathrm{pre}})$. During the query, the algorithm may want to access symbols of strings from $\mathcal{A}$. We do not store them explicitly (note that storing $\mathrm{mstr}(T)$ would not be enough), but instead perform the mapping on-the-fly.

*Construction algorithm* We start by building a data structure that supports LCE queries for suffixes of $T$. Using [50, Theorem 5.4], the construction takes $\mathcal{O}(n/\log_{\sigma} n)$ time, and the resulting data structure answers queries in $\mathcal{O}(1)$ time. Denote the length of the longest common prefix between suffixes $T[i \mathinner{.\,.} n]$ and $T[j \mathinner{.\,.} n]$ as $\mathrm{LCE}(i, j)$. Observe that for any $i, j \in [1 \mathinner{.\,.} n]$ such that $i \neq j$, the longest common prefix of $\mathrm{mstr}(T[i \mathinner{.\,.} n])$ and $\mathrm{mstr}(T[j \mathinner{.\,.} n])$ has length $\lfloor \frac{\mathrm{LCE}(i,j)}{\kappa} \rfloor$, which can be computed in $\mathcal{O}(1)$ time. We construct $\mathcal{T}$ by inserting elements of $\{\mathrm{mstr}(T[A[i] \mathinner{.\,.} n])\}_{i \in [1 \mathinner{.\,.} q]}$ in the order given by $A$. The construction proceeds as in the proof of Proposition 4.3 and takes $\mathcal{O}(q)$ time. Once the trie is constructed, we add the $\lceil\sqrt{n}\rceil$ dummy length-1 strings and augment $\mathcal{T}$ using [30, Theorem 1] in $\mathcal{O}(q + \sqrt{n})$ time. In total, the construction takes $\mathcal{O}(q + n/\log_{\sigma} n)$ time. $\qquad\square$

## 4.3 (Prefix) Rank and Selection Queries

We start with an implementation of rank and selection queries for larger alphabets.

**Lemma 4.5** (Belazzougui and Puglisi [9])**.** *For all integers $N \geq n \geq \sigma \geq 2$ and every string $S \in [0 \mathinner{.\,.} \sigma)^{\leq n}$, there exists a data structure of $\mathcal{O}(n\log\sigma)$ bits that answers rank queries in $\mathcal{O}(\log\log N)$ time and selection queries in $\mathcal{O}(1)$ time. Moreover, given a table precomputed in $\mathcal{O}(N)$ time (shareable across all instances with common parameter $N$) and the packed representation of $S$, the data structure can be constructed in $\mathcal{O}(\min(n, \sigma + n\log\sigma/\sqrt{\log N}))$ time using $\mathcal{O}(n\log\sigma)$ bits of space.*

*Proof.* If $\log^2\sigma \geq \log N$, we use the data structure of [9, Lemma C.2], which occupies $\mathcal{O}(n\log\sigma)$ bits, answers rank queries in $\mathcal{O}(\log\log n)$ time and selection queries in $\mathcal{O}(1)$ time, and can be

constructed in $\mathcal{O}(n)$ time using $\mathcal{O}(n \log \sigma)$ bits of space.[7] Otherwise, we use the data structure of [9, Lemma C.3], which occupies $\mathcal{O}(n \log \sigma)$ bits, answers rank queries in $\mathcal{O}(\log \log N)$ time and selection queries in $\mathcal{O}(1)$ time, and can be constructed in $\mathcal{O}(\sigma + n \log^2 \sigma / \log N)$ time using $\mathcal{O}(n \log \sigma)$ bits of space. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

The following proposition, instantiated with $h = \lceil \ell^{\epsilon/2} \rceil$, immediately yields Theorem 2.2.

**Proposition 4.6.** *For all integers $h, m, \ell, \sigma \in \mathbb{Z}_{\geq 1}$ satisfying $h \geq 2$ and $m \geq \sigma^\ell \geq 2$, and every string $W \in ([0 .. \sigma)^\ell)^{\leq m}$, there exists a data structure of size $\mathcal{O}(m \log_h(h\ell))$ that answers prefix rank queries in $\mathcal{O}(h \log \log m \log_h(h\ell))$ time and prefix selection queries in $\mathcal{O}(h \log_h(h\ell))$ time. Moreover, it can be constructed in $\mathcal{O}(m \min(\ell, \sqrt{\log m}) \log_h(h\ell))$ time using $\mathcal{O}(m \log_h(h\ell))$ space given the packed representation of $W$ and the parameter $h$.*

*Proof.* The data structure consists in the wavelet tree of $W$ and, when $h \leq \ell$, an instance constructed recursively for an auxiliary string $\tilde{W}$ defined below.

*Wavelet tree* Let $\Sigma = [0 .. \sigma)$ so that the alphabet of $W$ is $\Sigma^\ell$. The wavelet tree of $W$ [41] is the trie of $\Sigma^\ell$ with each internal node $v_X$ (representing a string $X \in \Sigma^{\leq \ell - 1}$) associated to a string $B_X[1 .. \mathsf{rank}_{W,X}(|W|)] \in \Sigma^*$ such that $B_X[r] = W[\mathsf{select}_{W,X}(r)][|X| + 1]$ for $r \in [1 .. \mathsf{rank}_{W,X}(|W|)]$. The strings $B_X$ are augmented with the component of Lemma 4.5 (for rank and selection queries) with parameter $N := m$.

*Recursive instance* We shall define $\tilde{W}$ as a string of length $|W|$ over the alphabet $\tilde{\Sigma}^{\tilde\ell}$, where $\tilde\ell := \lfloor \ell/h \rfloor$, $\tilde\sigma = \sigma^h$, and $\tilde\Sigma := [0 .. \tilde\sigma)$. We identify $\tilde\Sigma$ with $\Sigma^h$, treating each string in $\Sigma^h$ as the $h$-digit base-$\sigma$ representation of an integer in $\tilde\Sigma$. For every string $X \in \Sigma^*$, define $\tilde{X} \in \tilde\Sigma^*$ so that $|\tilde{X}| = \lfloor |X|/h \rfloor$ and $\tilde{X}[i] = X(h(i-1) .. hi]$ for $i \in [1 .. |\tilde{X}|]$. Moreover, we set $\tilde{W}[1 .. |W|]$ so that $\tilde{W}[j] = \widetilde{W[j]}$ for $j \in [1 .. |W|]$. Note that the recursive application of Proposition 4.6 to $\tilde{W}$ is possible because $2 \leq \tilde\sigma^{\tilde\ell} \leq \sigma^\ell \leq m$ and $\tilde\ell \geq 1$ hold when $h \leq \ell$.

*Data structure size* It is easy to see that, for a fixed length $d \in [0 .. \ell)$, the strings $B_X$ for $X \in \Sigma^d$ are of total length $m$. Across all $X \in \Sigma^{\leq \ell - 1}$, this sums up to $m\ell$, so the raw strings $B_X$ occupy $\mathcal{O}(m\ell \log \sigma) = \mathcal{O}(m \log m)$ bits. The augmentation of $B_X$ using Lemma 4.5 adds $\mathcal{O}((\sigma + |B_X|) \log \sigma)$ extra bits (we set $n := \max(\sigma, |B_X|) = \Theta(\sigma + |B_X|)$ to ensure $\sigma \leq n$), which sums up to $\mathcal{O}((\sigma^\ell + m\ell) \log \sigma) = \mathcal{O}(m \log m)$ bits, i.e., $\mathcal{O}(m)$ machine words. The recursion depth is $\mathcal{O}(\log_h(h\ell))$, so the overall size is $\mathcal{O}(m \log_h(h\ell))$.

*Answering queries* To handle any query concerning $X \in \Sigma^{\leq \ell}$, we compute auxiliary strings $\tilde{X}$ (as defined above) and $X' := X[1 .. |X| - (|X| \bmod h)]$ (obtained by expanding the letters in $\tilde{X}$ into length-$h$ strings).

Answering a prefix rank query $\mathsf{rank}_{W,X}(j)$, we traverse the path from $v_{X'}$ to $v_X$, maintaining a value $r$ such that $r = \mathsf{rank}_{W,Y}(j)$ holds while the algorithm visits $v_Y$. We initialize $r := j = \mathsf{rank}_{W,\varepsilon}(j)$ if $X' = \varepsilon$ and $r := \mathsf{rank}_{\tilde{W},\tilde{X}}(j)$ (computed recursively) otherwise; this is valid due to $\mathsf{rank}_{\tilde{W},\tilde{X}}(j) = \mathsf{rank}_{W,X'}(j)$. Upon entering a node $v_{Ya}$ from its parent $v_Y$, we set $r := \mathsf{rank}_{B_Y,a}(r)$ since $\mathsf{rank}_{W,Ya}(j) = \mathsf{rank}_{B_Y,a}(\mathsf{rank}_{W,Y}(j))$; see [41]. When reaching $v_X$, we return $r = \mathsf{rank}_{W,X}(j)$. The running time is $\mathcal{O}(h \log \log m)$ per recursive level, for a total of $\mathcal{O}(h \log \log m \cdot \log_h(h\ell))$.

---

[7]The statement of [9, Lemma C.2] does not bound the space consumption of the construction algorithm. Nevertheless, it is straightforward to implement the underlying construction procedure in $\mathcal{O}(n \log \sigma)$ bits of working space. The original algorithm scans the input sequence $S$ from left to right and, for each $a \in \Sigma$, builds an array $P_a[1 .. n_a]$ such that $n_a = \mathsf{rank}_{S,a}(|S|)$ and $P_a[r] = \mathsf{select}_{S,a}(r)$ for $r \in [1 .. n_a]$. The array $P_a[1 .. n_a]$ is then converted to the Elias–Fano representation: an array $A_a[1 .. n_a]$ with $A_a[r] = P_a[r] \bmod \sigma$ for $r \in [1 .. n_a]$ and a bit vector $V_a = \mathsf{unary}((\lfloor P_a[r]/\sigma \rfloor - \lfloor P_a[r-1]/\sigma \rfloor)_{r \in [1 .. n_a]})$, where we assume $P_a[0] = 0$ to streamline the formula. To achieve $\mathcal{O}(n \log \sigma)$ bits of working space, instead of storing $P_a$ explicitly, we convert $P_a$ to the Elias–Fano representation on the fly as subsequent positions are appended to $P_a$.

Answering a prefix selection query $\mathsf{select}_{W,X}(r)$, we traverse the path from $v_X$ to $v_{X'}$, maintaining a value $q$ such that $\mathsf{select}_{W,Y}(q) = \mathsf{select}_{W,X}(r)$ holds while the algorithm visits $v_Y$. We initialize $q := r$ and, upon entering a node $v_Y$ from its child $v_{Ya}$, we set $q := \mathsf{select}_{B_Y,a}(q)$ since $\mathsf{select}_{W,Ya}(r) = \mathsf{select}_{W,Y}(\mathsf{select}_{B_Y,a}(r))$; see [41]. When reaching $v_{X'}$, we return $q = \mathsf{select}_{W,\varepsilon}(q)$ if $X' = \varepsilon$ and $\mathsf{select}_{\widetilde{W},\widetilde{X}}(q)$ otherwise; this is valid due to $\mathsf{select}_{\widetilde{W},\widetilde{X}}(q) = \mathsf{select}_{W,X'}(q)$. The running time is $\mathcal{O}(h)$ per recursive level, for a total of $\mathcal{O}(h \cdot \log_h(h\ell))$.

*Construction algorithm* If $\ell \leq \sqrt{\log m}$, we use the original wavelet tree construction algorithm [41], which takes $\mathcal{O}(m\ell)$ time and $\mathcal{O}(m)$ space. Building the data structure of Lemma 4.5 for $B_X$ takes $\mathcal{O}(\sigma + |B_X|)$ time and $\mathcal{O}((\sigma + |B_X|)\log\sigma/\log m)$ space, which sums up to $\mathcal{O}(\sigma^\ell + m\ell) = \mathcal{O}(m\ell)$ time and $\mathcal{O}(m\ell\log\sigma/\log m) = \mathcal{O}(m)$ space across $X \in \Sigma^{\leq\ell-1}$ (due to $\ell\log\sigma \leq \log m$). Precomputing the table shared by all instances of Lemma 4.5 takes $\mathcal{O}(m)$ time and space. Considering all levels of recursion, we get $\mathcal{O}(m\ell)$ time (due to $\tilde{\ell} \leq \frac{1}{2}\ell$) and $\mathcal{O}(m\log_h(h\ell))$ space.

If $\ell > \sqrt{\log m}$, on the other hand, we apply the bit-parallel wavelet tree construction algorithm of [64, 3], which has been adapted to large alphabets in [50, Lemma 6.4]. Due to $\ell\log\sigma \leq \log m$, this procedure takes $\mathcal{O}(m\ell\log\sigma/\sqrt{\log m}+m\ell\log^2\sigma/\log m) = \mathcal{O}(m\sqrt{\log m})$ time and $\mathcal{O}(m)$ space. Building the data structure of Lemma 4.5 for $B_X$ takes $\mathcal{O}(\sigma + (\sigma + |B_X|)\log\sigma/\sqrt{\log m}) = \mathcal{O}(\sigma + |B_X|\log\sigma/\sqrt{\log m})$ time and $\mathcal{O}((\sigma + |B_X|)\log\sigma/\log m)$ space, which sums up to $\mathcal{O}(m\sqrt{\log m})$ time and $\mathcal{O}(m)$ space across $X \in \Sigma^{\leq\ell-1}$. Precomputing the table shared by all instances of Lemma 4.5 takes $\mathcal{O}(m)$ time and space. Considering all levels of recursion, we get a multiplicative overhead of $\mathcal{O}(\log_h(h\ell))$, for a total of $\mathcal{O}(m\sqrt{\log m}\log_h(h\ell))$ time and $\mathcal{O}(m\log_h(h\ell))$ space. $\square$

## 4.4 Range Counting and Selection

**Proposition 2.3.** *An array $A[1..m']$ of $m' \in [2..m]$ nonnegative integers satisfying $\sum_{i=1}^{m'} A[i] = \mathcal{O}(m\log m)$ can be preprocessed in $\mathcal{O}(m)$ time so that range counting and selection queries can be answered in $\mathcal{O}(\log\log m)$ time and $\mathcal{O}(1)$ time, respectively.*

*Proof.* We use the following definitions. Denote $h = \lfloor\log m\rfloor$. For any $k \geq 0$, by $P_k[1..m_k]$, where $m_k = \mathsf{rcount}_A(kh, m)$, we denote the array defined by $P_k[i] = \mathsf{rselect}_A(kh, i)$. Let $v \geq 0$. We define a bitvector $M_v[1..m_k]$, where $k = \lfloor\frac{v}{h}\rfloor$ as follows. For any $i \in [1..m_k]$, $M_v[i] = 1$ holds if and only if $A[P_k[i]] \geq v$. For any $k \geq 0$, we define the concatenation $M_k' = M_{kh}M_{kh+1}\cdots M_{(k+1)h-1}$. Let $k_{\max} = \max\{k \geq 0 : m_k > 0\}$. Since all elements of $A$ are nonnegative, and $\sum_{i=1}^{m'} A[i] \in \mathcal{O}(m\log m)$, we obtain $\max_{i\in[1..m']} A[i] \in \mathcal{O}(m\log m)$, and consequently, $k_{\max} = \lfloor\frac{1}{h}\max_{i\in[1..m']} A[i]\rfloor \in \mathcal{O}(m)$.

*Components* The data structure consists of two components:

1. First, for $k \in [0..k_{\max}]$, we store a plain representation of the sequence $P_k[1..m_k]$ using $\mathcal{O}(m_k)$ space. Each array is augmented with a static predecessor data structure. We use [30, Proposition 2], and hence achieve linear space and $\mathcal{O}(\log\log m)$ query time. Each $i \in [1..m']$ occurs in $\lceil\frac{A[i]+1}{h}\rceil$ arrays. Thus, $\sum_{k\geq 0} m_k = \sum_{i=1}^{m'}\lceil\frac{A[i]+1}{h}\rceil \leq 2m' + \sum_{i=1}^{m'}\lfloor\frac{A[i]}{h}\rfloor \leq 2m' + \frac{1}{h}\sum_{i=1}^{m'} A[i] \in \mathcal{O}(m)$ and hence we can store the arrays $P_k$ (including the associated predecessor data structures) using $\mathcal{O}((k_{\max} + 1) + \sum_{k\geq 0} m_k) \subseteq \mathcal{O}(m)$ space, so that we can access each array in $\mathcal{O}(1)$ time.

2. Second, for every $k \in [0..k_{\max}]$, we store the plain representation of bitvector $M_k'$, augmented using Theorem 2.1. By $|M_k'| = h \cdot m_k$, the total length of bitvectors $M_k'$ is $\sum_{k\geq 0} |M_k'| = h\sum_{k\geq 0} m_k \in \mathcal{O}(m\log m)$. All bitvectors $M_k'$ can thus be stored in $\mathcal{O}((k_{\max} + 1) + \frac{1}{\log m}\sum_{k\geq 0} |M_k'|) \subseteq \mathcal{O}(m)$ words of space, so that we can access each in $\mathcal{O}(1)$ time. For a bitvector of length $t$, the augmentation of Theorem 2.1 adds only $\mathcal{O}(\log m + t)$ bits of space, and hence does not increase the space usage.

In total, the data structure takes $\mathcal{O}(m)$ space.

*Implementation of queries*    Using the above two components, we answer range counting/selection queries on $A$ as follows. To compute $\mathsf{rcount}_A(v,j)$, we first let $k = \lfloor \frac{v}{h} \rfloor$. If $k > k_{\max}$, then we return $\mathsf{rcount}_A(v,j) = 0$. Otherwise, we observe that if $j' = |\{i \in [1 \mathinner{\ldotp\ldotp} m_k] : P_k[i] \leq j\}|$, then $\mathsf{rcount}_A(v,j) = \mathsf{rank}_{M_v,1}(j')$. Computing $j'$ using the predecessor data structure takes $\mathcal{O}(\log \log m)$ time, and then $\mathsf{rank}_{M_v,1}(j')$ is computed using the rank support data structure of the bitvector $M_k'$ as $\mathsf{rank}_{M_k',1}(j' + (v - kh)m_k) - \mathsf{rank}_{M_k',1}((v - kh)m_k)$ in $\mathcal{O}(1)$ time. To compute $\mathsf{rselect}_A(v,r)$, we observe that letting again $k = \lfloor \frac{v}{h} \rfloor$, it holds $\mathsf{rselect}_A(v,r) = P_k[\mathsf{select}_{M_v,1}(r)]$. The value $\mathsf{select}_{M_v,1}(r)$ is computed using the select support data structure of the bitvector $M_k'$ as $\mathsf{select}_{M_k',1}(\mathsf{rank}_{M_k',1}((v - kh)m_k) + r) - (v - kh)m_k$ in $\mathcal{O}(1)$ time.

*Construction algorithm*    We start by initializing $P_0[i] = i$ for $i \in [1 \mathinner{\ldotp\ldotp} m']$. For $k \in [1 \mathinner{\ldotp\ldotp} k_{\max}]$, the array $P_k$ is computed by iterating over $P_{k-1}$ and including only elements $P_{k-1}[i]$ satisfying $A[P_{k-1}[i]] \geq kh$. By $\sum_{k \geq 0} m_k \in \mathcal{O}(m)$, this takes $\mathcal{O}(m)$ time in total. We then augment all arrays $P_k$ with the predecessor data structures. Since the arrays are sorted, using [30, Proposition 2], the construction altogether again takes $\mathcal{O}(m)$ time. We then construct bitvectors $M_k'$ in the order of increasing $k \in [0 \mathinner{\ldotp\ldotp} k_{\max}]$. To build $M_k'$ we first scan $P_k$ and check if there exists $i \in [1 \mathinner{\ldotp\ldotp} m_k]$ such that $A[P_k[i]] < (k+1)h$.

1. If there is no such $i$, we set $M_k' := 1^{hm_k}$ in $\mathcal{O}(1 + \frac{1}{\log m} hm_k) = \mathcal{O}(m_k)$ time.
2. Otherwise, we scan again $P_k[1 \mathinner{\ldotp\ldotp} m_k]$ and prepare $h$ lists $L_0, L_1, \ldots, L_{h-1}$ such that $L_y$ contains all $i \in [1 \mathinner{\ldotp\ldotp} m_k]$ satisfying $A[P_k[i]] = kh + y$. Construction of all lists takes $\mathcal{O}(m_k + h)$ time. The bitvector $M_k'$ is then obtained as the concatenation of bitvectors $M_{kh}, M_{kh+1}, \ldots, M_{(k+1)h-1}$ computed in this order. We first initialize $M_{kh} := 1^{m_k}$ in $\mathcal{O}(1 + \frac{m_k}{\log m})$ time. The bitvector $M_{kh+y}$ for $y > 0$ is obtained by first copying the bitvector $M_{kh+y-1}$ in $\mathcal{O}(1 + \frac{m_k}{\log m})$ time, and then setting $M_{kh+y}[i] = 0$ for every position $i$ stored in $L_{y-1}$. The total length of all lists $L_y$ is bounded by $m_k$. Thus, the construction of $M_k'$ takes $\mathcal{O}(h + m_k + \frac{1}{\log m} hm_k) \subseteq \mathcal{O}(h + m_k)$ time.

To bound the total time spent constructing bitvectors $M_k'$, we consider two cases:

- $k \leq \frac{m}{h}$: The total time spent in the construction of bitvectors $M_k'$ for such $k$ is bounded by the sum $\sum_{k=0}^{\lfloor m/h \rfloor} \mathcal{O}(h + m_k) \subseteq \mathcal{O}(m + \sum_{k \geq 0} m_k) \subseteq \mathcal{O}(m)$.
- $k > \frac{m}{h}$: Let $k' = \lfloor \frac{m}{h} \rfloor + 1$. Note that for any $t$, it holds $m_{t+1} \leq m_t$. Moreover, whenever Case 2 above happens for some $t$, it holds $m_{t+1} < m_t$. Thus, Case 2 above can happen for $k > \frac{m}{h}$ only $m_{k'}$ times. Since for every $i \in [1 \mathinner{\ldotp\ldotp} m_{k'}]$ we have $A[P_{k'}[i]] \geq m$, by $\sum_{i \in [1 \mathinner{\ldotp\ldotp} m']} A[i] \in \mathcal{O}(m \log m)$ it holds $m_{k'} \in \mathcal{O}(\log m)$. The total time spend computing $M_k'$ for $k > \frac{m}{h}$ is thus bounded by $\mathcal{O}(m_{k'}(h + m_{k'}) + \sum_{k \geq k'} m_k) \subseteq \mathcal{O}(\log^2 m + \sum_{k \geq 0} m_k) \subseteq \mathcal{O}(m)$.

The total length of bitvectors $M_k'$ for $k \in [0 \mathinner{\ldotp\ldotp} k_{\max}]$, is $\sum_{k \in [0 \mathinner{\ldotp\ldotp} k_{\max}]} hm_k \in \mathcal{O}(hm)$. Thus, augmenting them all using Theorem 2.1 takes $\mathcal{O}((k_{\max} + 1) + \frac{1}{\log m} hm) \subseteq \mathcal{O}(m)$ time.    $\square$

## 5   SA and ISA Queries

Let $\epsilon \in (0,1)$ be any fixed constant and let $T \in [0 \mathinner{\ldotp\ldotp} \sigma)^n$, where $2 \leq \sigma < n^{1/7}$. In this section, we show how, given the packed representation of $T$, in $\mathcal{O}(n \min(1, \log \sigma / \sqrt{\log n}))$ time and using $\mathcal{O}(n / \log_\sigma n)$ working space, to construct a data structure of size $\mathcal{O}(n / \log_\sigma n)$ that answers SA and ISA queries in $\mathcal{O}(\log^\epsilon n)$ time. We also derive a general reduction depending on prefix rank and selection queries.

Let $\tau = \lfloor \mu \log_\sigma n \rfloor$, where $\mu$ is any positive constant smaller than $\frac{1}{6}$ such that $\tau \geq 1$ (such $\mu$

exists by $\sigma < n^{1/7}$), be fixed for the duration of this section. Throughout, we also use R as a shorthand for $\mathsf{R}(\tau, T)$.

**Definition 5.1.** Let $j \in [1 \mathinner{.\,.} n]$. We call position $j$ *periodic* if $j \in \mathsf{R}$. Otherwise, $j$ is *nonperiodic*.

**Organization** The structure and the query algorithm to compute $\mathrm{SA}[i]$ (resp. $\mathrm{ISA}[j]$), given any $i \in [1 \mathinner{.\,.} n]$ (resp. $j \in [1 \mathinner{.\,.} n]$), are different depending on whether $\mathrm{SA}[i]$ (resp. $j$) is periodic (Definition 5.1). Our description is thus split as follows. First (Section 5.1), we describe the set of data structures called collectively the index "core" that enables efficiently checking if $\mathrm{SA}[i] \in \mathsf{R}$ (resp. $j \in \mathsf{R}$); the core also contain some common components utilized by the remaining parts. In the following two parts (Sections 5.2 and 5.3), we describe structures handling each of the two cases. All ingredients are then put together in Section 5.4. Finally, we present our result in the general form (Section 5.5).

## 5.1 The Index Core

In this section, we present a data structure that, given any $j \in [1 \mathinner{.\,.} n]$ (resp. $i \in [1 \mathinner{.\,.} n]$), lets us in $\mathcal{O}(1)$ time determine if $j \in \mathsf{R}$ (resp. $\mathrm{SA}[i] \in \mathsf{R}$).

The section is organized as follows. First, we introduce the components of the data structure (Section 5.1.1). Next, we describe the query algorithms (Section 5.1.2). Finally, we show the construction algorithm (Section 5.1.3).

### 5.1.1 The Data Structure

**Definitions** Let $L_{\mathrm{range}}$ be a mapping from $X \in [0 \mathinner{.\,.} \sigma)^{\leq 3\tau-1} := \{\varepsilon\} \cup [0 \mathinner{.\,.} \sigma) \cup \ldots \cup [0 \mathinner{.\,.} \sigma)^{3\tau-1}$ to the pair of integers $(b, e) := (\mathrm{RangeBeg}(X, T), \mathrm{RangeEnd}(X, T))$. Let also $L_{\mathrm{per}}$ denote the mapping from $[0 \mathinner{.\,.} \sigma)^{3\tau-1}$ to $\mathbb{Z}_+$ such that every $X$ is mapped to $\mathrm{per}(X)$.

Let $B_{3\tau-1}[1 \mathinner{.\,.} n]$ be a bitvector defined such that $B_{3\tau-1}[i] = 1$ holds if and only if $i = n$, or $i < n$ and $X_{\mathrm{SA}[i]} \neq X_{\mathrm{SA}[i+1]}$, where $X_j = T[j \mathinner{.\,.} \min(n+1, j+3\tau-1))$ for every $j \in [1 \mathinner{.\,.} n]$.

Let $A_{\mathrm{short}}[1 \mathinner{.\,.} t]$ ($t = \mathsf{rank}_{B_{3\tau-1}, 1}(n)$) be defined by $A_{\mathrm{short}}[i] = X_{\mathrm{SA}[j]}$, where $j = \mathsf{select}_{B_{3\tau-1}, 1}(i)$.

**Components** The index core, denoted $\mathrm{C}_{\mathrm{SA}}(T)$, consists of five components:

1. The packed representation of $T$ using $\mathcal{O}(n/\log_\sigma n)$ space.
2. The lookup table $L_{\mathrm{range}}$. When accessing $L_{\mathrm{range}}$, strings $X \in [0 \mathinner{.\,.} \sigma)^{\leq 3\tau-1}$ are converted to small integers using the mapping $\mathrm{int}(X)$ defined in the proof of Proposition 4.4. By $\mathrm{int}(X) \in [0 \mathinner{.\,.} \sigma^{6\tau})$, $L_{\mathrm{range}}$ needs $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n^{6\mu}) = \mathcal{O}(n/\log_\sigma n)$ space.
3. The lookup table $L_{\mathrm{per}}$. Similarly as above, we utilize the mapping $\mathrm{int}(X)$. $L_{\mathrm{per}}$ thus also needs $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n/\log_\sigma n)$ space.
4. The bitvector $B_{3\tau-1}$ augmented using Theorem 2.1 for rank and selection queries. The augmented bitvector takes $\mathcal{O}(n/\log n)$ space.
5. The array $A_{\mathrm{short}}$. Every string $X \in \{A_{\mathrm{short}}[i]\}_{i \in [1 \mathinner{.\,.} t]}$ is encoded as $\mathrm{int}(X)$ using $6\tau \log \sigma = \mathcal{O}(\log n)$ bits. This implicitly encodes the length of the string and ensures that all strings are encoded using the same number of bits. By $\{A_{\mathrm{short}}[i]\}_{i \in [1 \mathinner{.\,.} t]} \subseteq [0 \mathinner{.\,.} \sigma)^{\leq 3\tau-1}$, we have $t = \mathcal{O}(n^{1/2})$, and hence the array $A_{\mathrm{short}}$ needs $\mathcal{O}(n^{1/2}) = \mathcal{O}(n/\log n)$ space.

In total, $\mathrm{C}_{\mathrm{SA}}(T)$ takes $\mathcal{O}(n/\log_\sigma n)$ space.

### 5.1.2 Navigation Primitives

**Proposition 5.2.** *Given* $\mathrm{C}_{\mathrm{SA}}(T)$, *for any* $j \in [1 \mathinner{.\,.} n]$ *we can in* $\mathcal{O}(1)$ *time determine if* $j \in \mathsf{R}$.

*Proof.* If $j > n - 3\tau + 2$, we return that $j \notin \mathsf{R}$ (Definition 2.4). Otherwise, we use the packed encoding of $T$ to extract $X = T[j \mathinner{.\,.} j + 3\tau - 1)$ in $\mathcal{O}(1)$ time and convert it to $x = \mathrm{int}(X)$. We then use the lookup table $L_{\mathrm{per}}$, to determine $p = \mathrm{per}(X)$, and return that $j \in \mathsf{R}$ if $p \leq \frac{1}{3}\tau$. $\qquad\square$

**Proposition 5.3.** *Given* $\mathrm{C}_{\mathrm{SA}}(T)$, *for any* $i \in [1 \mathinner{.\,.} n]$ *we can in* $\mathcal{O}(1)$ *time determine if* $\mathrm{SA}[i] \in \mathsf{R}$.

*Proof.* Given the position $i \in [1 \mathinner{.\,.} n]$, we first compute $y = \mathsf{rank}_{B_{3\tau-1},1}(i-1)$. The string $X = A_{\mathrm{short}}[y+1]$ is then a prefix of $T[\mathrm{SA}[i] \mathinner{.\,.} n]$. If $|X| < 3\tau - 1$, we must have $\mathrm{SA}[i] > n - 3\tau + 2$, and thus we return that $\mathrm{SA}[i] \in [1 \mathinner{.\,.} n] \setminus \mathsf{R}$ (see Definition 2.4). Otherwise (i.e., $|X| = 3\tau - 1$), using $L_{\mathrm{per}}$ we determine $p = \mathrm{per}(X)$ and return that $\mathrm{SA}[i] \in \mathsf{R}$ if $p \leq \frac{1}{3}\tau$. $\qquad\square$

### 5.1.3 Construction Algorithm

**Proposition 5.4.** *Given the packed representation of* $T \in [0 \mathinner{.\,.} \sigma)^n$, *we can construct* $\mathrm{C}_{\mathrm{SA}}(T)$ *in* $\mathcal{O}(n/\log_\sigma n)$ *time.*

*Proof.* To compute $L_{\mathrm{range}}$, we first compute for every $X \in [0 \mathinner{.\,.} \sigma)^\ell$ (where $\ell = 3\tau - 1$), its frequency $f_X := |\mathrm{Occ}(X, T)|$. Using the simple generalization of the algorithm described in [50, Section 6.1.2], this takes $\mathcal{O}(n/\log_\sigma n)$ time (note that the algorithm requires $\ell\sigma^{2\ell-1} = O(n/\log_\sigma n)$, which is satisfied here, since $2\ell - 1 < 6\mu\log_\sigma n$ and $\mu < \frac{1}{6}$). From the frequencies of $X \in [0 \mathinner{.\,.} \sigma)^{3\tau-1}$ we then compute the values of $f_X$ for all $X \in [0 \mathinner{.\,.} \sigma)^{<3\tau-1}$ by observing that unless $X$ is a nonempty suffix of $T$, it holds $f_X = \sum_{c \in [0 \mathinner{.\,.} \sigma)} f_{Xc}$, i.e., the frequency of each string shorter than $3\tau - 1$ is obtained in $\mathcal{O}(\sigma)$ time. If $X$ is a nonempty suffix of $T$ (which we can check in $\mathcal{O}(1)$ time), we additionally add one to the count. Since each string contributes exactly once to the frequency of another string, over all $X \in [0 \mathinner{.\,.} \sigma)^{<3\tau-1}$, this takes $\mathcal{O}(\sigma^{3\tau-1}) = \mathcal{O}(n/\log_\sigma n)$ time. Once $f_X$ is computed for all $X \in [0 \mathinner{.\,.} \sigma)^{\leq 3\tau-1}$, we compute $L_{\mathrm{range}}$ as follows. Denote $\Sigma = [0 \mathinner{.\,.} \sigma)$. Assume that $L_{\mathrm{range}}[\mathrm{int}(X)] = (b, e)$ holds for some $X \in [0 \mathinner{.\,.} \sigma)^{<3\tau-1}$. Then, for any $c \in \Sigma$, it holds $L_{\mathrm{range}}[\mathrm{int}(Xc)] = (e - x - f_{Xc}, e - x)$, where $x = \sum_{c' \in \Sigma, c' > c} f_{Xc'}$, e.g., for $\sigma = 2$, $L_{\mathrm{range}}[\mathrm{int}(X0)] = (e - f_{X1} - f_{X0}, e - f_{X1})$. We thus compute $L_{\mathrm{range}}[\mathrm{int}(X)]$ by initializing $L_{\mathrm{range}}[\mathrm{int}(\varepsilon)] = (0, n)$, and then enumerating all $X \in [0 \mathinner{.\,.} \sigma)^{\leq 3\tau-1}$ in the order of non-decreasing length (and, in case of ties, in the reverse lexicographical order). During the enumeration of strings of the form $Xc$, where $c \in \Sigma$, we maintain the sum $x = \sum_{c' \in \Sigma, c' > c} f_{Xc'}$. Then, using the above formula, the value of $L_{\mathrm{range}}[\mathrm{int}(Xc)]$ can be obtained in $\mathcal{O}(1)$ time. Over all $X$, the computation of $L_{\mathrm{range}}[\mathrm{int}(X)]$ thus takes $\mathcal{O}(\sigma^{3\tau-1}) = \mathcal{O}(n/\log_\sigma n)$ time.

To construct $L_{\mathrm{per}}$, we enumerate all $X \in [0 \mathinner{.\,.} \sigma)^{3\tau-1}$, and for each $X$ in $\mathcal{O}(\tau^2)$ time we compute $\mathrm{per}(X)$ by trying all $\ell \in [1 \mathinner{.\,.} 3\tau - 1]$. Initializing $L_{\mathrm{per}}$ takes $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n/\log_\sigma n)$. Over all $X \in [0 \mathinner{.\,.} \sigma)^{3\tau-1}$, we spend $\mathcal{O}(\sigma^{3\tau-1}\tau^2) = \mathcal{O}(n^{1/2}\log^2 n) = \mathcal{O}(n/\log_\sigma n)$ time.

We finish with the construction of $B_{3\tau-1}$ and $A_{\mathrm{short}}$. First, in $\mathcal{O}(n/\log n)$ time we initialize $B_{3\tau-1}$ to zeros. Next, we initialize temporary counters $k$ and $f$ to zero, and simulate a preorder traversal of the trie of $[0 \mathinner{.\,.} \sigma)^{3\tau-1}$. For each visited node with label $X$, we consider two cases:

- If $|X| < 3\tau - 1$, we check if $X$ is a suffix of $T$. If so, increment $k$ and $f$ by one, and report $X$.
- Otherwise (i.e., if $|X| = 3\tau - 1$), if $f_X > 0$, we increment $k$ by one, $f$ by $f_X$, and report $X$.

Each time some string $X$ is reported, we set $A_{\mathrm{short}}[k] = X$ and $B_{3\tau-1}[f] = 1$. The correctness of this procedure follows by noting that labels of nodes visited during the preorder traversal are lexicographically sorted. The traversal takes $\mathcal{O}(\sigma^{3\tau}) = \mathcal{O}(n^{3\mu}) = \mathcal{O}(n/\log n)$ time. Finally, using Theorem 2.1, in $\mathcal{O}(n/\log n)$ time we augment $B_{3\tau-1}$ with $\mathcal{O}(1)$-time rank and select queries. $\qquad\square$

## 5.2 The Nonperiodic Positions

In this section, we describe a data structure that, given any $j \in [1 \mathinner{.\,.} n]$ (resp. $i \in [1 \mathinner{.\,.} n]$) satisfying $j \in [1 \mathinner{.\,.} n] \setminus \mathsf{R}$ (resp. $\mathrm{SA}[i] \in [1 \mathinner{.\,.} n] \setminus \mathsf{R}$) computes $\mathrm{ISA}[j]$ (resp. $\mathrm{SA}[i]$) in $\mathcal{O}(\log^\epsilon n)$

time.

The section is organized as follows. First, we introduce the components of the data structure (Section 5.2.1). Next, we describe the query algorithms (Sections 5.2.2 and 5.2.3). Finally, we show the construction algorithm (Section 5.2.4).

### 5.2.1   The Data Structure

**Definitions**   We fix some $\tau$-synchronizing set $\mathsf{S}$ of $T$ obtained using Theorem 2.6 (recall, that $\tau = \lfloor \mu \log_\sigma n \rfloor$ is fixed for Section 5). We denote $n' = |\mathsf{S}| = \mathcal{O}(n/\tau)$. Let $(s_t^{\text{text}})_{t \in [1\,..\,n']}$ be the sequence containing the elements of $\mathsf{S}$ in sorted order, i.e., if $i < j$ then $s_i^{\text{text}} < s_j^{\text{text}}$. Let also $(s_t^{\text{lex}})_{t \in [1\,..\,n']}$ be the sequence containing elements of $\mathsf{S}$ sorted according to the lexicographical order of the corresponding suffixes, i.e., if $i < j$ then $T[s_i^{\text{lex}}\,..\,n] \prec T[s_j^{\text{lex}}\,..\,n]$. Let $W[1\,..\,n']$ be a sequence of length-$3\tau$ strings such that $W[i] = \overline{X_i}$, where $X_i = T^\infty[s_i^{\text{lex}} - \tau \,..\, s_i^{\text{lex}} + 2\tau)$.

For any $i \in [1\,..\,n{-}2\tau{+}1]$, we define $\text{succ}_{\mathsf{S}}(i) = \min\{j \in \mathsf{S} \cup \{n{-}2\tau{+}2\} : j \geq i\}$ and denote $\mathcal{D} := \{T[i\,..\,\text{succ}_{\mathsf{S}}(i) + 2\tau) : i \in [1\,..\,n{-}3\tau{+}2] \setminus \mathsf{R}\}$. Let $L_\mathcal{D}$ be a mapping from $[0\,..\,\sigma)^{3\tau-1}$ to $[0\,..\,\sigma)^{\leq 3\tau-1}$ such that for any $Y \in [0\,..\,\sigma)^{3\tau-1}$, if there exists $X \in \mathcal{D}$ that is a prefix of $Y$ (by the consistency of $\mathsf{S}$, there can be at most one such $X$), then $L_\mathcal{D}$ maps $Y$ to $X$. Otherwise (i.e., there is no such $X$), $L_\mathcal{D}$ maps $Y$ to $\varepsilon$. Let $L_{\text{rev}}$ be a mapping that for every string $X \in [0\,..\,\sigma)^{\leq 3\tau-1}$, returns the packed representation of $\overline{X}$.

Let $B_{\mathsf{S}}[1\,..\,n]$ be a bitvector defined so that $B_{\mathsf{S}}[i] = 1$ holds if and only if $i \in \mathsf{S}$.

Let $A_{\text{smap}}[1\,..\,n']$ be an array storing a permutation of $[1\,..\,n']$ such that $A_{\text{smap}}[i] = j$ holds if and only if $s_j^{\text{text}} = s_i^{\text{lex}}$. Let $A_{\text{smap}}^{-1}[1\,..\,n']$ be an array storing a permutation of $[1\,..\,n']$ such that $A_{\text{smap}}^{-1}[j] = i$ holds if and only if $s_j^{\text{text}} = s_i^{\text{lex}}$.

**Components**   The data structure to handle nonperiodic positions consists of seven components:

1. The index core $\mathrm{C}_{\text{SA}}(T)$ (Section 5.1.1). It takes $\mathcal{O}(n/\log_\sigma n)$ space.
2. The lookup table $L_{\text{rev}}$. When accessing $L_{\text{rev}}$, strings $X \in [0\,..\,\sigma)^{\leq 3\tau-1}$ are converted to $\text{int}(X)$. Thus, the mapping $L_{\text{rev}}$ needs $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n^{6\mu}) = \mathcal{O}(n/\log_\sigma n)$ space.
3. The lookup table $L_\mathcal{D}$. As above, $L_\mathcal{D}$ needs $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n/\log_\sigma n)$ space.
4. The bitvector $B_{\mathsf{S}}$ augmented using Theorem 2.1. It needs $\mathcal{O}(n/\log n)$ space.
5. The array $A_{\text{smap}}[1\,..\,n']$ in plain form, using $n' = \mathcal{O}(n/\log_\sigma n)$ words of space.
6. The array $A_{\text{smap}}^{-1}[1\,..\,n']$ in plain form, using $n' = \mathcal{O}(n/\log_\sigma n)$ words of space.
7. The data structure of Theorem 2.2 for the sequence $W[1\,..\,n']$. By $n' = \mathcal{O}(n/\log_\sigma n)$ and $\sigma^{3\tau} = \mathcal{O}(\sqrt{n}) = o(n/\log n)$, it needs $\mathcal{O}(n/\log_\sigma n)$ space.

In total, the data structure takes $\mathcal{O}(n/\log_\sigma n)$ space.

### 5.2.2   Implementation of ISA Queries

For any $j \in [1\,..\,n - 3\tau + 2] \setminus \mathsf{R}$, letting $X \in \mathcal{D}$ be a prefix of $T[j\,..\,n]$ (by [50, Lemma 6.1], $\mathcal{D}$ is prefix-free, and hence there is exactly one such $X$), we define

$$\text{Pos}(j) = \{j' \in [1\,..\,n] : \text{LCE}(j, j') \geq |X| \text{ and } T[j'\,..\,n] \preceq T[j\,..\,n]\},$$

and denote $\delta(j) := |\text{Pos}(j)|$.

**Lemma 5.5.** *Let $j \in [1\,..\,n{-}3\tau{+}2]\setminus\mathsf{R}$ and $X \in \mathcal{D}$ be a prefix of $T[j\,..\,n]$. Denote $\delta_{\text{text}} = |X|{-}2\tau$ and $b_X = \text{RangeBeg}(X, T)$. Then:*

1. *It holds $\text{ISA}[j] = b_X + \delta(j)$.*
2. *If $y \in [1\,..\,n']$ is such that $s_y^{\text{lex}} = j + \delta_{\text{text}}$, then $\delta(j) = \text{rank}_{W,\overline{X}}(y)$.*

*Proof.* 1. Observe that $j' \in \mathrm{Occ}(X, T)$ holds if and only if $\mathrm{LCE}(j, j') \geq |X|$. Thus, by definition of ISA, we have $\mathrm{ISA}[j] = \mathrm{RangeBeg}(X, T) + |\{j' \in \mathrm{Occ}(X, T) : T[j' \mathinner{.\,.} n] \preceq T[j \mathinner{.\,.} n]\}| = b_X + |\{j' \in [1 \mathinner{.\,.} n] : \mathrm{LCE}(j, j') \geq |X| \text{ and } T[j' \mathinner{.\,.} n] \preceq T[j \mathinner{.\,.} n]\}| = b_X + \delta(j)$.

2. Denote $s = j + \delta_{\mathrm{text}}$. By definition of $\mathcal{D}$, we have $s \in \mathsf{S}$. By the consistency of $\mathsf{S}$, there exists a bijection (given by the mapping $j' \mapsto j' + \delta_{\mathrm{text}}$) between positions $j' \in [1 \mathinner{.\,.} n] \setminus \mathsf{R}$ satisfying $T[j' \mathinner{.\,.} \mathrm{succ}_{\mathsf{S}}(j') + 2\tau) = X$ and $T[j' \mathinner{.\,.} n] \preceq T[j \mathinner{.\,.} n]$, and positions $s' \in \mathsf{S}$ such that $T^\infty[s' - \delta_{\mathrm{text}} \mathinner{.\,.} s' + 2\tau) = X$ and $T[s' \mathinner{.\,.} n] \preceq T[s \mathinner{.\,.} n]$. Thus, letting $y \in [1 \mathinner{.\,.} n']$ be such that $s_y^{\mathrm{lex}} = s$, we obtain that $\delta(j) = |\{i \in [1 \mathinner{.\,.} y] : T^\infty[s_i^{\mathrm{lex}} - \delta_{\mathrm{text}} \mathinner{.\,.} s_i^{\mathrm{lex}} + 2\tau) = X\}|$. Since we defined $W[i] = \overline{X_i}$, where $X_i = T^\infty[s_i^{\mathrm{lex}} - \tau \mathinner{.\,.} s_i^{\mathrm{lex}} + 2\tau)$, we conclude that $\delta(j) = \mathsf{rank}_{W, \overline{X}}(y)$. $\qquad\square$

**Proposition 5.6.** *Let $j \in [1 \mathinner{.\,.} n]$ be such that $j \in [1 \mathinner{.\,.} n] \setminus \mathsf{R}$. Given the data structure from Section 5.2.1 and the position $j$, we can compute $\mathrm{ISA}[j]$ in $\mathcal{O}(\log^\epsilon n)$ time.*

*Proof.* Given $j \in [1 \mathinner{.\,.} n] \setminus \mathsf{R}$, we compute $\mathrm{ISA}[j]$ as follows. If $j > n - 3\tau + 2$, then letting $X = T[j \mathinner{.\,.} n]$, in $\mathcal{O}(1)$ time we compute $(b_X, e_X) = (\mathrm{RangeBeg}(X, T), \mathrm{RangeEnd}(X, T))$ using the lookup table $L_{\mathrm{range}}$. By definition of the lexicographical order, we then have $\mathrm{SA}[b+1] = j$, and hence we return $\mathrm{ISA}[j] = b + 1$. Let us thus assume $j \leq n - 3\tau - 2$. By $j \notin \mathsf{R}$ and the density condition of $\mathsf{S}$ (see Definition 2.4), this implies that $\mathsf{S} \cap [j \mathinner{.\,.} j + \tau) \neq \emptyset$. In $\mathcal{O}(1)$ time we compute $x = \mathsf{rank}_{B_{\mathsf{S}}, 1}(j - 1)$. Then, in $\mathcal{O}(1)$ we compute $s = \mathsf{select}_{B_{\mathsf{S}}, 1}(x + 1) = \mathrm{succ}_{\mathsf{S}}(j) \in \mathsf{S}$. We then have $X = T[j \mathinner{.\,.} s + 2\tau) \in \mathcal{D}$, and in particular, $|X| = s + 2\tau - j$. In $\mathcal{O}(1)$ time we lookup $(b_X, e_X) = L_{\mathrm{range}}[\mathrm{int}(X)]$, i.e., $b_X = \mathrm{RangeBeg}(X, T)$. Letting $y = A_{\mathrm{smap}}^{-1}[x + 1]$, we then have $s_y^{\mathrm{lex}} = s = j + |X| - 2\tau$. By Lemma 5.5, it thus remains to determine $\mathsf{rank}_{W, \overline{X}}(y)$. In $\mathcal{O}(1)$ time we compute $\overline{X}$ using the lookup table $L_{\mathrm{rev}}$. In $\mathcal{O}(\log^\epsilon n)$ time, we then compute $\delta(j) = \mathsf{rank}_{W, \overline{X}}(y)$ using Theorem 2.2, and finally return $\mathrm{ISA}[j] = b_X + \delta(j)$. $\qquad\square$

### 5.2.3 Implementation of SA Queries

**Lemma 5.7.** *Let $i \in [1 \mathinner{.\,.} n]$ be such that $\mathrm{SA}[i] \in [1 \mathinner{.\,.} n - 3\tau + 2] \setminus \mathsf{R}$ and $X \in \mathcal{D}$ be a prefix of $T[\mathrm{SA}[i] \mathinner{.\,.} n]$. Denote $\delta_{\mathrm{text}} = |X| - 2\tau$ and $b_X = \mathrm{RangeBeg}(X, T)$. Then:*

  *1. It holds $i = b_X + \delta(\mathrm{SA}[i])$.*
  *2. If $y = \mathsf{select}_{W, \overline{X}}(i - b_X)$, then $s_y^{\mathrm{lex}} = \mathrm{SA}[i] + \delta_{\mathrm{text}}$.*

*Proof.* 1. Denote $j = \mathrm{SA}[i]$. By Lemma 5.5(1), $i = \mathrm{ISA}[j] = b_X + \delta(j) = b_X + \delta(\mathrm{SA}[i])$.

2. By the consistency of $\mathsf{S}$, we have $\mathrm{SA}[i] + \delta_{\mathrm{text}} \in \mathsf{S}$. Thus, there exists $y \in [1 \mathinner{.\,.} n']$ such that $s_y^{\mathrm{lex}} = \mathrm{SA}[i] + \delta_{\mathrm{text}}$. By Lemma 5.5(2) applied for $j = \mathrm{SA}[i]$, for any such $y$ it holds $\delta(\mathrm{SA}[i]) = \mathsf{rank}_{W, \overline{X}}(y)$. By (1), we thus have $i - b_X = \mathsf{rank}_{W, \overline{X}}(y)$. Since $X$ is a prefix of $T[\mathrm{SA}[i] \mathinner{.\,.} n]$, such $y$ must also satisfy $T[s_y^{\mathrm{lex}} - \delta_{\mathrm{text}} \mathinner{.\,.} s_y^{\mathrm{lex}} + 2\tau) = X$, or equivalently, $\overline{X}$ must be a prefix of $W[y]$. The only $y \in [1 \mathinner{.\,.} n']$ for which $\overline{X}$ is a prefix of $W[y]$ and that satisfies $\mathsf{rank}_{W, \overline{X}}(y) = i - b_X$, is $y = \mathsf{select}_{W, \overline{X}}(i - b_X)$. $\qquad\square$

**Proposition 5.8.** *Let $i \in [1 \mathinner{.\,.} n]$ be such that $\mathrm{SA}[i] \in [1 \mathinner{.\,.} n] \setminus \mathsf{R}$. Given the data structure from Section 5.2.1 and the index $i$, we can compute $\mathrm{SA}[i]$ in $\mathcal{O}(\log^\epsilon n)$ time.*

*Proof.* Given $i \in [1 \mathinner{.\,.} n]$ such that $\mathrm{SA}[i] \in [1 \mathinner{.\,.} n] \setminus \mathsf{R}$, we compute $\mathrm{SA}[i]$ as follows. First, we compute $y = \mathsf{rank}_{B_{3\tau-1}, 1}(i - 1)$. The string $Y = A_{\mathrm{short}}[y + 1]$ is then a prefix of $T[\mathrm{SA}[i] \mathinner{.\,.} n]$ of length $\min(3\tau - 1, n + 1 - i)$. If $|Y| < 3\tau - 1$, we therefore have $\mathrm{SA}[i] > n - 3\tau + 2$ and moreover, $\mathrm{SA}[i] + |Y| = n + 1$. Thus, we return $\mathrm{SA}[i] = n + 1 - |Y|$. Otherwise (i.e., $|Y| = 3\tau - 1$), using $L_{\mathcal{D}}$ on $Y$ we determine $x = \mathrm{int}(X)$, where $X \in \mathcal{D}$ is a prefix of $T[\mathrm{SA}[i] \mathinner{.\,.} n]$. In $\mathcal{O}(1)$ time we lookup $(b_X, e_X) = L_{\mathrm{range}}[x]$, i.e., $b_X = \mathrm{RangeBeg}(X, T)$. In $\mathcal{O}(\log^\epsilon n)$ time, we then compute $y = \mathsf{select}_{W, \overline{X}}(i - b_X)$ using Theorem 2.2 (the packed representation of $\overline{X}$ is obtained using the lookup table $L_{\mathrm{rev}}$ in $\mathcal{O}(1)$ time). By Lemma 5.7(2), we then have $\mathrm{SA}[i] = s_y^{\mathrm{lex}} - \delta_{\mathrm{text}}$, where

$\delta_{\text{text}} = |X| - 2\tau$. Using $B_\mathsf{S}$, in $\mathcal{O}(1)$ time we compute $j' = \mathsf{select}_{B_\mathsf{S},1}(A_{\text{smap}}[y])$. We then have $j' = s_y^{\text{lex}}$, and hence we return $\mathrm{SA}[i] = j' - \delta_{\text{text}}$. Altogether, the query takes $\mathcal{O}(\log^\epsilon n)$ time. $\quad\square$

### 5.2.4 Construction Algorithm

**Proposition 5.9.** *Given* $\mathrm{C}_{\mathrm{SA}}(T)$, *we can augment it into a data structure from Section 5.2.1 in* $\mathcal{O}(n \min(1, \log \sigma / \sqrt{\log n}))$ *time and using* $\mathcal{O}(n / \log_\sigma n)$ *working space.*

*Proof.* First, using Theorem 2.6, we construct a $\tau$-synchronizing set $\mathsf{S}$ of size $\mathcal{O}(n/\tau)$ in $\mathcal{O}(n/\tau) = \mathcal{O}(n/\log_\sigma n)$ time from a packed representation of $T$. The set $\mathsf{S}$ is returned as an array taking $\mathcal{O}(n/\log_\sigma n)$ space. Using this array, we initialize the bitvector $B_\mathsf{S}$ in $\mathcal{O}(n/\log_\sigma n)$ time. Augmenting $B_\mathsf{S}$ with Theorem 2.1 takes $\mathcal{O}(n/\log n)$ time.

Next, we construct the arrays $A_{\text{smap}}^{-1}$ and $A_{\text{smap}}$. We start by creating the sequence $(s_t^{\text{text}})_{t \in [1..n']}$ using select queries on $B_\mathsf{S}$. This takes $\mathcal{O}(n/\log_\sigma n)$ time. Then, given $(s_t^{\text{text}})_{t \in [1..n']}$, and the packed representation of $T$, by [50, Theorem 4.3], we compute the sequence $(s_t^{\text{lex}})_{t \in [1..n']}$ in $\mathcal{O}(n/\log_\sigma n)$ time. Given $(s_t^{\text{lex}})_{t \in [1..n']}$, we then easily obtain the arrays $A_{\text{smap}}^{-1}$ and $A_{\text{smap}}$: simply scan the sequence $(s_t^{\text{lex}})_{t \in [1..n']}$ and for each $i \in [1..n']$, let $j = \mathsf{rank}_{B_\mathsf{S},1}(s_i^{\text{lex}})$ and note that then $s_j^{\text{text}} = s_i^{\text{lex}}$ and hence we can set $A_{\text{smap}}^{-1}[j] = i$ and $A_{\text{smap}}[i] = j$.

Next, we initialize $L_{\text{rev}}$. In the RAM model, such array is easily initialized in $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n/\log_\sigma n)$ time. The sequence $W[1..n']$ is then obtained from $(s_t^{\text{lex}})_{t \in [1..n']}$ using $L_{\text{rev}}$ in $\mathcal{O}(n/\log_\sigma n)$ time. We then process $W$ using Theorem 2.2, which takes $\mathcal{O}(n \min(1, \log \sigma / \sqrt{\log n}))$ time and $\mathcal{O}(n/\log_\sigma n)$ working space.

Finally, to construct $L_\mathcal{D}$, we first compute a lookup table that for every $Z \in [0..\sigma)^{2\tau}$ tells whether $T[j..j+2\tau) = Z$ implies $j \in \mathsf{S}$ (by consistency of $\mathsf{S}$ this does not depend on $j$). Given the array containing the positions in $\mathsf{S}$ and the packed representation of $T$, this takes $\mathcal{O}(\sigma^{2\tau} + |\mathsf{S}|) = \mathcal{O}(n/\log_\sigma n)$ time. Given such lookup table, we iterate through every $Y \in [0..\sigma)^{3\tau-1}$ and in $\mathcal{O}(\tau)$ time we compute the shortest prefix $X$ of $Y$ whose length-$2\tau$ suffix is marked true in the above lookup table. If such $X$ exists, we have $X \in \mathcal{D}$. Accounting for the initialization of $L_\mathcal{D}$, over all $Y \in [0..\sigma)^{3\tau-1}$, this takes $\mathcal{O}(\sigma^{6\tau} + \sigma^{3\tau-1} \log_\sigma n) = \mathcal{O}(n/\log_\sigma n)$ time. $\quad\square$

## 5.3 The Periodic Positions

In this section, we describe a data structure that, given any $j \in [1..n]$ (resp. $i \in [1..n]$) satisfying $j \in \mathsf{R}$ (resp. $\mathrm{SA}[i] \in \mathsf{R}$) computes $\mathrm{ISA}[j]$ (resp. $\mathrm{SA}[i]$) in $\mathcal{O}(\log \log n)$ time.

The section is organized as follows. First, we present the toolbox of combinatorial properties for periodic positions (Section 5.3.1). Next, we introduce the components of the data structure (Section 5.3.2). We then show how using this structure to implement some basic navigational routines (Section 5.3.3). Next, we describe the query algorithms (Sections 5.3.4 and 5.3.5). Finally, we show the construction algorithm (Section 5.3.6).

### 5.3.1 Preliminaries

We start by introducing the definitions to express the properties utilized in our data structures. For any $j \in \mathsf{R}$, we define $\text{L-root}(j) = \min\{T[j+t..j+t+p) : t \in [0..p)\}$, where $p = \text{per}(T[j..j+3\tau-1))$. We denote $\text{Roots} = \{\text{L-root}(j) : j \in \mathsf{R}\}$. For any $j \in \mathsf{R}$, let $e(j) = \min\{j' \geq j : j' \notin \mathsf{R}\} + 3\tau - 2$.

**Lemma 5.10.** *Let* $j \in \mathsf{R}$ *and* $p = \text{per}(T[j..j+3\tau-1))$. *Then:*

    *1. If* $j+1 \in \mathsf{R}$ *then* $\text{per}(T[j+1..j+3\tau)) = p$,
    *2. It holds* $e(j) = j + p + \text{LCE}(j, j+p)$.

*Proof.* 1. Denote $P = T[j \mathinner{.\,.} j + 3\tau - 1)$, $P' = T[j+1 \mathinner{.\,.} j + 3\tau)$, and $p' = \mathrm{per}(P')$. Our goal is to show that $p' = p$. For a proof by contradiction, assume $p' \neq p$. By the assumption, $\mathrm{per}(P) = p$. Denote $Y = P'[1 \mathinner{.\,.} \tau]$, and note that since $P$ and $P'$ overlap by $3\tau - 2 \geq \tau$ symbols, $Y$ is a substring of $P$, and hence has periods $p$ and $p'$. Observe that we cannot have $p \mid p'$ since this would imply that $Y[1 \mathinner{.\,.} p']$ is not primitive which would contradict $p' = \mathrm{per}(P')$. Observe now that we have $p, p' \leq \frac{1}{3}\tau$. By the Weak Periodicity Lemma [29], we thus have that $Y$ has period $p'' = \gcd(p, p')$. By our assumptions, this implies $p'' < p'$ and $p'' \mid p'$. Thus, again we obtain that $Y[1 \mathinner{.\,.} p']$ is not primitive. Therefore, we must have $p' = p$.

2. Denote $j' = e(j) - 3\tau + 2$. By definition, we then have $[j \mathinner{.\,.} j'] \subseteq \mathsf{R}$ and $j' \notin \mathsf{R}$. By the above, for every $t \in [0 \mathinner{.\,.} j' - j)$, it holds $\mathrm{per}(T[j + t \mathinner{.\,.} j + t + 3\tau - 1)) = p$. Thus, for every $j'' \in [j \mathinner{.\,.} j' + 3\tau - 2 - p)$, we have $T[j''] = T[j'' + p]$, i.e., the substring $T[j \mathinner{.\,.} j' + 3\tau - 2)$ has period $p$, and thus $\mathrm{LCE}(j, j + p) \geq (j' + 3\tau - 2) - j - p$, or equivalently, $j + p + \mathrm{LCE}(j, j + p) \geq j' + 3\tau - 2 = e(j)$. To show that this lower bound on $j + p + \mathrm{LCE}(j, j + p)$ is tight, let us assume that $e(j) \leq n$ (otherwise, the claim follows immediately). Equivalently, we then have $j' + 3\tau - 2 = e(j) \leq n$ and to finish the proof, it remains to show $T[e(j)] \neq T[e(j) - p]$. Recall that $\mathrm{per}(T[j' - 1 \mathinner{.\,.} j' + 3\tau - 2)) = p$. Thus, $T[j' + 3\tau - 2] = T[e(j)] = T[e(j) - p] = T[j' + 3\tau - 2 - p]$ would imply that $\mathrm{per}(T[j' \mathinner{.\,.} j' + 3\tau - 1)) = p$, or equivalently, that $j' \in \mathsf{R}$, a contradiction. $\square$

Observe that by definition of L-root, letting $p = |\text{L-root}(j)|$, there exists $s \in [0 \mathinner{.\,.} p)$ such that $T[j + s \mathinner{.\,.} j + s + p) = \text{L-root}(j)$. Combining this with Lemma 5.10 implies that for every $j \in \mathsf{R}$, we can write $T[j \mathinner{.\,.} e(j)) = H'H^k H''$, where $H = \text{L-root}(j)$, and $H'$ (resp. $H''$) is a proper suffix (resp. prefix) of $H$. We call such factorization the *L-decomposition* of $T[j \mathinner{.\,.} e(j))$. Note that the L-decomposition is unique, since otherwise would contradict the synchronization property of primitive strings [24, Lemma 1.11]. We denote $\text{L-head}(j) = |H'|$, $\text{L-exp}(j) = k$, and $\text{L-tail}(j) = |H''|$. For $j \in \mathsf{R}$, we let $\mathrm{type}(j) = +1$ if $e(j) \leq n$ and $T[e(j)] \succ T[e(j) - p]$ (where $p = |\text{L-root}(j)|$), and $\mathrm{type}(j) = -1$ otherwise. For any $j \in \mathsf{R}$, we denote $e^{\mathrm{full}}(j) = e(j) - \text{L-tail}(j)$. Observe that $e^{\mathrm{full}}(j) = j + \text{L-head}(j) + \text{L-exp}(j) \cdot |\text{L-root}(j)|$.

We repeatedly refer to the following subsets of $\mathsf{R}$. First, denote $\mathsf{R}^- = \{j \in \mathsf{R} : \mathrm{type}(j) = -1\}$ and $\mathsf{R}^+ = \mathsf{R} \setminus \mathsf{R}^-$. For any $H \in \Sigma^+$ and any $s \in \mathbb{Z}_{\geq 0}$ we then let $\mathsf{R}_H = \{j \in \mathsf{R} : \text{L-root}(j) = H\}$, $\mathsf{R}_H^- = \mathsf{R}^- \cap \mathsf{R}_H$, $\mathsf{R}_H^+ = \mathsf{R}^+ \cap \mathsf{R}_H$, $\mathsf{R}_{s,H} = \{j \in \mathsf{R}_H : \text{L-head}(j) = s\}$, $\mathsf{R}_{s,H}^- = \mathsf{R}^- \cap \mathsf{R}_{s,H}$, and $\mathsf{R}_{s,H}^+ = \mathsf{R}^+ \cap \mathsf{R}_{s,H}$.

The following lemmas establish the key properties of periodic positions. First, we prove that the set of positions $\mathsf{R}_{s,H}$ occupies a contiguous block in SA and describe the structure of such block.

**Lemma 5.11.** *Let $j \in \mathsf{R}_{s,H}$. For any $j' \in [1 \mathinner{.\,.} n]$, $\mathrm{LCE}(j, j') \geq 3\tau - 1$ holds if and only if $j' \in \mathsf{R}_{s,H}$. Moreover, if $j' \in \mathsf{R}_{s,H}$ then, letting $t = e(j) - j$ and $t' = e(j') - j'$, it holds $\mathrm{LCE}(j, j') \geq \min(t, t')$ and:*

1. *If $\mathrm{type}(j) \neq \mathrm{type}(j')$, then $T[j \mathinner{.\,.} n] \prec T[j' \mathinner{.\,.} n]$ if and only if $\mathrm{type}(j) < \mathrm{type}(j')$,*
2. *If $\mathrm{type}(j) = \mathrm{type}(j') = -1$ and $t \neq t'$, then $T[j \mathinner{.\,.} n] \prec T[j' \mathinner{.\,.} n]$ if and only if $t < t'$,*
3. *If $\mathrm{type}(j) = \mathrm{type}(j') = +1$ and $t \neq t'$, then $T[j \mathinner{.\,.} n] \prec T[j' \mathinner{.\,.} n]$ if and only if $t > t'$,*
4. *If $\mathrm{type}(j) \neq \mathrm{type}(j')$ or $t \neq t'$, then $\mathrm{LCE}(j, j') = \min(t, t')$.*

*Proof.* Let $j' \in [1 \mathinner{.\,.} n]$ be such that $\mathrm{LCE}(j, j') \geq 3\tau - 1$. Denoting $p = \mathrm{per}(T[j \mathinner{.\,.} j + 3\tau - 1))$ and $p' = \mathrm{per}(T[j' \mathinner{.\,.} j' + 3\tau - 1))$ we then have $p' = p \leq \frac{1}{3}$. Thus, $j' \in \mathsf{R}$ and $\text{L-root}(j') = \min\{T[j' + t \mathinner{.\,.} j' + t + p') : t \in [0 \mathinner{.\,.} p')\} = \min\{T[j' + t \mathinner{.\,.} j' + t + p) : t \in [0 \mathinner{.\,.} p)\} = \min\{T[j + t \mathinner{.\,.} j + t + p) : t \in [0 \mathinner{.\,.} p)\} = H$. To show that $\text{L-head}(j') = s$, note that by $|H| \leq \tau$, the string $H'H^2$ (where $H'$ is a length-$s$ suffix of $H$) is a prefix of $T[j \mathinner{.\,.} j + 3\tau - 1) = T[j' \mathinner{.\,.} j' + 3\tau - 1)$. On the other hand, $\text{L-head}(j') = s'$ implies that $\widehat{H}'H^2$ (where $\widehat{H}'$ is a length-$s'$ suffix of $H$) is a prefix of $T[j' \mathinner{.\,.} j' + 3\tau - 1)$. Thus, by the synchronization property of primitive strings [24,

Lemma 1.11] applied to the two copies of $H$, we have $s' = s$, and consequently, $j' \in \mathsf{R}_{s,H}$. For the converse implication, assume $j' \in \mathsf{R}_{s,H}$. This implies that both $T[j \mathinner{.\,.} e(j))$ and $T[j' \mathinner{.\,.} e(j'))$ are prefixes of $H' \cdot H^\infty[1 \mathinner{.\,.})$ (where $H'$ is as above). Thus, by $e(j) - j$, $e(j') - j' \geq 3\tau - 1$, we obtain $\mathrm{LCE}(j, j') \geq 3\tau - 1$.

Let us now assume $j' \in \mathsf{R}_{s,H}$. Since, as noted above, both $T[j \mathinner{.\,.} e(j)) = T[j \mathinner{.\,.} j + t)$ and $T[j' \mathinner{.\,.} e(j')) = T[j' \mathinner{.\,.} j' + t')$ are prefixes of $H' \cdot H^\infty[1 \mathinner{.\,.})$, we have $\mathrm{LCE}(j, j') \geq \min(t, t')$.

1. Assume $\mathrm{type}(j) < \mathrm{type}(j')$. Let $Q = H' \cdot H^\infty[1 \mathinner{.\,.})$, where $H'$ is a length-$s$ suffix of $H$. We will prove $T[j \mathinner{.\,.} n] \prec Q \prec T[j' \mathinner{.\,.} n]$, which implies the claim. First, we note that $\mathrm{type}(j) = -1$ implies that either $e(j) = n + 1$, or $e(j) \leq n$ and $T[e(j)] \prec T[e(j) - |H|]$. In the first case, $T[j \mathinner{.\,.} e(j)) = T[j \mathinner{.\,.} n]$ is a proper prefix of $Q$ and hence $T[j \mathinner{.\,.} n] \prec Q$. In the second case, we have $T[j \mathinner{.\,.} e(j)) = T[j \mathinner{.\,.} j + t) = Q[1 \mathinner{.\,.} t]$ and $T[j + t] \prec T[j + t - |H|] = Q[1 + t - |H|] = Q[1 + t]$. Consequently, $T[j \mathinner{.\,.} n] \prec Q$. To show $Q \prec T[j' \mathinner{.\,.} n]$ we observe that $\mathrm{type}(j') = +1$ implies $e(j') \leq n$. Thus, we have $Q[1 \mathinner{.\,.} t'] = T[j' \mathinner{.\,.} e(j')) = T[j' \mathinner{.\,.} j' + t')$ and $Q[1 + t'] = Q[1 + t' - |H|] = T[j' + t' - |H|] \prec T[j' + t']$. Hence, we obtain $Q \prec T[j' \mathinner{.\,.} n]$. We have thus obtained $T[j \mathinner{.\,.} n] \prec Q \prec T[j' \mathinner{.\,.} n]$ which implies $T[j \mathinner{.\,.} n] \prec T[j' \mathinner{.\,.} n]$. The opposite implication follows easily by symmetry. More precisely, in a proof by contraposition, assuming $\mathrm{type}(j) \geq \mathrm{type}(j')$ we immediately obtain $\mathrm{type}(j) > \mathrm{type}(j')$ from the assumption. By the analogous argument as above we then have $T[j \mathinner{.\,.} n] \succ T[j' \mathinner{.\,.} n]$.

2. Assume $t < t'$. Similarly as above, we consider two cases for $e(j)$. If $e(j) = n + 1$, then by $t < t'$, the string $T[j \mathinner{.\,.} e(j)) = T[j \mathinner{.\,.} n]$ is a proper prefix of $T[j' \mathinner{.\,.} e(j')) = T[j' \mathinner{.\,.} j' + t')$ and hence $T[j \mathinner{.\,.} n] \prec T[j' \mathinner{.\,.} j' + t') \preceq T[j' \mathinner{.\,.} n]$. On the other hand, if $e(j) \leq n$, then we have $T[j \mathinner{.\,.} j + t) = T[j' \mathinner{.\,.} j' + t)$ and by $t < t'$, $T[j + t] \prec T[j + t - |H|] = T[j' + t - |H|] = T[j' + t]$. Hence, $T[j \mathinner{.\,.} n] \prec T[j' \mathinner{.\,.} n]$. The opposite implication follows by symmetry similarly as in Item 1.

3. Assume $t > t'$. By $\mathrm{type}(j') = +1$ we have $e(j') \leq n$. Thus, by $t > t'$, we have $T[j \mathinner{.\,.} j + t') = T[j' \mathinner{.\,.} j' + t')$ and $T[j + t'] = T[j + t' - |H|] = T[j' + t' - |H|] \prec T[j' + t']$. Hence, $T[j \mathinner{.\,.} n] \prec T[j' \mathinner{.\,.} n]$. The opposite implication follows by symmetry similarly as in Item 1.

4. By the earlier implication, $\mathrm{LCE}(j, j') \geq \min(t, t')$. Thus, it remains to show $\mathrm{LCE}(j, j') \leq \min(t, t')$. First, let $\mathrm{type}(j) \neq \mathrm{type}(j')$ and without the loss of generality let us assume $\mathrm{type}(j) < \mathrm{type}(j')$ (i.e., $\mathrm{type}(j) = -1$ and $\mathrm{type}(j') = +1$). Consider two cases:

- First, assume $t \leq t'$. Our goal is to prove $\mathrm{LCE}(j, j') \leq t$. If $j + t = n + 1$, then we immediately obtain the claim. Let us thus assume $j + t \leq n$. In the proof of Item 1 we showed that in this case $\mathrm{type}(j) = -1$ implies $T[j + t] \prec Q[1 + t]$. On the other hand, there we also proved that $\mathrm{type}(j') = +1$ implies $Q[1 \mathinner{.\,.} t'] = T[j' \mathinner{.\,.} j' + t')$ and $Q[1 + t'] \prec T[j' + t']$. By $t \leq t'$, we thus obtain $Q[1 + t] \preceq T[j' + t]$. Consequently, $T[j + t] \neq T[j' + t]$ and hence $\mathrm{LCE}(j, j') \leq t$.

- Let us now assume $t > t'$. Our goal is to prove $\mathrm{LCE}(j, j') \leq t'$. In the proof of Item 1 we showed that $\mathrm{type}(j) = -1$ implies that $T[j \mathinner{.\,.} j + t) = Q[1 \mathinner{.\,.} t]$. Thus, by $t > t'$ we have $T[j + t'] = Q[1 + t']$. On the other hand, in the proof of Item 1 we also proved that $\mathrm{type}(j') = +1$ implies $Q[1 + t'] \prec T[j' + t']$. Thus, we obtain $T[j + t'] \neq T[j' + t']$ and hence $\mathrm{LCE}(j, j') \leq t'$.

This concludes the proof of the claim in the case $\mathrm{type}(j) \neq \mathrm{type}(j')$. Let us thus assume $\mathrm{type}(j) = \mathrm{type}(j')$ and $t \neq t'$. First, consider the case $\mathrm{type}(j) = \mathrm{type}(j') = -1$ and assume without the loss of generality that $t < t'$ (to match the assumption in Item 2). Our goal is thus to show $\mathrm{LCE}(j, j') \leq t$. In the proof of Item 2, we showed that we then either have $T[j \mathinner{.\,.} j + t) = T[j \mathinner{.\,.} n]$ (in which case $\mathrm{LCE}(j, j') \leq n - j + 1 = t$), or $T[j \mathinner{.\,.} j + t) = T[j' \mathinner{.\,.} j' + t)$ and $T[j + t] \prec T[j' + t]$ (which also immediately implies $\mathrm{LCE}(j, j') \leq t$). Let us now consider the case $\mathrm{type}(j) = \mathrm{type}(j') = +1$ and assume without the loss of generality that $t > t'$ (to match the assumption in Item 3). Our goal is thus to show $\mathrm{LCE}(j, j') \leq t'$. In the proof of Item 3, we

showed that we then have $T[j \mathinner{.\,.} j + t') = T[j' \mathinner{.\,.} j' + t')$ and $T[j + t'] \prec T[j' + t']$. This implies $\mathrm{LCE}(j, j') \le t'$. □

The key to the efficient computation of SA and ISA values for periodic positions is processing of the elements of $\mathsf{R}$ in blocks (note that unlike in Lemma 5.11, which describes the structure of blocks in SA, here we mean blocks of positions in the text). The starting positions of these blocks are defined as $\mathsf{R}' := \{j \in \mathsf{R} : j - 1 \notin \mathsf{R}\}$. We also let $\mathsf{R}'^{-} = \mathsf{R}' \cap \mathsf{R}^{-}$, $\mathsf{R}'^{+} = \mathsf{R}' \cap \mathsf{R}^{+}$, $\mathsf{R}'^{-}_{H} = \mathsf{R}' \cap \mathsf{R}^{-}_{H}$, and $\mathsf{R}'^{+}_{H} = \mathsf{R}' \cap \mathsf{R}^{+}_{H}$ for any $H \in \Sigma^{+}$. The following lemma justifies this strategy.

**Lemma 5.12.** *For every $j \in \mathsf{R} \setminus \mathsf{R}'$ it holds:*

- L-root$(j - 1) = $ L-root$(j)$,
- $e(j - 1) = e(j)$,
- L-tail$(j - 1) = $ L-tail$(j)$,
- $e^{\mathrm{full}}(j - 1) = e^{\mathrm{full}}(j)$,
- $\mathrm{type}(j - 1) = \mathrm{type}(j)$.

*Proof.* Denote $p = \mathrm{per}(T[j-1 \mathinner{.\,.} j-1+3\tau-1))$. By Lemma 5.10(1), it holds $\mathrm{per}(T[j \mathinner{.\,.} j + 3\tau - 1)) = p$. By $p \le \frac{\tau}{3}$, we thus have $T[j-1 \mathinner{.\,.} j-1+p) = T[j-1+p \mathinner{.\,.} j-1+2p)$. Consequently, $\{T[j-1+t \mathinner{.\,.} j-1+t+p) : t \in [0 \mathinner{.\,.} p)\} = \{T[j+t \mathinner{.\,.} j+t+p) : t \in [0 \mathinner{.\,.} p)\}$, and hence L-root$(j-1) = $ L-root$(j)$.

Denote $p' = \mathrm{per}(T[j \mathinner{.\,.} j+3\tau-1))$. By Lemma 5.10(2), $e(j-1) = j-1+p+\mathrm{LCE}(j-1, j-1+p)$ and $e(j) = j+p'+\mathrm{LCE}(j, j+p')$. Thus, by $p = p'$ (following by the above) and $T[j-1] = T[j-1+p]$, we have $e(j-1) = j-1+p+\mathrm{LCE}(j-1, j-1+p) = j+p+\mathrm{LCE}(j, j+p) = j+p'+\mathrm{LCE}(j, j+p') = e(j)$.

Assume $T[j - 1 \mathinner{.\,.} e(j - 1)) = H'H^k H''$, where $H = $ L-root$(j-1)$, $|H'| = $ L-head$(j-1)$, and $|H''| = $ L-tail$(j-1)$. By $e(j-1) = e(j)$ and the uniqueness of L-decomposition, this implies that either $T[j \mathinner{.\,.} e(j)) = H'[2 \mathinner{.\,.} |H'|]H^k H''$ (if $|H'| > 0$) or $T[j \mathinner{.\,.} e(j)) = H[2 \mathinner{.\,.} |H|]H^{k-1}H''$ (otherwise) is the L-decomposition of $T[j \mathinner{.\,.} e(j))$. In both cases, L-tail$(j-1) = $ L-tail$(j) = |H''|$.

By the above two properties, $e^{\mathrm{full}}(j-1) = e(j-1) - $ L-tail$(j-1) = e(j) - $ L-tail$(j) = e^{\mathrm{full}}(j)$.

The last claim follows from the definition of type and equalities $e(j-1) = e(j)$ and $p = p'$. □

The above is complemented by the following results establishing the lower bound on the gap between blocks of positions in $\mathsf{R}$, and that a mapping from $j$ to $e^{\mathrm{full}}(j)$ establishes an injective mapping of blocks of positions in $\mathsf{R}$ to positions in $T$.

**Lemma 5.13.** *Let $j, j', j'' \in [1 \mathinner{.\,.} n]$ be such that $j, j'' \in \mathsf{R}$, $j' \notin \mathsf{R}$, and $j < j' < j''$. Then, it holds $e(j) \le j'' + \tau - 1$ and $j'' - j \ge 2\tau$.*

*Proof.* Let $r = \min\{i \in (j' \mathinner{.\,.} j''] : i \in \mathsf{R}\}$. Then, $r \in \mathsf{R}'$. Observe that by Lemma 5.10 (resp. by $r - 1 \notin \mathsf{R}$), it holds $\mathrm{per}(T[j \mathinner{.\,.} e(j))) \le \lfloor \frac{1}{3}\tau \rfloor$ (resp. $\mathrm{per}(T[r \mathinner{.\,.} e(r))) \le \lfloor \frac{1}{3}\tau \rfloor$), $e(j) - j \ge 3\tau - 1$ (resp. $e(r) - r \ge 3\tau - 1$), and the substring $T[j \mathinner{.\,.} e(j))$ (resp. $T[r \mathinner{.\,.} e(r))$) cannot be extended in $T$ to the right (resp. left) without changing its shortest period. By [54, Fact 2.2.4], the fragments $T[j \mathinner{.\,.} e(j))$ and $T[r \mathinner{.\,.} e(r))$ must therefore overlap by less than $2\lfloor \frac{1}{3}\tau \rfloor$ symbols. In other words, $e(j) - r < 2\lfloor \frac{1}{3}\tau \rfloor$. By $r \le j''$ we thus obtain $e(j) \le r + 2\lfloor \frac{1}{3}\tau \rfloor \le j'' + \tau - 1$, i.e., the first claim. Equivalently, we can state that $j'' \ge e(j) - \tau + 1$. By combining this with $e(j) - j \ge 3\tau - 1$, we then obtain $j'' \ge e(j) - \tau + 1 \ge j + 3\tau - 1 - \tau + 1 = j + 2\tau$, i.e., the second claim. □

**Lemma 5.14.** *For any $j, j' \in \mathsf{R}'$, $j \ne j'$ implies $e^{\mathrm{full}}(j) \ne e^{\mathrm{full}}(j')$.*

*Proof.* Assume without the loss of generality that $j < j'$. Then, $j' - 1 \notin \mathsf{R}$. By Lemma 5.13 applied for $j$, $j' - 1$, and $j'$ we obtain $e(j) \le j' + \tau - 1$. Consequently, $e^{\mathrm{full}}(j) \le e(j) \le j' + \tau - 1$. Let now $r' = \min\{t \in (j' \mathinner{.\,.} n] : t \notin \mathsf{R}\}$. We then have $e(j') = r' + 3\tau - 2$. Since for every $t \in \mathsf{R}$,

it holds $e(t) - e^{\text{full}}(t) = \text{L-tail}(t) = |\text{L-root}(t)| \leq \lfloor \frac{1}{3}\tau \rfloor$, we thus have $e^{\text{full}}(j') \geq e(j') - \lfloor \frac{1}{3}\tau \rfloor \geq r' + 3\tau - 2 - \lfloor \frac{1}{3}\tau \rfloor \geq j' + 2\tau - 1$. Combining this with the earlier upper bound on $e^{\text{full}}(j)$, we thus obtain $e^{\text{full}}(j) \leq j' + \tau - 1 < j' + 2\tau - 1 \leq e^{\text{full}}(j')$. In particular, $e^{\text{full}}(j) \neq e^{\text{full}}(j')$. $\qquad\square$

### 5.3.2 The Data Structure

**Definitions**  Let $q = |\mathsf{R}'^{-}|$ and let $(r_i^{\text{text}})_{i \in [1..q]}$ be a sequence containing all elements of $\mathsf{R}'^{-}$ in sorted order, i.e, for any $i, i' \in [1..q]$, $i < i'$ implies $r_i^{\text{text}} < r_{i'}^{\text{text}}$. Let $(r_i^{\text{lex}})_{i \in [1..q]}$ also be a sequence containing all elements $k \in \mathsf{R}'^{-}$, but sorted first according to L-root$(k)$ and in case of ties, by $T[e^{\text{full}}(k)\mathbin{..}n]$. Formally, for any $i, i' \in [1..q]$, $i < i'$ implies that L-root$(r_i^{\text{lex}}) \prec$ L-root$(r_{i'}^{\text{lex}})$, or L-root$(r_i^{\text{lex}}) =$ L-root$(r_{i'}^{\text{lex}})$ and $T[e^{\text{full}}(r_i^{\text{lex}})\mathbin{..}n] \prec T[e^{\text{full}}(r_{i'}^{\text{lex}})\mathbin{..}n]$. Note that by Lemma 5.14, the sequence $(r_i^{\text{lex}})_{i \in [1..q]}$ is well-defined. Based on $(r_i^{\text{lex}})_{i \in [1..q]}$ we define the sequence of integers $(\ell_i)_{i \in [1..q]}$ as $\ell_i = e^{\text{full}}(r_i^{\text{lex}}) - r_i^{\text{lex}}$.

Let $L_{\text{root}}$ denote the mapping from $[0\mathbin{..}\sigma)^{3\tau-1}$ to $\mathbb{N}^2$ such that for any $X \in [0\mathbin{..}\sigma)^{3\tau-1}$ satisfying $\text{per}(X) \leq \frac{1}{3}\tau$, $L_{\text{root}}$ maps $X$ to a pair $(s, p)$, where $p = \text{per}(X)$ and $s \in [0\mathbin{..}p)$ is such that $X[1+s\mathbin{..}1+s+p] = \min\{X[1+t\mathbin{..}1+t+p] : t \in [0\mathbin{..}p)\}$. We also define $L_{\text{minexp}} : [0\mathbin{..}\sigma)^{3\tau-1} \to [1\mathbin{..}n]$ as the mapping such that for every $X \in [0\mathbin{..}\sigma)^{3\tau-1}$ satisfying $\text{per}(X) \leq \frac{1}{3}\tau$, if we let $p = \text{per}(X)$, $H = \min\{X[1+t\mathbin{..}1+t+p] : t \in [0\mathbin{..}p)\}$ and $s \in [0\mathbin{..}p)$ be such that $X[1+s\mathbin{..}1+s+p] = H$, then assuming $\mathsf{R}_{s,H}^{-} \neq \emptyset$, $L_{\text{minexp}}$ maps $X$ to $\min\{\text{L-exp}(j) : j \in \mathsf{R}_{s,H}^{-}\}$. Let $L_{\text{runs}}$ be a mapping, such that for every $H \in [0\mathbin{..}\sigma)^{\leq\tau}$ and every $H' \in [0\mathbin{..}\sigma)^{\leq\tau}$, $L_{\text{runs}}$ maps the pair $(H, H')$ to $(b, e)$ defined by $b = |\{k \in \mathsf{R}'^{-} : \text{L-root}(k) \prec H$, or L-root$(k) = H$ and $T[e^{\text{full}}(k)\mathbin{..}n] \prec H'\}|$ and $e = b + |\{k \in \mathsf{R}_H'^{-} : H'$ is a prefix of $T[e^{\text{full}}(k)\mathbin{..}n]\}|$. Note that then the set $\{r_i^{\text{lex}} : i \in (b\mathbin{..}e]\}$ consists of all positions $k \in \mathsf{R}_H'^{-}$ such that $H'$ is a prefix of $T[e^{\text{full}}(k)\mathbin{..}n]$. In particular, every $(H, \varepsilon)$ maps to a pair $(b, e)$ such that $e = \sum_{H' \preceq H} |\mathsf{R}_{H'}'^{-}|$. For any $\ell > 0$, $H \in [0\mathbin{..}\sigma)^{+}$, and $s \in [0\mathbin{..}|H|)$, we define $\text{Pref}_\ell(s, H)$ as the length-$\ell$ prefix of $H' \cdot H^{\infty}[1\mathbin{..}]$, where $H'$ is a length-$s$ suffix of $H$. Let $L_{\text{pref}}$ denote the mapping that, given the pair $(H, s)$, where $H \in [0\mathbin{..}\sigma)^{\leq\tau}$, and $s \in [0\mathbin{..}|H|)$, returns the packed encoding of $\text{Pref}_{3\tau-1}(s, H)$.

Let $B_{\text{exp}}[1\mathbin{..}n]$ be a bitvector such that for every $i \in [1\mathbin{..}n]$, it holds $B_{\text{exp}}[i] = 0$ if and only if $\text{SA}[i] \in [1\mathbin{..}n] \setminus \mathsf{R}^{-}$, or $i < n$ and the positions $j = \text{SA}[i]$ and $j' = \text{SA}[i+1]$ satisfy $j, j' \in \mathsf{R}_{s,H}^{-}$ and $\text{L-exp}(j) = \text{L-exp}(j')$ for some $H \in \text{Roots}$ and $s \in [0\mathbin{..}|H|)$. Let $B_{\mathsf{R}'}[1\mathbin{..}n]$ be a bitvector defined such that $B_{\mathsf{R}'}[i] = 1$ holds if and only if $i \in \mathsf{R}'$.

Let $A_{\text{len}}[1\mathbin{..}q]$ by an array defined by $A_{\text{len}}[i] = \ell_i$. Let $A_{\text{rmap}}[1\mathbin{..}q]$ be an array containing a permutation of $[1\mathbin{..}q]$ such that $A_{\text{rmap}}[i] = i'$ holds if and only if $r_i^{\text{text}} = r_{i'}^{\text{lex}}$. By $A_{\text{rmap}}^{-1}[1\mathbin{..}q]$ we denote an array containing a permutation of $[1\mathbin{..}q]$ such that $A_{\text{rmap}}^{-1}[i'] = i$ holds if and only if $r_i^{\text{text}} = r_{i'}^{\text{lex}}$.

**Components**  The data structure consists of two parts. The first part, designed to compute $\text{SA}[i]$ (resp. $\text{ISA}[j]$) for $i \in [1\mathbin{..}n]$ (resp. $j \in [1\mathbin{..}n]$) satisfying $\text{SA}[i] \in \mathsf{R}^{-}$ (resp. $j \in \mathsf{R}^{-}$), consists of the following eleven components:

1. $\text{C}_{\text{SA}}(T)$ (Section 5.1.1). It takes $\mathcal{O}(n/\log_\sigma n)$ space.
2. The lookup table $L_{\text{root}}$. When accessing $L_{\text{root}}$, strings $X \in [0\mathbin{..}\sigma)^{3\tau-1}$ are converted to $\text{int}(X)$. Thus, $L_{\text{root}}$ needs $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n^{6\mu}) = \mathcal{O}(n/\log_\sigma n)$ space.
3. The lookup table $L_{\text{minexp}}$. As above, $L_{\text{minexp}}$ also needs $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n/\log_\sigma n)$ space.
4. The lookup table $L_{\text{runs}}$. When storing the mapping from the key $(H, H')$, we first concatenate $H$ and $H'$ and convert it to an integer $x = \text{int}(HH')$ in the range $[0\mathbin{..}\sigma^{6\tau})$. We then create a triple $(x, |H|)$ (this contains enough information to decode $H$ and $H'$) and injectively map it to a positive integer not exceeding $\sigma^{6\tau}\tau < \sigma^{6\tau}\log_\sigma n$. Thus, $L_{\text{runs}}$ can be stored using $\mathcal{O}(n^{6\mu}\log n) = \mathcal{O}(n/\log_\sigma n)$ space.
5. The lookup table $L_{\text{pref}}$. When storing the mapping, we convert the string $H$ to $\text{int}(H)$.

By $\text{int}(H) \in [0 \mathinner{.\,.} \sigma^{6\tau})$ and $|H| \leq \tau$, each pair $(\text{int}(H), s)$ can then be injectively encoded as an integer in the range of size $\sigma^{6\tau} \tau < \sigma^{6\tau} \log_\sigma n$ and hence $L_{\text{pref}}$ needs $\mathcal{O}(n^{6\mu} \log n) = \mathcal{O}(n/\log_\sigma n)$ space.

6. The bitvector $B_{\exp}$ augmented using Theorem 2.1. It needs $\mathcal{O}(n/\log n)$ space.

7. The bitvector $B_{\mathsf{R}'}$ augmented using Theorem 2.1. It needs $\mathcal{O}(n/\log n)$ space.

8. The array $A_{\text{len}}[1 \mathinner{.\,.} q]$ augmented with a structure from Proposition 2.3. To analyze its space usage, consider any $j_1, j_2, j_3 \in \mathsf{R}'$ such that $j_1 < j_2 < j_3$. Then, $j_2 - 1 \notin \mathsf{R}$ and $j_3 - 1 \notin \mathsf{R}$. By Lemma 5.13 applied first for $j_1, j_2 - 1$, and $j_2$ we have $e(j_1) \leq j_2 + \tau - 1$. Applying it again for $j_2, j_3 - 1$, and $j_3$, we obtain $j_3 - j_2 \geq 2\tau$, or equivalently, $j_2 \leq j_3 - 2\tau$. Combining the two inequalities, we thus obtain that $e(j_1) \leq j_3 - \tau - 1 < j_3$. This implies that each position of $T$ belongs to at most two intervals in the collection $\{[j \mathinner{.\,.} e(j)) : j \in \mathsf{R}'\}$, and consequently, $\sum_{i=1}^{q} \ell_i \leq 2n$. On the other hand, by Lemma 5.13, for every $j, j' \in \mathsf{R}'$, $j \neq j'$ implies $|j' - j| \geq 2\tau$. Thus, $q = \mathcal{O}(n/\tau) = \mathcal{O}(n/\log_\sigma n)$. The array $A$ augmented using Proposition 2.3 thus needs $\mathcal{O}(n/\log_\sigma n)$ space.

9. The array $A_{\text{rmap}}$ in plain form using $\mathcal{O}(1 + q) = \mathcal{O}(n/\log_\sigma n)$ space.

10. The array $A_{\text{rmap}}^{-1}$ in plain form using $\mathcal{O}(1 + q) = \mathcal{O}(n/\log_\sigma n)$ space

11. The $\mathcal{O}(n/\log_\sigma n)$-space data structure from [50, Theorem 5.4] that, given any $i, i' \in [1 \mathinner{.\,.} n]$, returns $\text{LCE}(i, i')$ in $\mathcal{O}(1)$ time.

The second part of the structure, designed to compute $\text{SA}[i]$ (resp. $\text{ISA}[j]$) for $i \in [1 \mathinner{.\,.} n]$ (resp. $j \in [1 \mathinner{.\,.} n]$) satisfying $\text{SA}[i] \in \mathsf{R}^+$ (resp. $j \in \mathsf{R}^+$), consists of the symmetric counterparts adapted according to Lemma 5.11.

In total, the data structure takes $\mathcal{O}(n/\log_\sigma n)$ space.

### 5.3.3 Navigation Primitives

**Proposition 5.15.** *Given the data structure from Section 5.3.2 and any position $j \in \mathsf{R}$, we can in $\mathcal{O}(1)$ time compute the values* L-root$(j)$, L-head$(j)$, L-exp$(j)$, L-tail$(j)$, *and* type$(j)$.

*Proof.* We first compute $x \in [0 \mathinner{.\,.} \sigma^{6\tau})$ such that $x = \text{int}(T[j \mathinner{.\,.} j + 3\tau - 1))$. Given the packed encoding of text $T$, such $x$ is obtained in $\mathcal{O}(1)$ time. We then look up $(s, p) = L_{\text{root}}[x]$, and in $\mathcal{O}(1)$ time obtain L-root$(j) = T[j+s \mathinner{.\,.} j+s+p)$ and L-head$(j) = s$. Next, we compute L-exp$(j)$ and L-tail$(j)$. For this we recall that by Lemma 5.10(2), it holds $e(j) = j + p + \text{LCE}(j, j + p)$. Thus, given $j$ and $p$, we can compute $e(j)$ in $\mathcal{O}(1)$ time. We then obtain L-exp$(j) = \lfloor \frac{e(j)-j-s}{p} \rfloor$ and L-tail$(j) = (e(j) - j - s) \bmod p$. Finally, to test if type$(j) = +1$, we check whether $e(j) \leq n$, and if so, whether $T[e(j)] \succ T[e(j) - p]$. $\qquad\square$

**Proposition 5.16.** *Let $i \in [1 \mathinner{.\,.} n]$ be such that $\text{SA}[i] \in \mathsf{R}$. Given the data structure from Section 5.3.2 and the index $i$, in $\mathcal{O}(1)$ time we can compute* L-root$(\text{SA}[i])$ *and* L-head$(\text{SA}[i])$.

*Proof.* We first compute $y = \mathsf{rank}_{B_{3\tau-1},1}(i - 1)$. The string $X = A_{\text{short}}[y + 1]$ is then a prefix of $T[\text{SA}[i] \mathinner{.\,.} n]$ of length $3\tau - 1$. Let $x = \text{int}(X)$. We then look up $(s, p) = L_{\text{root}}[x]$, and in $\mathcal{O}(1)$ time obtain L-root$(\text{SA}[i]) = X[1+s \mathinner{.\,.} 1+s+p)$ and L-head$(\text{SA}[i]) = s$. $\qquad\square$

### 5.3.4 Implementation of ISA Queries

For any $j \in \mathsf{R}$, we define

$$\text{Pos}(j) = \{j' \in [1 \mathinner{.\,.} n] : \text{LCE}(j, j') \geq 3\tau - 1 \text{ and } T[j' \mathinner{.\,.} n] \preceq T[j \mathinner{.\,.} n]\},$$

and denote $\delta(j) = |\text{Pos}(j)|$.

**Lemma 5.17.** *Let $j \in \mathsf{R}$ and $X = T[j \mathinner{.\,.} j + 3\tau - 1)$. Then,* $\text{ISA}[j] = \text{RangeBeg}(X, T) + \delta(j)$.

*Proof.* It suffices to observe that $j' \in \mathrm{Occ}(X, T)$ holds if and only if $\mathrm{LCE}(j, j') \geq 3\tau - 1$. Thus, it holds by definition of $\mathrm{ISA}[j]$ that $\mathrm{ISA}[j] = \mathrm{RangeBeg}(X, T) + |\{j' \in \mathrm{Occ}(X, T) : T[j'\mathinner{.\,.} n] \preceq T[j \mathinner{.\,.} n]\}| = \mathrm{RangeBeg}(X, T) + |\{j' \in [1 \mathinner{.\,.} n] : \mathrm{LCE}(j, j') \geq 3\tau - 1 \text{ and } T[j'\mathinner{.\,.} n] \preceq T[j \mathinner{.\,.} n]\}| = \mathrm{RangeBeg}(X, T) + \delta(j)$. $\qquad\square$

We focus on computing $\delta(j)$ for $j \in \mathsf{R}^-$. The elements of $\mathsf{R}^+$ are processed symmetrically (see the proof of Proposition 5.23). For any $H \in \mathrm{Roots}$, $s \in [0 \mathinner{.\,.} |H|)$, and $j \in \mathsf{R}^-_{s,H}$, we define $\mathrm{Pos}^{\mathsf{a}}(j) = \{j' \in \mathsf{R}^-_{s,H} : \mathrm{L\text{-}exp}(j') \leq \mathrm{L\text{-}exp}(j)\}$ and $\mathrm{Pos}^{\mathsf{s}}(j) = \{j' \in \mathsf{R}^-_{s,H} : \mathrm{L\text{-}exp}(j') = \mathrm{L\text{-}exp}(j) \text{ and } T[j' \mathinner{.\,.} n] \succ T[j \mathinner{.\,.} n]\}$. For any $j \in \mathsf{R}^-$, we denote $\delta^{\mathsf{a}}(j) = |\mathrm{Pos}^{\mathsf{a}}(j)|$ and $\delta^{\mathsf{s}}(j) = |\mathrm{Pos}^{\mathsf{s}}(j)|$.

**Lemma 5.18.** *For any $j \in \mathsf{R}^-$, it holds $\delta(j) = \delta^{\mathsf{a}}(j) - \delta^{\mathsf{s}}(j)$.*

*Proof.* We will prove that $\mathrm{Pos}^{\mathsf{a}}(j)$ is a disjoint union of $\mathrm{Pos}(j)$ and $\mathrm{Pos}^{\mathsf{s}}(j)$. This implies $\delta(j) + \delta^{\mathsf{s}}(j) = \delta^{\mathsf{a}}(j)$, and consequently, the equality in the claim.

By Lemma 5.11, letting $j \in \mathsf{R}^-_{s,H}$, we have $\mathrm{Pos}(j) = \{j' \in \mathsf{R}^-_{s,H} : T[j' \mathinner{.\,.} n] \preceq T[j \mathinner{.\,.} n]\}$, and moreover, if $j' \in \mathrm{Pos}(j)$, then $e(j') - j' \leq e(j) - j$. In particular, $\mathrm{L\text{-}exp}(j') = \lfloor \frac{e(j') - j' - s}{|H|} \rfloor \leq \lfloor \frac{e(j) - j - s}{|H|} \rfloor = \mathrm{L\text{-}exp}(j)$. Hence, $\mathrm{Pos}(j) \subseteq \mathrm{Pos}^{\mathsf{a}}(j)$. On the other hand, clearly $\mathrm{Pos}^{\mathsf{s}}(j) \subseteq \mathrm{Pos}^{\mathsf{a}}(j)$ and $\mathrm{Pos}^{\mathsf{s}}(j) \cap \mathrm{Pos}(j) = \emptyset$. Thus, to obtain the claim, it suffices to show that $\mathrm{Pos}^{\mathsf{a}}(j) \setminus \mathrm{Pos}^{\mathsf{s}}(j) \subseteq \mathrm{Pos}(j)$.

Let $j' \in \mathrm{Pos}^{\mathsf{a}}(j) \setminus \mathrm{Pos}^{\mathsf{s}}(j)$. Consider two cases. If $\mathrm{L\text{-}exp}(j') = \mathrm{L\text{-}exp}(j)$, then by definition of $\mathrm{Pos}^{\mathsf{s}}(j)$, it must hold $T[j' \mathinner{.\,.} n] \preceq T[j \mathinner{.\,.} n]$. Thus, we have $j' \in \mathrm{Pos}(j)$. Let us therefore assume $\mathrm{L\text{-}exp}(j') < \mathrm{L\text{-}exp}(j)$. Then, $e(j') - j' = s + \mathrm{L\text{-}exp}(j') \cdot |H| + \mathrm{L\text{-}tail}(j') < s + \mathrm{L\text{-}exp}(j') \cdot |H| + |H| \leq s + \mathrm{L\text{-}exp}(j) \cdot |H| \leq s + \mathrm{L\text{-}exp}(j) \cdot |H| + \mathrm{L\text{-}tail}(j) = e(j) - j$. By Lemma 5.11(2), this implies $T[j' \mathinner{.\,.} n] \prec T[j \mathinner{.\,.} n]$, and consequently, $j' \in \mathrm{Pos}(j)$. $\qquad\square$

**Computing $\delta^{\mathsf{a}}(j)$** We now describe the algorithm to compute $\delta^{\mathsf{a}}(j)$ for $j \in \mathsf{R}^-$.

**Proposition 5.19.** *Given the data structure from Section 5.3.2 and any $j \in \mathsf{R}^-$, in $\mathcal{O}(1)$ time we can compute $\delta^{\mathsf{a}}(j)$.*

*Proof.* Let $X = T[j \mathinner{.\,.} j + 3\tau - 1)$. First, using the lookup table $L_{\mathrm{range}}$, we compute $(b_X, e_X) = (\mathrm{RangeBeg}(X, T), \mathrm{RangeEnd}(X, T))$. Then, by Lemma 5.11, $\mathrm{SA}(b_X \mathinner{.\,.} e_X]$ contains all positions from $\mathsf{R}_{s,H}$, where $H = \mathrm{L\text{-}root}(j)$ and $s = \mathrm{L\text{-}head}(j)$. Next, using Proposition 5.15, we compute in $\mathcal{O}(1)$ time the value $k = \mathrm{L\text{-}exp}(j)$. Finally, we retrieve $k_{\min} = L_{\mathrm{minexp}}[\mathrm{int}(X)]$. Observe now that by Lemma 5.11, all positions in $\mathsf{R}^-_{s,H}$ occur in $\mathrm{SA}(b_X \mathinner{.\,.} e_X]$ before $\mathsf{R}^+_{s,H}$. Furthermore, by Lemma 5.11(2), $[k_{\min} \mathinner{.\,.} k] \subseteq \{\mathrm{L\text{-}exp}(j') : j' \in \mathsf{R}^-_{s,H}\}$ (for $k' \in (k_{\min} \mathinner{.\,.} k]$, we can take $j' = j + (k - k')|H|$). Thus, by the definition of $B_{\mathrm{exp}}$, we can finally return $\delta^{\mathsf{a}}(j) = \mathsf{select}_{B_{\mathrm{exp}}, 1}(\mathsf{rank}_{B_{\mathrm{exp}}, 1}(b_X) + (k - k_{\min}) + 1) - b_X$ in $\mathcal{O}(1)$ time. $\qquad\square$

**Computing $\delta^{\mathsf{s}}(j)$** Next, we describe the algorithm to compute $\delta^{\mathsf{s}}(j)$ for any position $j \in \mathsf{R}^-$.

**Lemma 5.20.** *Assume $i, j \in \mathsf{R}^-_H$ and let $\ell = e(i) - i - 3\tau + 2$. Then $|\mathrm{Pos}^{\mathsf{s}}(j) \cap [i \mathinner{.\,.} i + \ell]| \leq 1$. Moreover, $|\mathrm{Pos}^{\mathsf{s}}(j) \cap [i \mathinner{.\,.} i + \ell]| = 1$ if and only if $T[e^{\mathrm{full}}(i) \mathinner{.\,.} n] \succ T[e^{\mathrm{full}}(j) \mathinner{.\,.} n]$ and $e^{\mathrm{full}}(i) - i \geq e^{\mathrm{full}}(j) - j$.*

*Proof.* By Lemma 5.12, we have $[i \mathinner{.\,.} i + \ell] \subseteq \mathsf{R}^-_H$ with $e(i + \delta) = e(i)$ for every $\delta \in [0 \mathinner{.\,.} \ell]$. Moreover, by the uniqueness of L-decomposition, $\mathrm{L\text{-}tail}(i + \delta) = \mathrm{L\text{-}tail}(i)$. Together, these imply that $e^{\mathrm{full}}(i + \delta) = e^{\mathrm{full}}(i)$, and consequently $e^{\mathrm{full}}(i + \delta) - (i + \delta) = e^{\mathrm{full}}(i) - i - \delta$. It remains to observe that, letting $j \in \mathsf{R}^-_{s,H}$, for $j' \in \mathrm{Pos}^{\mathsf{s}}(j)$ it holds $e^{\mathrm{full}}(j') - j' = s + \mathrm{L\text{-}exp}(j') \cdot |H| = s + \mathrm{L\text{-}exp}(j) \cdot |H| = e^{\mathrm{full}}(j) - j$. Thus, $i + \delta \in \mathrm{Pos}^{\mathsf{s}}(j)$ implies $e^{\mathrm{full}}(i + \delta) - (i + \delta) = e^{\mathrm{full}}(i) - (i + \delta) = e^{\mathrm{full}}(j) - j$, or equivalently, $\delta = (e^{\mathrm{full}}(i) - i) - (e^{\mathrm{full}}(j) - j)$, and therefore, $|\mathrm{Pos}^{\mathsf{s}}(j) \cap [i \mathinner{.\,.} i + \ell]| \leq 1$.

For the second part, assume first that $i+\delta \in \mathrm{Pos}^{\mathsf{s}}(j)$ holds for some $\delta \in [0\mathinner{.\,.}\ell)$. Then, as noted above, we have $e^{\mathrm{full}}(j) - j = e^{\mathrm{full}}(i) - (i+\delta) \leq e^{\mathrm{full}}(i) - i$. Moreover, letting $j \in \mathsf{R}^{-}_{s,H}$, by definition of $\mathrm{Pos}^{\mathsf{s}}(j)$, we have $i+\delta \in \mathsf{R}^{-}_{s,H}$, $\mathrm{L\text{-}exp}(j) = \mathrm{L\text{-}exp}(i+\delta)$, and $T[i+\delta\mathinner{.\,.}n] \succ T[j\mathinner{.\,.}n]$. Therefore, we obtain that $T[i+\delta\mathinner{.\,.}e^{\mathrm{full}}(i+\delta)) = T[i+\delta\mathinner{.\,.}e^{\mathrm{full}}(i)) = T[j\mathinner{.\,.}e^{\mathrm{full}}(j)) = H'H^k$ (where $k = \mathrm{L\text{-}exp}(j)$ and $H'$ is the length-$s$ suffix of $H$), and consequently, $T[e^{\mathrm{full}}(i)\mathinner{.\,.}n] \succ T[e^{\mathrm{full}}(j)\mathinner{.\,.}n]$. To show the converse implication, assume $T[e^{\mathrm{full}}(i)\mathinner{.\,.}n] \succ T[e^{\mathrm{full}}(j)\mathinner{.\,.}n]$ and $e^{\mathrm{full}}(i) - i \geq e^{\mathrm{full}}(j) - j$. Let $\delta = (e^{\mathrm{full}}(i) - i) - (e^{\mathrm{full}}(j) - j)$. We will prove that $\delta \in [0\mathinner{.\,.}\ell)$ and $i + \delta \in \mathrm{Pos}^{\mathsf{s}}(j)$. Clearly $\delta \geq 0$. To show $\delta < \ell$, we first prove $e(i) - e^{\mathrm{full}}(i) \geq e(j) - e^{\mathrm{full}}(j)$. Suppose that $q = e(i) - e^{\mathrm{full}}(i) < e(j) - e^{\mathrm{full}}(j)$. By $i \in \mathsf{R}^{-}_{H}$, we then either have $e^{\mathrm{full}}(i) + q = n + 1$, or $e^{\mathrm{full}}(i) + q \leq n$ and $T[e^{\mathrm{full}}(i) + q] \prec T[e^{\mathrm{full}}(i) + q - |H|] = T[e^{\mathrm{full}}(j) + q - |H|] = T[e^{\mathrm{full}}(j) + q]$, both of which contradict $T[e^{\mathrm{full}}(i)\mathinner{.\,.}n] \succ T[e^{\mathrm{full}}(j)\mathinner{.\,.}n]$. Thus, $e(i) - e^{\mathrm{full}}(i) \geq e(j) - e^{\mathrm{full}}(j)$. This implies, $e(i) - (i+\delta) = (e^{\mathrm{full}}(i) - (i+\delta)) + (e(i) - e^{\mathrm{full}}(i)) = (e^{\mathrm{full}}(j) - j) + (e(i) - e^{\mathrm{full}}(i)) \geq (e^{\mathrm{full}}(j) - j) + (e(j) - e^{\mathrm{full}}(j)) = e(j) - j \geq 3\tau - 1$, or equivalently $\delta \leq e(i) - i - 3\tau + 1 < \ell$. To show $i + \delta \in \mathrm{Pos}^{\mathsf{s}}(j)$, it remains to observe that $e^{\mathrm{full}}(i+\delta) - (i+\delta) = e^{\mathrm{full}}(i) - (i+\delta) = e^{\mathrm{full}}(j) - j$ and $i+\delta, j \in \mathsf{R}^{-}_{H}$ imply $T[i+\delta\mathinner{.\,.}e^{\mathrm{full}}(i)) = T[j\mathinner{.\,.}e^{\mathrm{full}}(j))$. This in particular gives, letting $j \in \mathsf{R}_{s,H}$, that $i + \delta \in \mathsf{R}_{s,H}$ and $\mathrm{L\text{-}exp}(i+\delta) = \mathrm{L\text{-}exp}(j)$. Moreover, combining it with $T[e^{\mathrm{full}}(i)\mathinner{.\,.}n] \succ T[e^{\mathrm{full}}(j)\mathinner{.\,.}n]$ yields $T[i+\delta\mathinner{.\,.}n] \succ T[j\mathinner{.\,.}n]$. Finally, by Lemma 5.12, $\mathrm{type}(i+\delta) = \mathrm{type}(i) = -1$. Therefore, $i + \delta \in \mathrm{Pos}^{\mathsf{s}}(j)$. $\qquad\square$

**Proposition 5.21.** *Given the data structure from Section 5.3.2 and any $j \in \mathsf{R}^{-}$, in $\mathcal{O}(\log \log n)$ time we can compute $\delta^{\mathsf{s}}(j)$.*

*Proof.* Given $j \in \mathsf{R}^{-}$, we first compute $H = \mathrm{L\text{-}root}(j)$, $s = \mathrm{L\text{-}head}(j)$, and $k = \mathrm{L\text{-}exp}(j)$. By Proposition 5.15, this takes $\mathcal{O}(1)$ time. This lets us deduce $e^{\mathrm{full}}(j) = j + s + k|H|$. Then, we compute $i \in [1\mathinner{.\,.}q]$ satisfying $j \in [r^{\mathrm{text}}_i\mathinner{.\,.}e(r^{\mathrm{text}}_i) - 3\tau + 2)$, i.e., $j$ is in the maximal block of positions from $\mathsf{R}^{-}$ starting at position $r^{\mathrm{text}}_i$. Using $B_{\mathsf{R}'}$ we obtain $i = \mathrm{rank}_{B_{\mathsf{R}'},1}(j)$ in $\mathcal{O}(1)$ time. Observe now that, letting $j' = r^{\mathrm{text}}_i$, by $e^{\mathrm{full}}(j') = e^{\mathrm{full}}(j)$, we have $T[e^{\mathrm{full}}(j')\mathinner{.\,.}n] = T[e^{\mathrm{full}}(j)\mathinner{.\,.}n]$. Therefore, letting $x = A_{\mathrm{rmap}}[i]$ and $x' = \sum_{H' \preceq H} |\mathsf{R}'^{-}_{H'}|$ (obtained in $\mathcal{O}(1)$ time using $L_{\mathrm{runs}}$), by Lemma 5.20 we have $\delta^{\mathsf{s}}(j) = |\{i' \in (x\mathinner{.\,.}x'] : \bar{\ell}_{i'} \geq e^{\mathrm{full}}(j) - j\}| = \mathrm{rcount}_{A_{\mathrm{len}}}(e^{\mathrm{full}}(j) - j, x') - \mathrm{rcount}_{A_{\mathrm{len}}}(e^{\mathrm{full}}(j) - j, x)$, which we compute in $\mathcal{O}(\log \log n)$ time using the data structure from Proposition 2.3. $\qquad\square$

*Remark* 5.22. In Lemma 5.18, we presented an equation relating the sizes of $\mathrm{Pos}^{\mathsf{a}}(j)$ and $\mathrm{Pos}^{\mathsf{s}}(j)$, and the size of $\mathrm{Pos}(j)$, where $j \in \mathsf{R}^{-}$. In this formula, some positions are first counted as part of $\mathrm{Pos}^{\mathsf{a}}(j)$, and then canceled when subtracting the size of $\mathrm{Pos}^{\mathsf{s}}(j)$. To see the reason for this counterintuitive formula, let $J := \{j' \in \mathsf{R}^{-}_{s,H} : \mathrm{L\text{-}exp}(j') = \mathrm{L\text{-}exp}(j)\}$, where $s = \mathrm{L\text{-}head}(j)$ and $H = \mathrm{L\text{-}root}(j)$, and consider the problem of computing the size of $J' = \{j' \in J : T[j'\mathinner{.\,.}n] \succeq T[j\mathinner{.\,.}n]\}$. As shown in Lemma 5.20, to count such positions, it suffices to first align all $j'' \in \mathsf{R}'^{-}$ by the position $e^{\mathrm{full}}(j'')$, and then count those $j''$ that satisfy (1) $T[e^{\mathrm{full}}(j'')\mathinner{.\,.}n] \succeq T[e^{\mathrm{full}}(j)\mathinner{.\,.}n]$, and (2) $e^{\mathrm{full}}(j'') - j'' \geq e^{\mathrm{full}}(j) - j$. For every $j'' \in \mathsf{R}'^{-}$ satisfying these conditions, there exists exactly one $j' \in \mathsf{R}^{-}_{s,H}$ such that $[j''\mathinner{.\,.}j'] \subseteq \mathsf{R}$, $\mathrm{L\text{-}exp}(j') = \mathrm{L\text{-}exp}(j)$ and $T[j'\mathinner{.\,.}n] \succeq T[j\mathinner{.\,.}n]$, because for $j, j'' \in \mathsf{R}^{-}$, $T[e^{\mathrm{full}}(j'')\mathinner{.\,.}n] \succeq T[e^{\mathrm{full}}(j)\mathinner{.\,.}n]$ implies $e(j'') - e^{\mathrm{full}}(j'') \geq e(j) - e^{\mathrm{full}}(j)$ (Lemma 5.11). Thus, letting $\ell = e^{\mathrm{full}}(j) - j$, such $j'$ is given by $j' = e^{\mathrm{full}}(j'') - \ell$. In particular, such $j'$ satisfies $j' \in \mathsf{R}$ because $(e(j'') - e^{\mathrm{full}}(j'')) + \ell \geq e(j) - e^{\mathrm{full}}(j) + \ell = e(j) - j \geq 3\tau - 1$.

Consider now the problem of computing the size of $J'' = \{j' \in J : T[j'\mathinner{.\,.}n] \prec T[j\mathinner{.\,.}n]\}$ (defining $\mathrm{Pos}^{\mathsf{s}}(j)$ as $J''$ may seem like a simpler alternative to the current definition). Observe that the above method does not work for this problem. The reason for this is that position $j'' \in \mathsf{R}'^{-}$ satisfying $T[e^{\mathrm{full}}(j')\mathinner{.\,.}n] \prec T[e^{\mathrm{full}}(j)\mathinner{.\,.}n]$ does not necessarily imply that $e^{\mathrm{full}}(j'') - \ell \in \mathsf{R}$. This is because we may have $e(j'') - e^{\mathrm{full}}(j'') < e(j) - e^{\mathrm{full}}(j)$, which implies that it is possible that $(e(j'') - e^{\mathrm{full}}(j'')) + \ell < 3\tau - 1$. This motivates the current definition of $\mathrm{Pos}^{\mathsf{s}}(j)$.

**Summary** By combining all above results, we obtain the following algorithm to compute ISA$[j]$ for periodic positions.

**Proposition 5.23.** *Given the data structure from Section 5.3.2 and any $j \in \mathsf{R}$, in $\mathcal{O}(\log \log n)$ time we can compute* ISA$[j]$.

*Proof.* First, in $\mathcal{O}(1)$ time we compute $x = \mathrm{int}(X)$, where $X = T[j \mathbin{..} j + 3\tau - 1]$. In $\mathcal{O}(1)$ we then look up $(b_X, e_X) = L_{\mathrm{range}}[x]$. In particular, we have $b_X = \mathrm{RangeBeg}(X, T)$. Then, using Proposition 5.15 we determine type$(j)$. Depending on whether $j \in \mathsf{R}^-$ or $j \in \mathsf{R}^+$ we use either a combination of Propositions 5.19 and 5.21, or their symmetric counterparts (more precisely, if $j \in \mathsf{R}^+$, letting $s = \text{L-head}(j)$ and $H = \text{L-root}(j)$, we have $\delta^{\mathsf{a}}(j) = |\mathrm{Pos}^{\mathsf{a}}(j)|$ and $\delta^{\mathsf{s}}(j) = |\mathrm{Pos}^{\mathsf{s}}(j)|$, where $\mathrm{Pos}^{\mathsf{a}}(j) = \{j' \in \mathsf{R}^+_{s,H} : \text{L-exp}(j') \leq \text{L-exp}(j)\}$ and $\mathrm{Pos}^{\mathsf{s}}(j) = \{j' \in \mathsf{R}^+_{s,H} : \text{L-exp}(j) = \text{L-exp}(j)$ and $T[j' \mathbin{..} n] \prec T[j \mathbin{..} n]\}$), to compute $\delta^{\mathsf{a}}(j)$ and $\delta^{\mathsf{s}}(j)$ in $\mathcal{O}(1)$ and $\mathcal{O}(\log \log n)$ time, respectively. If $j \in \mathsf{R}^-$, then by Lemma 5.18 we have $\delta(j) = \delta^{\mathsf{a}}(j) - \delta^{\mathsf{s}}(j)$. Otherwise, by the counterpart of Lemma 5.18, $\delta(j) = (e_X - b_X) - (\delta^{\mathsf{a}}(j) - \delta^{\mathsf{s}}(j))$. Finally, we return ISA$[j] = b_X + \delta(j)$ as the answer. In total, the query takes $\mathcal{O}(\log \log n)$ time. $\qquad\square$

### 5.3.5 Implementation of SA Queries

We focus on positions $i \in [1 \mathbin{..} n]$ satisfying SA$[i] \in \mathsf{R}^-$. Positions satisfying SA$[i] \in \mathsf{R}^+$ are processed symmetrically (see the proof of Proposition 5.26). The algorithm to query SA$[i]$ for $i \in [1 \mathbin{..} n]$ satisfying SA$[i] \in \mathsf{R}^-$ proceeds in two steps. First, we compute L-exp(SA$[i]$) and $\delta^{\mathsf{s}}(\text{SA}[i])$. In the second steps, these values are used to compute SA$[i]$.

**Computing L-exp(SA$[i]$) and $\delta^{\mathsf{s}}(\text{SA}[i])$** We now describe the first step during the computation of SA$[i]$ for $i \in [1 \mathbin{..} n]$ satisfying SA$[i] \in \mathsf{R}$.

**Proposition 5.24.** *Let $i \in [1 \mathbin{..} n]$ be such that SA$[i] \in \mathsf{R}$. Given the data structure from Section 5.3.2 and the index $i$, in $\mathcal{O}(1)$ time we can check if* type(SA$[i]$) $= -1$, *and if so, return* L-exp(SA$[i]$) *and* $\delta^{\mathsf{s}}(\text{SA}[i])$.

*Proof.* To check if type(SA$[i]$) $= -1$, we first compute $y = \mathsf{rank}_{B_{3\tau-1},1}(i - 1)$. The string $X = A_{\mathrm{short}}[y + 1]$ is then a prefix of $T[\text{SA}[i] \mathbin{..} n]$ of length $3\tau - 1$. Let $x = \mathrm{int}(X)$. In $\mathcal{O}(1)$ time we then look up $(b_X, e_X) = L_{\mathrm{range}}[x]$. By Lemma 5.11 we then have type(SA$[i]$) $= -1$ if and only if $B_{\mathrm{exp}}[i \mathbin{..} e_X]$ contains a bit with value 1. This can be checked in $\mathcal{O}(1)$ time by checking if $\mathsf{rank}_{B_{\mathrm{exp}},1}(e_X) > \mathsf{rank}_{B_{\mathrm{exp}},1}(i - 1)$. Let us assume type(SA$[i]$) $= -1$. To compute L-exp(SA$[i]$), we first in $\mathcal{O}(1)$ retrieve $k_{\min} = L_{\mathrm{minexp}}[x]$, and then compute L-exp(SA$[i]$) $= k_{\min} + (\mathsf{rank}_{B_{\mathrm{exp}},1}(i - 1) - \mathsf{rank}_{B_{\mathrm{exp}},1}(b_X))$. Then, $\delta^{\mathsf{a}}(\text{SA}[i])$ can be computed in $\mathcal{O}(1)$ time as $\delta^{\mathsf{a}}(\text{SA}[i]) = \mathsf{select}_{B_{\mathrm{exp}},1}(\mathsf{rank}_{B_{\mathrm{exp}},1}(i - 1) + 1) - b_X$. Finally, by applying Lemma 5.17 and Lemma 5.18 for $j = \text{SA}[i]$, it holds $i - b_X = \delta^{\mathsf{a}}(\text{SA}[i]) - \delta^{\mathsf{s}}(\text{SA}[i])$. Thus, we obtain $\delta^{\mathsf{s}}(\text{SA}[i]) = b_X + \delta^{\mathsf{a}}(\text{SA}[i]) - i$. $\qquad\square$

**Computing SA$[i]$** We now describe the algorithm to complete the computation of SA$[i]$ for any $i \in [1 \mathbin{..} n]$ such that SA$[i] \in \mathsf{R}^-$.

**Proposition 5.25.** *In $\mathcal{O}(n / \log_\sigma n)$ time, we can augment the structure of Proposition 5.16 so that, given any $i \in [1 \mathbin{..} n]$ such that SA$[i] \in \mathsf{R}^-$, along with* L-exp(SA$[i]$) *and* $\delta^{\mathsf{s}}(\text{SA}[i])$, *we can compute* SA$[i]$ *in $\mathcal{O}(\log \log n)$ time.*

*Proof.* First, we compute $H = \text{L-root}(\text{SA}[i])$ and L-head(SA$[i]$) in $\mathcal{O}(1)$ time using Proposition 5.16. This lets us deduce that $e^{\mathrm{full}}(\text{SA}[i]) - \text{SA}[i] = \ell$, where $\ell = \text{L-head}(\text{SA}[i]) + \text{L-exp}(\text{SA}[i])|H|$. Let $x = \sum_{H' \preceq H} |\mathsf{R}'^-_{H'}|$ (obtained using $L_{\mathrm{runs}}$ in $\mathcal{O}(1)$ time). Next, we compute

$\delta = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, x)$. Using the structure from Proposition 2.3, this takes $\mathcal{O}(\log\log n)$ time. Let $k = \delta - \delta^{\mathsf{s}}(\mathrm{SA}[i])$. We then compute the position $p \in [1\mathbin{..}q]$ of the $k$th leftmost element in $A_{\mathrm{len}}$ that is greater or equal than $\ell$. Using Proposition 2.3, we compute $p = \mathsf{rselect}_{A_{\mathrm{len}}}(\ell, k)$ in $\mathcal{O}(1)$ time. By Lemma 5.14 and Lemma 5.20, we then have $e^{\mathrm{full}}(r_p^{\mathrm{lex}}) = e^{\mathrm{full}}(\mathrm{SA}[i])$. By combining Lemma 5.12 and Lemma 5.14, for any $j', j'' \in \mathsf{R}$ such that $j' < j''$ and $e^{\mathrm{full}}(j') = e^{\mathrm{full}}(j'')$, it holds $[j'\mathbin{..}j''] \subseteq \mathsf{R}$, i.e., $j'$ and $j''$ must belong to the same contiguous block of positions from $\mathsf{R}$. Since $r_p^{\mathrm{lex}} \in \mathsf{R}'$, we thus have $\mathrm{SA}[i] \in [r_p^{\mathrm{lex}}\mathbin{..}e(r_p^{\mathrm{lex}}) - 3\tau + 2) \subseteq \mathsf{R}_H^-$. In $\mathcal{O}(1)$ time we obtain $p' = A_{\mathrm{rmap}}^{-1}[p]$ and $j := \mathsf{select}_{B_{R'},1}(p') = r_p^{\mathrm{lex}}$. Observe now that in the block $[j\mathbin{..}e(j) - 3\tau + 2)$ there is at most one element with given values of L-exp and L-head, and we already have values L-exp($\mathrm{SA}[i]$) and L-head($\mathrm{SA}[i]$). We thus proceed as follows. First, we compute $e(j)$. For this, we recall that by Lemma 5.10(2), it holds $e(j) = j + |H| + \mathrm{LCE}(j, j + |H|)$. Thus, given $j$ and $|H|$, we can compute $e(j)$ in $\mathcal{O}(1)$ time. We then in $\mathcal{O}(1)$ time compute $s = \mathrm{L\text{-}head}(j)$ using the lookup table $L_{\mathrm{root}}$. This lets us determine $e^{\mathrm{full}}(j) = e(j) - ((e(j) - j - s) \bmod |H|)$. In $\mathcal{O}(1)$ time we then obtain $\mathrm{SA}[i] = e^{\mathrm{full}}(j) - \mathrm{L\text{-}head}(\mathrm{SA}[i]) - \mathrm{L\text{-}exp}(\mathrm{SA}[i])|H|$. In total, the query takes $\mathcal{O}(\log\log n)$ time. $\qquad\square$

**Summary**  By combining all above results, we obtain the following algorithm to compute $\mathrm{SA}[i]$ for periodic positions.

**Proposition 5.26.** *Let $i \in [1\mathbin{..}n]$ be such that $\mathrm{SA}[i] \in \mathsf{R}$. Given the data structure from Section 5.3.2 and the index $i$, in $\mathcal{O}(\log\log n)$ time we can compute $\mathrm{SA}[i]$.*

*Proof.* First, using Proposition 5.24, in $\mathcal{O}(1)$ time we compute type($\mathrm{SA}[i]$). Depending on whether $\mathrm{SA}[i] \in \mathsf{R}^-$ or $\mathrm{SA}[i] \in \mathsf{R}^+$, we use either a combination of Propositions 5.24 and 5.25 or their symmetric counterparts (see the proof of Proposition 5.23), to first compute L-exp($\mathrm{SA}[i]$) and $\delta^{\mathsf{s}}(\mathrm{SA}[i])$ in $\mathcal{O}(1)$ time, and then $\mathrm{SA}[i]$ in $\mathcal{O}(\log\log n)$ time. $\qquad\square$

### 5.3.6 Construction Algorithm

**Proposition 5.27.** *Given $\mathrm{C_{SA}}(T)$, we can in $\mathcal{O}(n/\log_\sigma n)$ time augment it into a data structure from Section 5.3.2.*

*Proof.* Due to a large number of components, as well as dependency of some components on others, we present the description in separate paragraphs, in the order in which it occurs.

*Construction of $L_{\mathrm{root}}$*  To compute $L_{\mathrm{root}}$, we observe that, given $X \in [0\mathbin{..}\sigma)^{3\tau - 1}$, we can check in $\mathcal{O}(\tau^2)$ time if $\mathrm{per}(X) \le \frac{1}{3}\tau$, and if so, determine the value $L_{\mathrm{root}}[\mathrm{int}(X)] = (s, p)$. To compute $\mathrm{per}(X)$, we try all $\ell \in [1\mathbin{..}\lfloor\frac{\tau}{3}\rfloor]$ until we find that $\ell$ is a period of $X$, or that there is no such $\ell$. Assuming $p := \mathrm{per}(X) \le \frac{1}{3}\tau$, finding $s \in [0\mathbin{..}p)$ satisfying $X[1+s\mathbin{..}1+s+p] = \min\{X[t\mathbin{..}t+p] : t \in [1\mathbin{..}p]\}$ also takes $\mathcal{O}(\tau^2)$ time. Initializing $L_{\mathrm{root}}$ takes $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n/\log_\sigma n)$. Over all $X \in [0\mathbin{..}\sigma)^{3\tau - 1}$, we spend $\mathcal{O}(\sigma^{3\tau - 1}\tau^2) = \mathcal{O}(n^{1/2}\log^2 n) = \mathcal{O}(n/\log_\sigma n)$ time.

*Construction of the structure for* LCE *queries*  By [50, Theorem 5.4], the data structure for LCE queries on $T$ can be constructed from the packed representation of $T$ in $\mathcal{O}(n/\log_\sigma n)$ time.

*Construction of $B_{\mathrm{exp}}$*  To simplify the notation, for the duration of this proof, we denote $E := B_{\mathrm{exp}}$. We use the following definitions. For any $H \in \mathrm{Roots}$ and $s \in [0\mathbin{..}|H|)$, let $E_{s,H}^-$ denote the block of $E$ corresponding to suffixes starting in $\mathsf{R}_{s,H}^-$, i.e., $E_{s,H}^- = E(b\mathbin{..}e]$, where $(b\mathbin{..}e] \subseteq [1\mathbin{..}n]$ is such that $\mathsf{R}_{s,H}^- = \{\mathrm{SA}[i] : i \in (b\mathbin{..}e]\}$ (such $(b\mathbin{..}e]$ exists by Lemma 5.11(1)). Finally, let $\mathrm{unary}(x) := \mathtt{0}^x\mathtt{1}$ denote the unary encoding of an integer $x \ge 0$, and let $\mathrm{unary}^+(x)$ be $\mathrm{unary}(x)$ with the first symbol removed (in particular, $\mathrm{unary}^+(0)$ is the empty string). If $(a_i)_{i \in [1\mathbin{..}k]}$ is a sequence of non-negative integers, we define $\mathrm{unary}((a_i)_{i \in [1\mathbin{..}k]}) := \bigodot_{i=1}^k \mathrm{unary}(a_i)$,

where $\odot$ denotes concatenation. Analogously, $\mathrm{unary}^+((a_i)_{i\in[1..k]}) := \odot_{i=1}^k \mathrm{unary}^+(a_i)$. The definitions of unary and $\mathrm{unary}^+$ are naturally extended to infinite sequences $(a_i)_{i\in[1..\infty)}$.

Let $\alpha < 1$ be a positive constant. We first show an algorithm that, given the set of positions $\mathsf{R}'^-_H$ (where $H \in \mathrm{Roots}$) as input, computes all bitvectors $E^-_{0,H}, \ldots, E^-_{|H|-1,H}$ in $\mathcal{O}(|\mathsf{R}'^-_H| + |\mathsf{R}^-_H|/\log n + n^\alpha)$ time. For any $s \in [0..|H|)$ and $k \geq 0$, denote $e_{s,k,H} = |\{j' \in \mathsf{R}^-_{s,H} : \text{L-exp}(j') = k\}|$. We start by observing that by Lemma 5.11(2), $E^-_{s,H} = \mathrm{unary}^+((e_{s,k,H})_{k\in[0..\infty)})$. The values $e_{s,k,H}$ can be efficiently determined based on the following observation. First, note that if $j \in \mathsf{R}'^-_H$, then $[j..e(j) - 3\tau + 2) \subseteq \mathsf{R}^-_H$, and $j - 1, e(j) - 3\tau + 2 \notin \mathsf{R}$, i.e., the block of positions in $\mathsf{R}^-_H$ is maximal. By Lemma 5.12, for any $j' \in [j..e(j) - 3\tau + 2)$, it holds $e(j') = e(j)$. Thus, for any $j' \in [j..e(j) - 3\tau + 2)$, we have $\text{L-exp}(j') = \lfloor\frac{e-j'}{|H|}\rfloor$ and $\text{L-head}(j') = (e - j') \bmod |H|$, where $e = e(j) - \text{L-tail}(j)$. With this in mind, for any $j \in \mathsf{R}'^-_H$, we let $\mathcal{I}_j = (3\tau - 2 - t .. s + k|H|]$, where $s = \text{L-head}(j)$, $k = \text{L-exp}(j)$, and $t = \text{L-tail}(j)$. By the above discussion, for any $s \in [0..|H|)$ and $k \geq 0$, we have $e_{s,k,H} = |\{j \in \mathsf{R}'^-_H : s + k|H| \in \mathcal{I}_j\}|$. The algorithm consists of three steps:

1. First, we compute the string $\mathrm{unary}((e_{0,k,H})_{k=0}^{k_{\max}})$, where $k_{\max} = \max\{\text{L-exp}(j') : j' \in \mathsf{R}^-_H\}$. We start by computing $k_{\max}$. For this we observe that $k_{\max} = \max\{\text{L-exp}(j') : j' \in \mathsf{R}'^-_H\}$. Thus, using Proposition 5.15, we can compute $k_{\max}$ in $\mathcal{O}(|\mathsf{R}'^-_H|)$ time. To compute $\mathrm{unary}((e_{0,k,H})_{k\in[0..k_{\max}]})$, we generate the sequence of "events" from $\mathsf{R}'^-_H$, sort them, and then output $\mathrm{unary}((e_{0,k,H})_{k\in[0..k_{\max}]})$ left-to-right. More precisely, let $m = |\mathsf{R}'^-_H|$, and let $(p_i, v_i)_{i\in[0..2m]}$ be a sequence containing the multiset $\{(0,0), (k_{\max} + 1, 0)\} \cup \{(\lceil\min\mathcal{I}_j/|H|\rceil, +1) : j \in \mathsf{R}'^-_H\} \cup \{(\lfloor\max\mathcal{I}_j/|H|\rfloor + 1, -1) : j \in \mathsf{R}'^-_H\}$ such that for any $i \in [1..2m]$, it holds $p_{i-1} \leq p_i$. To compute the sequence $(p_i, v_i)_{i\in[0..2m]}$, we observe that, given $j \in \mathsf{R}'^-_H$, we can compute $\mathcal{I}_j$ in $\mathcal{O}(1)$ time using Proposition 5.15. Thus, in $\mathcal{O}(m)$ time we can generate all pairs in the above multiset. We then sort the pairs by the first element. Using $\lceil 1/\alpha\rceil$-round radix sort, this takes $\mathcal{O}(m + n^\alpha)$ time. Consequently, we can compute $(p_i, v_i)_{i\in[0..2m]}$ in $\mathcal{O}(|\mathsf{R}'^-_H| + n^\alpha)$ time. Given the sequence $(p_i, v_i)_{i\in[0..2m]}$, we compute $\mathrm{unary}((e_{0,k,H})_{k\in[0..k_{\max}]})$ as follows. First, we initialize the output bitvector to the empty string and set $v = 0$. We then iterate through $i = 1, \ldots, 2m$. For every $i$, we first append $p_i - p_{i-1}$ copies of the string $\mathrm{unary}(v)$ to the output string. We then add $v_i$ to $v$. To efficiently append multiple copies of $\mathrm{unary}(v)$ to the output, we first precompute (in $\mathcal{O}(\log^2 n) = \mathcal{O}(n^\alpha)$ time) the prefix of length $\log n$ of the string $\mathrm{unary}(x)^\infty[1..)$ for every $x \in [0..\log n)$. This way, we can append $\mathrm{unary}(v)^\ell$ to the output in $\mathcal{O}(1 + (v+1)\ell/\log n)$ time. Consequently, the construction of $\mathrm{unary}((e_{0,k,H})_{k\in[0..k_{\max}]})$ takes $\mathcal{O}(|\mathsf{R}'^-_H| + |\mathrm{unary}((e_{0,k,H})_{k\in[0..k_{\max}]})|/\log n + n^\alpha) = \mathcal{O}(|\mathsf{R}'^-_H| + |E^-_{0,H}|/\log n + n^\alpha) = \mathcal{O}(|\mathsf{R}'^-_H| + |\mathsf{R}^-_H|/\log n + n^\alpha)$ time. To show the first upper bound, observe that $k_{\max} \leq |E^-_{0,H}| + \mathcal{O}(\tau/|H|)$. Thus, $|\mathrm{unary}((e_{0,k,H})_{k\in[0..k_{\max}]})| = |\mathrm{unary}^+((e_{0,k,H})_{k\in[0..\infty)})| + k_{\max} + 1 = |E^-_{0,H}| + k_{\max} + 1 \leq 2|E^-_{0,H}| + \mathcal{O}(\log n)$. The second upper bound follows by observing that $|E^-_{0,H}| + \cdots + |E^-_{|H|-1,H}| = |\mathsf{R}^-_H|$.

2. The second step of the algorithm is to compute the strings $\mathrm{unary}((e_{s,k,H})_{k\in[0..k_{\max}]})$ for $s \in [1..|H|)$. For any $s \in [1..|H|)$, let $(q_i^{(s)}, p_i^{(s)}, v_i^{(s)})_{i\in[0..m_s]}$ denote the sequence containing all the elements $(q, p, v)$ of the multiset $\{(q, 0, 0) : q \in [1..|H|)\} \cup \{(q, k_{\max} + 1, 0) : q \in [1..|H|)\} \cup \{(\min\mathcal{I}_j \bmod |H|, \lfloor\min\mathcal{I}_j/|H|\rfloor, +1) : j \in \mathsf{R}'^-_H\} \cup \{((\max\mathcal{I}_j + 1) \bmod |H|, \lfloor(\max\mathcal{I}_j + 1)/|H|\rfloor, -1) : j \in \mathsf{R}'^-_H\}$ that satisfy $q = s$, and for any $i \in [1..m_s]$, it holds $p_{i-1}^{(s)} \leq p_i^{(s)}$ (note that the elements of this multiset satisfying $q = 0$ are not included in any sequence). To compute the sequences $(q_i^{(s)}, p_i^{(s)}, v_i^{(s)})_{i\in[0..m_s]}$ for all $s \in [1..|H|)$, we first enumerate all triples in the above multiset. Using Proposition 5.15, this takes $\mathcal{O}(m)$ time. We then sort the triples lexicographically. Using $\lceil 1/\alpha\rceil$-round radix sort, this takes $\mathcal{O}(m + n^\alpha)$ time. This yields all sequences concatenated together. It is easy to discard unused elements, and to detect boundaries between lists with a single scan. Consequently, we can construct all sequences in $\mathcal{O}(|\mathsf{R}'^-_H| + n^\alpha)$ time. Given the above sequences, we can compute the strings $\mathrm{unary}((e_{s,k,H})_{k\in[0..k_{\max}]})$ for $s \in [1..|H|)$ as follows.

The algorithm computes the strings in the order of increasing $s$. More precisely, given the string $U := \mathrm{unary}((e_{s-1,k,H})_{k\in[0\,..\,k_{\max}]})$ and the sequence $(q_i^{(s)}, p_i^{(s)}, v_i^{(s)})_{i\in[0\,..\,m_s]}$ (where $s \in [1\,..\,|H|)$), we compute the string $\mathrm{unary}((e_{s,k,H})_{k\in[0\,..\,k_{\max}]})$ in $\mathcal{O}(m_s + |U|/\log n)$ time as follows. First, we initialize the output bitvector to the empty string, and set $v = 0$ and $y = 0$. We then iterate through $i = 1, \ldots, m_s$. For every $i$, we first check if $p_i^{(s)} > p_{i-1}^{(s)}$. If yes, we perform the following three steps. First, find the position $y'$ of the $p_i^{(s)}$th 1-bit in $U$. Second, append the substring $U(y\,..\,y']$ to the output, except we first prepend it with $v$ zeros (if $v \geq 0$) or discard its first $-v$ bits (if $v < 0$). Finally, we set $y = y'$ and $v = 0$. Then (regardless of whether $p_i^{(s)} > p_{i-1}^{(s)}$), we add $v_i^{(s)}$ to $v$. To efficiently compute $y'$ we observe that the arguments of the consecutive select queries are increasing. We can thus precompute in $\mathcal{O}(n^\alpha)$ time a lookup table such that the computation of $y'$ takes $\mathcal{O}(1 + (y' - y)/\log n)$ time (these lookup tables can be shared among algorithms for different $s$). Note that for any $s \in [0\,..\,|H|)$, we have $k_{\max} \leq |E_{s,H}^-| + \mathcal{O}(\tau/|H|)$. Thus, $|U| \leq 2|E_{s-1,H}^-| + \mathcal{O}(\log n)$, and hence the algorithm runs in $\mathcal{O}(m_s + |E_{s-1,H}^-|/\log n)$ time. Consequently, by $m_0 + \cdots + m_{|H|-1} \leq 2|\mathsf{R}_H'^-| + 2|H|$ and $|E_{0,H}^-| + \cdots + |E_{|H|-1,H}^-| = |\mathsf{R}_H^-|$, over all $s \in [1\,..\,|H|)$, we spend $\mathcal{O}(|\mathsf{R}_H'^-| + |\mathsf{R}_H^-|/\log n + n^\alpha)$ time.

3. The third and final step of the algorithm is to convert the string $\mathrm{unary}((e_{s,k,H})_{k\in[0\,..\,k_{\max}]})$ into $\mathrm{unary}^+((e_{s,k,H})_{k\in[0\,..\,k_{\max}]}) = E_{s,H}^-$ for every $s \in [0\,..\,|H|)$. Let us fix some $s \in [0\,..\,|H|)$. Observe that to implement the conversion, it suffices to remove the first bit, as well as every bit following a 1-bit in $\mathrm{unary}((e_{s,k,H})_{k\in[0\,..\,k_{\max}]})$. In the RAM model, such local operation is easy implemented in $\mathcal{O}(1 + |\mathrm{unary}((e_{s,k,H})_{k\in[0\,..\,k_{\max}]})|/\log n)$ time after a $\mathcal{O}(n^\alpha)$-time preprocessing (we do the preprocessing once for all $s \in [0\,..\,|H|)$). As observed above, $|\mathrm{unary}((e_{s,k,H})_{k\in[0\,..\,k_{\max}]})| \leq 2|E_{s,H}^-| + \mathcal{O}(\log n)$. Thus, the total time to perform the conversion for all $s$ is $\mathcal{O}(|R_H^-|/\log n + n^\alpha)$.

Using the above algorithm, we construct $E$ as follows. We start by computing the set $\{(\mathrm{int}(\text{L-root}(j)), j)\}_{j\in\mathsf{R}'^-}$. For this, observe that for every $\tau$-synchronizing set $\mathsf{P}$ of $T$, by the density condition (see also [50, Section 6.1.2]), $i \in \mathsf{R}'$ implies that either $i = 1$ or $i > 1$ and $i-1 \in \mathsf{P}$. In particular, $|\mathsf{R}'^-| \leq |\mathsf{R}'| \leq 1 + |\mathsf{P}|$. We thus proceed as follows. First, using [50, Theorem 8.11] in $\mathcal{O}(n/\log_\sigma n)$ time we construct any $\tau$-synchronizing set $\mathsf{P}$ of $T$ of size $\mathcal{O}(n/\tau)$. Then, using the above observation together with Proposition 5.15, we enumerate the set $\{(\mathrm{int}(\text{L-root}(j)), j)\}_{j\in\mathsf{R}'^-}$ in $\mathcal{O}(1 + |\mathsf{P}|) = \mathcal{O}(n/\log_\sigma n)$ time. We then discard $\mathsf{P}$. Using $\lceil 1/\alpha \rceil$-round radix sort we then sort in $\mathcal{O}(|\mathsf{R}'^-| + n^\alpha) = \mathcal{O}(n/\log_\sigma n + n^\alpha)$ time the set of pairs by the first coordinate. This yields the representation of sets $\mathsf{R}_H'^-$ for all $H \in \mathrm{Roots}$. For each $H \in \mathrm{Roots}$, we then use the above algorithm to compute bitvectors $E_{0,H}^-, \ldots, E_{|H|-1,H}^-$ in $\mathcal{O}(|\mathsf{R}_H'^-| + |\mathsf{R}_H^-|/\log n + n^\alpha)$ time. By $\mathrm{Roots} \subseteq [0\,..\,\sigma)^{\leq\tau}$, over all $H$, this takes $\mathcal{O}(|\mathsf{R}'^-| + |R^-|/\log n + n^{\alpha+\mu})$ time (recall that $\tau = \lfloor \mu \log_\sigma n \rfloor$ and $\mu < \frac{1}{6}$). Choosing $\alpha < 1 - \mu$ results in $\mathcal{O}(n/\log_\sigma n)$ total time. When bitvectors $E_{s,H}^-$ are computed for all $H \in \mathrm{Roots}$ and $s \in [0\,..\,|H|)$, we initialize $E$ to the string $0^n$ in $\mathcal{O}(n/\log n)$ time, and then "paste" all the non-empty bitvectors $E_{s,H}^-$ into their correct positions. Given $H \in \mathrm{Roots}$ and $s \in [0\,..\,|H|)$, we first compute in $\mathcal{O}(\log n)$ time the corresponding string $X \in [0\,..\,\sigma)^{3\tau-1}$, and then compute the position to paste $E_{s,H}^-$ using the lookup table $L_{\mathrm{range}}$. Over all $H \in \mathrm{Roots}$ and $s \in [0\,..\,|H|)$, this takes $\mathcal{O}(n^\mu \log^2 n + n/\log n) = \mathcal{O}(n/\log_\sigma n)$ time. Thus, altogether, constructing $E$ and augmenting it using Theorem 2.1 takes $\mathcal{O}(n/\log_\sigma n)$ time.

*Construction of $L_{\mathrm{minexp}}$* Observe that in the above algorithm, if $i$ is the position of the leftmost 0-bit in $\mathrm{unary}((e_{s,k,H})_{k\in[0\,..\,k_{\max}]})$, then $\min\{\text{L-exp}(j) : j \in \mathsf{R}_{s,H}^-\} = i - 1$. Given the packed representation of $\mathrm{unary}((e_{s,k,H})_{k\in[0\,..\,k_{\max}]})$, the position $i$ can be easily found in $\mathcal{O}(1 + |\mathrm{unary}((e_{s,k,H})_{k\in[0\,..\,k_{\max}]})|/\log n)$ time. Thus, accounting for the computation of $X \in [0\,..\,\sigma)^{3\tau-1}$ corresponding to the choice of $H \in \mathrm{Roots}$ and $s \in [0\,..\,|H|)$, we can initialize $L_{\mathrm{minexp}}$ in $\mathcal{O}(n/\log n + n^\mu \log^2 n) = \mathcal{O}(n/\log_\sigma n)$ time.

*Construction of $B_{\mathsf{R}'}$*  As seen above, we can enumerate $\mathsf{R}'$, and thereby compute $B_{\mathsf{R}'}$, in $\mathcal{O}(n/\log_\sigma n)$ time. Augmenting $B_{\mathsf{R}'}$ with Theorem 2.1 takes $\mathcal{O}(n/\log n)$ time.

*Construction of $A_{\mathrm{rmap}}$*  Since for any $j \in \mathsf{R}$, we can in $\mathcal{O}(1)$ compute L-root$(j)$, $e(j)$, $s = $ L-head$(j)$, and $k = $ L-exp$(j)$, in $\mathcal{O}(n/\log_\sigma n)$ time we can also enumerate all $j \in \mathsf{R}'^-$. The key challenge is computing the sequence $(r_i^{\mathrm{lex}})_{i\in[1..q]}$. By the density condition, for every $\tau$-synchronizing set $\mathsf{P}$ of $T$, it holds that if $j \in \mathsf{R}$, then $e(j) - 2\tau + 1 \in \mathsf{P}$ (for a proof, simply compare the claims of Lemma 5.10 and [50, Fact 3.2]). This lets us compute $(r_i^{\mathrm{lex}})_{i\in[1..q]}$ as follows. First, using [50, Theorem 8.11], in $\mathcal{O}(n/\log_\sigma n)$ time we construct any $\tau$-synchronizing set $\mathsf{P}$ of $T$ of size $\mathcal{O}(n/\tau)$. The set $\mathsf{P}$ is returned as an array of size $|\mathsf{P}|$. In $\mathcal{O}(n/\log_\sigma n)$ time we then create a bitvector $B_{\mathsf{P}}[1..n]$ such that $B_{\mathsf{P}}[i] = 1$ holds if and only if $i \in \mathsf{P}$. In $\mathcal{O}(n/\log n)$ time we augment $B_{\mathsf{P}}$ using Theorem 2.1. Let $(p_t^{\mathrm{text}})_{t\in[1..|\mathsf{P}|]}$ denote a sequence containing elements of $\mathsf{P}$ in increasing order and let $(p_t^{\mathrm{lex}})_{t\in[1..|\mathsf{P}|]}$ denote a sequence containing $\mathsf{P}$ sorted according to the lexicographical order of the corresponding suffixes in $T$, i.e., such that for any $i, i' \in [1..|\mathsf{P}|]$, $i < i'$ implies $T[p_i^{\mathrm{lex}}..n] \prec T[p_{i'}^{\mathrm{lex}}..n]$. Given the array containing $\mathsf{P}$, we compute the sequence $(p_t^{\mathrm{lex}})_{t\in[1..|\mathsf{P}|]}$ in $\mathcal{O}(n/\log_\sigma n)$ time using [50, Theorem 4.3]. Let $\mathrm{ISA}_{\mathsf{P}}[1..|\mathsf{P}|]$ be an array storing a permutation of $[1..|\mathsf{P}|]$ such that $\mathrm{ISA}_{\mathsf{P}}[j] = i$ holds if and only if $p_j^{\mathrm{text}} = p_i^{\mathrm{lex}}$. Using the sequence $(p_t^{\mathrm{lex}})_{t\in[1..|\mathsf{P}|]}$ and the bitvector $B_{\mathsf{P}}$, we compute $\mathrm{ISA}_{\mathsf{P}}$ in $\mathcal{O}(|\mathsf{P}|) = \mathcal{O}(n/\log_\sigma n)$ time: For every $i \in [1..|\mathsf{P}|]$, we first compute $j = \mathsf{rank}_{B_{\mathsf{P}},1}(p_i^{\mathrm{lex}})$ and then set $\mathrm{ISA}_{\mathsf{P}}[j] = i$. Next, for each $j \in \mathsf{R}'^-$, letting $H = $ L-root$(j)$ and $j_{\mathsf{P}} = \mathsf{rank}_{B_{\mathsf{P}},1}(e(j) - 2\tau + 1)$, we form a tuple $(\mathrm{int}(H), e(j) - e^{\mathrm{full}}(j), \mathrm{ISA}_{\mathsf{P}}[j_{\mathsf{P}}], j)$. Observe, that $X \prec X'$ holds if and only if $\mathrm{int}(X) < \mathrm{int}(X')$. Let $j, j' \in \mathsf{R}'^-_H$. Note that since both $T[e^{\mathrm{full}}(j)..e(j))$ and $T[e^{\mathrm{full}}(j')..e(j'))$ are prefixes of $H$, by definition of $\mathsf{R}^-$, $e(j) - e^{\mathrm{full}}(j) < e(j') - e^{\mathrm{full}}(j')$ implies $T[e^{\mathrm{full}}(j)..n] \prec T[e^{\mathrm{full}}(j')..n]$. If $e(j) - e^{\mathrm{full}}(j) = e(j') - e^{\mathrm{full}}(j')$, then $T[e(j) - 2\tau + 1..e^{\mathrm{full}}(j)) = T[e(j') - 2\tau + 1..e^{\mathrm{full}}(j'))$, and consequently, $T[e^{\mathrm{full}}(j)..n] \prec T[e^{\mathrm{full}}(j')..n]$ holds if and only if $\mathrm{ISA}_{\mathsf{P}}[j_{\mathsf{P}}] < \mathrm{ISA}_{\mathsf{P}}[j'_{\mathsf{P}}]$. We have thus shown that sorting the tuples lexicographically yields a sequence $(r_i^{\mathrm{lex}})_{i\in[1..q]}$ on the fourth coordinate. Given $j \in \mathsf{R}'^-$, we can compute the corresponding tuple in $\mathcal{O}(1)$ time. Thus, since all its elements are integers in the range $[1..n]$, using LSD radix-sort, we can compute $(r_i^{\mathrm{lex}})_{i\in[1..q]}$ in $\mathcal{O}(n/\log_\sigma n)$ time. With a single scan of $(r_i^{\mathrm{lex}})_{i\in[1..q]}$ and the help of rank queries on $B_{\mathsf{R}'}$ we can then compute table $A_{\mathrm{rmap}}$ in $\mathcal{O}(n/\log_\sigma n)$ time.

*Construction of $A_{\mathrm{rmap}}^{-1}$*  Given $A_{\mathrm{rmap}}$, we can compute $A_{\mathrm{rmap}}^{-1}$ in $\mathcal{O}(q) = \mathcal{O}(n/\log_\sigma n)$ time, since these two arrays are inverses of each other.

*Construction of $L_{\mathrm{runs}}$*  In $\mathcal{O}(\sigma^\tau + |\mathsf{R}'^-|)$ time we perform a synchronized enumeration of all $H \in [0..\sigma)^{\leq\tau}$ in lexicographical order and the L-root values (obtained using Proposition 5.15) for positions in the sequence $(r_i^{\mathrm{lex}})_{i\in[1..q]}$. This lets us obtain the pair $(b_H, e_H)$ satisfying $\{r_i^{\mathrm{lex}} : i \in (b_H..e_H]\} = \mathsf{R}'^-_H$ for every $H \in [0..\sigma)^{\leq\tau}$ satisfying $\mathsf{R}'^-_H \neq \emptyset$. For each such $H$, we then enumerate all $H' \in [0..\sigma)^{\leq\tau}$ and for each we find corresponding subrange of $(b_H..e_H]$ in $\mathcal{O}(\tau \log n)$ time using binary search. Overall, the initialization of $L_{\mathrm{runs}}$ takes $\mathcal{O}(\sigma^{6\tau}\tau + |\mathsf{R}'^-| + \sigma^{2\tau}\tau \log n) = \mathcal{O}(n/\log_\sigma n)$ time.

*Construction of $L_{\mathrm{pref}}$*  To construct $L_{\mathrm{pref}}$, we enumerate all possible $H \in [0..\sigma)^{\leq\tau}$. For each $H$, we try all $s \in [0..|H|)$, and for each we construct the string $\mathrm{Pref}_{3\tau-1}(s, H)$ in $\mathcal{O}(\tau)$ time. Over all $H$, and including the initialization of $L_{\mathrm{pref}}$, this takes $\mathcal{O}(\sigma^{6\tau}\tau + \sigma^\tau\tau^2) = \mathcal{O}(n^{6\mu}\log n) = \mathcal{O}(n/\log_\sigma n)$ time.

*Construction of range counting/selection for $A$*  From $(r_i^{\mathrm{lex}})_{i\in[1..q]}$ we construct in $\mathcal{O}(n/\log_\sigma n)$ time the sequence $(\ell_i)_{i\in[1..q]}$, and then build the array $A_{\mathrm{len}}[1..q]$ and augment it with a range counting/selection data structure. Using Proposition 2.3, by $q = \mathcal{O}(n/\log_\sigma n)$ and $\sum_{i=1}^q A_{\mathrm{len}}[i] = \mathcal{O}(n)$, this takes $\mathcal{O}(n/\log_\sigma n)$ time.

*Construction of the remaining components*   After the above components are constructed, we then analogously construct their symmetric counterparts (adapted according to Lemma 5.11).  □

## 5.4   The Final Data Structure

In this section, we put together Sections 5.1 to 5.3 to obtain a data structure that, given any $j \in [1 \mathinner{\ldotp\ldotp} n]$ (resp. $i \in [1 \mathinner{\ldotp\ldotp} n]$) computes $\mathrm{ISA}[j]$ (resp. $\mathrm{SA}[i]$) in $\mathcal{O}(\log^\epsilon n)$ time.

The section is organized as follows. First, we introduce the components of the data structure (Section 5.4.1). Next, we describe the query algorithms (Sections 5.4.2 and 5.4.3). Finally, we show the construction algorithm (Section 5.4.4).

### 5.4.1   The Data Structure

The data structure consists of two components:

1. The structure from Section 5.2.1 (used to handle nonperiodic positions).
2. The structure from Section 5.3.2 (used to handle periodic positions).

In total, the data structure needs $\mathcal{O}(n/\log_\sigma n)$ space.

### 5.4.2   Implementation of ISA Queries

**Proposition 5.28.** *Given the data structure from Section 5.4.1 and any $j \in [1 \mathinner{\ldotp\ldotp} n]$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\mathrm{ISA}[j]$.*

*Proof.* First, we use Proposition 5.2 to check in $\mathcal{O}(1)$ time if $j \in \mathsf{R}$. Depending on whether $j \in \mathsf{R}$ or not, we use Proposition 5.6 or Proposition 5.23 to compute $\mathrm{ISA}[j]$ in $\mathcal{O}(\log^\epsilon n)$ or $\mathcal{O}(\log\log n)$ time (respectively).  □

### 5.4.3   Implementation of SA Queries

**Proposition 5.29.** *Given the data structure from Section 5.4.1 and any $i \in [1 \mathinner{\ldotp\ldotp} n]$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\mathrm{SA}[i]$.*

*Proof.* First, we use Proposition 5.3 to check in $\mathcal{O}(1)$ time if $\mathrm{SA}[i] \in \mathsf{R}$. Depending on whether $\mathrm{SA}[i] \in \mathsf{R}$ or not, we use Proposition 5.8 or Proposition 5.26 to compute $\mathrm{SA}[i]$ in $\mathcal{O}(\log^\epsilon n)$ or $\mathcal{O}(\log\log n)$ time (respectively).  □

### 5.4.4   Construction Algorithm

**Proposition 5.30.** *Given the packed representation of $T \in [0 \mathinner{\ldotp\ldotp} \sigma)^n$, we can construct the data structure from Section 5.4.1 in $\mathcal{O}(n\min(1, \log\sigma/\sqrt{\log n}))$ time and $\mathcal{O}(n/\log_\sigma n)$ working space.*

*Proof.* First, from a packed representation of $T$, we construct $\mathrm{C_{SA}}(T)$ in $\mathcal{O}(n/\log_\sigma n)$ time using Proposition 5.4. Then, using Propositions 5.9 and 5.27, we augment $\mathrm{C_{SA}}(T)$ into the two components of the structure from Section 5.4.1 in $\mathcal{O}(n\min(1, \log\sigma/\sqrt{\log n}))$ and $\mathcal{O}(n/\log_\sigma n)$ time (respectively) and using $\mathcal{O}(n/\log_\sigma n)$ working space.  □

## 5.5   Summary

By combining Propositions 5.28 to 5.30 we obtain the following final result of this section.

**Theorem 5.31.** *Given any constant $\epsilon \in (0,1)$ and the packed representation of a text $T \in [0 \mathinner{\ldotp\ldotp} \sigma)^n$ with $2 \le \sigma < n^{1/7}$, in $\mathcal{O}(n\min(1, \log\sigma/\sqrt{\log n}))$ time and $\mathcal{O}(n/\log_\sigma n)$ working space we can construct a data structure of size $\mathcal{O}(n/\log_\sigma n)$ that:*

- *Given any $i \in [1 \mathinner{.\,.} n]$ returns $\mathrm{SA}[i]$ in $\mathcal{O}(\log^\epsilon n)$ time,*
- *Given any $j \in [1 \mathinner{.\,.} n]$ returns $\mathrm{ISA}[j]$ in $\mathcal{O}(\log^\epsilon n)$ time.*

We also immediately obtain the following more general result.

**Theorem 5.32.** *Consider a data structure answering prefix rank and selection queries that, for any string of length $m$ over alphabet $[0 \mathinner{.\,.} \sigma)^\ell$, achieves the following complexities:*

1. *Space usage $S(m, \ell, \sigma)$,*
2. *Preprocessing time $P_t(m, \ell, \sigma)$,*
3. *Preprocessing space $P_s(m, \ell, \sigma)$,*
4. *Query time $Q(m, \ell, \sigma)$.*

*For every $T \in [0 \mathinner{.\,.} \sigma)^n$ with $2 \leq \sigma < n^{1/7}$, there exist $m = \mathcal{O}(n/\log_\sigma n)$ and $\ell = \mathcal{O}(\log_\sigma n)$ such that, given the packed representation of $T$, we can in $\mathcal{O}(n/\log_\sigma n + P_t(m, \ell, \sigma))$ time and $\mathcal{O}(n/\log_\sigma n + P_s(m, \ell, \sigma))$ working space build a structure of size $\mathcal{O}(n/\log_\sigma n + S(m, \ell, \sigma))$ that:*

- *Given any $i \in [1 \mathinner{.\,.} n]$ returns $\mathrm{SA}[i]$ in $\mathcal{O}(\log \log n + Q(m, \ell, \sigma))$ time,*
- *Given any $j \in [1 \mathinner{.\,.} n]$ returns $\mathrm{ISA}[j]$ in $\mathcal{O}(\log \log n + Q(m, \ell, \sigma))$ time.*

# 6 Pattern Matching Queries

Let $\epsilon \in (0, 1)$ be any fixed constant and let $T \in [0 \mathinner{.\,.} \sigma)^n$, where $2 \leq \sigma < n^{1/7}$. In this section we show how, given the packed representation of $T$, to construct in $\mathcal{O}(n \min(1, \log \sigma / \sqrt{\log n}))$ time and $\mathcal{O}(n/\log_\sigma n)$ working space a structure of size $\mathcal{O}(n/\log_\sigma n)$ that, given the packed representation of a pattern $P \in [0 \mathinner{.\,.} \sigma)^m$, returns $\mathrm{RangeBeg}(P, T)$ and $\mathrm{RangeEnd}(P, T)$ in $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$ time. We also derive a general reduction depending on prefix rank and selection queries.

As in Section 5, we let $\tau = \lfloor \mu \log_\sigma n \rfloor$, where $\mu$ is some positive constant smaller than $\frac{1}{6}$ such that $\tau \geq 1$, be fixed for the duration of this section. Throughout, we also use $\mathsf{R}$ as a shorthand for $\mathsf{R}(\tau, T)$.

**Definition 6.1.** Let $P \in [0 \mathinner{.\,.} \sigma)^m$. We call pattern $P$ *periodic* if it holds that $m \geq 3\tau - 1$ and $\mathrm{per}(P[1 \mathinner{.\,.} 3\tau-1]) \leq \frac{1}{3}\tau$. Otherwise, $P$ is *nonperiodic*.

**Organization** The structure and the query algorithm for a pattern $P$ are different depending on whether $P$ is periodic (Definition 6.1). Our description is thus split as follows. First (Section 6.1), we describe the set of data structures called collectively the index "core" that enables efficiently checking if $P$ is periodic (it is also used to handle very short patterns and contains some common components utilized by the remaining parts). In the following two parts (Sections 6.2 and 6.3), we describe structures handling each of the two cases. All ingredients are then put together in Section 6.4. Finally, we present our result in the general form (Section 6.5).

## 6.1 The Index Core

In this section, we present a data structure that, given a packed representation of any pattern $P \in [0 \mathinner{.\,.} \sigma)^m$, lets us in $\mathcal{O}(1)$ time check if $P$ is periodic. It also let us compute $(\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T))$ if $m < 3\tau - 1$.

The section is organized as follows. First, we introduce the components of the data structure (Section 6.1.1). We then show how using this structure to implement the periodicity check (Section 6.1.2). Next, we describe the query algorithm for short patterns (Section 6.1.3). Finally, we show the construction algorithm (Section 6.1.4).

### 6.1.1 The Data Structure

The index core, denoted $C_{PM}(T)$ consists of the following subset of components of $C_{SA}(T)$:

1. The packed representation of $T$ using $\mathcal{O}(n/\log_\sigma n)$ space.
2. The lookup table $L_{\text{range}}$ using $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n/\log_\sigma n)$ space.
3. The lookup table $L_{\text{per}}$ using $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n/\log_\sigma n)$ space.

In total, $C_{PM}(T)$ needs $\mathcal{O}(n/\log_\sigma n)$ space.

### 6.1.2 Navigation Primitives

**Proposition 6.2.** *Given $C_{PM}(T)$ and a packed representation of $P \in [0 \mathinner{.\,.} \sigma)^m$, we can in $\mathcal{O}(1)$ time determine whether $P$ is periodic (Definition 6.1).*

*Proof.* If $m < 3\tau - 1$, we return false. Otherwise, in $\mathcal{O}(1)$ time we compute $x = \text{int}(X)$, where $X = P[1 \mathinner{.\,.} 3\tau{-}1]$. We then look up $p = L_{\text{per}}[x]$ and return true if and only if $p \le \frac{1}{3}\tau$. $\qquad\square$

### 6.1.3 Implementation of Queries

**Proposition 6.3.** *Let $P \in [0 \mathinner{.\,.} \sigma)^m$ be a pattern satisfying $m < 3\tau - 1$. Given $C_{PM}(T)$ and the packed representation of $P$, in $\mathcal{O}(1)$ time we can compute $(\text{RangeBeg}(P,T), \text{RangeEnd}(P,T))$.*

*Proof.* Using $L_{\text{range}}$ on $P$, we immediately obtain and return $(\text{RangeBeg}(P,T), \text{RangeEnd}(P,T))$ in $\mathcal{O}(1)$ time. $\qquad\square$

### 6.1.4 Construction Algorithm

**Proposition 6.4.** *Given the packed representation of $T \in [0 \mathinner{.\,.} \sigma)^n$, we can construct $C_{PM}(T)$ in $\mathcal{O}(n/\log_\sigma n)$ time.*

*Proof.* Since $C_{PM}(T)$ contains a subset of components of $C_{SA}(T)$, this follows by Proposition 5.4. $\qquad\square$

## 6.2 The Nonperiodic Patterns

In this section, we describe a data structure that, given a packed representation of any nonperiodic pattern $P \in [0 \mathinner{.\,.} \sigma)^m$ (see Definition 6.1), computes $(\text{RangeBeg}(P,T), \text{RangeEnd}(P,T))$ in $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$ time.

The section is organized as follows. First, we introduce the components of the data structure (Section 6.2.1). Next, we describe the query algorithm (Section 6.2.2). Finally, we show the construction algorithm (Section 6.2.3).

### 6.2.1 The Data Structure

**Definitions** Let $S$ be a $\tau$-synchronizing set, as defined in Section 5.2.1. Let $A_S[1 \mathinner{.\,.} n']$ be an array defined by $A_S[i] = s_i^{\text{lex}}$ (where $(s_t^{\text{lex}})_{t \in [1 \mathinner{.\,.} n']}$ is a sequence as defined in Section 5.2.1).

**Components** The data structure to handle nonperiodic patterns consists of three components:

1. The index core $C_{PM}(T)$ (Section 6.1.1) using $\mathcal{O}(n/\log_\sigma n)$ space.
2. The data structure from Section 5.2.1 using $\mathcal{O}(n/\log_\sigma n)$ space.
3. The data structure from Proposition 4.4 for the array $A_S[1 \mathinner{.\,.} n']$. By $n' = \mathcal{O}(n/\log_\sigma n)$ and Proposition 4.4, it needs $\mathcal{O}(n/\log_\sigma n)$ space.

In total, the data structure takes $\mathcal{O}(n/\log_\sigma n)$ space.

### 6.2.2 Implementation of Queries

**Lemma 6.5.** *Let $P \in [0 \mathinner{.\,.} \sigma)^m$ be a nonperiodic pattern satisfying $m \geq 3\tau - 1$, and let $X \in \mathcal{D}$ be a prefix of $P$. Denote $\delta_{\text{text}} = |X| - 2\tau$ and $P' = P(\delta_{\text{text}} \mathinner{.\,.} m]$. Let $(b_{\text{pre}}, e_{\text{pre}})$ be such that $b_{\text{pre}} = |\{i \in [1 \mathinner{.\,.} n'] : T[s_i^{\text{lex}} \mathinner{.\,.} n] \prec P'\}|$ and $(b_{\text{pre}} \mathinner{.\,.} e_{\text{pre}}] = \{i \in [1 \mathinner{.\,.} n'] : P' \text{ is a prefix of } T[s_i^{\text{lex}} \mathinner{.\,.} n]\}$. Then, it holds*

$$(\text{RangeBeg}(P, T), \text{RangeEnd}(P, T)) = (b_X + \delta_1, b_X + \delta_2),$$

*where $b_X = \text{RangeBeg}(X, T)$, $\delta_1 = \mathsf{rank}_{W, \overline{X}}(b_{\text{pre}})$, and $\delta_2 = \mathsf{rank}_{W, \overline{X}}(e_{\text{pre}})$.*

*Proof.* Observe that by the consistency of $\mathsf{S}$ and $X \in \mathcal{D}$, $j \in \text{Occ}(X, T)$ implies $j + \delta_{\text{text}} \in \mathsf{S}$. Thus, $\text{Occ}(X, T) = \{s - \delta_{\text{text}} : s \in \mathsf{S} \text{ and } s - \delta_{\text{text}} \in \text{Occ}(X, T)\}$. Note also that if $S_1$ is a prefix of $S_2$ then $\text{RangeBeg}(S_2, T) = \text{RangeBeg}(S_1, T) + |\{j \in \text{Occ}(S_1, T) : T[j \mathinner{.\,.} n] \prec S_2\}|$. Together with the definition of $b_{\text{pre}}$, this implies

$$
\begin{aligned}
\text{RangeBeg}(P, T) &= \text{RangeBeg}(X, T) + |\{j \in \text{Occ}(X, T) : T[j \mathinner{.\,.} n] \prec P\}| \\
&= b_X + |\{s - \delta_{\text{text}} : s \in \mathsf{S}, s - \delta_{\text{text}} \in \text{Occ}(X, T), \text{ and } T[s - \delta_{\text{text}} \mathinner{.\,.} n] \prec P\}| \\
&= b_X + |\{s \in \mathsf{S} : s - \delta_{\text{text}} \in \text{Occ}(X, T) \text{ and } T[s - \delta_{\text{text}} \mathinner{.\,.} n] \prec P\}| \\
&= b_X + |\{s \in \mathsf{S} : s - \delta_{\text{text}} \in \text{Occ}(X, T) \text{ and } T[s \mathinner{.\,.} n] \prec P'\}| \\
&= b_X + |\{i \in [1 \mathinner{.\,.} n'] : s_i^{\text{lex}} - \delta_{\text{text}} \in \text{Occ}(X, T) \text{ and } T[s_i^{\text{lex}} \mathinner{.\,.} n] \prec P'\}| \\
&= b_X + |\{i \in [1 \mathinner{.\,.} n'] : s_i^{\text{lex}} - \delta_{\text{text}} \in \text{Occ}(X, T) \text{ and } i \leq b_{\text{pre}}\}| \\
&= b_X + |\{i \in [1 \mathinner{.\,.} b_{\text{pre}}] : s_i^{\text{lex}} - \delta_{\text{text}} \in \text{Occ}(X, T)\}| \\
&= b_X + \mathsf{rank}_{W, \overline{X}}(b_{\text{pre}}) \\
&= b_X + \delta_1,
\end{aligned}
$$

where the second-to-last equality follows by $W[i] = \overline{X_i}$, where $X_i = T^\infty[s_i^{\text{lex}} - \tau \mathinner{.\,.} s_i^{\text{lex}} + 2\tau)$, since $s_i^{\text{lex}} - \delta_{\text{text}} \in \text{Occ}(X, T)$ holds if and only if $X$ if a suffix of $X_i$ (i.e., if $\overline{X}$ is a prefix of $\overline{X_i}$).

Next, we show that $|\text{Occ}(P, T)| = \delta_2 - \delta_1$. We start by observing that (similarly as above, except applied to $P$) by the consistency of $\mathsf{S}$ and $X \in \mathcal{D}$ being a prefix of $P$, $j \in \text{Occ}(P, T)$ implies $j + \delta_{\text{text}} \in \mathsf{S}$. Thus, $\text{Occ}(P, T) = \{s - \delta_{\text{text}} : s \in \mathsf{S} \text{ and } s - \delta_{\text{text}} \in \text{Occ}(P, T)\}$ and hence,

$$
\begin{aligned}
|\text{Occ}(P, T)| &= |\{s - \delta_{\text{text}} : s \in \mathsf{S}, s - \delta_{\text{text}} \in \text{Occ}(P, T)\}| \\
&= |\{s \in \mathsf{S} : s - \delta_{\text{text}} \in \text{Occ}(P, T)\}| \\
&= |\{i \in [1 \mathinner{.\,.} n'] : s_i^{\text{lex}} - \delta_{\text{text}} \in \text{Occ}(P, T)\}| \\
&= |\{i \in [1 \mathinner{.\,.} n'] : s_i^{\text{lex}} - \delta_{\text{text}} \in \text{Occ}(X, T) \text{ and } s_i^{\text{lex}} \in \text{Occ}(P', T)\}| \\
&= |\{i \in [1 \mathinner{.\,.} n'] : s_i^{\text{lex}} - \delta_{\text{text}} \in \text{Occ}(X, T) \text{ and } b_{\text{pre}} < i \leq e_{\text{pre}}\}| \\
&= |\{i \in (b_{\text{pre}} \mathinner{.\,.} e_{\text{pre}}] : s_i^{\text{lex}} - \delta_{\text{text}} \in \text{Occ}(X, T)\}| \\
&= \mathsf{rank}_{W, \overline{X}}(b_{\text{pre}}) - \mathsf{rank}_{W, \overline{X}}(e_{\text{pre}}) \\
&= \delta_1 - \delta_2.
\end{aligned}
$$

Combining the above with the earlier equality, we obtain $\text{RangeEnd}(P, T) = \text{RangeBeg}(P, T) + |\text{Occ}(P, T)| = b_X + \delta_2$, i.e., the second part of the claim. $\qquad\square$

*Remark* 6.6. Note that since the range $(b_{\text{pre}} \mathinner{.\,.} e_{\text{pre}}]$ is well-defined even if $e_{\text{pre}} - b_{\text{pre}} = 0$, the above lemma holds even if $|\text{Occ}(P, T)| = 0$.

**Proposition 6.7.** *Let $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$ be a nonperiodic pattern satisfying $m \geq 3\tau - 1$. Given the data structure from Section 6.2.1 and the packed representation of $P$, in $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$ time we can compute $(\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T))$.*

*Proof.* Let $Y = P[1 \mathinner{\ldotp\ldotp} 3\tau - 1]$. First, using the lookup table $L_{\mathrm{range}}$ on $Y$, in $\mathcal{O}(1)$ time we compute $(b_Y, e_Y) = (\mathrm{RangeBeg}(Y, T), \mathrm{RangeEnd}(Y, T))$. If $e_Y - b_Y = 0$, then $\mathrm{Occ}(Y, T) = \emptyset$, and it is easy to see that then we have $\mathrm{RangeBeg}(P, T) = \mathrm{RangeBeg}(Y, T)$ and $\mathrm{RangeEnd}(P, T) = \mathrm{RangeEnd}(Y, T)$. We thus return $(b_Y, e_Y)$. Let us thus assume $b_Y \neq e_Y$, i.e., $\mathrm{Occ}(Y, T) \neq \emptyset$. Together with $\mathrm{per}(Y) > \frac{1}{3}\tau$, this implies (see Section 5.2.1) that there exists a unique prefix $X \in \mathcal{D}$ of $P$. Using $L_{\mathcal{D}}$ on $Y$ in $\mathcal{O}(1)$ time we compute the prefix $X \in \mathcal{D}$ of $P$. Let $\delta = |X| - 2\tau$. Using again the lookup table $L_{\mathrm{range}}$, in $\mathcal{O}(1)$ time we compute $(b_X, e_X) = (\mathrm{RangeBeg}(X, T), \mathrm{RangeEnd}(X, T))$. Using Proposition 4.4, we then compute in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time the pair $(b_{\mathrm{pre}}, e_{\mathrm{pre}})$ for the pattern $P' := P(\delta \mathinner{\ldotp\ldotp} m]$. By Lemma 6.5, we then return $(\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T)) = (b_X + \mathsf{rank}_{W, \overline{X}}(b_{\mathrm{pre}}), b_X + \mathsf{rank}_{W, \overline{X}}(e_{\mathrm{pre}}))$, with the two prefix rank queries implemented using Theorem 2.2, in $\mathcal{O}(\log^\epsilon n)$ time each (the string $\overline{X}$ is obtained using the lookup table $L_{\mathrm{rev}}$). Altogether, the query time is $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$. $\square$

### 6.2.3 Construction Algorithm

**Proposition 6.8.** *Given $\mathrm{C}_{\mathrm{PM}}(T)$, we can in $\mathcal{O}(n \min(1, \log\sigma/\sqrt{\log n}))$ time and in $\mathcal{O}(n/\log_\sigma n)$ working space augment it into a data structure from Section 6.2.1.*

*Proof.* First, we combine Propositions 5.4 and 5.9 (recall that the packed representation of $T$ is a component of $\mathrm{C}_{\mathrm{PM}}(T)$) to construct the structure from Section 5.2.1 in $\mathcal{O}(n \min(1, \log\sigma/\sqrt{\log n}))$ time and using $\mathcal{O}(n/\log_\sigma n)$ working space. In particular, this constructs $(s_i^{\mathrm{lex}})_{i \in [1 \mathinner{\ldotp\ldotp} n']}$. We thus initialize $A_{\mathsf{S}}[i] = s_i^{\mathrm{lex}}$ for $i \in [1 \mathinner{\ldotp\ldotp} n']$ and in $\mathcal{O}(n/\log_\sigma n)$ time and $\mathcal{O}(n/\log_\sigma n)$ working space construct the data structure from Proposition 4.4. The overall runtime is $\mathcal{O}(n \min(1, \log\sigma/\sqrt{\log n}))$. The working space never exceed $\mathcal{O}(n/\log_\sigma n)$ words. $\square$

## 6.3 The Periodic Patterns

In this section, we describe a data structure that, given a packed representation of any periodic pattern $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$ (see Definition 6.1), computes $(\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T))$ in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time.

The section is organized as follows. First, we present the toolbox of combinatorial properties for periodic patterns (Section 6.3.1). Next, we introduce the components of the data structure (Section 6.3.2). We then show how using this structure to implement some basic navigational routines (Section 6.3.3). Next, we describe the query algorithm (Section 6.3.4). Finally, we show the construction algorithm (Section 6.3.5).

### 6.3.1 Preliminaries

Let $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$ be a periodic pattern (see Definition 6.1). We define $\mathrm{L\text{-}root}(P) = \min\{P[1 + t \mathinner{\ldotp\ldotp} 1 + t + p) : t \in [0 \mathinner{\ldotp\ldotp} p)\}$, where $p = \mathrm{per}(P[1 \mathinner{\ldotp\ldotp} 3\tau - 1])$. Let $H = \mathrm{L\text{-}root}(P)$. We define $e(P) = 1 + p + \mathrm{lcp}(P[1 \mathinner{\ldotp\ldotp} m], P[1 + p \mathinner{\ldotp\ldotp} m])$, where $p = |H|$. By definition, there exists $s \in [0 \mathinner{\ldotp\ldotp} p)$ such that $P[1 + s \mathinner{\ldotp\ldotp} 1 + s + p) = H$. Thus, we can write $P[1 \mathinner{\ldotp\ldotp} e(P)] = H'H^kH''$, where $H'$ (resp. $H''$) is a proper suffix (resp. prefix) of $H$. By $e(P) \geq 3\tau$ and $|H| \leq \tau$, such decomposition is unique (see also Section 5.3.1). We denote $\mathrm{L\text{-}head}(P) = |H'|$, $\mathrm{L\text{-}exp}(P) = k$, and $\mathrm{L\text{-}tail}(P) = |H''|$. We also let $e^{\mathrm{full}}(P) = e(P) - \mathrm{L\text{-}tail}(P)$. We define $\mathrm{type}(P) = +1$ if $e(P) \leq m$ and $P[e(P)] \succ P[e(P) - p]$ (where $p = |\mathrm{L\text{-}root}(P)|$), and $\mathrm{type}(P) = -1$ otherwise.

**Lemma 6.9.** *Let $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$ be a periodic pattern and let $s = \text{L-head}(P)$ and $H = \text{L-root}(P)$. For any $j \in [1 \mathinner{\ldotp\ldotp} n]$, $\text{lcp}(P, T[j \mathinner{\ldotp\ldotp} n]) \geq 3\tau - 1$ holds if and only if $j \in \mathsf{R}_{s,H}$. Moreover, if $j \in \mathsf{R}_{s,H}$ then, letting $t = e(P) - 1$ and $t' = e(j) - j$, it holds $\text{lcp}(P, T[j \mathinner{\ldotp\ldotp} n]) \geq \min(t, t')$ and:*

 1. *If $\text{type}(P) \neq \text{type}(j)$, then $P \prec T[j \mathinner{\ldotp\ldotp} n]$ if and only if $\text{type}(P) < \text{type}(j)$,*
 2. *If $\text{type}(P) = \text{type}(j) = -1$ and $t \neq t'$, then $P \prec T[j \mathinner{\ldotp\ldotp} n]$ if and only if $t < t'$,*
 3. *If $\text{type}(P) = \text{type}(j) = +1$ and $t \neq t'$, then $P \prec T[j \mathinner{\ldotp\ldotp} n]$ if and only if $t > t'$,*
 4. *If $\text{type}(P) \neq \text{type}(j)$ or $t \neq t'$, then $P \neq T[j \mathinner{\ldotp\ldotp} n]$ and $\text{lcp}(P, T[j \mathinner{\ldotp\ldotp} n]) = \min(t, t')$.*

*Proof.* Let $j \in [1 \mathinner{\ldotp\ldotp} n]$ be such that $\text{lcp}(P, T[j \mathinner{\ldotp\ldotp} n]) \geq 3\tau - 1$. Denoting $p = \text{per}(P[1 \mathinner{\ldotp\ldotp} 3\tau - 1])$ and $p' = \text{per}(T[j \mathinner{\ldotp\ldotp} j + 3\tau - 1))$ we then have $p' = p \leq \frac{1}{3}\tau$. Thus, $j \in \mathsf{R}$. Moreover, this implies $\text{L-root}(j) = \min\{T[j + \delta \mathinner{\ldotp\ldotp} j + \delta + p') : \delta \in [0 \mathinner{\ldotp\ldotp} p')\} = \min\{T[j + \delta \mathinner{\ldotp\ldotp} j + \delta + p) : \delta \in [0 \mathinner{\ldotp\ldotp} p)\} = \min\{P[1 + \delta \mathinner{\ldotp\ldotp} 1 + \delta + p) : \delta \in [0 \mathinner{\ldotp\ldotp} p)\} = H$. To show that $\text{L-head}(j) = s$, note that by $|H| \leq \tau$, the string $H'H^2$ (where $H'$ is a length-$s$ suffix of $H$) is a prefix of $P[1 \mathinner{\ldotp\ldotp} 3\tau - 1] = T[j \mathinner{\ldotp\ldotp} j + 3\tau - 1)$. On the other hand, $\text{L-head}(j) = s'$ implies that $\widehat{H}'H^2$ (where $\widehat{H}'$ is a length-$s'$ suffix of $H$) is a prefix of $T[j \mathinner{\ldotp\ldotp} j + 3\tau - 1)$. Thus, by the synchronization property of primitive strings [24, Lemma 1.11] applied to the two copies of $H$, we have $s' = s$, and hence, $j \in \mathsf{R}_{s,H}$. For the converse implication, assume $j \in \mathsf{R}_{s,H}$. This implies that both $P[1 \mathinner{\ldotp\ldotp} e(P))$ and $T[j \mathinner{\ldotp\ldotp} e(j))$ are prefixes of $H' \cdot H^{\infty}[1 \mathinner{\ldotp\ldotp})$ (where $H'$ is as above). Thus, by $e(P) - 1$, $e(j) - j \geq 3\tau - 1$, we obtain $\text{lcp}(P, T[j \mathinner{\ldotp\ldotp} n]) \geq 3\tau - 1$.

Let us now assume $j \in \mathsf{R}_{s,H}$. Since, as noted above, both $P[1 \mathinner{\ldotp\ldotp} e(P)) = P[1 \mathinner{\ldotp\ldotp} t]$ and $T[j \mathinner{\ldotp\ldotp} e(j)) = T[j \mathinner{\ldotp\ldotp} j + t']$ are prefixes of $H' \cdot H^{\infty}[1 \mathinner{\ldotp\ldotp})$, we have $\text{lcp}(P, T[j \mathinner{\ldotp\ldotp} n]) \geq \min(t, t')$.

1. Let $Q = H' \cdot H^{\infty}[1 \mathinner{\ldotp\ldotp})$, where $H'$ is a length-$s$ suffix of $H$. In the proof of Lemma 5.11, it is shown that $\text{type}(j) = -1$ implies $T[j \mathinner{\ldotp\ldotp} n] \prec Q$, and $\text{type}(j) = +1$ implies $Q \prec T[j \mathinner{\ldotp\ldotp} n]$. We now prove an analogous fact for $P$. We first note that $\text{type}(P) = -1$ implies that either $e(P) = m + 1$, or $e(P) \leq m$ and $P[e(P)] \prec P[e(P) - |H|]$. In the first case, $P[1 \mathinner{\ldotp\ldotp} e(P)) = P$ is a proper prefix of $Q$ and hence $P \prec Q$. In the second case, we have $P[1 \mathinner{\ldotp\ldotp} t] = Q[1 \mathinner{\ldotp\ldotp} t]$ and $P[1 + t] \prec P[1 + t - |H|] = Q[1 + t - |H|] = Q[1 + t]$. Consequently, $P \prec Q$. If $\text{type}(P) = +1$ holds, then $e(P) \leq m$. Thus, we have $Q[1 \mathinner{\ldotp\ldotp} t] = P[1 \mathinner{\ldotp\ldotp} t]$ and $Q[1 + t] = Q[1 + t - |H|] = P[1 + t - |H|] \prec P[1 + t]$. Hence, we obtain $Q \prec P$. We are now ready to prove the claim. Assume first that $\text{type}(P) < \text{type}(j)$. By the above we then have $P \prec Q \prec T[j \mathinner{\ldotp\ldotp} n]$. The opposite implication is proved by contraposition. Assume $\text{type}(P) > \text{type}(j)$. By the above we then have $T[j \mathinner{\ldotp\ldotp} n] \prec Q \prec P$.

2. Assume $t < t'$. If $e(P) = m + 1$, then $P[1 \mathinner{\ldotp\ldotp} t] = P[1 \mathinner{\ldotp\ldotp} e(P)) = P$ is proper prefix of $T[j \mathinner{\ldotp\ldotp} j + t'] = T[j \mathinner{\ldotp\ldotp} e(j))$, and hence $P \prec T[j \mathinner{\ldotp\ldotp} e(j)) \preceq T[j \mathinner{\ldotp\ldotp} n]$. If $e(P) \leq m$, then we have $P[1 \mathinner{\ldotp\ldotp} t] = T[j \mathinner{\ldotp\ldotp} j + t)$ and by $t < t'$, $P[1 + t] \prec P[1 + t - |H|] = T[j + t - |H|] = T[j + t]$. Thus, we also obtain $P \prec T[j \mathinner{\ldotp\ldotp} n]$. The opposite implication is proved by contraposition. Assume $t > t'$. If $e(j) = n + 1$, then by $t > t'$, the string $T[j \mathinner{\ldotp\ldotp} j + t'] = T[j \mathinner{\ldotp\ldotp} e(j)) = T[j \mathinner{\ldotp\ldotp} n]$ is a proper prefix of $P[1 \mathinner{\ldotp\ldotp} t] = P[1 \mathinner{\ldotp\ldotp} e(P))$, and hence $T[j \mathinner{\ldotp\ldotp} n] \prec P[1 \mathinner{\ldotp\ldotp} e(P)) \preceq P$. If $e(j) \leq n$, then we have $T[j \mathinner{\ldotp\ldotp} j + t'] = P[1 \mathinner{\ldotp\ldotp} t']$ and by $t > t'$, $T[j + t'] \prec T[j + t' - |H|] = P[1 + t' - |H|] = P[1 + t']$. Consequently, we also obtain $T[j \mathinner{\ldotp\ldotp} n] \prec P$.

3. Assume $t > t'$. By $\text{type}(j) = +1$, we have $e(j) \leq n$. Thus, by $t > t'$ we have $P[1 \mathinner{\ldotp\ldotp} t'] = T[j \mathinner{\ldotp\ldotp} j + t')$ and $P[1 + t'] = P[1 + t' - |H|] = T[j + t' - |H|] \prec T[j + t']$. Consequently, $P \prec T[j \mathinner{\ldotp\ldotp} n]$. The opposite implication is proved by contraposition. Assume $t < t'$. By $\text{type}(P) = +1$, we have $e(P) \leq m$. Thus, by $t < t'$ we have $T[j \mathinner{\ldotp\ldotp} j + t) = P[1 \mathinner{\ldotp\ldotp} t]$ and $T[j + t] = T[j + t - |H|] = P[1 + t - |H|] \prec P[1 + t]$. Consequently, we obtain $T[j \mathinner{\ldotp\ldotp} n] \prec P$.

4. By the earlier implication, $\text{lcp}(P, T[j \mathinner{\ldotp\ldotp} n]) \geq \min(t, t')$. Thus, it remains to show that $P \neq T[j \mathinner{\ldotp\ldotp} n]$ and $\text{lcp}(P, T[j \mathinner{\ldotp\ldotp} n]) \leq \min(t, t')$. First, let us assume $\text{type}(P) < \text{type}(j)$ (i.e., $\text{type}(P) = -1$ and $\text{type}(j) = +1$). Consider two cases:

- First, assume $t \leq t'$. Our goal is to prove $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t$. First, recall from the proof of Lemma 5.11(1) that $\mathrm{type}(j) = +1$ implies $j + t' \leq n$, $Q[1 \mathinner{.\,.} t'] = T[j \mathinner{.\,.} j + t')$, and $Q[1 + t'] \prec T[j + t']$. Consider now two subcases. If $e(P) = m + 1$, then $t = m$, and hence $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq m = t$. By $t' \leq n - j$ we then also have $|P| = t < t' + 1 \leq n - j + 1 = |T[j \mathinner{.\,.} n]|$. Thus, $P \neq T[j \mathinner{.\,.} n]$. Let us thus assume $e(P) \leq m$. In the proof of Item 1 we showed that in this case $\mathrm{type}(P) = -1$ implies $P[1 + t] \prec Q[1 + t]$. On the other hand, as noted above, $\mathrm{type}(j) = +1$ implies $Q[1 \mathinner{.\,.} t'] = T[j \mathinner{.\,.} j + t')$, and $Q[1 + t'] \prec T[j + t']$. By $t \leq t'$ we thus have $Q[1 + t] \preceq T[j + t]$. Consequently, $P[1 + t] \neq T[j + t]$. This immediately implies $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t$.
- Let us now assume $t > t'$. Our goal is to prove $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t'$. In the proof of Item 1 we showed that $\mathrm{type}(P) = -1$ implies $P[1 \mathinner{.\,.} t] = Q[1 \mathinner{.\,.} t]$. Thus, by $t > t'$ we have $P[1 + t'] = Q[1 + t']$. On the other hand, in the proof of Lemma 5.11(1) we showed that $\mathrm{type}(j) = +1$ implies $Q[1 + t'] \prec T[j + t']$. Thus, we obtain $P[1 + t'] \neq T[j + t']$. This immediately implies $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t'$.

Assume now $\mathrm{type}(P) > \mathrm{type}(j)$ (i.e., $\mathrm{type}(P) = +1$ and $\mathrm{type}(j) = -1$). Consider two cases:

- First, assume $t < t'$. Our goal is to prove $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t$. In the proof of Lemma 5.11(1) we showed that $\mathrm{type}(j) = -1$ implies $T[j \mathinner{.\,.} j + t') = Q[1 \mathinner{.\,.} t']$. Thus, by $t < t'$ we have $T[j + t] = Q[1 + t]$. On the other hand, in the proof of Item 1 we showed that $\mathrm{type}(P) = +1$ implies $Q[1 + t] \prec P[1 + t]$. Thus, we obtain $T[j + t] \neq P[1 + t]$. This immediately implies $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t$.
- Let us now assume $t \geq t'$. Our goal is to prove $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t'$. First, recall from the proof of Item 1 that $\mathrm{type}(P) = +1$ implies $t + 1 = e(P) \leq m$, $Q[1 \mathinner{.\,.} t] = P[1 \mathinner{.\,.} t]$, and $Q[1 + t] \prec P[1 + t]$. Consider now two subcases. If $e(j) = n + 1$, then $j + t' = n + 1$ (or equivalently, $t' = n - j + 1$) and hence $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq |T[j \mathinner{.\,.} n]| = t'$. By $t + 1 \leq m$ we then also have $|T[j \mathinner{.\,.} n]| = t' < t + 1 \leq m = |P|$. Thus, $P \neq T[j \mathinner{.\,.} n]$. Let us thus assume $e(j) \leq n$. In the proof of Lemma 5.11(1) we showed that in this case $\mathrm{type}(j) = -1$ implies $T[j + t'] \prec Q[1 + t']$. On the other hand, as noted above, $\mathrm{type}(P) = +1$ implies $Q[1 \mathinner{.\,.} t] = P[1 \mathinner{.\,.} t]$ and $Q[1 + t] \prec P[1 + t]$. By $t \geq t'$ we thus have $Q[1 + t'] \preceq P[1 + t']$. Consequently, $T[j + t'] \neq P[1 + t']$. This immediately implies $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t'$.

This concludes the proof of the claim if $\mathrm{type}(P) \neq \mathrm{type}(j)$. Let us now assume $\mathrm{type}(P) = \mathrm{type}(j) = -1$ and $t \neq t'$. Consider two cases:

- First, assume $t < t'$. Our goal is to prove $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t$. In the proof of Item 2 we showed that either it holds $P[1 \mathinner{.\,.} t] = P$ (in which case $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq |P| = t$ and $|P| = t < t' = e(j) - j \leq n + 1 - j = |T[j \mathinner{.\,.} n]|$ which in turn implies $P \neq T[j \mathinner{.\,.} n]$), or $P[1 + t] \prec T[j + t]$ (which also implies $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t$).
- Let us now assume $t > t'$. Our goal is to prove $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t'$. In the proof of Item 2, we showed that either it holds $T[j \mathinner{.\,.} j + t') = T[j \mathinner{.\,.} n]$ (in which case $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq n - j + 1 = t'$ and $|T[j \mathinner{.\,.} n]| = n - j + 1 = t' < t = e(P) - 1 \leq |P|$ which in turn implies $P \neq T[j \mathinner{.\,.} n]$), or $T[j + t'] \prec P[1 + t']$ (which also implies $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t'$).

Let us now assume $\mathrm{type}(P) = \mathrm{type}(j) = +1$ and $t \neq t'$. Consider two cases:

- First, assume $t < t'$. Our goal is to prove $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t$. In the proof of Item 3 we showed that $T[j + t] \prec P[1 + t]$. This immediately implies the claims.
- Let us now assume $t > t'$. Our goal is to prove $P \neq T[j \mathinner{.\,.} n]$ and $\mathrm{lcp}(P, T[j \mathinner{.\,.} n]) \leq t'$. In the proof of Item 3 we showed that $P[1 + t'] \prec T[j + t']$. This immediately implies the claims. $\qquad\square$

**Lemma 6.10.** *Let $P \in [0 \mathinner{.\,.} \sigma)^m$ be a periodic pattern. For every $S \in [0 \mathinner{.\,.} \sigma)^+$, $\mathrm{lcp}(P, S) \geq 3\tau - 1$ implies that $S$ is periodic, and that it holds $\mathrm{L\text{-}root}(S) = \mathrm{L\text{-}root}(P)$ and $\mathrm{L\text{-}head}(S) = \mathrm{L\text{-}head}(P)$.*

*Proof.* Denote $X = P[1 \mathinner{.\,.} 3\tau - 1]$. Letting $p := \mathrm{per}(X)$ we then have $p \leq \frac{1}{3}\tau$. By $\mathrm{lcp}(P, S) \geq 3\tau - 1$, $X$ is thus a prefix of $S$ and hence $\mathrm{per}(S[1 \mathinner{.\,.} 3\tau - 1]) = p \leq \frac{1}{3}\tau$, i.e., $S$ is periodic. Moreover, we then have $\mathrm{L\text{-}root}(S) = \min\{S[1 + t \mathinner{.\,.} 1 + t + p) : t \in [0 \mathinner{.\,.} p)\} = \min\{X[1 + t \mathinner{.\,.} 1 + t + p) : t \in [0 \mathinner{.\,.} p)\} = \min\{P[1 + t \mathinner{.\,.} 1 + t + p) : t \in [0 \mathinner{.\,.} p)\} = \mathrm{L\text{-}root}(P)$. To show the last claim, denote $s = \mathrm{L\text{-}head}(P)$ and $s' = \mathrm{L\text{-}head}(S)$. Then, letting $H = \mathrm{L\text{-}root}(P) = \mathrm{L\text{-}root}(S)$, the string $H'H^2$ (resp. $\widehat{H}'H^2$) is a prefix of $P$ (resp. $S$), where $H'$ (resp. $\widehat{H}'$) is a length-$s$ (resp. length-$s'$) suffix of $H$. Note, however, that $s, s' < |H| = p \leq \frac{1}{3}\tau$ and $|X| \geq \tau \geq 3|H|$. This implies that $H'H^2$ and $\widehat{H}'H^2$ are both prefixes of $X$. By the synchronization property of primitive strings [24, Lemma 1.11], this implies $|H'| = |\widehat{H}'|$. Thus, we obtain $\mathrm{L\text{-}head}(P) = s = |H'| = |\widehat{H}'| = s' = \mathrm{L\text{-}head}(S)$. $\square$

**Lemma 6.11.** *Let $P \in [0 \mathinner{.\,.} \sigma)^+$ be a periodic pattern satisfying $e(P) \leq |P|$. Then:*

1. *For every $S \in [0 \mathinner{.\,.} \sigma)^+$, $\mathrm{lcp}(P, S) \geq e(P)$ (in particular, $P$ being a prefix of $S$) implies that $S$ is periodic and it holds:*

   - $e(S) = e(P)$,
   - $\mathrm{L\text{-}tail}(S) = \mathrm{L\text{-}tail}(P)$,
   - $e^{\mathrm{full}}(S) = e^{\mathrm{full}}(P)$,
   - $\mathrm{L\text{-}exp}(S) = \mathrm{L\text{-}exp}(P)$,
   - $\mathrm{type}(S) = \mathrm{type}(P)$.

2. *If $j \in \mathrm{Occ}(P, T)$, then $j \in \mathsf{R}$ and it holds:*

   - $e(j) - j = e(P) - 1$,
   - $\mathrm{L\text{-}tail}(j) = \mathrm{L\text{-}tail}(P)$,
   - $e^{\mathrm{full}}(j) - j = e^{\mathrm{full}}(P) - 1$,
   - $\mathrm{L\text{-}exp}(j) = \mathrm{L\text{-}exp}(P)$,
   - $\mathrm{type}(j) = \mathrm{type}(P)$.

*Proof.* Denote $X = P[1 \mathinner{.\,.} 3\tau - 1]$, $H = \mathrm{L\text{-}root}(P)$, $s = \mathrm{L\text{-}head}(P)$, and $p = \mathrm{per}(X) = |H| \leq \frac{1}{3}\tau$.

1. First, observe that by definition, $e(P) = 1 + p + \mathrm{lcp}(P, P[1 + p \mathinner{.\,.} |P|]) > |X|$. Thus, $\mathrm{lcp}(P, S) \geq e(P)$ implies that $X$ is a prefix of $S$, and hence $S$ is periodic. By Lemma 6.10, we then also have $\mathrm{L\text{-}root}(S) = H$ and $\mathrm{L\text{-}head}(S) = s$. To show $e(S) = e(P)$, observe that by $e(P) \leq |P|$ and the definition of $e(P)$, we have $P[e(P)] \neq P[e(P) - p]$. Consequently, $\mathrm{lcp}(P, S) \geq e(P)$ yields $\mathrm{lcp}(P, P[1 + p \mathinner{.\,.} |P|]) = \mathrm{lcp}(S, S[1 + p \mathinner{.\,.} |S|])$. Combining this with $|\mathrm{L\text{-}root}(S)| = p$ we thus obtain $e(S) = 1 + p + \mathrm{lcp}(S, S[1 + p \mathinner{.\,.} |S|]) = 1 + p + \mathrm{lcp}(P, P[1 + p \mathinner{.\,.} |P|]) = e(P)$. By $\mathrm{L\text{-}head}(S) = s$ we then also obtain $\mathrm{L\text{-}tail}(S) = (e(S) - 1 - \mathrm{L\text{-}head}(S)) \bmod |\mathrm{L\text{-}root}(S)| = (e(P) - 1 - s) \bmod |H| = \mathrm{L\text{-}tail}(P)$, and consequently $e^{\mathrm{full}}(S) = e(S) - \mathrm{L\text{-}tail}(S) = e(P) - \mathrm{L\text{-}tail}(P) = e^{\mathrm{full}}(P)$. We then also have $\mathrm{L\text{-}exp}(S) = \lfloor \frac{e(S) - 1 - \mathrm{L\text{-}head}(S)}{|\mathrm{L\text{-}root}(S)|} \rfloor = \lfloor \frac{e(P) - 1 - s}{|H|} \rfloor = \mathrm{L\text{-}exp}(P)$. Finally, by $e(P) \leq |P|$ and $\mathrm{lcp}(P, S) \geq e(P)$, we then also have $S[e(S)] = P[e(P)]$. Consequently, $S[e(S)] \prec S[e(S) - p]$ holds if and only of $P[e(P)] \prec P[e(P) - p]$. Therefore, $\mathrm{type}(S) = \mathrm{type}(P)$.

2. We start by noting that $j \in \mathrm{Occ}(P, T)$ implies $j \in \mathrm{Occ}(X, T)$. Thus, $\mathrm{per}(T[j \mathinner{.\,.} j + 3\tau - 1]) = \mathrm{per}(X) = p \leq \frac{1}{3}\tau$, and hence $j \in \mathsf{R}$. By Lemma 6.9, we then also have $\mathrm{L\text{-}root}(j) = H$ and $\mathrm{L\text{-}head}(j) = s$. To show $e(j) - j = e(P) - 1$, denote $S = T[j \mathinner{.\,.} n]$. Since $P$ is a prefix of $S$, by Item 1, it follows that $e(S) = e(P)$. By $\mathrm{L\text{-}root}(S) = \mathrm{L\text{-}root}(P)$ (Lemma 6.10) and the definition of $e(S)$, we thus have $1 + p + \mathrm{lcp}(S, S[1 + p \mathinner{.\,.} |S|]) = e(S) = e(P)$, or equivalently, $\mathrm{lcp}(S, S[1 + p \mathinner{.\,.} |S|]) = e(P) - p - 1$. Since $\mathrm{lcp}(S, S[1 + p \mathinner{.\,.} |S|]) = \mathrm{LCE}(j, j + p)$, we thus obtain $p + \mathrm{LCE}(j, j + p) = e(P) - 1$. It remains to note that for $j \in \mathsf{R}$, by Lemma 5.10(2), $e(j) - j = p + \mathrm{LCE}(j, j + p)$. Therefore, we have $e(j) - j = e(P) - 1$. Combining this with

43

L-root$(j) = H$ and L-head$(j) = s$ yields L-tail$(j) = (e(j) - j - \text{L-head}(j)) \bmod |\text{L-root}(j)| = (e(P) - 1 - s) \bmod |H| = \text{L-tail}(P)$, $e^{\text{full}}(j) - j = e(j) - j - \text{L-tail}(j) = e(P) - 1 - \text{L-tail}(P) = e^{\text{full}}(P) - 1$, and L-exp$(j) = \lfloor \frac{e(j)-j-s}{|\text{L-root}(j)|} \rfloor = \lfloor \frac{e(P)-1-s}{|H|} \rfloor = \text{L-exp}(P)$. Finally, by $e(P) \leq |P|$ we have $e(j) \leq n$ and $T[e(j)] = P[e(P)]$. Consequently, $T[e(j)] \prec T[e(j) - p]$ holds if and only if $P[e(P)] \prec P[e(P) - p]$. Therefore, type$(j) =$ type$(P)$. $\qquad\square$

### 6.3.2 The Data Structure

**Definitions** Let $q = |\mathsf{R}'^{-}|$. Recall (Section 5.3.2), that $(r_i^{\text{lex}})_{i \in [1..q]}$ denotes the sequence containing all positions $j \in \mathsf{R}'^{-}$ sorted first by L-root$(j)$, and in case of ties, by $T[e^{\text{full}}(j) .. n]$. Recall also that Roots $= \{\text{L-root}(j) : j \in \mathsf{R}\}$. For any string $H \in$ Roots, let pow$(H) = H^{\infty}[1 .. |H| \lceil \frac{\tau}{|H|} \rceil]$. This function satisfies the following properties:

- The set $\{\text{pow}(H) : H \in \text{Roots}\}$ is prefix-free.
- For any $X, Y \in$ Roots, $X \prec Y$ implies pow$(X) \prec$ pow$(Y)$.

For a proof, consider $X, Y \in$ Roots such that $X \prec Y$. By [54, Fact 9.1.6], it holds $X \preceq \text{pow}(X) \prec X^{\infty}[1 ..) \prec Y \preceq \text{pow}(Y)$. Since $|Y| < \tau \leq |\text{pow}(X)|$, the set $\{\text{pow}(X), \text{pow}(Y)\}$ is prefix-free.

We define $\mathsf{Z} = \{e^{\text{full}}(j) - |\text{pow}(\text{L-root}(j))| : j \in \mathsf{R}'^{-}\}$. We also define an array $A_{\mathsf{Z}}[1 .. q]$ so that, for any $i \in [1 .. q]$, $A_{\mathsf{Z}}[i] = e^{\text{full}}(j) - |\text{pow}(H_i)|$, where $j = r_i^{\text{lex}}$ and $H_i = \text{L-root}(r_i^{\text{lex}})$. Note that $\{A_{\mathsf{Z}}[i] : i \in [1 .. q]\} = \mathsf{Z}$. Observe also that $T[A_{\mathsf{Z}}[i] .. n] = \text{pow}(H_i) \cdot T[e^{\text{full}}(j) .. n]$. Together with the properties of the pow function and with the definition of $(r_i^{\text{lex}})_{i \in [1..q]}$, this implies that the positions in $A_{\mathsf{Z}}$ are sorted according to the lexicographic order of the corresponding suffixes of $T$, i.e., $i < i'$ implies $T[A_{\mathsf{Z}}[i] .. n] \prec T[A_{\mathsf{Z}}[i'] .. n]$.

**Components** The data structure to handle periodic patterns consists of two parts. The first part (designed to handle periodic patterns $P$ satisfying type$(P) = -1$) consists of three components:

1. The index core $\mathsf{C}_{\text{PM}}(T)$ (Section 6.1.1) using $\mathcal{O}(n/\log_\sigma n)$ space.
2. The first part of the structure from Section 5.3.2 using $\mathcal{O}(n/\log_\sigma n)$ space.
3. The data structure from Proposition 4.4 for the array $A_{\mathsf{Z}}[1 .. q]$. By $q = \mathcal{O}(n/\log_\sigma n)$ and Proposition 4.4, it needs $\mathcal{O}(n/\log_\sigma n)$ space.

The second part of the structure (to handle $P$ satisfying type$(P) = +1$) consists of the symmetric counterparts of the above components adapted according to Lemma 6.9.

In total, the data structure takes $\mathcal{O}(n/\log_\sigma n)$ space.

### 6.3.3 Navigation Primitives

**Proposition 6.12.** *Let $P \in [0 .. \sigma)^m$ be a periodic pattern. Given the data structure from Section 6.3.2 and the packed representation of $P$, we can in $\mathcal{O}(1 + m/\log_\sigma n)$ time compute* L-root$(P)$, L-head$(P)$, L-exp$(P)$, L-tail$(P)$, *and* type$(P)$.

*Proof.* We first compute $x \in [0 .. \sigma^{6\tau})$ such that $x = \text{int}(P[1 .. 3\tau - 1])$. Given the packed encoding of $P$, such $x$ is obtained in $\mathcal{O}(1)$ time. We then look up $(s, p) = L_{\text{root}}[x]$, and in $\mathcal{O}(1)$ time obtain L-root$(P) = P[1+s .. 1+s+p]$ and L-head$(P) = s$. Next, we compute L-exp$(P)$ and L-tail$(P)$. For this, we first determine the length $\ell$ of the longest common prefix of $P$ and $P(p .. m]$. Using the packed representation of $P$, we can do this in $\mathcal{O}(1 + m/\log_\sigma n)$ time (see, e.g., [50, Proposition 2.3]). Consequently, we obtain $e(P) = 1 + p + \ell$, L-exp$(P) = \lfloor \frac{e(P)-1-s}{p} \rfloor$, and L-tail$(P) = (e(P) - 1 - s) \bmod p$. Finally, to test if type$(P) = +1$, we check whether $e(P) \leq m$, and if so, whether $P[e(P)] \succ P[e(P) - p]$. $\qquad\square$

### 6.3.4 Implementation of Queries

**Overview**  The query algorithm is derived in two steps. First, we establish how, given the structure from Section 6.3.2 and a packed representation of any periodic pattern $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$ to compute $|\mathrm{Occ}(P, T)|$ in $\mathcal{O}(m / \log_\sigma n + \log \log n)$ time. This culminates in Proposition 6.18. We then show how to extend this algorithm to instead return $(\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T))$ in the same time complexity, culminating in Proposition 6.24. The reason for this two-step approach is explained in Remark 6.25.

**Computing** $|\mathrm{Occ}(P, T)|$  Let $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$ be a periodic pattern. Denote $s = \mathrm{L\text{-}head}(P)$ and $H = \mathrm{L\text{-}root}(P)$. We define $\mathrm{Occ}^{\mathsf{a}}(P, T) = \{j \in \mathsf{R}_{s,H} \cap \mathrm{Occ}(P, T) : \mathrm{L\text{-}exp}(j) > \mathrm{L\text{-}exp}(P)\}$ and $\mathrm{Occ}^{\mathsf{s}}(P, T) = \{j \in \mathsf{R}_{s,H} \cap \mathrm{Occ}(P, T) : \mathrm{L\text{-}exp}(j) = \mathrm{L\text{-}exp}(P)\}$.

**Lemma 6.13.** *For any periodic pattern $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$, the set $\mathrm{Occ}(P, T)$ is a disjoint union of $\mathrm{Occ}^{\mathsf{a}}(P, T)$ and $\mathrm{Occ}^{\mathsf{s}}(P, T)$.*

*Proof.* By definition, $\mathrm{Occ}^{\mathsf{a}}(P, T) \cap \mathrm{Occ}^{\mathsf{s}}(P, T) = \emptyset$ and $\mathrm{Occ}^{\mathsf{a}}(P, T) \cup \mathrm{Occ}^{\mathsf{s}}(P, T) \subseteq \mathrm{Occ}(P, T)$. Thus, it suffices to show $\mathrm{Occ}(P, T) \subseteq \mathrm{Occ}^{\mathsf{a}}(P, T) \cup \mathrm{Occ}^{\mathsf{s}}(P, T)$. Assume $j \in \mathrm{Occ}(P, T)$. By $m \geq 3\tau - 1$, this implies $\mathrm{lcp}(T[j \mathinner{\ldotp\ldotp} n], P) \geq 3\tau - 1$. Thus, by Lemma 6.9, it holds $j \in \mathsf{R}_{s,H}$, where $s = \mathrm{L\text{-}head}(P)$ and $H = \mathrm{L\text{-}root}(P)$. To obtain $j \in \mathrm{Occ}^{\mathsf{a}}(P, T) \cup \mathrm{Occ}^{\mathsf{s}}(P, T)$ it remains to show $\mathrm{L\text{-}exp}(j) \geq \mathrm{L\text{-}exp}(P)$. First, note that for any $t \in [1 \mathinner{\ldotp\ldotp} m]$, $j \in \mathrm{Occ}(P, T)$ implies $\mathrm{LCE}(j, j+t) \geq \mathrm{lcp}(P[1 \mathinner{\ldotp\ldotp} m], P[1+t \mathinner{\ldotp\ldotp} m])$. In particular, letting $p = |H|$, by definition of $e(P)$ and Lemma 5.10(2), we have $e(j) - j = p + \mathrm{LCE}(j, j+p) \geq p + \mathrm{lcp}(P[1 \mathinner{\ldotp\ldotp} m], P[1+p \mathinner{\ldotp\ldotp} m]) = e(P) - 1$. Consequently, $\mathrm{L\text{-}exp}(j) = \lfloor \frac{e(j) - j - s}{p} \rfloor \geq \lfloor \frac{e(P) - 1 - s}{p} \rfloor = \mathrm{L\text{-}exp}(P)$. $\qquad\square$

By the above lemma, if $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$ is periodic, then $\mathrm{Occ}(P, T) \subseteq \mathsf{R}$. We focus on computing sizes of sets $\mathrm{Occ}^{\mathsf{a}-}(P, T) := \mathrm{Occ}^{\mathsf{a}}(P, T) \cap \mathsf{R}^-$ and $\mathrm{Occ}^{\mathsf{s}-}(P, T) := \mathrm{Occ}^{\mathsf{s}}(P, T) \cap \mathsf{R}^-$. The sizes of the sets $\mathrm{Occ}^{\mathsf{a}+}(P, T) := \mathrm{Occ}^{\mathsf{a}}(P, T) \cap \mathsf{R}^+$ and $\mathrm{Occ}^{\mathsf{s}+}(P, T) := \mathrm{Occ}^{\mathsf{s}}(P, T) \cap \mathsf{R}^+$ are computed analogously (see Proposition 6.18).

We now describe the algorithm to compute $|\mathrm{Occ}^{\mathsf{a}-}(P, T)|$ for any periodic pattern $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$.

**Lemma 6.14.** *Assume that $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$ is periodic. If $e(P) \leq m$, then it holds $\mathrm{Occ}^{\mathsf{a}-}(P, T) = \emptyset$. Otherwise, it holds $\mathrm{Occ}^{\mathsf{a}-}(P, T) = \{j \in \mathsf{R}_{s,H}^- : \mathrm{L\text{-}exp}(j) > \mathrm{L\text{-}exp}(P)\}$, where $s = \mathrm{L\text{-}head}(P)$ and $H = \mathrm{L\text{-}root}(P)$.*

*Proof.* Let $e(P) \leq m$. Denote $k = \mathrm{L\text{-}exp}(P)$. Suppose $\mathrm{Occ}^{\mathsf{a}-}(P, T) \neq \emptyset$, and let $j \in \mathrm{Occ}^{\mathsf{a}-}(P, T)$. By definition, $s + k|H| \leq e(P) - 1 < s + (k+1)|H|$ and $P[e(P)] \neq P[e(P) - |H|]$. On the other hand, by $j \in \mathsf{R}_{s,H}$ and $\mathrm{L\text{-}exp}(j) > k$, the string $H'H^{k+1}$ (where $H'$ is a length-$s$ suffix of $H$) is a prefix of $T[j \mathinner{\ldotp\ldotp} n]$. Thus, we have $T[j+e(P)-1] = T[j+e(P)-1-|H|] = P[e(P)-|H|] \neq P[e(P)]$. This implies $j \notin \mathrm{Occ}(P, T)$, contradicting $j \in \mathrm{Occ}^{\mathsf{a}-}(P, T)$. Thus, $\mathrm{Occ}^{\mathsf{a}-}(P, T) = \emptyset$.

Let $e(P) > m$. The inclusion $\mathrm{Occ}^{\mathsf{a}-}(P, T) \subseteq \{j \in \mathsf{R}_{s,H}^- : \mathrm{L\text{-}exp}(j) > \mathrm{L\text{-}exp}(P)\}$ follows by definition. To show the opposite inclusion, let $j \in \mathsf{R}_{s,H}^-$ be such that $\mathrm{L\text{-}exp}(j) > \mathrm{L\text{-}exp}(P)$. Denote $k = \mathrm{L\text{-}exp}(P)$. Then, $P = H'H^kH''$, where $|H'| = s$, and $H'$ (resp. $H''$) is a suffix (resp. prefix) of $H$. Thus, $P$ is a prefix of $H'H^{k+1}$. The latter string, on the other hand, is by $\mathrm{L\text{-}exp}(j) \geq k + 1$ and $j \in \mathsf{R}_{s,H}$, a prefix of $T[j \mathinner{\ldotp\ldotp} n]$. Thus, $j \in \mathrm{Occ}(P, T)$. By $j \in \mathsf{R}_{s,H}^-$ and $\mathrm{L\text{-}exp}(j) > \mathrm{L\text{-}exp}(P)$, we therefore also have $j \in \mathrm{Occ}^{\mathsf{a}-}(P, T)$. $\qquad\square$

**Proposition 6.15.** *Let $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$ be a periodic pattern. Given the data structure from Section 6.3.2 and the packed representation of $P$, we can compute $|\mathrm{Occ}^{\mathsf{a}-}(P, T)|$ in $\mathcal{O}(1 + m / \log_\sigma n)$ time.*

*Proof.* First, using Proposition 6.12, we compute $s = \text{L-head}(P)$, $H = \text{L-root}(P)$, $k = \text{L-exp}(P)$, and $t = \text{L-tail}(P)$ in $\mathcal{O}(1 + m/\log_\sigma n)$ time. This lets us determine $e(P) = 1 + s + k|H| + t$. If $e(P) \leq m$, then by Lemma 6.14, we return $|\text{Occ}^{\mathsf{a}-}(P,T)| = 0$. Otherwise, using the array $L_{\text{range}}$, we compute in $\mathcal{O}(1)$ time a pair of integers $b, e$ such that $\text{SA}(b \mathinner{\ldotp\ldotp} e)$ contains the starting positions of all suffixes of $T$ prefixed with $X = P[1 \mathinner{\ldotp\ldotp} 3\tau - 1]$. Equivalently, by Lemma 5.11 (see also the implementation of queries in Proposition 5.19), $\text{SA}(b \mathinner{\ldotp\ldotp} e)$ contains all positions from $\mathsf{R}_{s,H}$. If $b = e$, then it holds $\mathsf{R}_{s,H} = \emptyset$, and thus we return $|\text{Occ}^{\mathsf{a}-}(P,T)| = 0$. Let us thus assume $b < e$. Our goal now is to determine the subrange of $\text{SA}(b \mathinner{\ldotp\ldotp} e)$ containing all positions in $\{j \in \mathsf{R}^-_{s,H} : \text{L-exp}(j) > \text{L-exp}(P)\}$ (these positions form a subrange by Lemma 5.11). For that, we first compute $d = \mathsf{rank}_{B_{\exp},1}(e) - \mathsf{rank}_{B_{\exp},1}(b)$ in $\mathcal{O}(1)$ time. If $d = 0$, then $\mathsf{R}^-_{s,H} = \emptyset$, and hence we return $|\text{Occ}^{\mathsf{a}-}(P,T)| = 0$. Otherwise, we retrieve $k_{\min} = L_{\text{minexp}}[\text{int}(X)]$ in $\mathcal{O}(1)$ time. Then, letting $k_{\max} = k_{\min} + d - 1$, we have $k_{\min} \leq k_{\max}$ and $[k_{\min} \mathinner{\ldotp\ldotp} k_{\max}] = \{\text{L-exp}(j) : j \in \mathsf{R}^-_{s,H}\}$ (see the proof of Proposition 5.19). If $k \geq k_{\max}$, by Lemma 6.14, we return $|\text{Occ}^{\mathsf{a}-}(P,T)| = 0$. Otherwise, we have two cases. Let $p = \mathsf{rank}_{B_{\exp},1}(b)$. If $k < k_{\min}$, then we return $|\text{Occ}^{\mathsf{a}-}(P,T)| = |\mathsf{R}^-_{s,H}| = \mathsf{select}_{B_{\exp},1}(p + d) - b$. Otherwise (i.e., $k \geq k_{\min}$), we return $|\text{Occ}^{\mathsf{a}-}(P,T)| = \mathsf{select}_{B_{\exp},1}(p + d) - \mathsf{select}_{B_{\exp},1}(p + k - k_{\min} + 1)$. In total, the query takes $\mathcal{O}(1 + m/\log_\sigma n)$ time. $\square$

Next, we now describe the algorithm compute $|\text{Occ}^{\mathsf{s}-}(P,T)|$ for any periodic pattern $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$.

**Lemma 6.16.** *Let $P \in [0 \mathinner{\ldotp\ldotp} \sigma)^m$ be a periodic pattern. Denote $H = \text{L-root}(P)$. Assume $i \in \mathsf{R}^-_H$ and let $\ell = e(i) - i - 3\tau + 2$. Then, $|\text{Occ}^{\mathsf{s}-}(P,T) \cap [i \mathinner{\ldotp\ldotp} i + \ell]| \leq 1$. Moreover, $|\text{Occ}^{\mathsf{s}-}(P,T) \cap [i \mathinner{\ldotp\ldotp} i + \ell]| = 1$ holds if and only if $P[e^{\text{full}}(P) \mathinner{\ldotp\ldotp} m]$ is a prefix of $T[e^{\text{full}}(i) \mathinner{\ldotp\ldotp} n]$ and $e^{\text{full}}(i) - i \geq e^{\text{full}}(P) - 1$.*

*Proof.* As observed in the proof of Lemma 5.20, $[i \mathinner{\ldotp\ldotp} i + \ell) \subseteq \mathsf{R}^-_H$, and for any $\delta \in [0 \mathinner{\ldotp\ldotp} \ell)$, it holds $e(i + \delta) = e(i)$, $\text{L-tail}(i + \delta) = \text{L-tail}(i)$, and consequently, $e^{\text{full}}(i + \delta) = e^{\text{full}}(i)$ and $e^{\text{full}}(i + \delta) - (i + \delta) = e^{\text{full}}(i) - i - \delta$. Moreover, by definition of $\text{Occ}^{\mathsf{s}-}(P,T)$, letting $\text{L-head}(P) = s$, for any $j \in \text{Occ}^{\mathsf{s}-}(P,T)$ it holds $e^{\text{full}}(j) - j = s + \text{L-exp}(j) \cdot |H| = s + \text{L-exp}(P) \cdot |H| = e^{\text{full}}(P) - 1$. Thus, $i + \delta \in \text{Occ}^{\mathsf{s}-}(P,T)$ implies $e^{\text{full}}(i + \delta) - (i + \delta) = e^{\text{full}}(i) - (i + \delta) = e^{\text{full}}(P) - 1$, or equivalently, $\delta = (e^{\text{full}}(i) - i) - (e^{\text{full}}(P) - 1)$, and therefore, $|\text{Occ}^{\mathsf{s}-}(P,T) \cap [i \mathinner{\ldotp\ldotp} i + \ell]| \leq 1$.

For the second part, assume first that $i + \delta \in \text{Occ}^{\mathsf{s}-}(P,T)$ holds for some $\delta \in [0 \mathinner{\ldotp\ldotp} \ell]$. Then, as noted above, we have $e^{\text{full}}(P) - 1 = e^{\text{full}}(i) - (i + \delta) \leq e^{\text{full}}(i) - i$. Moreover, letting $\text{L-head}(P) = s$, by definition of $\text{Occ}^{\mathsf{s}-}(P,T)$, we have $i + \delta \in \mathsf{R}^-_{s,H}$, $\text{L-exp}(P) = \text{L-exp}(i + \delta)$, and $T[i + \delta \mathinner{\ldotp\ldotp} i + \delta + m) = P$. Therefore, we obtain that $T[i + \delta \mathinner{\ldotp\ldotp} e^{\text{full}}(i + \delta)) = T[i + \delta \mathinner{\ldotp\ldotp} e^{\text{full}}(i)) = P[1 \mathinner{\ldotp\ldotp} e^{\text{full}}(P)) = H'H^k$ (where $k = \text{L-exp}(P)$ and $H'$ is the length-$s$ suffix of $H$), and consequently, $P[e^{\text{full}}(P) \mathinner{\ldotp\ldotp} m]$ is a prefix of $T[e^{\text{full}}(i) \mathinner{\ldotp\ldotp} n]$. To show the converse implication, assume that $P[e^{\text{full}}(P) \mathinner{\ldotp\ldotp} m]$ is a prefix of $T[e^{\text{full}}(i) \mathinner{\ldotp\ldotp} n]$ and $e^{\text{full}}(i) - i \geq e^{\text{full}}(P) - 1$. Let $\delta = (e^{\text{full}}(i) - i) - (e^{\text{full}}(P) - 1)$. We will prove that $\delta \in [0 \mathinner{\ldotp\ldotp} \ell)$ and $i + \delta \in \text{Occ}^{\mathsf{s}-}(P,T)$. Clearly $\delta \geq 0$. To show $\delta < \ell$, we first prove $e(i) - e^{\text{full}}(i) \geq e(P) - e^{\text{full}}(P)$. Suppose that $q = e(i) - e^{\text{full}}(i) < e(P) - e^{\text{full}}(P)$. By $i \in \mathsf{R}^-_H$, we then either have $e^{\text{full}}(i) + q = n + 1$, or $e^{\text{full}}(i) + q \leq n$ and $T[e^{\text{full}}(i) + q] \neq T[e^{\text{full}}(i) + q - |H|] = P[e^{\text{full}}(P) + q - |H|] = P[e^{\text{full}}(P) + q]$, both of which contradict that $P[e^{\text{full}}(P) \mathinner{\ldotp\ldotp} m]$ is a prefix of $T[e^{\text{full}}(i) \mathinner{\ldotp\ldotp} n]$. Thus, $e(i) - e^{\text{full}}(i) \geq e(P) - e^{\text{full}}(P)$. This implies, $e(i) - (i + \delta) = (e^{\text{full}}(i) - (i + \delta)) + (e(i) - e^{\text{full}}(i)) = (e^{\text{full}}(P) - 1) + (e(i) - e^{\text{full}}(i)) \geq (e^{\text{full}}(P) - 1) + (e(P) - e^{\text{full}}(P)) = e(P) - 1 \geq 3\tau - 1$, or equivalently $\delta \leq e(i) - i - 3\tau + 1 < \ell$. To show $i + \delta \in \text{Occ}^{\mathsf{s}-}(P,T)$, it remains to observe that $e^{\text{full}}(i + \delta) - (i + \delta) = e^{\text{full}}(i) - (i + \delta) = e^{\text{full}}(P) - 1$ and $\text{L-root}(i + \delta) = \text{L-root}(i) = H = \text{L-root}(P)$ (following from Lemma 5.12) imply $T[i + \delta \mathinner{\ldotp\ldotp} e^{\text{full}}(i)) = P[1 \mathinner{\ldotp\ldotp} e^{\text{full}}(P))$. This in particular gives, letting $\text{L-head}(P) = s$, that $i + \delta \in \mathsf{R}_{s,H}$ and $\text{L-exp}(i + \delta) = \text{L-exp}(P)$. Moreover, combining it with $P[e^{\text{full}}(P) \mathinner{\ldotp\ldotp} m]$ being a prefix of $T[e^{\text{full}}(i) \mathinner{\ldotp\ldotp} n]$ yields $T[i + \delta \mathinner{\ldotp\ldotp} i + \delta + m) = P$. Finally, by Lemma 5.12, $\text{type}(i + \delta) = \text{type}(i) = -1$. Therefore, $i + \delta \in \text{Occ}^{\mathsf{s}-}(P,T)$. $\square$

**Proposition 6.17.** *Let $P \in [0\mathinner{.\,.}\sigma)^m$ be a periodic pattern. Given the data structure from Section 6.3.2 and the packed representation of $P$, we can compute $|\mathrm{Occ}^{\mathsf{s}-}(P,T)|$ in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time.*

*Proof.* First, using Proposition 6.12, we compute $s = \text{L-head}(P)$, $H = \text{L-root}(P)$, and $k = \text{L-exp}(P)$ in $\mathcal{O}(1 + m/\log_\sigma n)$ time. This lets us determine $e^{\mathrm{full}}(P) = 1 + s + k|H|$ and $P' := P[e^{\mathrm{full}}(P) - |\mathrm{pow}(H)|\mathinner{.\,.}m]$. Then, using Proposition 4.4, we compute in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time a range $(b_{\mathrm{pre}}\mathinner{.\,.}e_{\mathrm{pre}}] = \{i \in [1\mathinner{.\,.}q] : P' \text{ is a prefix of } T[A_{\mathsf{Z}}[i]\mathinner{.\,.}n]\}$. Observe that the set $\{r_i^{\mathrm{lex}} : i \in (b_{\mathrm{pre}}\mathinner{.\,.}e_{\mathrm{pre}}]\}$ consists of all positions $j \in \mathsf{R}'^{-}_H$ such that $P[e^{\mathrm{full}}(P)\mathinner{.\,.}m]$ is a prefix of $T[e^{\mathrm{full}}(j)\mathinner{.\,.}n]$. Thus, by Lemma 6.16, we have $|\mathrm{Occ}^{\mathsf{s}-}(P,T)| = |\{i \in (b_{\mathrm{pre}}\mathinner{.\,.}e_{\mathrm{pre}}] : e^{\mathrm{full}}(r_i^{\mathrm{lex}}) - r_i^{\mathrm{lex}} \geq e^{\mathrm{full}}(P) - 1\}|$, which we compute in $\mathcal{O}(\log\log n)$ time using the range counting structure as $\mathsf{rcount}_{A_{\mathrm{len}}}(e^{\mathrm{full}}(P) - 1, e_{\mathrm{pre}}) - \mathsf{rcount}_{A_{\mathrm{len}}}(e^{\mathrm{full}}(P) - 1, b_{\mathrm{pre}})$ (recall, that $A_{\mathrm{len}}[i] = e^{\mathrm{full}}(r_i^{\mathrm{lex}}) - r_i^{\mathrm{lex}}$; see Section 5.3.2). $\square$

By combining all above results, we obtain the following algorithm to compute $|\mathrm{Occ}(P,T)|$ for any periodic pattern $P$.

**Proposition 6.18.** *Let $P \in [0\mathinner{.\,.}\sigma)^m$ be a periodic pattern. Given the data structure from Section 6.3.2 and the packed representation of $P$, we can compute $|\mathrm{Occ}(P,T)|$ in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time.*

*Proof.* Given a packed representation of a periodic pattern $P$, we compute $|\mathrm{Occ}(P,T)| = |\mathrm{Occ}^{\mathsf{a}-}(P,T)| + |\mathrm{Occ}^{\mathsf{s}-}(P,T)| + |\mathrm{Occ}^{\mathsf{a}+}(P,T)| + |\mathrm{Occ}^{\mathsf{s}+}(P,T)|$ using Propositions 6.15 and 6.17 and their symmetric counterparts (adapted according to Lemma 6.9). The total time is $\mathcal{O}(m/\log_\sigma n + \log\log n)$. $\square$

**Generalizing the Query Algorithm**  We now show how to generalize the above algorithms to compute $|\mathrm{Occ}(P,T)|$, to instead return $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$.

For any periodic pattern $P \in [0\mathinner{.\,.}\sigma)^m$ we define

$$\mathrm{Pos}(P,T) = \{j \in [1\mathinner{.\,.}n] : \mathrm{lcp}(T[j\mathinner{.\,.}n], P) \geq 3\tau - 1 \text{ and } T[j\mathinner{.\,.}n] \prec P\},$$

and denote $\delta(P,T) = |\mathrm{Pos}(P,T)|$.

**Lemma 6.19.** *Let $P \in [0\mathinner{.\,.}\sigma)^m$ be a periodic pattern and let $X = P[1\mathinner{.\,.}3\tau-1]$. Then, it holds $\mathrm{RangeBeg}(P,T) = \mathrm{RangeBeg}(X,T) + \delta(P,T)$.*

*Proof.* It suffices to observe that $j \in \mathrm{Occ}(X,T)$ holds if and only if $\mathrm{lcp}(T[j\mathinner{.\,.}n], P) \geq 3\tau - 1$. Thus, it holds by definition of $\mathrm{RangeBeg}(P,T)$ that $\mathrm{RangeBeg}(P,T) = \mathrm{RangeBeg}(X,T) + |\{j \in \mathrm{Occ}(X,T) : T[j\mathinner{.\,.}n] \prec P\}| = \mathrm{RangeBeg}(X,T) + |\{j \in [1\mathinner{.\,.}n] : \mathrm{lcp}(T[j\mathinner{.\,.}n], P) \geq 3\tau - 1 \text{ and } T[j\mathinner{.\,.}n] \prec P\}| = \mathrm{RangeBeg}(X,T) + \delta(P,T)$. $\square$

We focus on computing $\delta(P,T)$ for $P$ satisfying $\mathrm{type}(P) = -1$ (the structure for $P$ satisfying $\mathrm{type}(P) = +1$ is symmetric; see the proof of Proposition 6.24). We define $\mathrm{Pos}^{\mathsf{a}}(P,T) = \{j \in \mathsf{R}^-_{s,H} : \text{L-exp}(j) \leq \text{L-exp}(P)\}$ and $\mathrm{Pos}^{\mathsf{s}}(P,T) = \{j \in \mathsf{R}^-_{s,H} : \text{L-exp}(j) = \text{L-exp}(P) \text{ and } T[j\mathinner{.\,.}n] \succeq P\}$, where $s = \text{L-head}(P)$ and $H = \text{L-root}(P)$. We denote $\delta^{\mathsf{a}}(P,T) = |\mathrm{Pos}^{\mathsf{a}}(P,T)|$ and $\delta^{\mathsf{s}}(P,T) = |\mathrm{Pos}^{\mathsf{s}}(P,T)|$.

**Lemma 6.20.** *For any periodic pattern $P \in [0\mathinner{.\,.}\sigma)^m$ that satisfies $\mathrm{type}(P) = -1$, it holds $\delta(P,T) = \delta^{\mathsf{a}}(P,T) - \delta^{\mathsf{s}}(P,T)$.*

*Proof.* We will prove that $\mathrm{Pos}^{\mathsf{a}}(P,T)$ is a disjoint union of $\mathrm{Pos}(P,T)$ and $\mathrm{Pos}^{\mathsf{s}}(P,T)$. This implies $\delta(P,T) + \delta^{\mathsf{s}}(P,T) = \delta^{\mathsf{a}}(P,T)$, and consequently, the equality in the claim.

Denote $s = \text{L-head}(P)$ and $H = \text{L-root}(P)$. By Lemma 6.9, letting $j \in \mathsf{R}_{s,H}^{-}$, we have $\mathrm{Pos}(P,T) = \{j \in \mathsf{R}_{s,H}^{-} : T[j\mathinner{\ldotp\ldotp}n] \prec P\}$, and moreover, if $j \in \mathrm{Pos}(P,T)$, then $e(j) - j \leq e(P) - 1$. In particular, $\text{L-exp}(j) = \lfloor \frac{e(j)-j-s}{|H|} \rfloor \leq \lfloor \frac{e(P)-1-s}{|H|} \rfloor = \text{L-exp}(P)$. Hence, $\mathrm{Pos}(P,T) \subseteq \mathrm{Pos}^{\mathsf{a}}(P,T)$. On the other hand, clearly $\mathrm{Pos}^{\mathsf{s}}(P,T) \subseteq \mathrm{Pos}^{\mathsf{a}}(P,T)$ and $\mathrm{Pos}^{\mathsf{s}}(P,T) \cap \mathrm{Pos}(P,T) = \emptyset$. Thus, to obtain the claim, it suffices to show that $\mathrm{Pos}^{\mathsf{a}}(P,T) \setminus \mathrm{Pos}^{\mathsf{s}}(P,T) \subseteq \mathrm{Pos}(P,T)$.

Let $j \in \mathrm{Pos}^{\mathsf{a}}(P,T) \setminus \mathrm{Pos}^{\mathsf{s}}(P,T)$. Consider two cases. If $\text{L-exp}(j) = \text{L-exp}(P)$, then by definition of $\mathrm{Pos}^{\mathsf{s}}(P,T)$, it must hold $T[j\mathinner{\ldotp\ldotp}n] \prec P$. Thus, we have $j \in \mathrm{Pos}(P,T)$. Let us therefore assume $\text{L-exp}(j) < \text{L-exp}(P)$. Then, $e(j) - j = s + \text{L-exp}(j) \cdot |H| + \text{L-tail}(j) < s + \text{L-exp}(j) \cdot |H| + |H| \leq s + \text{L-exp}(P) \cdot |H| \leq s + \text{L-exp}(P) \cdot |H| + \text{L-tail}(P) = e(P) - 1$. By Lemma 6.9(2) and Lemma 6.9(4), this implies $T[j\mathinner{\ldotp\ldotp}n] \prec P$, and consequently, $j \in \mathrm{Pos}(P,T)$. $\square$

We now describe how, given any periodic pattern $P \in [0\mathinner{\ldotp\ldotp}\sigma)^m$ that satisfies $\mathrm{type}(P) = -1$, to compute $\delta^{\mathsf{a}}(P,T)$.

**Proposition 6.21.** *Let $P \in [0\mathinner{\ldotp\ldotp}\sigma)^m$ be a periodic pattern satisfying $\mathrm{type}(P) = -1$. Given the data structure from Section 6.3.2 and the packed representation of $P$, we can in $\mathcal{O}(1 + m/\log_\sigma n)$ time compute $\delta^{\mathsf{a}}(P,T)$.*

*Proof.* First, using Proposition 6.12, we compute $H = \text{L-root}(P)$ and $k = \text{L-exp}(P)$ in $\mathcal{O}(1 + m/\log_\sigma n)$ time. Then, using $L_{\mathrm{range}}$, we compute in $\mathcal{O}(1)$ time a pair of integers $b, e$ such that $\mathrm{SA}(b\mathinner{\ldotp\ldotp}e]$ contains the starting positions of all suffixes of $T$ prefixed with $X = P[1\mathinner{\ldotp\ldotp}3\tau-1]$. Equivalently, by Lemma 6.9, $\mathrm{SA}(b\mathinner{\ldotp\ldotp}e]$ contains all positions from $\mathsf{R}_{s,H}$, where $s = \text{L-head}(P)$. If $b = e$, then it holds $\mathsf{R}_{s,H} = \emptyset$, and thus we return $\delta^{\mathsf{a}}(P,T) = 0$. Let us thus assume $b < e$. Our goal now is to determine the subrange of $\mathrm{SA}(b\mathinner{\ldotp\ldotp}e]$ containing all positions in $\{j \in \mathsf{R}_{s,H}^{-} : \text{L-exp}(j) \leq \text{L-exp}(P)\}$ (these positions form a subrange by Lemma 5.11). For that, we first compute $d = \mathsf{rank}_{B_{\exp},1}(e) - \mathsf{rank}_{B_{\exp},1}(b)$ in $\mathcal{O}(1)$ time. If $d = 0$, then $\mathsf{R}_{s,H}^{-} = \emptyset$, and hence we return $\delta^{\mathsf{a}}(P,T) = 0$. Otherwise, we retrieve $k_{\min} = L_{\mathrm{minexp}}[\mathrm{int}(X)]$ in $\mathcal{O}(1)$ time. Then, letting $k_{\max} = k_{\min} + d - 1$, we have $k_{\min} \leq k_{\max}$ and $[k_{\min}\mathinner{\ldotp\ldotp}k_{\max}] = \{\text{L-exp}(j) : j \in \mathsf{R}_{s,H}^{-}\}$ (see the proof of Proposition 5.19). If $k < k_{\min}$, we return $\delta^{\mathsf{a}}(P,T) = 0$. Otherwise, we have two cases. Let $p = \mathsf{rank}_{B_{\exp},1}(b)$. If $k \geq k_{\max}$, then we return $\delta^{\mathsf{a}}(P,T) = |\mathsf{R}_{s,H}^{-}| = \mathsf{select}_{B_{\exp},1}(p+d) - b$. Otherwise (i.e., $k < k_{\max}$), we return $\delta^{\mathsf{a}}(P,T) = \mathsf{select}_{B_{\exp},1}(p + k - k_{\min} + 1) - b$. In total, the query takes $\mathcal{O}(1 + m/\log_\sigma n)$ time. $\square$

We now describe how, given any periodic pattern $P \in [0\mathinner{\ldotp\ldotp}\sigma)^m$ that satisfies $\mathrm{type}(P) = -1$, to compute $\delta^{\mathsf{s}}(P,T)$.

**Lemma 6.22.** *Let $P \in [0\mathinner{\ldotp\ldotp}\sigma)^m$ be a periodic pattern that satisfies $\mathrm{type}(P) = -1$. Denote $H = \text{L-root}(P)$. Assume $i \in \mathsf{R}_{H}^{-}$ and let $\ell = e(i) - i - 3\tau + 2$. Then, we have $|\mathrm{Pos}^{\mathsf{s}}(P,T) \cap [i\mathinner{\ldotp\ldotp}i+\ell]| \leq 1$. Moreover, $|\mathrm{Pos}^{\mathsf{s}}(P,T) \cap [i\mathinner{\ldotp\ldotp}i+\ell]| = 1$ holds if and only if $T[e^{\mathrm{full}}(i)\mathinner{\ldotp\ldotp}n] \succeq P[e^{\mathrm{full}}(P)\mathinner{\ldotp\ldotp}m]$ and $e^{\mathrm{full}}(i) - i \geq e^{\mathrm{full}}(P) - 1$.*

*Proof.* In the proof of Lemma 6.16, it is shown that $[i\mathinner{\ldotp\ldotp}i+\ell] \subseteq \mathsf{R}_{H}^{-}$, and for any $\delta \in [0\mathinner{\ldotp\ldotp}\ell]$, it holds $e^{\mathrm{full}}(i+\delta) - (i+\delta) = e^{\mathrm{full}}(i) - i - \delta$. By definition of $\mathrm{Pos}^{\mathsf{s}}(P,T)$, letting $s = \text{L-head}(P)$, for any $j \in \mathrm{Pos}^{\mathsf{s}}(P,T)$ it holds $e^{\mathrm{full}}(j) - j = s + \text{L-exp}(j) \cdot |H| = s + \text{L-exp}(P) \cdot |H| = e^{\mathrm{full}}(P) - 1$. Thus, $i + \delta \in \mathrm{Pos}^{\mathsf{s}}(P,T)$ implies $e^{\mathrm{full}}(i+\delta) - (i+\delta) = e^{\mathrm{full}}(i) - (i+\delta) = e^{\mathrm{full}}(P) - 1$, or equivalently, $\delta = (e^{\mathrm{full}}(i) - i) - (e^{\mathrm{full}}(P) - 1)$, and therefore, $|\mathrm{Pos}^{\mathsf{s}}(P,T) \cap [i\mathinner{\ldotp\ldotp}i+\ell]| \leq 1$.

For the second part, assume first that $i + \delta \in \mathrm{Pos}^{\mathsf{s}}(P,T)$ holds for some $\delta \in [0\mathinner{\ldotp\ldotp}\ell]$. Then, as noted above, we have $e^{\mathrm{full}}(P) - 1 = e^{\mathrm{full}}(i) - (i+\delta) \leq e^{\mathrm{full}}(i) - i$. Moreover, letting $\text{L-head}(P) = s$, by definition of $\mathrm{Pos}^{\mathsf{s}}(P,T)$, we have $i + \delta \in \mathsf{R}_{s,H}^{-}$, $\text{L-exp}(P) = \text{L-exp}(i+\delta)$, and $T[i+\delta\mathinner{\ldotp\ldotp}n] \succeq P$.

Therefore, we obtain that $T[i+\delta \mathinner{.\,.} e^{\mathrm{full}}(i+\delta)) = T[i+\delta \mathinner{.\,.} e^{\mathrm{full}}(i)) = P[1 \mathinner{.\,.} e^{\mathrm{full}}(P)) = H'H^k$ (where $k = \mathrm{L\text{-}exp}(P)$ and $H'$ is the length-$s$ suffix of $H$), and consequently, $T[e^{\mathrm{full}}(i) \mathinner{.\,.} n] \succeq P[e^{\mathrm{full}}(P) \mathinner{.\,.} m]$. To show the converse implication, assume that $T[e^{\mathrm{full}}(i) \mathinner{.\,.} n] \succeq P[e^{\mathrm{full}}(P) \mathinner{.\,.} m]$ and $e^{\mathrm{full}}(i) - i \geq e^{\mathrm{full}}(P) - 1$. Let $\delta = (e^{\mathrm{full}}(i) - i) - (e^{\mathrm{full}}(P) - 1)$. We will prove that $\delta \in [0 \mathinner{.\,.} \ell)$ and $i+\delta \in \mathrm{Pos}^{\mathsf{s}}(P,T)$. Clearly $\delta \geq 0$. To show $\delta < \ell$, we first prove $e(i) - e^{\mathrm{full}}(i) \geq e(P) - e^{\mathrm{full}}(P)$. Suppose that $q = e(i) - e^{\mathrm{full}}(i) < e(P) - e^{\mathrm{full}}(P)$. By $i \in \mathsf{R}_H^-$, we then either have $e^{\mathrm{full}}(i) + q = n+1$, or $e^{\mathrm{full}}(i) + q \leq n$ and $T[e^{\mathrm{full}}(i)+q] \prec T[e^{\mathrm{full}}(i)+q-|H|] = P[e^{\mathrm{full}}(P)+q-|H|] = P[e^{\mathrm{full}}(P)+q]$, both of which contradict $T[e^{\mathrm{full}}(i) \mathinner{.\,.} n] \succeq P[e^{\mathrm{full}}(P) \mathinner{.\,.} m]$. Thus, $e(i) - e^{\mathrm{full}}(i) \geq e(P) - e^{\mathrm{full}}(P)$. This implies, $e(i) - (i+\delta) = (e^{\mathrm{full}}(i) - (i+\delta)) + (e(i) - e^{\mathrm{full}}(i)) = (e^{\mathrm{full}}(P) - 1) + (e(i) - e^{\mathrm{full}}(i)) \geq (e^{\mathrm{full}}(P) - 1) + (e(P) - e^{\mathrm{full}}(P)) = e(P) - 1 \geq 3\tau - 1$, or equivalently $\delta \leq e(i) - i - 3\tau + 1 < \ell$. To show $i+\delta \in \mathrm{Pos}^{\mathsf{s}}(P,T)$, it remains to observe that $e^{\mathrm{full}}(i+\delta) - (i+\delta) = e^{\mathrm{full}}(i) - (i+\delta) = e^{\mathrm{full}}(P) - 1$ and $\mathrm{L\text{-}root}(i+\delta) = \mathrm{L\text{-}root}(i) = H = \mathrm{L\text{-}root}(P)$ (following from Lemma 5.12) imply $T[i+\delta \mathinner{.\,.} e^{\mathrm{full}}(i)) = P[1 \mathinner{.\,.} e^{\mathrm{full}}(P))$. This in particular gives, letting $\mathrm{L\text{-}head}(P) = s$, that $i+\delta \in \mathsf{R}_{s,H}$ and $\mathrm{L\text{-}exp}(i+\delta) = \mathrm{L\text{-}exp}(P)$. Moreover, combining it with $T[e^{\mathrm{full}}(i) \mathinner{.\,.} n] \succeq P[e^{\mathrm{full}}(P) \mathinner{.\,.} m]$ yields $T[i+\delta \mathinner{.\,.} n] \succeq P$. Finally, by Lemma 5.12, $\mathrm{type}(i+\delta) = \mathrm{type}(i) = -1$. Therefore, $i+\delta \in \mathrm{Pos}^{\mathsf{s}}(P,T)$. $\qquad\square$

**Proposition 6.23.** *Let $P \in [0 \mathinner{.\,.} \sigma)^m$ be a periodic pattern satisfying* $\mathrm{type}(P) = -1$. *Given the data structure from Section 6.3.2 and the packed representation of $P$, we can in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time compute* $\delta^{\mathsf{s}}(P,T)$.

*Proof.* First, using Proposition 6.12, we compute $s = \mathrm{L\text{-}head}(P)$, $H = \mathrm{L\text{-}root}(P)$, and $k = \mathrm{L\text{-}exp}(P)$ in $\mathcal{O}(1 + m/\log_\sigma n)$ time. This lets us determine $e^{\mathrm{full}}(P) = 1+s+k|H|$ and $P' := P[e^{\mathrm{full}}(P) - |\mathrm{pow}(H)| \mathinner{.\,.} m]$. Then, using Proposition 4.4, we compute in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time a value $x = |\{i \in [1 \mathinner{.\,.} q] : T[A_{\mathsf{Z}}[i] \mathinner{.\,.} n] \prec P'\}|$. Then, letting $x' = \sum_{H' \preceq H} |\mathsf{R}_{H'}^{\prime-}|$ (obtained from $L_{\mathrm{runs}}$ in $\mathcal{O}(1)$ time as explained in the proof of Proposition 5.21), by definition of $A_{\mathsf{Z}}$ and properties of function pow (see the proof of Proposition 6.17), the set $\{r_i^{\mathrm{lex}} : i \in (x \mathinner{.\,.} x']\}$ (where $r_i^{\mathrm{lex}}$ is defined as in the proof of Proposition 6.17) consists of all positions $j \in \mathsf{R}_H^{\prime-}$ satisfying $T[e^{\mathrm{full}}(j) \mathinner{.\,.} n] \succeq P[e^{\mathrm{full}}(P) \mathinner{.\,.} m]$. Thus, by Lemma 6.22, it holds $\delta^{\mathsf{s}}(P,T) = |\mathrm{Pos}^{\mathsf{s}}(P,T)| = |\{i \in (x \mathinner{.\,.} x'] : \ell_i \geq e^{\mathrm{full}}(P) - 1\}|$ (where $\ell_i$ is defined as in Proposition 6.17), which we compute in $\mathcal{O}(\log\log n)$ time using the range counting structure as $\mathsf{rcount}_{A_{\mathrm{len}}}(e^{\mathrm{full}}(P) - 1, x') - \mathsf{rcount}_{A_{\mathrm{len}}}(e^{\mathrm{full}}(P) - 1, x)$. $\qquad\square$

By combining the above results, we obtain the algorithm to efficiently compute the pair $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$ for periodic patterns.

**Proposition 6.24.** *Let $P \in [0 \mathinner{.\,.} \sigma)^m$ be a periodic pattern. Given the data structure from Section 6.3.2 and the packed representation of $P$, we can in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time compute the pair* $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$.

*Proof.* First, using Proposition 6.18 in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time we compute $|\mathrm{Occ}(P,T)|$. Next, using the lookup table $L_{\mathrm{range}}$, in $\mathcal{O}(1)$ time we compute $(b_X, e_X) = (\mathrm{RangeBeg}(X,T), \mathrm{RangeEnd}(X,T))$, where $X = P[1 \mathinner{.\,.} 3\tau - 1]$. Then, in $\mathcal{O}(1 + m/\log_\sigma n)$ time using Proposition 6.12 we determine $\mathrm{type}(P)$. Depending on whether $\mathrm{type}(P) = -1$ or $\mathrm{type}(P) = +1$, we use either a combination of Propositions 6.21 and 6.23, or their symmetric counterparts (more precisely, if $\mathrm{type}(P) = +1$, we have $\delta^{\mathsf{a}}(P,T) = |\mathrm{Pos}^{\mathsf{a}}(P,T)|$ and $\delta^{\mathsf{s}}(P,T) = |\mathrm{Pos}^{\mathsf{s}}(P,T)|$, where $\mathrm{Pos}^{\mathsf{a}}(P,T) = \{j \in \mathsf{R}_{s,H}^+ : \mathrm{L\text{-}exp}(j) \leq \mathrm{L\text{-}exp}(P)\}$ and $\mathrm{Pos}^{\mathsf{s}}(P,T) = \{j \in \mathsf{R}_{s,H}^+ : \mathrm{L\text{-}exp}(j) = \mathrm{L\text{-}exp}(j)$ and $T[j \mathinner{.\,.} n] \prec P\})$, to compute $\delta^{\mathsf{a}}(P,T)$ and $\delta^{\mathsf{s}}(P,T)$ in $\mathcal{O}(1 + m/\log_\sigma n)$ and $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time, respectively. If $\mathrm{type}(P) = -1$, then by Lemma 6.20 we have $\delta(P,T) = \delta^{\mathsf{a}}(P,T) - \delta^{\mathsf{s}}(P,T)$. Otherwise, by the counterpart of Lemma 6.20, $\delta(P,T) = (e_X - b_X) - (\delta^{\mathsf{a}}(P,T) - \delta^{\mathsf{s}}(P,T))$. Finally, we return $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T)) =$

$(b_X + \delta(P,T), b_X + \delta(P,T) + |\mathrm{Occ}(P,T)|)$ (see Lemma 6.19) as the answer. In total, the query takes $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time. $\qquad\square$

*Remark* 6.25. Note the subtle difference in the type of symmetry used during the computation of $|\mathrm{Pos}(P,T)|$ and $|\mathrm{Occ}(P,T)|$. When computing $\delta(P,T) = |\mathrm{Pos}(P,T)|$, by Lemma 6.9 we have $\mathrm{Pos}(P,T) \subseteq \mathsf{R}^-$ for any $P$ satisfying $\mathrm{type}(P) = -1$ (and $\mathrm{Pos}(P,T) \subseteq \mathsf{R}^+$ for $P$ satisfying $\mathrm{type}(P) = +1$). However, when computing $|\mathrm{Occ}(P,T)|$ for $P$ satisfying $\mathrm{type}(P) = -1$, it is possible that $\mathrm{Occ}(P,T) \cap \mathsf{R}^- \neq \emptyset$ and $\mathrm{Occ}(P,T) \cap \mathsf{R}^+ \neq \emptyset$. Consequently, during the computation of $|\mathrm{Occ}(P,T)|$, we partition the output set $\mathrm{Occ}^{\mathsf{a}}(P,T)$ (resp. $\mathrm{Occ}^{\mathsf{s}}(P,T)$) into two subsets $\mathrm{Occ}^{\mathsf{a}-}(P,T)$ and $\mathrm{Occ}^{\mathsf{a}+}(P,T)$ (resp. $\mathrm{Occ}^{\mathsf{s}-}(P,T)$ and $\mathrm{Occ}^{\mathsf{s}+}(P,T)$), but the computation is always performed regardless of $\mathrm{type}(P)$, leading to two queries for each periodic pattern $P$. During the computation of $\delta(P,T)$, on the other hand, the computation is performed separately for $P$ satisfying $\mathrm{type}(P) = -1$ and $P$ satisfying $\mathrm{type}(P) = +1$, without the need to partition $\mathrm{Pos}(P,T)$ within each case, leading to a single query but only on the appropriate structure depending on $\mathrm{type}(P)$. This is the reason for why the seemingly related computation of $|\mathrm{Pos}(P,T)|$ and $|\mathrm{Occ}(P,T)|$ is (unlike for nonperiodic patterns; see Section 6.2.2) described separately.

### 6.3.5 Construction Algorithm

**Proposition 6.26.** *Given* $\mathrm{C}_{\mathrm{PM}}(T)$, *we can in* $\mathcal{O}(n/\log_\sigma n)$ *time augment it into a data structure from Section 6.3.2.*

*Proof.* First, we combine Propositions 5.4 and 5.27 (recall that the packed representation of $T$ is a component of $\mathrm{C}_{\mathrm{PM}}(T)$) to construct the data structure from Section 5.3.2 in $\mathcal{O}(n/\log_\sigma n)$ time. In particular, this constructs $(r_i^{\mathrm{lex}})_{i \in [1 \mathinner{.\,.} q]}$. Using Proposition 5.15, we can now compute $A_{\mathsf{Z}}[i]$ for any $i \in [1 \mathinner{.\,.} q]$ in $\mathcal{O}(1)$ time. Then, in $\mathcal{O}(n/\log_\sigma n)$ time, we construct the data structure from Proposition 4.4.

After the above components are constructed, we then analogously construct their symmetric counterparts (adapted according to Lemma 6.9). $\qquad\square$

## 6.4 The Final Data Structure

In this section, we put together Sections 6.1 to 6.3 to obtain a data structure that, given a packed representation of any pattern $P \in [0 \mathinner{.\,.} \sigma)^m$, computes $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$ in $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$ time.

The section is organized as follows. First, we introduce the components of the data structure (Section 6.4.1). Next, we describe the query algorithms (Section 6.4.2). Finally, we show the construction algorithm (Section 6.4.3).

### 6.4.1 The Data Structure

The data structure consists of two components:

1. The structure from Section 6.2.1 (used to handle nonperiodic patterns).
2. The structure from Section 6.3.2 (used to handle periodic patterns).

In total, the data structure takes $\mathcal{O}(n/\log_\sigma n)$ space.

### 6.4.2 Implementation of Queries

**Proposition 6.27.** *Given the data structure from Section 6.4.1 and the packed representation of any $P \in [0 \mathinner{.\,.} \sigma)^m$, we can in $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$ time compute the pair $(\mathrm{RangeBeg}(P,T)$, $\mathrm{RangeEnd}(P,T))$.*

*Proof.* First, using Proposition 6.2, in $\mathcal{O}(1)$ time we check if $P$ is periodic. If so, we obtain $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$ in $\mathcal{O}(m/\log_\sigma n + \log\log n)$ time using Proposition 6.24. Otherwise (i.e., if $P$ is not periodic), we consider two cases, depending on whether it holds $m < 3\tau - 1$. If so, then we obtain $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$ in $\mathcal{O}(1)$ time using Proposition 6.3. Otherwise, we obtain $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$ in $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$ time using Proposition 6.7. $\qquad\square$

### 6.4.3 Construction Algorithm

**Proposition 6.28.** *Given the packed representation of $T \in [0 \mathinner{.\,.} \sigma)^n$, we can construct the data structure from Section 6.4.1 in $\mathcal{O}(n\min(1, \log\sigma/\sqrt{\log n}))$ time and $\mathcal{O}(n/\log_\sigma n)$ working space.*

*Proof.* First, from a packed representation of $T$, we construct $\mathrm{C_{PM}}(T)$ in $\mathcal{O}(n/\log_\sigma n)$ time using Proposition 6.4. Then, using Propositions 6.8 and 6.26, we augment $\mathrm{C_{PM}}(T)$ into the two components of the structure from Section 6.4.1 in $\mathcal{O}(n\min(1, \log\sigma/\sqrt{\log n}))$ and $\mathcal{O}(n/\log_\sigma n)$ time (respectively) and using $\mathcal{O}(n/\log_\sigma n)$ working space. The overall runtime is thus $\mathcal{O}(n\min(1, \log\sigma/\sqrt{\log n}))$. $\qquad\square$

## 6.5 Summary

By combining Proposition 6.27 and Proposition 6.28 we obtain the following final result of this section.

**Theorem 6.29.** *Given any constant $\epsilon \in (0,1)$ and the packed representation of a text $T \in [0 \mathinner{.\,.} \sigma)^n$ with $2 \le \sigma < n^{1/7}$, we can in $\mathcal{O}(n\min(1, \log\sigma/\sqrt{\log n}))$ time and $\mathcal{O}(n/\log_\sigma n)$ working space construct a data structure of size $\mathcal{O}(n/\log_\sigma n)$ that, given the packed representation of any $P \in [0 \mathinner{.\,.} \sigma)^m$, returns the pair $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$ (and hence, in particular, the value $|\mathrm{Occ}(P,T)|$) in $\mathcal{O}(m/\log_\sigma n + \log^\epsilon n)$ time.*

By combining the above result with Theorem 5.31, we moreover obtain the following result.

**Theorem 6.30.** *Given any constant $\epsilon \in (0,1)$ and the packed representation of a text $T \in [0 \mathinner{.\,.} \sigma)^n$ with $2 \le \sigma < n^{1/7}$, we can in $\mathcal{O}(n\min(1, \log\sigma/\sqrt{\log n}))$ time and $\mathcal{O}(n/\log_\sigma n)$ working space construct a data structure of size $\mathcal{O}(n/\log_\sigma n)$ that, given the packed representation of any $P \in [0 \mathinner{.\,.} \sigma)^m$, returns the set $\mathrm{Occ}(P,T)$ in $\mathcal{O}(m/\log_\sigma n + (|\mathrm{Occ}(P,T)| + 1)\log^\epsilon n)$ time.*

By observing that the dominating operations in the above index are prefix rank and selection queries, we obtain the following more general result.

**Theorem 6.31.** *Consider a data structure answering prefix rank and selection queries that, for any string of length $m$ over alphabet $[0 \mathinner{.\,.} \sigma)^\ell$, achieves the following complexities:*

1. *Space usage $S(m, \ell, \sigma)$,*
2. *Preprocessing time $P_t(m, \ell, \sigma)$,*
3. *Preprocessing space $P_s(m, \ell, \sigma)$,*
4. *Query time $Q(m, \ell, \sigma)$.*

| Operation | Description |
|---|---|
| isleaf($v$) | Return true if and only if $v$ is a leaf |
| index($v$) | Any position $j \in \mathrm{Occ}(\mathrm{str}(v), T)$ |
| findleaf($j$) | The leaf $v$ satisfying $\mathrm{str}(v) = T[j\mathbin{..}n]$ |
| count($v$) | The number of leaves in the subtree rooted in $v$ |
| sdepth($v$) | The string-depth of node $v$, i.e., $|\mathrm{str}(v)|$ |
| parent($v$) | The parent of $v \neq \mathrm{root}(\mathcal{T}_{\mathrm{st}})$ |
| firstchild($v$) | The leftmost child of $v$, or $\perp$ if $v$ is a leaf |
| lastchild($v$) | The rightmost child of $v$, or $\perp$ if $v$ is a leaf |
| rightsibling($v$) | The right sibling of $v$, or $\perp$ if there is no such node |
| leftsibling($v$) | The left sibling of $v$, or $\perp$ if there is no such node |
| slink($v$) | A node $v'$ satisfying $\mathrm{str}(v') = \mathrm{str}(v)[2\mathbin{..}|\mathrm{str}(v)|]$ |
| slink($v, i$) | A node $v'$ satisfying $\mathrm{str}(v') = \mathrm{str}(v)[i{+}1\mathbin{..}|\mathrm{str}(v)|]$, i.e., iterated slink |
| wlink($v, c$) | A node $v'$ satisfying $\mathrm{str}(v') = c \cdot \mathrm{str}(v)$, or $\perp$ if there is no such node [8] |
| child($v, c$) | A child $v'$ of $v$ satisfying $\mathrm{str}(v')[|\mathrm{str}(v)|{+}1] = c$, or $\perp$ if there is no such node |
| pred($v, c$) | A node child($v, c'$), where $c' = \max\{c'' \in [0\mathbin{..}c) : \mathrm{child}(v, c'') \neq \perp\}$ (or $\perp$) |
| letter($v, i$) | The $i$th leftmost character of $\mathrm{str}(v)$ |
| WA($v, d$) | The most shallow ancestor of $v$ satisfying $\mathrm{sdepth}(v) \geq d$ |
| LCA($u, v$) | The lowest common ancestor of nodes $u$ and $v$ |
| isancestor($u, v$) | Return true if and only if $u$ is an ancestor of $v$ |

**Table 1:** Operations on suffix tree $\mathcal{T}_{\mathrm{st}}$ supported by our data structure.

*For every $T \in [0\mathbin{..}\sigma)^n$ with $2 \leq \sigma < n^{1/7}$, there exist $m = \mathcal{O}(n/\log_\sigma n)$ and $\ell = \mathcal{O}(\log_\sigma n)$ such that, given the packed representation of $T$, we can in $\mathcal{O}(n/\log_\sigma n + P_t(m, \ell, \sigma))$ time and $\mathcal{O}(n/\log_\sigma n + P_s(m, \ell, \sigma))$ working space build a structure of size $\mathcal{O}(n/\log_\sigma n + S(m, \ell, \sigma))$ that, given the packed representation of any $P \in [0\mathbin{..}\sigma)^p$, performs the following queries:*

- *Return $(\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T))$ in $\mathcal{O}(p/\log_\sigma n + \log\log n + Q(m, \ell, \sigma))$ time,*
- *Return $\mathrm{Occ}(P, T)$ in $\mathcal{O}(p/\log_\sigma n + (|\mathrm{Occ}(P, T)| + 1)(\log\log n + Q(m, \ell, \sigma)))$ time.*

# 7 Suffix Tree Queries

Let $\epsilon \in (0, 1)$ be any fixed constant and let $T \in [0\mathbin{..}\sigma)^n$, where $2 \leq \sigma < n^{1/7}$. Let $\mathcal{T}_{\mathrm{st}}$ denote the suffix tree of $T$, i.e., a compact trie of the set $\{T[1\mathbin{..}n], T[2\mathbin{..}n], \ldots, T[n]\}$. In this section, we show how given the packed representation of $T$, to construct in $\mathcal{O}(n\min(1, \log\sigma/\sqrt{\log n}))$ time and $\mathcal{O}(n/\log_\sigma n)$ working space a representation of $\mathcal{T}_{\mathrm{st}}$ occupying $\mathcal{O}(n/\log_\sigma n)$ space, and supporting each of the operations listed in Table 1 in $\mathcal{O}(\log^\epsilon n)$ time. [9] We also derive a general reduction depending on prefix rank and selection queries.

---

[8]Our data structure supports also a slightly stronger operation wlink$'(v, c)$ (see Proposition 7.63), that returns a node $v'$ satisfying $\mathrm{repr}(v') = (\mathrm{RangeBeg}(c \cdot \mathrm{str}(v), T), \mathrm{RangeEnd}(c \cdot \mathrm{str}(v), T))$, if such node exists. This generalizes wlink$(v, c)$, since wlink$(v, c) \neq \perp$ holds if and only if wlink$'(v, c) \neq \perp$ and $\mathrm{sdepth}(\mathrm{wlink}'(v, c)) = \mathrm{sdepth}(v) + 1$. Therefore, we can use wlink$'(v, c)$ to compute wlink$(v, c)$. Note, however, that it is possible that wlink$(v, c) = \perp$ and yet wlink$'(v, c) \neq \perp$. In that case, there exists a node corresponding to wlink$(v, c)$ in the suffix *trie* of $T$, but in $\mathcal{T}_{\mathrm{st}}$ this node is not explicit.

[9]Similarly as in prior CST implementations [79, 31, 77, 34, 14, 16], the time complexity of some operations is actually $\mathcal{O}(1)$. We also note that some prior CST implementations (e.g., [79, 77, 34, 16]) support two additional operations called *tree depth* and *tree level ancestor* which are analogous to sdepth($v$) and WA($v, d$) but with distance to the root defined by the number of ancestor nodes rather than the total length of edge labels.

As in Sections 5 and 6, we let $\tau = \lfloor \mu \log_\sigma n \rfloor$, where $\mu$ is any positive constant smaller than $\frac{1}{6}$ such that $\tau \geq 1$, be fixed for the duration of this section. Throughout, we also use R as a shorthand for $\mathsf{R}(\tau, T)$.

**Definition 7.1.** Let $v$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$. The node $v$ is said to be *periodic* if $\mathrm{str}(v)$ is periodic (Definition 6.1). Otherwise, $v$ is *nonperiodic*.

**Representation of a Node**  For any explicit node $v$ of $\mathcal{T}_{\mathrm{st}}$ we denote $\mathrm{Occ}(v) := \mathrm{Occ}(\mathrm{str}(v), T)$. In our data structure we represent each explicit node $v$ of $\mathcal{T}_{\mathrm{st}}$ in one of two ways:

- A pair $(j, \ell)$, where $j \in \mathrm{Occ}(v)$ (i.e., $j$ is the starting position of some occurrence of $\mathrm{str}(v)$ in $T$) and $\ell = \mathrm{sdepth}(v)$.
- A pair $(\mathrm{lrank}(v), \mathrm{rrank}(v))$. Note that since $v$ is a node of suffix tree, in this special case we have $(\mathrm{lrank}(v), \mathrm{rrank}(v)) = (\mathrm{RangeBeg}(\mathrm{str}(v), T), \mathrm{RangeEnd}(\mathrm{str}(v), T))$. Thus, letting $(b, e) = (\mathrm{lrank}(v), \mathrm{rrank}(v))$, we then have $\{\mathrm{SA}[i]\}_{i \in (b..e]} = \mathrm{Occ}(v)$. Note also that $b < e$.

In most cases, the latter representation leads to a more convenient implementation. Thus, we adopt it as a default and denote $\mathrm{repr}(v) := (\mathrm{lrank}(v), \mathrm{rrank}(v))$ (while using the first one mostly as a temporary internal representation). We also define $\mathrm{repr}(\bot) = (0, 0)$.

**Organization**  The structure and query algorithms for a node $v$ are different depending on whether $v$ is periodic (Definition 7.1). Our description is thus split as follows. First (Section 7.1), we describe the set of data structures called collectively the index "core" that enables efficiently checking if $v$ is periodic (it is also used to perform operations on nodes with very small depth and contains some common components utilized by the remaining parts). In the following two parts (Sections 7.2 and 7.3), we describe structures handling each of the two cases. All ingredients are then put together in Section 7.4. Finally, we present our result in the general form (Section 7.5).

## 7.1  The Index Core

In this section, we describe a data structure used to check in $\mathcal{O}(1)$ time if a given node is periodic. It also lets us perform operations concerning nodes at depth smaller than $3\tau - 1$ in $\mathcal{O}(1)$ time.

The section is organized as follows. First, we introduce the components of the data structure (Section 7.1.1). We then show how using this structure to implement some basic navigational routines (Section 7.1.2). Next, we describe the query algorithms for the fundamental operations (Sections 7.1.3 to 7.1.6). Finally, we show the construction algorithm (Section 7.1.7).

### 7.1.1  The Data Structure

**Definitions**  For any $k \geq 1$, let $\mathcal{S}_k := \{S \in [0 \mathinner{.\,.} \sigma)^k : S \text{ occurs in } T\}$ denote the set of length-$k$ substrings of $T$. Let $\mathcal{T}_{3\tau-1}$ denote the compact trie of $\mathcal{S}_{3\tau-1}$.

**Components**  The index core, denoted $\mathrm{C}_{\mathrm{ST}}(T)$, consists of two components:

1. The index core $\mathrm{C}_{\mathrm{SA}}(T)$ (Section 5.1.1). It takes $\mathcal{O}(n/\log_\sigma n)$ space.
2. The compact trie $\mathcal{T}_{3\tau-1}$. All nodes of $\mathcal{T}_{3\tau-1}$ are stored in an array and pointers to nodes are implemented as indexes to this array. Each node $v$ of $\mathcal{T}_{3\tau-1}$ stores the string $\mathrm{str}(v)$ encoded as an integer $\mathrm{int}(\mathrm{str}(v))$, the pointer $\mathrm{parent}(v)$, the value $\mathrm{sdepth}(v)$, and the doubly linked list containing pointers to all children of $v$, in ascending order of the first letter on the connecting edge. Since each node $v$ of $\mathcal{T}_{3\tau-1}$ corresponds to a unique string $S \in [0 \mathinner{.\,.} \sigma)^{\leq 3\tau-1}$, in total $\mathcal{T}_{3\tau-1}$ needs $\mathcal{O}(\sigma^{3\tau}) = \mathcal{O}(\sqrt{n})$ space. The trie $\mathcal{T}_{3\tau-1}$ is augmented with the following structures:

(a) A linear-space data structure answering the LCA queries in $\mathcal{T}_{3\tau-1}$ in $\mathcal{O}(1)$ time [10]. By the above bound, the data structure uses $\mathcal{O}(\sqrt{n})$ space.

(b) A lookup table $L_{\text{child}}$ that for each edge of $\mathcal{T}_{3\tau-1}$ connecting a node $v$ to its parent $p$ and labeled with a string starting with the character $c$, maps the pair $(i_p, c)$ to $i_v$, where $i_p$ and $i_v$ are pointers to $p$ and $v$. $\mathcal{T}_{3\tau-1}$ has less than $2\sigma^{3\tau-1}$ nodes and thus $i_v < 2\sigma^{3\tau-1}$. On the other hand, $c \in [0 \mathinner{.\,.} \sigma)$. Thus, each pair $(i_v, c)$ can be (in $\mathcal{O}(1)$ time) injectively mapped to an integer not exceeding $2\sigma^{3\tau} = \mathcal{O}(\sqrt{n})$ and hence $L_{\text{child}}$ needs $\mathcal{O}(n/\log_\sigma n)$ space.

(c) A lookup table $L_{\text{WA}}$ that for every node $v$ of $\mathcal{T}_{3\tau-1}$ and every $d \in [0 \mathinner{.\,.} 3\tau-1)$, maps the pair $(i_v, d)$ to $i_u$, where $u = \text{WA}(v, d)$ and $i_v$ (resp. $i_u$) is the pointer to $v$ (resp. $u$). Since $i_v < 2\sigma^{3\tau-1}$ and $\tau = \mathcal{O}(\log n)$, each pair $(i_v, d)$ can be injectively mapped to an integer not exceeding $\mathcal{O}(\sqrt{n}\log n)$ and hence $L_{\text{WA}}$ needs $\mathcal{O}(n/\log_\sigma n)$ space.

(d) An array storing the pointers to leaves of $\mathcal{T}_{3\tau-1}$ in the left-to-right order. Since the number of leaves is $\mathcal{O}(\sigma^{3\tau-1})$, the array needs $\mathcal{O}(n/\log_\sigma n)$ space.

In total, $\text{C}_{\text{ST}}(T)$ takes $\mathcal{O}(n/\log_\sigma n)$ space.

*Remark* 7.2. Note that $\mathcal{T}_{3\tau-1}$ corresponds to $\mathcal{T}_{\text{st}}$ truncated at depth $3\tau - 1$. The key reason motivating this definition is that the pair $(b, e) = (\text{RangeBeg}(\text{str}(v), T), \text{RangeEnd}(\text{str}(v), T))$ for every node $v$ of $\mathcal{T}_{3\tau-1}$ at depth $3\tau - 1$ that corresponds to an implicit node of $\mathcal{T}_{\text{st}}$ (in the middle of an edge connecting some node $v'$ of $\mathcal{T}_{\text{st}}$ to one of its children $v''$) satisfies $(b, e) = (\text{RangeBeg}(\text{str}(v''), T), \text{RangeEnd}(\text{str}(v''), T))$. In all our uses, this is sufficient, and the value $\text{sdepth}(v'')$ is never needed.

### 7.1.2 Navigation Primitives

**Mapping from $\mathcal{T}_{\text{st}}$ to $\mathcal{T}_{3\tau-1}$**  For any explicit node $v$ of $\mathcal{T}_{\text{st}}$, we define $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{3\tau-1}}(v)$ as the deepest explicit node $u$ of $\mathcal{T}_{3\tau-1}$ such that $\text{str}(u)$ is a prefix of $\text{str}(v)$.

**Lemma 7.3.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$ and $u = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{3\tau-1}}(v)$. Then, $\text{sdepth}(v) \geq 3\tau - 1$ holds if and only if $\text{sdepth}(u) = 3\tau - 1$. Moreover,*

1. *If $\text{sdepth}(v) \geq 3\tau - 1$ then $\text{str}(u) = \text{str}(v)[1 \mathinner{.\,.} 3\tau{-}1]$.*
2. *Otherwise (i.e., if $\text{sdepth}(v) < 3\tau - 1$), $\text{str}(u) = \text{str}(v)$.*

*Proof.* 1. If $\text{sdepth}(v) \geq 3\tau - 1$ then $\text{str}(v)[1 \mathinner{.\,.} 3\tau{-}1] \in \mathcal{S}_{3\tau-1}$. Therefore, by definition of $\mathcal{T}_{3\tau-1}$, there exists a node $u'$ in $\mathcal{T}_{3\tau-1}$ satisfying $\text{str}(u') = \text{str}(v)[1 \mathinner{.\,.} 3\tau{-}1]$. Since $\text{str}(u')$ is a prefix of $\text{str}(v)$ and $u'$ is a leaf of $\mathcal{T}_{3\tau-1}$, we thus have $u' = u$, and hence $\text{str}(u) = \text{str}(v)[1 \mathinner{.\,.} 3\tau{-}1]$.

2. Let $\text{sdepth}(v) < 3\tau - 1$ and $X = \text{str}(v)$. Since $v$ is explicit, there exists distinct $c, c' \in \Sigma$ such that $Xc$ and $Xc'$ occur in $T$. By $|X| < 3\tau - 1$, $\mathcal{T}_{3\tau-1}$ therefore has an explicit node $u'$ satisfying $\text{str}(u') = X$. By $\text{sdepth}(u') = \text{sdepth}(v)$, we thus have $u' = u$ and hence $\text{str}(u) = \text{str}(v)$.

The equivalence follows immediately from the two items. $\qquad\square$

**Lemma 7.4.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$. Let $i_1 = \text{lrank}(v) + 1$, $i_2 = \text{rrank}(v)$, $y_1 = \text{rank}_{B_{3\tau-1}, 1}(i_1 - 1) + 1$, $y_2 = \text{rank}_{B_{3\tau-1}, 1}(i_2 - 1) + 1$, $u_1$ (resp. $u_2$) be the $y_1$th (resp. $y_2$th) leftmost leaf of $\mathcal{T}_{3\tau-1}$, and $u = \text{LCA}(u_1, u_2)$. Then, $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{3\tau-1}}(v) = u$.*

*Proof.* By definition of $\mathcal{T}_{3\tau-1}$ and $A_{\text{short}}$ (Section 5.1), if $\widehat{u}$ is the $k$th leftmost leaf of $\mathcal{T}_{3\tau-1}$, then $\text{str}(\widehat{u}) = A_{\text{short}}[k]$. Thus, $\text{str}(u_1) = A_{\text{short}}[y_1]$ and $\text{str}(u_2) = A_{\text{short}}[y_2]$. Denote $Q = \text{str}(v)$ and consider two cases:

- Let $\text{sdepth}(v) \geq 3\tau - 1$. Denote $X = Q[1 \mathinner{.\,.} 3\tau{-}1]$. By $i_1 = \text{lrank}(v) + 1$ and $i_2 = \text{rrank}(v)$, we then have $\text{SA}[i_1], \text{SA}[i_2] \in \text{Occ}(Q, T) \subseteq \text{Occ}(X, T)$. By definition of $B_{3\tau-1}$ and $A_{\text{short}}$, positions $y_1 = \text{rank}_{B_{3\tau-1}, 1}(i_1 - 1) + 1$ and $y_2 = \text{rank}_{B_{3\tau-1}, 1}(i_2 - 1) + 1$ then satisfy $A_{\text{short}}[y_1] =$

$A_{\text{short}}[y_2] = X$. Thus, by the above observation, $\text{str}(u_1) = \text{str}(u_2) = X$ and hence, by Observation 4.2, $\text{str}(u) = X = \text{str}(v)[1 \mathinner{.\,.} 3\tau - 1]$. Consequently, by Lemma 7.3(1) and since all nodes of $\mathcal{T}_{3\tau-1}$ have different value of str, this yields $\text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(v) = u$.

- Let us now assume $\text{sdepth}(v) < 3\tau - 1$. Let $v_1$ (resp. $v_2$) be the $i_1$th (resp. $i_2$th) leftmost leaf of $\mathcal{T}_{\text{st}}$. Then, $\text{str}(v_1) = T[\text{SA}[i_1] \mathinner{.\,.} n]$ and $\text{str}(v_2) = T[\text{SA}[i_2] \mathinner{.\,.} n]$. By $i_1 = \text{lrank}(v) + 1$ and $i_2 = \text{rrank}(v)$ we have $v = \text{LCA}(v_1, v_2)$. Thus, by Observation 4.2, $\text{lcp}(T[\text{SA}[i_1] \mathinner{.\,.} n], T[\text{SA}[i_2] \mathinner{.\,.} n]) = \text{sdepth}(v) = |Q|$. Observe now that:

  - By definition of $A_{\text{short}}$ and $B_{3\tau-1}$, the string $A_{\text{short}}[y_1]$ (resp. $A_{\text{short}}[y_2]$) is a prefix of $T[\text{SA}[i_1] \mathinner{.\,.} n]$ (resp. $T[\text{SA}[i_2] \mathinner{.\,.} n]$). Thus, it holds $\text{lcp}(A_{\text{short}}[y_1], A_{\text{short}}[y_2]) \leq \text{lcp}(T[\text{SA}[i_1] \mathinner{.\,.} n], T[\text{SA}[i_2] \mathinner{.\,.} n]) = |Q|$.
  - On the other hand, since $Q$ is a prefix of $\text{str}(v_1) = T[\text{SA}[i_1] \mathinner{.\,.} n]$ and $\text{str}(v_2) = T[\text{SA}[i_2] \mathinner{.\,.} n]$, and it holds $|A_{\text{short}}[y_1]| = \min(3\tau - 1, n - \text{SA}[i_1] + 1)$, $|A_{\text{short}}[y_2]| = \min(3\tau - 1, n - \text{SA}[i_2] + 1)$, and $|Q| < 3\tau - 1$, we obtain that $Q$ is a prefix of $A_{\text{short}}[y_1]$ and $A_{\text{short}}[y_2]$. Thus, $\text{lcp}(A_{\text{short}}[y_1], A_{\text{short}}[y_2]) \geq |Q|$.

  We thus proved that $Q$ is a prefix of $A_{\text{short}}[y_1]$ and $A_{\text{short}}[y_2]$, and $\text{lcp}(A_{\text{short}}[y_1], A_{\text{short}}[y_2]) = |Q|$. Thus, since $\text{str}(u_1) = A_{\text{short}}[y_1]$ and $\text{str}(u_2) = A_{\text{short}}[y_2]$, we obtain from Observation 4.2, that $u = \text{LCA}(u_1, u_2)$ satisfies $\text{str}(u) = Q$. By Lemma 7.3(2) and since all nodes of $\mathcal{T}_{3\tau-1}$ have different value of str, this yields $\text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(v) = u$. $\qquad\square$

**Proposition 7.5.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$. Given $\text{C}_{\text{ST}}(T)$ and $\text{repr}(v)$, in $\mathcal{O}(1)$ time we can compute the pointer to the node $\text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(v)$.*

*Proof.* Denote $(b, e) = \text{repr}(v)$, $i_1 = b + 1$, and $i_2 = e$. First, in $\mathcal{O}(1)$ time we compute $y_1 = \text{rank}_{B_{3\tau-1},1}(i_1 - 1) + 1$ and $y_2 = \text{rank}_{B_{3\tau-1},1}(i_2 - 1) + 1$. In $\mathcal{O}(1)$ time we then retrieve the $y_1$th and $y_2$th leftmost leaves $u_1$ and $u_2$ of $\mathcal{T}_{3\tau-1}$ (respectively). Finally, using the LCA structure for $\mathcal{T}_{3\tau-1}$, we compute in $\mathcal{O}(1)$ time the pointer to node $u = \text{LCA}(u_1, u_2)$ of $\mathcal{T}_{3\tau-1}$. By Lemma 7.4, we then have $\text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(v) = u$. $\qquad\square$

**Proposition 7.6.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$. Given $\text{C}_{\text{ST}}(T)$ and $\text{repr}(v)$, we can in $\mathcal{O}(1)$ time check if $v$ is periodic. If $v$ is not periodic, then in $\mathcal{O}(1)$ we can additionally determine if it holds $\text{sdepth}(v) < 3\tau - 1$.*

*Proof.* First, using Proposition 7.5, in $\mathcal{O}(1)$ time we compute the pointer to $u = \text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(v)$. If $\text{sdepth}(u) = 3\tau - 1$ then by Lemma 7.3 it holds $\text{sdepth}(v) \geq 3\tau - 1$, and we can in $\mathcal{O}(1)$ time determine if $v$ is periodic by checking if $\text{per}(X) \leq \frac{1}{3}\tau$ for $X = \text{str}(v)[1 \mathinner{.\,.} 3\tau-1] = \text{str}(u)$ (stored with $u$) using the lookup table $L_{\text{per}}$. If $\text{sdepth}(u) < 3\tau - 1$, then by Lemma 7.3, we have $\text{sdepth}(v) < 3\tau - 1$, and hence $v$ is nonperiodic. Note that in the above algorithm, whenever $v$ is nonperiodic, we always know if $\text{sdepth}(v) < 3\tau - 1$. Thus, we can additionally return this information at no extra cost. Each of the steps takes $\mathcal{O}(1)$ time. $\qquad\square$

### 7.1.3 Implementation of $\text{LCA}(u, v)$

**Lemma 7.7.** *Let $v_1$ and $v_2$ be explicit nodes of $\mathcal{T}_{\text{st}}$. Then,*

$$\text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(\text{LCA}(v_1, v_2)) = \text{LCA}(\text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(v_1), \text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(v_2)).$$

*Proof.* Let $u_1 = \text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(v_1)$, $u_2 = \text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(v_2)$, $v = \text{LCA}(v_1, v_2)$, and $u = \text{LCA}(u_1, u_2)$. Then, the claim is that $\text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{3\tau-1}}(v) = u$. Denote $\ell = \text{lcp}(\text{str}(v_1), \text{str}(v_2))$ and recall that by Observation 4.2, we have $\text{sdepth}(v) = \ell$. We consider two cases:

- First, assume $\mathrm{sdepth}(v) \geq 3\tau - 1$. By $\mathrm{sdepth}(v_1) \geq \mathrm{sdepth}(v) \geq 3\tau - 1$, we obtain from Lemma 7.3(1) and Observation 4.2 that $\mathrm{str}(u_1) = \mathrm{str}(v_1)[1\mathinner{.\,.}3\tau-1] = \mathrm{str}(v)[1\mathinner{.\,.}3\tau-1]$. Analogously, $\mathrm{str}(u_2) = \mathrm{str}(v)[1\mathinner{.\,.}3\tau-1]$, and consequently, $\mathrm{str}(u) = \mathrm{str}(v)[1\mathinner{.\,.}3\tau-1]$. Since all nodes in $\mathcal{T}_{3\tau-1}$ have different values of str, by Lemma 7.3(1), this implies $u = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v)$.
- Let us now assume $\mathrm{sdepth}(v) < 3\tau - 1$. We will show that then $\mathrm{str}(u) = \mathrm{str}(v)$. Since all nodes in $\mathcal{T}_{3\tau-1}$ have different values of str, by Lemma 7.3(2), this immediately implies $u = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v)$. We first show that $\mathrm{sdepth}(u_1) \geq \ell$. Consider two cases. If $\mathrm{sdepth}(v_1) \geq 3\tau - 1$, then by Lemma 7.3(1), $\mathrm{str}(u_1) = \mathrm{str}(v_1)[1\mathinner{.\,.}3\tau - 1]$, i.e., $\mathrm{sdepth}(u_1) = 3\tau - 1 > \mathrm{sdepth}(v) = \ell$. Otherwise, by Lemma 7.3(2), it holds $\mathrm{str}(u_1) = \mathrm{str}(v_1)$, and thus also $\mathrm{sdepth}(u_1) = \mathrm{sdepth}(v_1) \geq \mathrm{sdepth}(v) = \ell$. By the analogous argument, $\mathrm{sdepth}(u_2) \geq \ell$. Recall now that, by definition, $\mathrm{str}(u_1)$ (resp. $\mathrm{str}(u_2)$) is a prefix of $\mathrm{str}(v_1)$ (resp. $\mathrm{str}(v_2)$). Thus, $\mathrm{str}(u_1)[1\mathinner{.\,.}\ell] = \mathrm{str}(v_1)[1\mathinner{.\,.}\ell] = \mathrm{str}(v_2)[1\mathinner{.\,.}\ell] = \mathrm{str}(u_2)[1\mathinner{.\,.}\ell]$. Denoting $\ell' = \mathrm{lcp}(\mathrm{str}(u_1), \mathrm{str}(u_2))$, we therefore have $\ell' \geq \ell$. On the other hand, $\mathrm{str}(u_1)$ (resp. $\mathrm{str}(u_2)$) being a prefix of $\mathrm{str}(v_1)$ (resp. $\mathrm{str}(v_2)$), implies $\ell' \leq \ell$. Consequently, $\ell' = \ell$. By Observation 4.2, we therefore obtain $\mathrm{str}(u) = \mathrm{str}(\mathrm{LCA}(u_1, u_2)) = \mathrm{str}(u_1)[1\mathinner{.\,.}\ell'] = \mathrm{str}(v_1)[1\mathinner{.\,.}\ell] = \mathrm{str}(\mathrm{LCA}(v_1, v_2)) = \mathrm{str}(v)$. $\qquad\square$

**Proposition 7.8.** *Let $v_1$ and $v_2$ be explicit nodes of $\mathcal{T}_{\mathrm{st}}$. Given $\mathrm{C}_{\mathrm{ST}}(T)$ and the pairs $\mathrm{repr}(v_1)$ and $\mathrm{repr}(v_2)$, we can in $\mathcal{O}(1)$ time check if $\mathrm{sdepth}(\mathrm{LCA}(v_1, v_2)) \geq 3\tau - 1$. If so, in $\mathcal{O}(1)$ time we can additionally determine if $\mathrm{LCA}(v_1, v_2)$ is periodic. Otherwise (i.e., if $\mathrm{sdepth}(\mathrm{LCA}(v_1, v_2)) < 3\tau-1$) in $\mathcal{O}(1)$ time we can compute $\mathrm{repr}(\mathrm{LCA}(v_1, v_2))$.*

*Proof.* Denote $v = \mathrm{LCA}(v_1, v_2)$. First, using Proposition 7.5, in $\mathcal{O}(1)$ time we compute pointers to $u_1 = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v_1)$ and $u_2 = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v_2)$ of $\mathcal{T}_{3\tau-1}$. Then, using the LCA structure for $\mathcal{T}_{3\tau-1}$, we compute in $\mathcal{O}(1)$ time the pointer to node $u = \mathrm{LCA}(u_1, u_2)$ of $\mathcal{T}_{3\tau-1}$. By Lemma 7.7, we now have $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v) = u$. If $\mathrm{sdepth}(u) = 3\tau - 1$, by Lemma 7.3 it holds $\mathrm{sdepth}(v) \geq 3\tau - 1$ and $\mathrm{str}(u) = \mathrm{str}(v)[1\mathinner{.\,.}3\tau-1]$, and thus we can in $\mathcal{O}(1)$ determine if $v$ is periodic by checking if $\mathrm{per}(X) \leq \frac{1}{3}\tau$ for $X = \mathrm{str}(u)$ (stored with $u$) using the lookup table $L_{\mathrm{per}}$. Otherwise (i.e., if $\mathrm{sdepth}(u) < 3\tau - 1$), by Lemma 7.3 we have $\mathrm{sdepth}(v) < 3\tau - 1$ and $\mathrm{str}(u) = \mathrm{str}(v)$. Thus, we return that $v$ is nonperiodic and in $\mathcal{O}(1)$ time we obtain the pair $\mathrm{repr}(v) = (\mathrm{RangeBeg}(\mathrm{str}(v), T), \mathrm{RangeEnd}(\mathrm{str}(v), T)) = (\mathrm{RangeBeg}(\mathrm{str}(u), T), \mathrm{RangeEnd}(\mathrm{str}(u), T))$ using the lookup table $L_{\mathrm{range}}$ on $\mathrm{str}(u)$. Each of the steps takes $\mathcal{O}(1)$ time. $\qquad\square$

### 7.1.4 Implementation of $\mathrm{child}(v, c)$

**Lemma 7.9.** *Let $v$ be an explicit internal node of $\mathcal{T}_{\mathrm{st}}$ satisfying $\mathrm{sdepth}(v) < 3\tau - 1$. Let $u = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v)$. For any $c \in [0\mathinner{.\,.}\sigma)$, $\mathrm{child}(v, c) = \bot$ holds if and only if $\mathrm{child}(u, c) = \bot$. Moreover, if $\mathrm{child}(v, c) \neq \bot$ then, letting $u' = \mathrm{child}(u, c)$, it holds*

$$\mathrm{repr}(\mathrm{child}(v, c)) = (\mathrm{RangeBeg}(\mathrm{str}(u'), T), \mathrm{RangeEnd}(\mathrm{str}(u'), T)).$$

*Proof.* By Lemma 7.3(2), $u$ satisfies $\mathrm{str}(u) = \mathrm{str}(v)$. Thus, since $\mathcal{T}_{3\tau-1}$ is a compact trie of substrings of $T$ of length $3\tau - 1$, we immediately obtain that for any $c \in [0\mathinner{.\,.}\sigma)$, $\mathrm{child}(v, c) \neq \bot$ if and only if $\mathrm{child}(u, c)$. Let us assume for some $c \in [0\mathinner{.\,.}\sigma)$, it holds $\mathrm{child}(v, c) = v' \neq \bot$. If $\mathrm{sdepth}(v') \leq 3\tau - 1$, then by definition of $\mathcal{T}_{3\tau-1}$, the node $u' = \mathrm{child}(u, c)$ must satisfy $\mathrm{str}(v') = \mathrm{str}(u')$. This implies the claim immediately. Otherwise ($\mathrm{sdepth}(v') > 3\tau - 1$), $u'$ satisfies $\mathrm{sdepth}(u') = 3\tau - 1$, and corresponds to the implicit node of $\mathcal{T}_{\mathrm{st}}$ on the edge connecting $v$ to $v'$. By definition of suffix tree, however, letting $S$ be such that $\mathrm{str}(v)S = \mathrm{str}(v')[1\mathinner{.\,.}3\tau-1]$, we have $(\mathrm{RangeBeg}(\mathrm{str}(v'), T), \mathrm{RangeEnd}(\mathrm{str}(v'), T)) = (\mathrm{RangeBeg}(\mathrm{str}(v)S, T), \mathrm{RangeEnd}(\mathrm{str}(v)S, T)) = (\mathrm{RangeBeg}(\mathrm{str}(u'), T), \mathrm{RangeEnd}(\mathrm{str}(u'), T))$, which by definition of repr implies the claim. $\qquad\square$

**Proposition 7.10.** *Let $v$ be an explicit internal node of $\mathcal{T}_{\text{st}}$ satisfying $\text{sdepth}(v) < 3\tau - 1$. Given $\text{C}_{\text{ST}}(T)$, $\text{repr}(v)$, and $c \in [0 \mathinner{.\,.} \sigma)$, we can in $\mathcal{O}(1)$ time compute $\text{repr}(\text{child}(v, c))$.*

*Proof.* First, using Proposition 7.5, in $\mathcal{O}(1)$ time we compute a pointer to $u = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{3\tau-1}}(v)$. Using the lookup table $L_{\text{child}}$, in $\mathcal{O}(1)$ time we check if $\text{child}(u, c) = \bot$. If so, then by Lemma 7.9, it holds $\text{child}(v, c) = \bot$ and we return $\text{repr}(\text{child}(v, c)) = (0, 0)$. Otherwise (i.e., $\text{child}(u, c) \neq \bot$), we obtain a pointer to $u' = \text{child}(u, c)$. By Lemma 7.9, we then have $\text{repr}(\text{child}(v, c)) = (\text{RangeBeg}(\text{str}(u'), T), \text{RangeEnd}(\text{str}(u'), T))$, which we obtain using the lookup table $L_{\text{range}}$ on $\text{str}(u')$. Each of the steps takes $\mathcal{O}(1)$ time. $\qquad\square$

### 7.1.5 Implementation of $\text{pred}(v, c)$

**Proposition 7.11.** *Let $v$ be an explicit internal node of $\mathcal{T}_{\text{st}}$ satisfying $\text{sdepth}(v) < 3\tau - 1$. Given $\text{C}_{\text{ST}}(T)$, $\text{repr}(v)$, and $c \in [0 \mathinner{.\,.} \sigma)$, we can in $\mathcal{O}(1)$ time compute $\text{RangeBeg}(\text{str}(v)c, T)$.*

*Proof.* First, using Proposition 7.5 we compute a pointer to $u = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{3\tau-1}}(v)$. By $\text{sdepth}(v) < 3\tau - 1$ and Lemma 7.3(2), node $u$ satisfies $\text{str}(u) = \text{str}(v)$. Thus, we have $\text{RangeBeg}(\text{str}(v)c, T) = \text{RangeBeg}(\text{str}(u)c, T)$. Next, we compute $Y = \text{str}(u)c$ (recall, that $\text{str}(u)$ is stored with $u$). Using the lookup table $L_{\text{range}}$, we then compute and return $\text{RangeBeg}(Y, T)$. Each of the steps takes $\mathcal{O}(1)$ time. $\qquad\square$

### 7.1.6 Implementation of $\text{WA}(v, d)$

**Lemma 7.12.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$ and $d$ be such that $0 \leq d \leq |\text{str}(v)|$ and $d < 3\tau - 1$. Then, letting $u = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{3\tau-1}}(v)$ and $u' = \text{WA}(u, d)$, it holds*

$$\text{repr}(\text{WA}(v, d)) = (\text{RangeBeg}(\text{str}(u'), T), \text{RangeEnd}(\text{str}(u'), T)).$$

*Proof.* Denote $v' = \text{WA}(v, d)$. We consider two cases:

- First, assume $\text{sdepth}(v) \geq 3\tau - 1$. By Lemma 7.3(1), we then have $\text{str}(u) = \text{str}(v)[1 \mathinner{.\,.} 3\tau-1]$. Therefore, utilizing one of the assumptions about $d$, we have $d < 3\tau - 1 = \text{sdepth}(u)$, i.e., $u'$ is well-defined (see Section 4.1). Moreover, this implies that for any ancestor $\bar{v}$ of $v$ at depth at most $3\tau - 1$, there exist a corresponding ancestor $\bar{u}$ of $u$ and there exists a one-to-one mapping between ancestors of $\bar{v}$ in $\mathcal{T}_{\text{st}}$ and ancestors of $\bar{u}$ in $\mathcal{T}_{3\tau-1}$ (with corresponding nodes having equal root-to-node labels). Therefore, if $\text{sdepth}(v') \leq 3\tau - 1$ then $\text{str}(u') = \text{str}(v')$ and the claim follows. Otherwise ($\text{sdepth}(v') > 3\tau - 1$), by $d < 3\tau - 1$, we must have $u' = u$ and $u'$ then corresponds to the implicit node on the edge connecting $v'$ to $\text{parent}(v')$. This implies $\text{repr}(v') = (\text{RangeBeg}(\text{str}(u'), T), \text{RangeEnd}(\text{str}(u'), T))$.
- Let us now assume $\text{sdepth}(v) < 3\tau - 1$. By Lemma 7.3(2), we then have $\text{str}(u) = \text{str}(v)$. In particular, utilizing one of the assumptions on $d$, we have $d \leq \text{sdepth}(v) = \text{sdepth}(u)$, i.e., $u'$ is well-defined (see Section 4.1). Moreover, this implies that there is a one-to-one correspondence between ancestors of $v$ in $\mathcal{T}_{\text{st}}$ and ancestors or $u$ in $\mathcal{T}_{3\tau-1}$. In particular, $\text{str}(v') = \text{str}(u')$, which implies the claim. $\qquad\square$

**Proposition 7.13.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$. Given $\text{C}_{\text{ST}}(T)$, $\text{repr}(v)$, and an integer $d$ satisfying $0 \leq d \leq |\text{str}(v)|$ and $d < 3\tau - 1$, in $\mathcal{O}(1)$ time we can compute $\text{repr}(\text{WA}(v, d))$.*

*Proof.* First, using Proposition 7.5, we compute a pointer to node $u = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{3\tau-1}}(v)$. Then, using the lookup table $L_{\text{WA}}$, in $\mathcal{O}(1)$ time we obtain the pointer to $u' = \text{WA}(u, d)$. By Lemma 7.12, we then have $\text{repr}(\text{WA}(v, d)) = (\text{RangeBeg}(\text{str}(u'), T), \text{RangeEnd}(\text{str}(u'), T))$, which is obtained using the lookup table $L_{\text{range}}$ on $\text{str}(u')$. Each of the steps takes $\mathcal{O}(1)$ time. $\qquad\square$

### 7.1.7 Construction Algorithm

**Proposition 7.14.** *Given the packed representation of $T \in [0 \ldots \sigma)^n$, we can construct $\mathrm{C}_{\mathrm{ST}}(T)$ in $\mathcal{O}(n/\log_\sigma n)$ time.*

*Proof.* First, in $\mathcal{O}(n/\log_\sigma n)$ time we construct $\mathrm{C}_{\mathrm{SA}}(T)$ using Proposition 5.4. Note that during the construction, we compute the frequency $f_X = |\mathrm{Occ}(X,T)|$ for every $X \in [0 \ldots \sigma)^{\leq 3\tau-1}$ (note that by definition of $\mathrm{Occ}(X,T)$ (Section 2), we have $f_X = n$ for the empty string $X = \varepsilon$).

Next, we construct the trie $\mathcal{T}_{3\tau-1}$ and the associated data structures. Observe that for every $X \in [0 \ldots \sigma)^{\leq 3\tau-1}$, the trie $\mathcal{T}_{3\tau-1}$ contains an explicit node $v$ satisfying $\mathrm{str}(v) = X$ if and only if $f_X > 0$, and either $|X| = 3\tau - 1$ or $|X| < 3\tau - 1$ and there exist distinct $c, c' \in [0 \ldots \sigma)$ such that $f_{Xc} > 0$ and $f_{Xc'} > 0$. [10] Thus, given any $X$, we can in $\mathcal{O}(\sigma)$ time check if there exists a node of $\mathcal{T}_{3\tau-1}$ corresponding to $X$. Moreover, if such $v$ exists and $|X| > 0$ then to find $X'$ satisfying $\mathrm{str}(\mathrm{parent}(v)) = X'$, it suffices to compute the longest prefix $X'$ of $X$ such that $L_{\mathrm{range}}$ for $X'$ is different from $L_{\mathrm{range}}$ for $X$. Thus, such $X'$ can be computed in $\mathcal{O}(\tau) = \mathcal{O}(\log n)$ time. Using the above observations, we construct $\mathcal{T}_{3\tau-1}$ as follows. We maintain a lookup table $L_{\mathrm{node}}$ that for any $X \in [0 \ldots \sigma)^{\leq 3\tau-1}$ maps the integer $\mathrm{int}(X)$ to a pointer to the node $v$ of $\mathcal{T}_{3\tau-1}$ satisfying $\mathrm{str}(v) = X$ if such $v$ exists. By $\mathrm{int}(X) \in [0 \ldots \sigma^{6\tau})$, the table needs $\mathcal{O}(\sigma^{6\tau}) = \mathcal{O}(n/\log_\sigma n)$ space and its initialization takes $\mathcal{O}(n/\log_\sigma n)$ time. During the construction, nodes are stored in a dynamic array with amortized $\mathcal{O}(1)$-time insertion at the end, and pointers are implemented as indexes of this array. We enumerate all $X \in [0 \ldots \sigma)^{\leq 3\tau-1}$ in the order of non-decreasing length, and in case of ties, in lexicographical order. For each $X$, using the above method in $\mathcal{O}(\sigma)$ time we check whether there should be a node in $\mathcal{T}_{3\tau-1}$ satisfying $\mathrm{str}(v) = X$. If so, we create a new node $v$, add it to the array of nodes, and update the lookup table $L_{\mathrm{node}}$. Associated with $v$ we store the string $X$ encoded as $\mathrm{int}(X)$ and the length $|X| = \mathrm{sdepth}(v)$. If $|X| > 0$, in $\mathcal{O}(\log n)$ time we then compute the longest prefix $X'$ of $X$ for which $(\mathrm{RangeBeg}(X',T), \mathrm{RangeEnd}(X',T)) \neq (\mathrm{RangeBeg}(X,T), \mathrm{RangeEnd}(X,T))$ (utilizing the lookup table $L_{\mathrm{range}}$), and then using $L_{\mathrm{node}}$ obtain $v'$ satisfying $\mathrm{str}(v') = X'$. We then set $\mathrm{parent}(v) = v'$ and add $v$ to the list of children of $v'$, updating also the links between children of $v'$. Over all $X \in [0 \ldots \sigma)^{\leq 3\tau-1}$, the construction takes $\mathcal{O}(\sigma^{3\tau-1}(\sigma + \log n)) = \mathcal{O}(n/\log_\sigma n)$ time. After constructing $\mathcal{T}_{3\tau-1}$, we augment it with auxiliary structures as follows:

(a) In $\mathcal{O}(\sigma^{3\tau-1}) = \mathcal{O}(n/\log_\sigma n)$ time we preprocess $\mathcal{T}_{3\tau-1}$ for $\mathcal{O}(1)$-time LCA queries using the structure from [10].

(b) Next, we perform a traversal of $\mathcal{T}_{3\tau-1}$. For each node $v$ different from the root, we obtain the pointer to $p = \mathrm{parent}(v)$ and $c = \mathrm{str}(v)[|\mathrm{str}(p)| + 1]$. We then injectively map $(i_p, c)$ (where $i_p$ is the pointer to $p$) to an integer $x$ not exceeding $2\sigma^{3\tau} = \mathcal{O}(\sqrt{n})$ and set $L_{\mathrm{child}}[x] := i_v$, where $i_v$ is the pointer to $v$. The construction takes $\mathcal{O}(\sqrt{n}) = \mathcal{O}(n/\log_\sigma n)$ time.

(c) Next, starting from each node $v$ of $\mathcal{T}_{3\tau-1}$ we compute the pointer to $\mathrm{WA}(v,d)$ for every $d \in [0 \ldots 3\tau-1)$. It suffices to perform one traversal towards the roots and thus this takes $\mathcal{O}(\log n)$ time per node (in total for all $d$). For each computed node $v'$, we map the pair $(i_v, d)$ (where $i_v$ is the pointer to $v$) to an integer $x$ not exceeding $\mathcal{O}(\sqrt{n}\log n)$ and set $L_{\mathrm{WA}}[x] := i_{v'}$, where $i_{v'}$ is the pointer to $v'$. Including the initialization of $L_{\mathrm{WA}}$, the construction takes $\mathcal{O}(\sqrt{n}\log n) = \mathcal{O}(n/\log_\sigma n)$ time.

(d) Finally, we perform the in-order traversal of the tree, collecting the leaves of $\mathcal{T}_{3\tau-1}$ in an array. By the bound on the number of nodes, this takes $\mathcal{O}(n/\log_\sigma n)$ time. $\qquad\square$

---

[10]Note, that this holds also for $X = \varepsilon$ because we defined $f_\varepsilon = n$ and assumed in Section 2 that $T$ contains at least two distinct symbols.

## 7.2 The Nonperiodic Nodes

In this section, we describe a data structure used to perform operations on nonperiodic nodes (see Definition 7.1) in $\mathcal{O}(\log^\epsilon n)$ time.

The section is organized as follows. First, we introduce the components of the data structure (Section 7.2.1). We then show how using this structure to implement some basic navigational routines (Section 7.2.2). Next, we describe the query algorithms for the fundamental operations (Sections 7.2.3 to 7.2.6). Finally, we show the construction algorithm (Section 7.2.7).

### 7.2.1 The Data Structure

**Definitions**  Let $\mathsf{S}$ be a $\tau$-synchronizing set, as defined in Section 5.2.1. Recall (Section 6.2.1) that $A_{\mathsf{S}}[1\mathbin{..}n']$ is an array defined by $A_{\mathsf{S}}[i] = s_i^{\text{lex}}$. Let $\mathcal{T}_{\mathsf{S}}$ denote the compact trie of the set $\{T[i\mathbin{..}n] : i \in \mathsf{S}\}$.

**Components**  The data structure to handle nonperiodic nodes consists of three components:

1. The index core $\mathrm{C}_{\text{ST}}(T)$ (Section 7.1.1). It takes $\mathcal{O}(n/\log_\sigma n)$ space.
2. The data structure from Section 5.2.1 using $\mathcal{O}(n/\log_\sigma n)$ space.
3. The compact trie $\mathcal{T}_{\mathsf{S}}$ represented as in Proposition 4.3 (i.e., for the array $A_{\mathsf{S}}[1\mathbin{..}n']$ defined above). By $n' = \mathcal{O}(n/\log_\sigma n)$ and Proposition 4.3, it needs $\mathcal{O}(n/\log_\sigma n)$ space.

In total, the data structure takes $\mathcal{O}(n/\log_\sigma n)$ space.

### 7.2.2 Navigation Primitives

**Mapping from $\mathcal{T}_{\text{st}}$ to $\mathcal{T}_{\mathsf{S}}$**  For any explicit nonperiodic node $v$ of $\mathcal{T}_{\text{st}}$ satisfying $\text{sdepth}(v) \geq 3\tau - 1$, we define $\text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{\mathsf{S}}}(v) = u$ as a node of $\mathcal{T}_{\mathsf{S}}$ satisfying $\text{str}(v) = X[1\mathbin{..}\delta_{\text{text}}] \cdot \text{str}(u)$, where $X \in \mathcal{D}$ is a prefix of $\text{str}(v)$ and $\delta_{\text{text}} = |X| - 2\tau$ (such $X$ exists and is unique, since for $Y = \text{str}(v)[1\mathbin{..}3\tau-1]$ it holds $\text{Occ}(Y,T) \neq \emptyset$ and $\text{per}(Y) > \frac{1}{3}\tau$; see Section 5.2).

**Lemma 7.15.** *Let $v$ be an explicit nonperiodic node of $\mathcal{T}_{\text{st}}$ satisfying* $\text{sdepth}(v) \geq 3\tau - 1$.

1. *The node* $\text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{\mathsf{S}}}(v)$ *is well-defined.*
2. *Let $X \in \mathcal{D}$ be a prefix of $\text{str}(v)$, $b_X = \text{RangeBeg}(X,T)$, $i_1 = \text{lrank}(v) + 1$, $i_2 = \text{rrank}(v)$, $y_1 = \text{select}_{W,\overline{X}}(i_1 - b_X)$, $y_2 = \text{select}_{W,\overline{X}}(i_2 - b_X)$, $u_1$ (resp. $u_2$) be the $y_1$th (resp. $y_2$th) leftmost leaf of $\mathcal{T}_{\mathsf{S}}$, and $u = \text{LCA}(u_1,u_2)$. Then,* $\text{map}_{\mathcal{T}_{\text{st}},\mathcal{T}_{\mathsf{S}}}(v) = u$.

*Proof.* 1. Let $X \in \mathcal{D}$ be a prefix of $\text{str}(v)$ and let $\delta_{\text{text}} = |X| - 2\tau$. If $v$ is a leaf of $\mathcal{T}_{\text{st}}$, then for $i \in \text{Occ}(\text{str}(v), T)$, it holds that $X$ is a prefix of $T[i\mathbin{..}n]$. Thus, by the consistency of $\mathsf{S}$, $i+\delta_{\text{text}} \in \mathsf{S}$, and consequently, there exist $u$ in $\mathcal{T}_{\mathsf{S}}$ such that $\text{str}(v) = X[1\mathbin{..}\delta_{\text{text}}]\cdot\text{str}(u)$. Otherwise (i.e., $v$ is an internal node), consider any two different leaves $v_1$ and $v_2$ in the subtree rooted in $v$ such that $v = \text{LCA}(v_1,v_2)$. Let $i_1 \in \text{Occ}(\text{str}(v_1), T)$ and $i_2 \in \text{Occ}(\text{str}(v_2), T)$. Then, $\text{str}(v)$ is a prefix of both $T[i_1\mathbin{..}n]$ and $T[i_2\mathbin{..}n]$. Since $X$ is a prefix of $\text{str}(v)$, $X$ is therefore also a prefix of $T[i_1\mathbin{..}n]$ and $T[i_2\mathbin{..}n]$. Thus, again by the consistency of $\mathsf{S}$, we have $i_1 + \delta_{\text{text}}, i_2 + \delta_{\text{text}} \in \mathsf{S}$. Consequently, there exist nodes $u_1$ and $u_2$ in $\mathcal{T}_{\mathsf{S}}$ satisfying $\text{str}(v_1) = X[1\mathbin{..}\delta_{\text{text}}] \cdot \text{str}(u_1)$ and $\text{str}(v_2) = X[1\mathbin{..}\delta_{\text{text}}]\cdot\text{str}(u_2)$. By Observation 4.2 applied to $v_1$ and $v_2$, for $\ell = \text{lcp}(\text{str}(v_1), \text{str}(v_2))$ it holds $\text{str}(v) = \text{str}(v_1)[1\mathbin{..}\ell]$. On the other hand, applying Observation 4.2 to $u_1$ and $u_2$ implies that for $\ell' = \text{lcp}(\text{str}(u_1), \text{str}(u_2))$ and $u = \text{LCA}(u_1,u_2)$, it holds $\text{str}(u) = \text{str}(u_1)[1\mathbin{..}\ell']$. Finally, by $\delta_{\text{text}} < |X|$, we have $\ell = \text{lcp}(\text{str}(v_1), \text{str}(v_2)) = \text{lcp}(X[1\mathbin{..}\delta_{\text{text}}]\cdot\text{str}(u_1), X[1\mathbin{..}\delta_{\text{text}}]\cdot\text{str}(u_2)) = \delta_{\text{text}} + \text{lcp}(\text{str}(u_1), \text{str}(u_2)) = \delta_{\text{text}} + \ell'$. Thus,

$$\text{str}(v) = \text{str}(v_1)[1\mathbin{..}\ell]$$

$$= X[1 \mathinner{.\,.} \delta_{\text{text}}] \cdot \text{str}(u_1)[1 \mathinner{.\,.} \ell - \delta_{\text{text}}]$$
$$= X[1 \mathinner{.\,.} \delta_{\text{text}}] \cdot \text{str}(u_1)[1 \mathinner{.\,.} \ell']$$
$$= X[1 \mathinner{.\,.} \delta_{\text{text}}] \cdot \text{str}(u),$$

i.e., there exists $u$ in $\mathcal{T}_{\mathsf{S}}$ satisfying $\text{str}(v) = X[1 \mathinner{.\,.} \delta_{\text{text}}] \cdot \text{str}(u)$, i.e., $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\mathsf{S}}}(v)$ is well-defined.

2. Let $\delta_{\text{text}} = |X| - 2\tau$. To see that $y_1$ and $y_2$ in the definition are well-defined (i.e., that $i_1 - b_X, i_2 - b_X \in [1 \mathinner{.\,.} \text{rank}_{W, \overline{X}}(n')]$), recall first that (similarly as in the proof of Lemmas 6.5 and 5.5) by consistency of $\mathsf{S}$, there exists a bijection (given by $j \mapsto j + \delta_{\text{text}}$) between $\text{Occ}(X, T)$ and positions $s \in \mathsf{S}$ such that $T^{\infty}[s - \delta_{\text{text}} \mathinner{.\,.} s + 2\tau) = X$. In particular, by definition of $W[1 \mathinner{.\,.} n']$, this implies $\text{rank}_{W, \overline{X}}(n') = |\text{Occ}(X, T)| = e_X - b_X$, where $e_X = \text{RangeEnd}(X, T)$. Observe now that in $\mathcal{T}_{\text{st}}$, for any node $v$, it holds $\text{lrank}(v) = \text{RangeBeg}(\text{str}(v), T)$ and $\text{rrank}(v) = \text{RangeEnd}(\text{str}(v), T)$ (this property does not hold, e.g., in $\mathcal{T}_{\mathsf{S}}$). Thus, since $X$ is a prefix of $\text{str}(v)$, we have $b_X < i_1 \leq i_2 \leq e_X$. Combining with the above, we thus obtain $1 \leq i_1 - b_X \leq i_2 - b_X \leq \text{rank}_{W, \overline{X}}(n')$.

Let $v_1$ (resp. $v_2$) be the $i_1$th (resp. $i_2$th) leftmost leaf of $\mathcal{T}_{\text{st}}$. Denote $\text{str}(v) = Q$. By $i_1, i_2 \in (b_X \mathinner{.\,.} e_X]$, the string $X$ is a prefix of $\text{str}(v_1)$ and $\text{str}(v_2)$. Since $v = \text{LCA}(v_1, v_2)$, $X$ is therefore also a prefix of $\text{str}(v)$. To show $\text{str}(v) = X[1 \mathinner{.\,.} \delta_{\text{text}}] \cdot \text{str}(u)$, it thus suffices to show $\text{str}(u) = Q(\delta_{\text{text}} \mathinner{.\,.} |Q|]$. To this end, we will prove that $Q(\delta_{\text{text}} \mathinner{.\,.} |Q|]$ is a prefix of $\text{str}(u_1)$ and $\text{str}(u_2)$, and that it holds $\text{lcp}(\text{str}(u_1), \text{str}(u_1)) = |Q| - \delta_{\text{text}}$. By $u = \text{LCA}(u_1, u_2)$, this immediately implies the claim. Let $i \in (b_X \mathinner{.\,.} e_X]$. By Lemma 5.5 for $j = \text{SA}[i]$, if $s_y^{\text{lex}} = \text{SA}[i] + \delta_{\text{text}}$ then $\delta(\text{SA}[i]) = \text{rank}_{W, \overline{X}}(y)$. Since by definition of $b_X$ it holds $\delta(\text{SA}[i]) = i - b_X$, we obtain $i - b_X = \text{rank}_{W, \overline{X}}(y)$. Since $y$ also satisfies $T[s_y^{\text{lex}} - \delta_{\text{text}} \mathinner{.\,.} s_y^{\text{lex}} + 2\tau) = X$, it must hold $y = \text{select}_{W, \overline{X}}(i - b_X)$. For such $y$ we have $s_y^{\text{lex}} = \text{SA}[i] + \delta_{\text{text}}$. Applied for $i_1$ and $i_2$, we obtain $s_{y_1}^{\text{lex}} = \text{SA}[i_1] + \delta_{\text{text}}$ and $s_{y_2}^{\text{lex}} = \text{SA}[i_2] + \delta_{\text{text}}$. Recall now that the sequence $(s_i^{\text{lex}})_{i \in [1 \mathinner{.\,.} n']}$ contain the positions in $\mathsf{S}$ sorted according to the lexicographical order of the corresponding suffixes of $T$. This implies that the $y_1$th (resp. $y_2$th) leftmost leaf $u_1$ (resp. $u_2$) of $\mathcal{T}_{\mathsf{S}}$ satisfies $\text{str}(u_1) = T[s_{y_1}^{\text{lex}} \mathinner{.\,.} n] = T[\text{SA}[i_1] + \delta_{\text{text}} \mathinner{.\,.} n]$ (resp. $\text{str}(u_2) = T[s_{y_2}^{\text{lex}} \mathinner{.\,.} n] = T[\text{SA}[i_2] + \delta_{\text{text}} \mathinner{.\,.} n]$). Since all suffixes of $T$ with starting positions in $\text{SA}(\text{lrank}(v) \mathinner{.\,.} \text{rrank}(v)]$ have $Q$ as a prefix, and we clearly have $i_1, i_2 \in (\text{lrank}(v) \mathinner{.\,.} \text{rrank}(v)]$, we immediately obtain that $Q(\delta_{\text{text}} \mathinner{.\,.} |Q|]$ is a prefix of both $T[\text{SA}[i_1] + \delta_{\text{text}} \mathinner{.\,.} n] = \text{str}(u_1)$ and $T[\text{SA}[i_2] + \delta_{\text{text}} \mathinner{.\,.} n] = \text{str}(u_2)$. To show the second claim, we first note that we have $\text{lcp}(T[\text{SA}[i_1] \mathinner{.\,.} n], T[\text{SA}[i_2] \mathinner{.\,.} n]) = \text{lcp}(\text{str}(v_1), \text{str}(v_2)) = |\text{str}(\text{LCA}(v_1, v_2))| = |\text{str}(v)| = |Q|$. Together with $\delta_{\text{text}} \leq |X| \leq |Q|$, this implies $\text{lcp}(\text{str}(u_1), \text{str}(u_2)) = \text{lcp}(T[\text{SA}[i_1] + \delta_{\text{text}} \mathinner{.\,.} n], T[\text{SA}[i_2] + \delta_{\text{text}} \mathinner{.\,.} n]) = \text{lcp}(T[\text{SA}[i_1] \mathinner{.\,.} n], T[\text{SA}[i_2] \mathinner{.\,.} n]) - \delta_{\text{text}} = |Q| - \delta_{\text{text}}$. As noticed earlier, these two facts yield $\text{str}(u) = Q(\delta_{\text{text}} \mathinner{.\,.} |Q|]$, and consequently $\text{str}(v) = X[1 \mathinner{.\,.} \delta_{\text{text}}] \cdot \text{str}(u)$. Thus, $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\mathsf{S}}}(v) = u$. $\qquad\square$

**Proposition 7.16.** *Let $v$ be an explicit nonperiodic node of $\mathcal{T}_{\text{st}}$ satisfying $\text{sdepth}(v) \geq 3\tau - 1$. Given the data structure from Section 7.2.1 and the pair $\text{repr}(v)$, we can in $\mathcal{O}(\log^{\epsilon} n)$ time compute the pointer to $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\mathsf{S}}}(v)$.*

*Proof.* Denote $(b, e) = \text{repr}(v)$. First, using Proposition 7.5, in $\mathcal{O}(1)$ time we compute pointer to $u = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{3\tau-1}}(v)$. By Lemma 7.3(1), we have $\text{str}(u) = \text{str}(v)[1 \mathinner{.\,.} 3\tau - 1]$. Letting $Y = \text{str}(u)$, we then have $\text{per}(Y) > \frac{1}{3}\tau$ and $\text{Occ}(Y, T) \neq \emptyset$. This implies (see Section 5.2.1) that there exists a unique prefix $X \in \mathcal{D}$ of $\text{str}(v)$. Using $L_{\mathcal{D}}$ on $Y$, in $\mathcal{O}(1)$ time we obtain $X$. Using the lookup table $L_{\text{range}}$ (stored as part of $C_{\text{ST}}(T)$; see Section 7.1), in $\mathcal{O}(1)$ time we compute $b_X = \text{RangeBeg}(X, T)$. Using the lookup table $L_{\text{rev}}$ stored in the structure from Section 7.2.1, we then obtain $\overline{X}$. Next, letting $i_1 = b + 1$ and $i_2 = e$ (recall that $\text{repr}(v) = (\text{lrank}(v), \text{rrank}(v))$), using Theorem 2.2 in $\mathcal{O}(\log^{\epsilon} n)$ time we compute $y_1 = \text{select}_{W, \overline{X}}(i_1 - b_X)$ and $y_2 = \text{select}_{W, \overline{X}}(i_2 - b_X)$. Then, using Proposition 4.3 in $\mathcal{O}(1)$ time we compute the pointers to the $y_1$th and $y_2$th leftmost leaves $u_1$ and $u_2$ (respectively) of $\mathcal{T}_{\mathsf{S}}$. Then, again using Proposition 4.3, in $\mathcal{O}(1)$ time we compute and return the pointer to $u = \text{LCA}(u_1, u_2)$. By Lemma 7.15(2), it holds $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\mathsf{S}}}(v) = u$. $\qquad\square$

**Mapping from $\mathcal{T}_\mathsf{S}$ to $\mathcal{T}_\mathrm{st}$** For any string $X \in [0 \mathinner{.\,.} \sigma)^{\leq 3\tau-1}$ and any node $u$ of the trie $\mathcal{T}_\mathsf{S}$, we define $\mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u) = (b_X + \delta_1, b_X + \delta_2)$, where $b_X = \mathrm{RangeBeg}(X, T)$, $z_1 = \mathrm{lrank}(u)$, $z_2 = \mathrm{rrank}(u)$, $\delta_1 = \mathsf{rank}_{W,\overline{X}}(z_1)$, and $\delta_2 = \mathsf{rank}_{W,\overline{X}}(z_2)$.

*Remark* 7.17. Note that the mapping from $\mathcal{T}_\mathrm{st}$ to $\mathcal{T}_\mathsf{S}$ is not necessarily injective, and hence it may not have an inverse. To perform the mapping from $\mathcal{T}_\mathsf{S}$ to $\mathcal{T}_\mathrm{st}$, we will use the above function. Note, however, that although the pair $\mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u)$ is always defined, not for *every* $X$ and $u$ it yields $\mathrm{repr}(v)$ for some node $v$ of $\mathcal{T}_\mathrm{st}$. Below we show a simple but useful condition where it does. In the following sections we show more subtle uses of $\mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u)$ (see, e.g., Remark 7.23).

**Lemma 7.18.** *Let $v$ be an explicit nonperiodic node of $\mathcal{T}_\mathrm{st}$ satisfying $\mathrm{sdepth}(v) \geq 3\tau - 1$ and let $u = \mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v)$. Then, letting $X \in \mathcal{D}$ be a prefix of $\mathrm{str}(v)$, it holds $\mathrm{repr}(v) = \mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u)$.*

*Proof.* First, recall that $\mathcal{D} \subseteq [0 \mathinner{.\,.} \sigma)^{\leq 3\tau-1}$ (Section 5.2.1). Thus, $\mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u)$ is well-defined. Denote $b_X = \mathrm{RangeBeg}(X, T)$, $Q = \mathrm{str}(v)$, $\delta_\mathrm{text} = |X| - 2\tau$, $Q_\mathrm{suf} = Q(\delta_\mathrm{text} \mathinner{.\,.} |Q|]$. Note that $\mathrm{str}(u) = Q_\mathrm{suf}$. Since $\{s_i^\mathrm{lex}\}_{i \in [1 \mathinner{.\,.} n']} = \mathsf{S}$ and $(T[s_i^\mathrm{lex} \mathinner{.\,.} n])_{i \in [1 \mathinner{.\,.} n']}$ is lexicographically sorted, it holds by definition of $\mathcal{T}_\mathsf{S}$ that $\mathrm{lrank}(u) = |\{i \in [1 \mathinner{.\,.} n'] : T[s_i^\mathrm{lex} \mathinner{.\,.} n] \prec Q_\mathrm{suf}\}|$ and $(\mathrm{lrank}(u) \mathinner{.\,.} \mathrm{rrank}(u)] = \{i \in [1 \mathinner{.\,.} n'] : Q_\mathrm{suf} \text{ is a prefix of } T[s_i^\mathrm{lex} \mathinner{.\,.} n]\}$ (in particular, we have $\{s_i^\mathrm{lex}\}_{i \in (\mathrm{lrank}(u) \mathinner{.\,.} \mathrm{rrank}(u)]} = \mathrm{Occ}(Q_\mathrm{suf}, T)$). Therefore, letting $\delta_1 = \mathsf{rank}_{W,\overline{X}}(\mathrm{lrank}(u))$ and $\delta_2 = \mathsf{rank}_{W,\overline{X}}(\mathrm{rrank}(u))$, by Lemma 6.5, it holds $\mathrm{repr}(v) = (\mathrm{RangeBeg}(Q, T), \mathrm{RangeEnd}(Q, T)) = (b_X + \delta_1, b_X + \delta_2) = \mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u)$. $\square$

**Proposition 7.19.** *Let $u$ be a node of $\mathcal{T}_\mathsf{S}$. Given the data structure from Section 7.2.1, a pointer to $u$, and the value $\mathrm{int}(X)$ for some $X \in [0 \mathinner{.\,.} \sigma)^{\leq 3\tau-1}$, we can in $\mathcal{O}(\log^\epsilon n)$ time compute the pair $\mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u)$.*

*Proof.* First, using the $L_\mathrm{range}$ lookup table (stored as part of $\mathrm{C_{ST}}(T)$; see Section 7.1), we compute $b_X = \mathrm{RangeBeg}(X, T)$. Using the lookup table $L_\mathrm{rev}$ stored in the structure from Section 7.2.1, we compute $\overline{X}$. In $\mathcal{O}(1)$ we obtain $z_1 = \mathrm{lrank}(u)$ and $z_2 = \mathrm{rrank}(u)$ (Proposition 4.3). Finally, using Theorem 2.2, in $\mathcal{O}(\log^\epsilon n)$ time we compute $\delta_1 = \mathsf{rank}_{W,\overline{X}}(z_1)$ and $\delta_2 = \mathsf{rank}_{W,\overline{X}}(z_2)$, and return $\mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u) = (b_X + \delta_1, b_X + \delta_2)$. $\square$

### 7.2.3 Implementation of $\mathrm{LCA}(u, v)$

**Lemma 7.20.** *Let $v_1$ and $v_2$ be explicit nodes of $\mathcal{T}_\mathrm{st}$ such that $\mathrm{LCA}(v_1, v_2)$ is nonperiodic and it holds $\mathrm{sdepth}(\mathrm{LCA}(v_1, v_2)) \geq 3\tau - 1$. Then, $v_1$ and $v_2$ are nonperiodic and it holds $\mathrm{sdepth}(v_1) \geq 3\tau - 1$ and $\mathrm{sdepth}(v_2) \geq 3\tau - 1$. Moreover,*

$$\mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(\mathrm{LCA}(v_1, v_2)) = \mathrm{LCA}(\mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v_1), \mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v_2)).$$

*Proof.* Denote $v = \mathrm{LCA}(v_1, v_2)$ and $Y = \mathrm{str}(v)[1 \mathinner{.\,.} 3\tau-1]$. By the assumption, we have $\mathrm{per}(Y) > \frac{1}{3}\tau$. Since by definition of $v$, the string $Y$ is a prefix of $\mathrm{str}(v_1)$ and $\mathrm{str}(v_2)$, we thus obtain that $v_1$ and $v_2$ are nonperiodic and it holds $\mathrm{sdepth}(v_1) \geq 3\tau - 1$ and $\mathrm{sdepth}(v_2) \geq 3\tau - 1$. Thus, $u_1 = \mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v_1)$ and $u_2 = \mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v_2)$ are well-defined (see Section 7.2.2).

Let $u = \mathrm{LCA}(u_1, u_2)$, $\ell' = \mathrm{sdepth}(u)$, and $\ell = \mathrm{sdepth}(v)$. By Observation 4.2, we have $\ell = \mathrm{lcp}(\mathrm{str}(v_1), \mathrm{str}(v_2))$, $\ell' = \mathrm{lcp}(\mathrm{str}(u_1), \mathrm{str}(u_2))$, $\mathrm{str}(v) = \mathrm{str}(v_1)[1 \mathinner{.\,.} \ell]$, and $\mathrm{str}(u) = \mathrm{str}(u_1)[1 \mathinner{.\,.} \ell']$. Let now $X \in \mathcal{D}$ be a prefix of $\mathrm{str}(v)$ (such $X$ exists and is unique since $\mathrm{per}(Y) > \frac{1}{3}\tau$ and since $\mathrm{str}(v)$ being a substring of $T$ implies $\mathrm{Occ}(Y, T) \neq \emptyset$; see also Section 5.2.1). By definition of $\mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v_1)$ and $\mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v_2)$, we have $\mathrm{str}(v_1) = X[1 \mathinner{.\,.} \delta_\mathrm{text}] \cdot \mathrm{str}(u_1)$ and $\mathrm{str}(v_2) = X[1 \mathinner{.\,.} \delta_\mathrm{text}] \cdot \mathrm{str}(u_2)$, where $\delta_\mathrm{text} = |X| - 2\tau$. This implies $\ell = \delta_\mathrm{text} + \ell'$, and consequently, $\mathrm{str}(v) = \mathrm{str}(v_1)[1 \mathinner{.\,.} \ell] = X[1 \mathinner{.\,.} \delta_\mathrm{text}] \cdot \mathrm{str}(u_1)[1 \mathinner{.\,.} \ell - \delta_\mathrm{text}] = X[1 \mathinner{.\,.} \delta_\mathrm{text}] \cdot \mathrm{str}(u_1)[1 \mathinner{.\,.} \ell'] = X[1 \mathinner{.\,.} \delta_\mathrm{text}] \cdot \mathrm{str}(u)$. Since $v$ is an explicit node of $\mathcal{T}_\mathrm{st}$, and no two nodes of $\mathcal{T}_\mathsf{S}$ have the same value of $\mathrm{str}$, we therefore obtain $\mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v) = u$. $\square$

**Proposition 7.21.** *Let $v_1$ and $v_2$ be explicit nodes of $\mathcal{T}_{\mathrm{st}}$ such that $\mathrm{LCA}(v_1, v_2)$ is nonperiodic and satisfies $\mathrm{sdepth}(\mathrm{LCA}(v_1, v_2)) \geq 3\tau - 1$. Given the data structure from Section 7.2.1 and the pairs $\mathrm{repr}(v_1)$ and $\mathrm{repr}(v_2)$, we can in $\mathcal{O}(\log^\epsilon n)$ time compute $\mathrm{repr}(\mathrm{LCA}(v_1, v_2))$.*

*Proof.* Denote $v = \mathrm{LCA}(v_1, v_2)$. By Lemma 7.20, $v_1$ and $v_2$ are nonperiodic and it holds $\mathrm{sdepth}(v_1) \geq 3\tau - 1$ and $\mathrm{sdepth}(v_2) \geq 3\tau - 1$. Thus, $u_1 = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v_1)$ and $u_2 = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v_2)$ are well-defined (see Section 7.2.2). Using Proposition 7.16, in $\mathcal{O}(\log^\epsilon n)$ time we compute pointers to $u_1$ and $u_2$. Next, using the representation of $\mathcal{T}_{\mathsf{S}}$ stored as part of the structure in Section 7.2.1, and Proposition 4.3, in $\mathcal{O}(1)$ time we compute a pointer to $u = \mathrm{LCA}(u_1, u_2)$. By Lemma 7.20, it holds $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v) = u$. We exploit this connection to compute $\mathrm{repr}(v)$. Since $v$ is nonperiodic, letting $Y = \mathrm{str}(v)[1\mathinner{.\,.}3\tau - 1]$, it holds $\mathrm{per}(Y) > \frac{1}{3}\tau$. Since $\mathrm{str}(v)$ is a substring of $T$, we have $\mathrm{Occ}(Y, T) \neq \emptyset$. Together with $\mathrm{per}(Y) > \frac{1}{3}\tau$, this implies (see Section 5.2.1) that there exists a unique prefix $X \in \mathcal{D}$ of $\mathrm{str}(v)$. We compute $X$ as follows. First, using Proposition 7.5, in $\mathcal{O}(1)$ time we compute pointers to $u_1' = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v_1)$ and $u_2' = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v_2)$ of $\mathcal{T}_{3\tau-1}$. Then, using the LCA structure for $\mathcal{T}_{3\tau-1}$, we compute in $\mathcal{O}(1)$ time the pointer to node $u' = \mathrm{LCA}(u_1', u_2')$ of $\mathcal{T}_{3\tau-1}$. By Lemma 7.7, we now have $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v) = u'$. Moreover, by $\mathrm{sdepth}(v) \geq 3\tau - 1$ and Lemma 7.3(1), $\mathrm{str}(u') = Y$. Using $L_{\mathcal{D}}$ on $Y$, in $\mathcal{O}(1)$ time we thus obtain $X$. Using Proposition 7.19, in $\mathcal{O}(\log^\epsilon n)$ time, we then compute the pair $(b, e) = \mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{S}}}(X, u)$. As noted above, $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v) = u$. By Lemma 7.18, we therefore have $\mathrm{repr}(v) = (b, e)$. $\qquad\square$

### 7.2.4 Implementation of $\mathrm{child}(v, c)$

**Lemma 7.22.** *Let $c \in [0\mathinner{.\,.}\sigma)$ and $v$ be an explicit nonperiodic internal node of $\mathcal{T}_{\mathrm{st}}$ satisfying $\mathrm{sdepth}(v) \geq 3\tau - 1$. Let $u = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v)$. If $\mathrm{child}(u, c) = \bot$ then $\mathrm{child}(v, c) = \bot$. Otherwise, letting $u' = \mathrm{child}(u, c)$, it holds*

$$\mathrm{repr}(\mathrm{child}(v, c)) = \begin{cases} (b, e) & \text{if } b \neq e, \\ (0, 0) & \text{otherwise,} \end{cases}$$

*where $(b, e) = \mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{S}}}(X, u')$ and $X \in \mathcal{D}$ is a prefix of $\mathrm{str}(v)$.*

*Proof.* Denote $P = \mathrm{str}(v)c$, $\delta_{\mathrm{text}} = |X| - 2\tau$, and $P' = P(\delta_{\mathrm{text}}\mathinner{.\,.}|P|]$. Observe that since $\mathrm{str}(v)$ is nonperiodic and it holds $\mathrm{sdepth}(v) \geq 3\tau - 1$, $P$ is also nonperiodic and satisfies $|P| \geq 3\tau - 1$. Let $(b_{\mathrm{pre}}, e_{\mathrm{pre}})$ be such that $b_{\mathrm{pre}} = |\{i \in [1\mathinner{.\,.}n'] : T[s_i^{\mathrm{lex}}\mathinner{.\,.}n] \prec P'\}|$ and $(b_{\mathrm{pre}}\mathinner{.\,.}e_{\mathrm{pre}}] = \{i \in [1\mathinner{.\,.}n'] : P' \text{ is a prefix of } T[s_i^{\mathrm{lex}}\mathinner{.\,.}n]\}$. Recall now that $u$ satisfies $\mathrm{str}(u) = \mathrm{str}(v)(\delta_{\mathrm{text}}\mathinner{.\,.}|\mathrm{str}(v)|]$, or equivalently, $\mathrm{str}(u)c = P'$. By definition of $\mathcal{T}_{\mathsf{S}}$ and $\mathrm{child}(u, c)$, we thus obtain that $\mathrm{child}(u, c) = \bot$ implies $e_{\mathrm{pre}} - b_{\mathrm{pre}} = 0$. Consequently, by Lemma 6.5, it holds $|\mathrm{Occ}(P, T)| = \mathrm{RangeEnd}(P, T) - \mathrm{RangeBeg}(P, T) = (b_X + \mathsf{rank}_{W, \overline{X}}(e_{\mathrm{pre}})) - (b_X + \mathsf{rank}_{W, \overline{X}}(b_{\mathrm{pre}})) = \mathsf{rank}_{W, \overline{X}}(b_{\mathrm{pre}}) - \mathsf{rank}_{W, \overline{X}}(b_{\mathrm{pre}}) = 0$, and thus $\mathrm{child}(v, c) = \bot$.

Let us now assume $\mathrm{child}(u, c) = u' \neq \bot$. By definition of $\mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{S}}}(X, u')$, we then have $(b, e) = (b_X + \mathsf{rank}_{W, \overline{X}}(\mathrm{lrank}(u')), b_X + \mathsf{rank}_{W, \overline{X}}(\mathrm{rrank}(u')))$, where $b_X = \mathrm{RangeBeg}(X, T)$. By definition of $\mathcal{T}_{\mathsf{S}}$ and $\mathrm{child}(u, c)$, however, we also have $b_{\mathrm{pre}} = \mathrm{lrank}(u')$ and $e_{\mathrm{pre}} = \mathrm{rrank}(u')$. Thus, by Lemma 6.5,

$$
\begin{aligned}
(\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T)) &= (b_X + \mathsf{rank}_{W, \overline{X}}(b_{\mathrm{pre}}), b_X + \mathsf{rank}_{W, \overline{X}}(e_{\mathrm{pre}})) \\
&= (b_X + \mathsf{rank}_{W, \overline{X}}(\mathrm{lrank}(u')), b_X + \mathsf{rank}_{W, \overline{X}}(\mathrm{rrank}(u'))) \\
&= (b, e).
\end{aligned}
$$

By the above, if $b \neq e$, then $\mathrm{Occ}(P, T) \neq \emptyset$. This implies $\mathrm{child}(v, c) \neq \bot$ and $\mathrm{repr}(\mathrm{child}(v, c)) = (\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T))$. We thus indeed have $\mathrm{repr}(\mathrm{child}(v, c)) = (b, e)$. Otherwise

(i.e., if $b = e$), by the above we have $\mathrm{Occ}(P, T) = \emptyset$. This implies $\mathrm{child}(v, c) = \bot$ and hence indeed we also have $\mathrm{repr}(\mathrm{child}(v, c)) = (0, 0)$. □

*Remark* 7.23. Note that even though in the above result we have $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_\mathsf{S}}(v) = u$ and $\mathrm{child}(u, c)$ contains information used to determine $\mathrm{child}(v, c)$, it does not necessarily hold that $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_\mathsf{S}}(\mathrm{child}(v, c)) = \mathrm{child}(u, c)$. The procedure is nevertheless correct, because such one-to-one correspondence is not required. The details of this mapping, however, become relevant for the $\mathrm{WA}(v, d)$ operation and are explained in detail in the proof of Lemma 7.27.

**Proposition 7.24.** *Let $v$ be an explicit nonperiodic internal node of $\mathcal{T}_{\mathrm{st}}$ satisfying $\mathrm{sdepth}(v) \geq 3\tau - 1$. Given the data structure from Section 7.2.1, $\mathrm{repr}(v)$, and $c \in [0 \mathinner{.\,.} \sigma)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\mathrm{repr}(\mathrm{child}(v, c))$.*

*Proof.* First, using Proposition 7.16, in $\mathcal{O}(\log^\epsilon n)$ time we compute a pointer to $u = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_\mathsf{S}}(v)$. Then, using the representation of $\mathcal{T}_\mathsf{S}$ stored as part of the structure in Section 7.2.1, and Proposition 4.3, in $\mathcal{O}(\log \log n)$ time we check if $\mathrm{child}(u, c) = \bot$. If so, by Lemma 7.22 we have $\mathrm{child}(v, c) = \bot$, and thus we return $\mathrm{repr}(\mathrm{child}(v, c)) = (0, 0)$. Otherwise ($\mathrm{child}(u, c) \neq \bot$), we obtain a pointer to $u' = \mathrm{child}(u, c)$. Next, using Proposition 7.5 in $\mathcal{O}(1)$ time we compute a pointer to $u'' = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v)$. By Lemma 7.3(1), $\mathrm{sdepth}(v) \geq 3\tau - 1$ implies $\mathrm{sdepth}(u'') = 3\tau - 1$ and $\mathrm{str}(u'') = \mathrm{str}(v)[1 \mathinner{.\,.} 3\tau{-}1]$. We obtain $Y = \mathrm{str}(u'')$ (stored with $u''$) in $\mathcal{O}(1)$ time. Using $L_\mathcal{D}$ on $Y$, in $\mathcal{O}(1)$ time we then compute a prefix $X \in \mathcal{D}$ of $Y$ (see Section 7.2.2). Finally, using Proposition 7.19, in $\mathcal{O}(\log^\epsilon n)$ time we compute the pair $(b, e) = \mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u')$. If $b = e$ then by Lemma 7.22 we return $\mathrm{repr}(\mathrm{child}(v, c)) = (0, 0)$. Otherwise, we return $\mathrm{repr}(\mathrm{child}(v, c)) = (b, e)$. □

### 7.2.5 Implementation of $\mathrm{pred}(v, c)$

**Lemma 7.25.** *Let $c \in [0 \mathinner{.\,.} \sigma)$ and $v$ be an explicit nonperiodic internal node of $\mathcal{T}_{\mathrm{st}}$ satisfying $\mathrm{sdepth}(v) \geq 3\tau - 1$. Let $u = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_\mathsf{S}}(v)$. If $\mathrm{pred}(u, c) = \bot$ then $\mathrm{RangeBeg}(\mathrm{str}(v)c, T) = \mathrm{RangeBeg}(\mathrm{str}(v), T)$. Otherwise, letting $u' = \mathrm{pred}(u, c)$, it holds*

$$\mathrm{RangeBeg}(\mathrm{str}(v)c, T) = e,$$

*where $(b, e) = \mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u')$ and $X \in \mathcal{D}$ is a prefix of $\mathrm{str}(v)$.*

*Proof.* We start by characterizing $\mathrm{RangeBeg}(\mathrm{str}(v), T)$ using Lemma 6.5 for pattern $\mathrm{str}(v)$. First, note that $v$ is nonperiodic and it holds $\mathrm{sdepth}(v) \geq 3\tau - 1$. On the other hand, by definition, we have $\mathrm{str}(u) = \mathrm{str}(v)(\delta_{\mathrm{text}} \mathinner{.\,.} |\mathrm{str}(v)|]$, where $\delta_{\mathrm{text}} = |X| - 2\tau$. Finally, by definition of $\mathcal{T}_\mathsf{S}$, we have $|\{i \in [1 \mathinner{.\,.} n'] : T[s_i^{\mathrm{lex}} \mathinner{.\,.} n] \prec \mathrm{str}(v)(\delta_{\mathrm{text}} \mathinner{.\,.} |\mathrm{str}(v)|]\}| = \mathrm{lrank}(u)$. Thus, by Lemma 6.5, we have $\mathrm{RangeBeg}(\mathrm{str}(v), T) = b_X + \mathrm{rank}_{W, \overline{X}}(\mathrm{lrank}(u))$, where $b_X = \mathrm{RangeBeg}(X, T)$.

Denote $P = \mathrm{str}(v)c$ and $P' = P(\delta_{\mathrm{text}} \mathinner{.\,.} |P|]$. Since $\mathrm{str}(v)$ is nonperiodic and satisfies $|\mathrm{str}(v)| \geq 3\tau - 1$, $P$ is also nonperiodic and it holds $|P| \geq 3\tau - 1$. Note also that by $\mathrm{str}(u) = \mathrm{str}(v)(\delta_{\mathrm{text}} \mathinner{.\,.} |\mathrm{str}(v)|]$, we have $\mathrm{str}(u)c = P'$.

Let us first assume $\mathrm{pred}(u, c) = \bot$. By definition, this implies that $|\{i \in [1 \mathinner{.\,.} n'] : T[s_i^{\mathrm{lex}} \mathinner{.\,.} n] \prec \mathrm{str}(u)c\}| = \mathrm{lrank}(u)$. Equivalently, by $\mathrm{str}(u)c = P'$, $|\{i \in [1 \mathinner{.\,.} n'] : T[s_i^{\mathrm{lex}} \mathinner{.\,.} n] \prec P'\}| = \mathrm{lrank}(u)$. By Lemma 6.5 for pattern $P = \mathrm{str}(v)c$ we thus obtain $\mathrm{RangeBeg}(\mathrm{str}(v)c, T) = b_X + \mathrm{rank}_{W, \overline{X}}(\mathrm{lrank}(u))$. Since above we also established that $\mathrm{RangeBeg}(\mathrm{str}(v), T) = b_X + \mathrm{rank}_{W, \overline{X}}(\mathrm{lrank}(u))$, we have thus proved that $\mathrm{pred}(u, c) = \bot$ implies $\mathrm{RangeBeg}(\mathrm{str}(v)c, T) = \mathrm{RangeBeg}(\mathrm{str}(v), T)$.

Let us now assume $\mathrm{pred}(u, c) = u' \neq \bot$. Observe that by definition of $\mathrm{pred}(u, c)$, this implies $|\{i \in [1 \mathinner{.\,.} n'] : T[s_i^{\mathrm{lex}} \mathinner{.\,.} n] \prec \mathrm{str}(u)c\}| = \mathrm{rrank}(u')$. By Lemma 6.5 applied for pattern $P = \mathrm{str}(v)c$,

we thus obtain $\mathrm{RangeBeg}(\mathrm{str}(v)c, T) = b_X + \mathsf{rank}_{W,\overline{X}}(\mathrm{rrank}(u'))$. On the other hand, observe that by definition of $(b, e) = \mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u')$, we have $e = b_X + \mathsf{rank}_{W,\overline{X}}(\mathrm{rrank}(u'))$. We thus obtain $\mathrm{RangeBeg}(\mathrm{str}(v)c, T) = e$. $\qquad\square$

**Proposition 7.26.** *Let $v$ be an explicit nonperiodic internal node of $\mathcal{T}_\mathrm{st}$ satisfying $\mathrm{sdepth}(v) \geq 3\tau - 1$. Given the data structure from Section 7.2.1, $\mathrm{repr}(v)$, and $c \in [0 \mathinner{.\,.} \sigma)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\mathrm{RangeBeg}(\mathrm{str}(v)c, T)$.*

*Proof.* First, using Proposition 7.16, in $\mathcal{O}(\log^\epsilon n)$ time we compute a pointer to $u = \mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v)$. Then, using the representation of $\mathcal{T}_\mathsf{S}$ stored as part of the structure in Section 7.2.1, and Proposition 4.3, in $\mathcal{O}(\log\log n)$ time we check if $\mathrm{pred}(u, c) = \bot$. If so, by Lemma 7.25 we have $\mathrm{RangeBeg}(\mathrm{str}(v)c, T) = \mathrm{RangeBeg}(\mathrm{str}(v), T)$ and hence we return $\mathrm{RangeBeg}(\mathrm{str}(v), T)$ (given as input) as the result. Otherwise ($\mathrm{pred}(u, c) \neq \bot$), we obtain a pointer to $u' = \mathrm{pred}(u, c)$. Next, using Proposition 7.5 in $\mathcal{O}(1)$ time we compute a pointer to $u'' = \mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_{3\tau-1}}(v)$. By Lemma 7.3(1), $\mathrm{sdepth}(v) \geq 3\tau - 1$ implies $\mathrm{sdepth}(u'') = 3\tau - 1$ and $\mathrm{str}(u'') = \mathrm{str}(v)[1 \mathinner{.\,.} 3\tau{-}1]$. We obtain $Y = \mathrm{str}(u'')$ (stored with $u''$) in $\mathcal{O}(1)$ time. Using $L_\mathcal{D}$ on $Y$, in $\mathcal{O}(1)$ time we then compute a prefix $X \in \mathcal{D}$ of $Y$ (see Section 7.2.2). Finally, using Proposition 7.19, in $\mathcal{O}(\log^\epsilon n)$ time we compute the pair $(b, e) = \mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u')$. By Lemma 7.25, it holds $\mathrm{RangeBeg}(\mathrm{str}(v)c, T) = e$. We thus return $e$ as the answer. $\qquad\square$

### 7.2.6 Implementation of $\mathrm{WA}(v, d)$

**Lemma 7.27.** *Let $v$ be an explicit nonperiodic node of $\mathcal{T}_\mathrm{st}$ and $d$ be an integer satisfying $3\tau - 1 \leq d \leq |\mathrm{str}(v)|$. Then, letting $u = \mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v)$, it holds*

$$\mathrm{repr}(\mathrm{WA}(v, d)) = \mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, \widehat{u}),$$

*where $X \in \mathcal{D}$ is a prefix of $\mathrm{str}(v)$, $\delta_\mathrm{text} = |X| - 2\tau$, and $\widehat{u} = \mathrm{WA}(u, d - \delta_\mathrm{text})$.*

*Proof.* Denote $f^{(0)}(x) = x$ and $f^{(i)}(x) = f(f^{(i-1)}(x))$ for $i \in \mathbb{Z}_+$. Let

$$\mathcal{V} := \{\mathrm{parent}^{(i)}(v) : i \in \mathbb{Z}_{\geq 0} \text{ and } \mathrm{sdepth}(\mathrm{parent}^{(i)}(v)) \geq |X|\} \text{ and}$$
$$\mathcal{U} := \{\mathrm{parent}^{(i)}(u) : i \in \mathbb{Z}_{\geq 0} \text{ and } \mathrm{sdepth}(\mathrm{parent}^{(i)}(u)) \geq 2\tau\}$$

For any $v' \in \mathcal{V}$, the node $u' = \mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v')$ satisfies $\mathrm{str}(v') = X[1 \mathinner{.\,.} \delta_\mathrm{text}] \cdot \mathrm{str}(u')$. In particular, $\mathrm{str}(u) = \mathrm{str}(v)(\delta_\mathrm{text} \mathinner{.\,.} |\mathrm{str}(v)|]$. Since for any $v' \in \mathcal{V}$, $\mathrm{str}(v') = \mathrm{str}(v)[1 \mathinner{.\,.} |\mathrm{str}(v')|]$, we thus obtain that for $u' = \mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v')$ it holds $\mathrm{str}(u') = \mathrm{str}(v')(\delta_\mathrm{text} \mathinner{.\,.} |\mathrm{str}(v')|] = \mathrm{str}(v)(\delta_\mathrm{text} \mathinner{.\,.} |\mathrm{str}(v')|] = \mathrm{str}(u)[1 \mathinner{.\,.} |\mathrm{str}(u')|]$, i.e., $u'$ is an ancestor of $u$. Moreover, $\mathrm{sdepth}(u') = |\mathrm{str}(v')| - \delta_\mathrm{text} \geq |X| - \delta_\mathrm{text} = 2\tau$. Consequently, $\mathcal{U}' := \{\mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v') : v' \in \mathcal{V}\}$ satisfies $\mathcal{U}' \subseteq \mathcal{U}$. Note also, that $v' \neq v''$ implies $\mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v') \neq \mathrm{map}_{\mathcal{T}_\mathrm{st}, \mathcal{T}_\mathsf{S}}(v'')$.

For any $u' \in \mathcal{U}$, denote $(s(u'), t(u')) = \mathrm{pseudoinv}_{\mathcal{T}_\mathsf{S}}(X, u')$. We prove the following property of $\mathcal{U}'$. Let $w, w' \in \mathcal{U}$ be such that $w' = \mathrm{parent}(w)$. We claim, that $(s(w), t(w)) \neq (s(w'), t(w'))$ implies $w' \in \mathcal{U}'$. Denote $Q' = \mathrm{str}(w')$ and $Q = X[1 \mathinner{.\,.} \delta_\mathrm{text}] \cdot Q'$. The proof consists of three steps:

- First, we show that it holds $\{\mathrm{SA}[i]\}_{i \in (s' \mathinner{.\,.} t']} = \mathrm{Occ}(Q, T)$, where $s' = s(w')$ and $t' = t(w')$. By the above discussion, we have $\mathrm{str}(v) = X[1 \mathinner{.\,.} \delta_\mathrm{text}] \cdot \mathrm{str}(u)$. Thus, $X(\delta_\mathrm{text} \mathinner{.\,.} |X|]$ is a prefix of $\mathrm{str}(u)$. On the other hand, by $w' \in \mathcal{U}$, $\mathrm{str}(w')$ is a prefix of $\mathrm{str}(u)$ and we have $|\mathrm{str}(w')| \geq 2\tau$. Consequently, $X(\delta_\mathrm{text} \mathinner{.\,.} |X|]$ is a prefix of $Q'$. As noted in the proof of Lemma 7.18, by the consistency of $\mathsf{S}$ we then have $\{s_i^\mathrm{lex}\}_{i \in (\mathrm{lrank}(w') \mathinner{.\,.} \mathrm{rrank}(w')]} = \mathrm{Occ}(Q', T)$ and consequently $\{\mathrm{SA}[i]\}_{i \in (s' \mathinner{.\,.} t']} = \mathrm{Occ}(X[1 \mathinner{.\,.} \delta_\mathrm{text}] \cdot Q', T) = \mathrm{Occ}(Q, T)$.

- Second, we prove that there exists a node $v'$ in $\mathcal{T}_{\mathrm{st}}$ such that $\mathrm{str}(v') = Q$. First, note that $\{\mathrm{SA}[i]\}_{i \in (s'..t']} = \mathrm{Occ}(Q, T)$ already implies that there exists some node $v'$ of $\mathcal{T}_{\mathrm{st}}$ such that $\mathrm{repr}(v') = (s', t')$ and $Q$ is a prefix of $\mathrm{str}(v')$. It thus remains to show that $\mathrm{str}(v') = Q$. For this, it suffices to show that there exists $c, c' \in [0..\sigma)$ such that $c \neq c'$, $\mathrm{Occ}(Qc, T) \neq \emptyset$, and $\mathrm{Occ}(Qc', T) \neq \emptyset$. Observe that by $(s(w'), t(w')) \neq (s(w), t(w))$, there exists a child $w'' \neq w$ of $w'$ such that $t(w'') > s(w'')$. This yields an occurrence of $\mathrm{str}(w')$ preceded in $T$ with $X[1..\delta_{\mathrm{text}}]$. The same holds for $\mathrm{str}(w)$, since $s(w') \leq s(u) < t(u) \leq t(w')$. In other words, for $c = \mathrm{str}(w)[|Q'| + 1]$ and $c' = \mathrm{str}(w'')[|Q'| + 1]$ we have $c \neq c'$, $\mathrm{Occ}(Qc, T) \neq \emptyset$, and $\mathrm{Occ}(Qc', T) \neq \emptyset$. This concludes the proof of $\mathrm{str}(v') = Q$.
- Finally, recall that by definition, the node $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v')$ satisfies $\mathrm{str}(v') = X[1..\delta_{\mathrm{text}}] \cdot \mathrm{str}(\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v'))$. Thus, by $\mathrm{str}(v') = Q = X[1..\delta_{\mathrm{text}}] \cdot Q'$ and $\mathrm{str}(w') = Q'$, we must have $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v') = w'$. This implies $w' \in \mathcal{U}'$.

We are now ready to prove the main claim. Let $v' = \mathrm{WA}(v, d)$ and $v'' = \mathrm{parent}(v')$. We then have $\mathrm{sdepth}(v'') < d \leq \mathrm{sdepth}(v')$. Moreover, by $|X| \leq 3\tau - 1 \leq d$, we have $v' \in \mathcal{V}$. Let $u' = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v')$. Then, $u' \in \mathcal{U}'$. By the above discussion, we also have $d - \delta_{\mathrm{text}} \leq \mathrm{sdepth}(\widehat{u}) \leq \mathrm{sdepth}(u')$. By $3\tau - 1 \leq d$ this implies $2\tau = |X| - \delta_{\mathrm{text}} \leq 3\tau - 1 - \delta_{\mathrm{text}} \leq d - \delta_{\mathrm{text}} \leq \mathrm{sdepth}(\widehat{u})$, i.e., $\widehat{u} \in \mathcal{U}$. Let $k \in \mathbb{Z}_{\geq 0}$ be such that $\widehat{u} = \mathrm{parent}^{(k)}(u')$. This implies that $\mathrm{parent}^{(i)}(u') \notin \mathcal{U}'$ holds for $i \in [1..k]$, since otherwise it would contradict $v' = \mathrm{WA}(v, d)$. If $k = 0$ then we trivially have $(s(u'), t(u')) = (s(\widehat{u}), t(\widehat{u}))$. Otherwise, by (the contraposition of) the above property of $\mathcal{U}'$ we have

$$
\begin{aligned}
(s(u'), t(u')) &= (s(\mathrm{parent}(u')), t(\mathrm{parent}(u'))) \\
&= \ldots \\
&= (s(\mathrm{parent}^{(k)}(u')), t(\mathrm{parent}^{(k)}(u'))) \\
&= (s(\widehat{u}), t(\widehat{u})).
\end{aligned}
$$

By Lemma 7.18, we obtain $\mathrm{repr}(\mathrm{WA}(v, d)) = \mathrm{repr}(v') = \mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{S}}}(X, u') = (s(u'), t(u')) = (s(\widehat{u}), t(\widehat{u})) = \mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{S}}}(X, \widehat{u})$. $\qquad\square$

**Proposition 7.28.** *Let $v$ be an explicit nonperiodic node of $\mathcal{T}_{\mathrm{st}}$ satisfying $3\tau - 1 \leq |\mathrm{str}(v)|$. Given the data structure from Section 7.2.1, $\mathrm{repr}(v)$, and an integer $d$ satisfying $3\tau - 1 \leq d \leq |\mathrm{str}(v)|$, in $\mathcal{O}(\log^{\epsilon} n)$ time we can compute $\mathrm{repr}(\mathrm{WA}(v, d))$.*

*Proof.* First, using Proposition 7.16, in $\mathcal{O}(\log^{\epsilon} n)$ time we compute a pointer to $u = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v)$. Next, using Proposition 7.5 in $\mathcal{O}(1)$ time we compute a pointer to $u' = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{3\tau-1}}(v)$. By Lemma 7.3(1), $\mathrm{sdepth}(v) \geq 3\tau - 1$ implies $\mathrm{sdepth}(u') = 3\tau - 1$ and $\mathrm{str}(u') = \mathrm{str}(v)[1..3\tau-1]$. We obtain $Y = \mathrm{str}(u')$ (stored with $u'$) in $\mathcal{O}(1)$ time. Using $L_{\mathcal{D}}$ on $Y$, in $\mathcal{O}(1)$ time we then compute a prefix $X \in \mathcal{D}$ of $Y$ (see Section 7.2.2). Let $\delta_{\mathrm{text}} = |X| - 2\tau$. Finally, using the representation of $\mathcal{T}_{\mathsf{S}}$ stored as part of the structure in Section 7.2.1, and Proposition 4.3, in $\mathcal{O}(\log \log n)$ time we compute a pointer to $\widehat{u} = \mathrm{WA}(u, d - \delta_{\mathrm{text}})$. Using Proposition 7.19, in $\mathcal{O}(\log^{\epsilon} n)$ time we then compute and return $\mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{S}}}(X, \widehat{u})$, which by Lemma 7.27 is equal to $\mathrm{repr}(\mathrm{WA}(v, d))$. $\qquad\square$

### 7.2.7 Construction Algorithm

**Proposition 7.29.** *Given $\mathrm{C}_{\mathrm{ST}}(T)$, we can in $\mathcal{O}(n \min(1, \log \sigma / \sqrt{\log n}))$ time and $\mathcal{O}(n / \log_{\sigma} n)$ working space augment it into a data structure from Section 7.2.1.*

*Proof.* First, we combine Propositions 5.4 and 5.9 (recall that the packed representation of $T$ is a component of $\mathrm{C}_{\mathrm{ST}}(T)$) to construct in $\mathcal{O}(n \min(1, \log \sigma / \sqrt{\log n}))$ time and using $\mathcal{O}(n / \log_{\sigma} n)$

65

working space the data structure from Section 5.2.1. In particular, this constructs $(s_i^{\mathrm{lex}})_{i \in [1 \mathrel{..} n']}$. We then initialize $A_{\mathsf{S}}[i] = s_i^{\mathrm{lex}}$ for $i \in [1 \mathrel{..} n']$ and in $\mathcal{O}(n / \log_\sigma n)$ time construct $\mathcal{T}_{\mathsf{S}}$ represented using Proposition 4.3. $\qquad\square$

## 7.3 The Periodic Nodes

In this section, we describe a data structure used to perform operations on periodic nodes (see Definition 7.1) in $\mathcal{O}(\log \log n)$ time.

The section is organized as follows. First, we introduce the components of the data structure (Section 7.3.1). We then show how using this structure to implement some basic navigational routines (Section 7.3.2). Next, we describe the query algorithms for the fundamental operations (Sections 7.3.3 to 7.3.6). Finally, we show the construction algorithm (Section 7.3.7).

### 7.3.1 The Data Structure

**Definitions**  Let $v$ be a periodic node of $\mathcal{T}_{\mathrm{st}}$. We define $\mathrm{L\text{-}root}(v) := \mathrm{L\text{-}root}(\mathrm{str}(v))$, $e(v) := e(\mathrm{str}(v))$, $\mathrm{L\text{-}head}(v) := \mathrm{L\text{-}head}(\mathrm{str}(v))$, $\mathrm{L\text{-}exp}(v) := \mathrm{L\text{-}exp}(\mathrm{str}(v))$, $\mathrm{L\text{-}tail}(v) := \mathrm{L\text{-}tail}(\mathrm{str}(v))$, $e^{\mathrm{full}}(v) := e^{\mathrm{full}}(\mathrm{str}(v))$, and $\mathrm{type}(v) := \mathrm{type}(\mathrm{str}(v))$. Let $q = |\mathsf{R}'^{-}|$. Recall (Section 5.3.2) that $(r_i^{\mathrm{lex}})_{i \in [1 \mathrel{..} q]}$ is a sequence containing all elements $k \in \mathsf{R}'^{-}$ sorted first according to $\mathrm{L\text{-}root}(k)$ and in case of ties, by $T[e^{\mathrm{full}}(k) \mathrel{..} n]$. Recall also (Section 6.3.2) that $\mathsf{Z} = \{e^{\mathrm{full}}(j) - |\mathrm{pow}(\mathrm{L\text{-}root}(j))| : j \in \mathsf{R}'^{-}\}$ and $A_{\mathsf{Z}}[1 \mathrel{..} q]$ is an array defined by $A_{\mathsf{Z}}[i] = e^{\mathrm{full}}(r_i^{\mathrm{lex}}) - |\mathrm{pow}(H_i)|$, where $H_i = \mathrm{L\text{-}root}(r_i^{\mathrm{lex}})$ and $\mathrm{pow}(H_i) = H_i^\infty[1 \mathrel{..} |H_i| \lceil \frac{\tau}{|H_i|} \rceil]$. Let $\mathcal{T}_{\mathsf{Z}}$ denote the compact trie of the set $\{T[i \mathrel{..} n] : i \in \mathsf{Z}\}$.

**Components**  The data structure consists of two parts. The first part consists of the following three components:

1. The index core $\mathrm{C}_{\mathrm{ST}}(T)$ (Section 7.1). It takes $\mathcal{O}(n / \log_\sigma n)$ space.
2. The first part of the structure from Section 5.3.2 using $\mathcal{O}(n / \log_\sigma n)$ space.
3. The compact trie $\mathcal{T}_{\mathsf{Z}}$ represented as in Proposition 4.3 (i.e., for the array $A_{\mathsf{Z}}[1 \mathrel{..} q]$ defined as above). By $q = \mathcal{O}(n / \log_\sigma n)$ and Proposition 4.3, it needs $\mathcal{O}(n / \log_\sigma n)$ space.

The second part of the structure consists of the symmetric counterparts of the above components adapted according to Lemma 5.11 (see also Section 5.3.2)

In total, the data structure takes $\mathcal{O}(n / \log_\sigma n)$ space.

### 7.3.2 Navigation Primitives

**Mapping from $\mathcal{T}_{\mathrm{st}}$ to $\mathcal{T}_{\mathsf{Z}}$**  For any periodic explicit node $v$ of $\mathcal{T}_{\mathrm{st}}$ satisfying $e(v) \le |\mathrm{str}(v)|$ and $\mathrm{type}(v) = -1$, we define $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v) = u$ as a node of $\mathcal{T}_{\mathsf{Z}}$ satisfying $\mathrm{str}(u) = \mathrm{pow}(H) \cdot \mathrm{str}(v)[e^{\mathrm{full}}(v) \mathrel{..} |\mathrm{str}(v)|]$, where $H = \mathrm{L\text{-}root}(v)$.

**Lemma 7.30.** *Let $v$ be a periodic explicit node $v$ of $\mathcal{T}_{\mathrm{st}}$ such that $e(v) \le |\mathrm{str}(v)|$ and $\mathrm{type}(v) = -1$.*

1. *The node $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v)$ is well-defined.*
2. *Let $i_1 = \mathrm{lrank}(v) + 1$, $i_2 = \mathrm{rrank}(v)$, $y_1$ and $y_2$ be such that $e^{\mathrm{full}}(r_{y_1}^{\mathrm{lex}}) = e^{\mathrm{full}}(\mathrm{SA}[i_1])$ and $e^{\mathrm{full}}(r_{y_2}^{\mathrm{lex}}) = e^{\mathrm{full}}(\mathrm{SA}[i_2])$ (respectively), $u_1$ and $u_2$ be the $y_1$th and $y_2$th leftmost leaf of $\mathcal{T}_{\mathsf{Z}}$ (respectively), and $u = \mathrm{LCA}(u_1, u_2)$. Then, $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v) = u$.*

*Proof.* Denote $s = \mathrm{L\text{-}head}(v)$, $H = \mathrm{L\text{-}root}(v)$, $p = |H|$, $Q = \mathrm{str}(v)$, and $Q_{\mathrm{suf}} = Q[e^{\mathrm{full}}(Q) \mathrel{..} |Q|]$. Note that by $e^{\mathrm{full}}(Q) \le e(Q) \le |Q|$, it holds $Q_{\mathrm{suf}} \ne \varepsilon$.

1. We start by observing that by Lemma 6.9 and Lemma 6.11(2), for every $i \in \mathrm{Occ}(Q, T)$, it holds $i \in \mathsf{R}_{s,H}^{-}$ and $e^{\mathrm{full}}(i) - i = e^{\mathrm{full}}(Q) - 1$. Note that this implies $e^{\mathrm{full}}(i) \in \mathrm{Occ}(Q_{\mathrm{suf}}, T)$. To show that $\mathcal{T}_{\mathsf{Z}}$ has a node $u$ satisfying $\mathrm{str}(u) = \mathrm{pow}(H) \cdot Q_{\mathrm{suf}}$, consider two cases:

- Assume that $v$ is a leaf. Let $i \in \mathrm{Occ}(\mathrm{str}(v), T)$. By the above, $i \in \mathsf{R}_H^-$. Let $j$ be the smallest integer such that $[j \mathinner{\ldotp\ldotp} i] \subseteq \mathsf{R}$. It holds $j \in \mathsf{R}'$ and moreover, by Lemma 5.12, $j \in \mathsf{R}_H^-$ and $e^{\mathrm{full}}(j) = e^{\mathrm{full}}(i)$. Thus, by $e^{\mathrm{full}}(i) \in \mathrm{Occ}(Q_{\mathrm{suf}}, T)$ (see above), we have $e^{\mathrm{full}}(j) \in \mathrm{Occ}(Q_{\mathrm{suf}}, T)$. Finally, by $i \in \mathrm{Occ}(Q, T)$ and $|Q| = n - i + 1$, we have $|Q_{\mathrm{suf}}| = |Q| - e^{\mathrm{full}}(Q) + 1 = |Q| - (e^{\mathrm{full}}(Q) - 1) = (n - i + 1) - (e^{\mathrm{full}}(i) - i) = n - e^{\mathrm{full}}(i) + 1 = n - e^{\mathrm{full}}(j) + 1$. We have thus shown that there exists $j \in \mathsf{R}_H'^-$ such that $e^{\mathrm{full}}(j) \in \mathrm{Occ}(Q_{\mathrm{suf}}, T)$ and $n - e^{\mathrm{full}}(j) + 1 = |Q_{\mathrm{suf}}|$. By definition of $A_{\mathsf{Z}}[1 \mathinner{\ldotp\ldotp} q]$ (see Section 6.3.2) this implies that there exists a leaf $u$ of $\mathcal{T}_{\mathsf{Z}}$ such that $\mathrm{str}(u) = \mathrm{pow}(H) \cdot T[e^{\mathrm{full}}(j) \mathinner{\ldotp\ldotp} n] = \mathrm{pow}(H) \cdot Q_{\mathrm{suf}}$.

- Assume that $v$ is an internal node. Consider the leftmost and the rightmost leaves $v_1$ and $v_2$ (respectively) in the subtree rooted in $v$. Letting $i_1 \in \mathrm{Occ}(\mathrm{str}(v_1), T)$ and $i_2 \in \mathrm{Occ}(\mathrm{str}(v_2), T)$, we have $i_1 \neq i_2$ and $i_1, i_2 \in \mathrm{Occ}(Q, T)$. Thus, $i_1, i_2 \in \mathsf{R}_H^-$ and $e^{\mathrm{full}}(i_1) - i_1 = e^{\mathrm{full}}(Q) - 1 = e^{\mathrm{full}}(i_2) - i_2$. Therefore, $e^{\mathrm{full}}(i_1) \neq e^{\mathrm{full}}(i_2)$ and, by Lemma 5.12, $i_1$ and $i_2$ are in different maximal contiguous blocks of positions from $\mathsf{R}$, i.e., letting $j_1$ (resp. $j_2$) be the smallest integer such that $[j_1 \mathinner{\ldotp\ldotp} i_1] \subseteq \mathsf{R}$ (resp. $[j_2 \mathinner{\ldotp\ldotp} i_2] \subseteq \mathsf{R}$), we have $j_1, j_2 \in \mathsf{R}'$ and $j_1 \neq j_2$. By Lemma 5.12, it then holds $j_1, j_2 \in \mathsf{R}_H^-$, $e^{\mathrm{full}}(j_1) = e^{\mathrm{full}}(i_1)$, and $e^{\mathrm{full}}(j_2) = e^{\mathrm{full}}(i_2)$. Thus, by $e^{\mathrm{full}}(i_1), e^{\mathrm{full}}(i_2) \in \mathrm{Occ}(Q_{\mathrm{suf}}, T)$ (following from $i_1, i_2 \in \mathrm{Occ}(Q, T)$), we obtain $e^{\mathrm{full}}(j_1), e^{\mathrm{full}}(j_2) \in \mathrm{Occ}(Q_{\mathrm{suf}}, T)$. Next, we show $\mathrm{LCE}(e^{\mathrm{full}}(j_1), e^{\mathrm{full}}(j_2)) = |Q_{\mathrm{suf}}|$. As noted earlier, $e^{\mathrm{full}}(i_1) - i_1 = e^{\mathrm{full}}(i_2) - i_2 = e^{\mathrm{full}}(Q) - 1$. Thus, by $\mathrm{LCE}(i_1, i_2) = |\mathrm{str}(\mathrm{LCA}(v_1, v_2))| = |\mathrm{str}(v)| = |Q|$ and $e^{\mathrm{full}}(Q) - 1 \leq e(Q) - 1 < |Q|$, we have $|Q| = \mathrm{LCE}(i_1, i_2) = e^{\mathrm{full}}(Q) - 1 + \mathrm{LCE}(e^{\mathrm{full}}(i_1), e^{\mathrm{full}}(i_2))$. Equivalently, $\mathrm{LCE}(e^{\mathrm{full}}(i_1), e^{\mathrm{full}}(i_2)) = |Q| - e^{\mathrm{full}}(Q) + 1 = |Q_{\mathrm{suf}}|$, which by $e^{\mathrm{full}}(j_1) = e^{\mathrm{full}}(i_1)$ and $e^{\mathrm{full}}(j_2) = e^{\mathrm{full}}(i_2)$ yields $\mathrm{LCE}(e^{\mathrm{full}}(j_1), e^{\mathrm{full}}(j_2)) = |Q_{\mathrm{suf}}|$. We have thus shown that there exist distinct positions $j_1, j_2 \in \mathsf{R}_H'^-$ satisfying $e^{\mathrm{full}}(j_1), e^{\mathrm{full}}(j_2) \in \mathrm{Occ}(Q_{\mathrm{suf}}, T)$ and $\mathrm{LCE}(e^{\mathrm{full}}(j_1), e^{\mathrm{full}}(j_1)) = |Q_{\mathrm{suf}}|$. By definition of $A_{\mathsf{Z}}[1 \mathinner{\ldotp\ldotp} q]$, this implies that there exists leaves $u_1$ and $u_2$ of $\mathcal{T}_{\mathsf{Z}}$ such that $\mathrm{str}(u_1) = \mathrm{pow}(H) \cdot T[e^{\mathrm{full}}(j_1) \mathinner{\ldotp\ldotp} n]$ and $\mathrm{str}(u_2) = \mathrm{pow}(H) \cdot T[e^{\mathrm{full}}(j_2) \mathinner{\ldotp\ldotp} n]$ (see the proof of Proposition 6.17) and consequently, by Observation 4.2, the node $u = \mathrm{LCA}(u_1, u_2)$ satisfies $\mathrm{str}(u) = \mathrm{pow}(H) \cdot Q_{\mathrm{suf}}$.

2. To show that $y_1$ and $y_2$ are well-defined note that for every $i \in \mathsf{R}^-$, letting $j$ be the smallest integer satisfying $[j \mathinner{\ldotp\ldotp} i] \subseteq \mathsf{R}$, we have $j \in \mathsf{R}'^-$ and (by Lemma 5.12) $e^{\mathrm{full}}(j) = e^{\mathrm{full}}(i)$. Consequently, since $\{r_i^{\mathrm{lex}}\}_{i \in [1 \mathinner{\ldotp\ldotp} q]} = \mathsf{R}'^-$, taking $y \in [1 \mathinner{\ldotp\ldotp} q]$ such that $r_y^{\mathrm{lex}} = j$, it holds $e^{\mathrm{full}}(r_y^{\mathrm{lex}}) = e^{\mathrm{full}}(i)$. Therefore, by $\mathrm{SA}[i_1], \mathrm{SA}[i_2] \in \mathrm{Occ}(Q, T) \subseteq \mathsf{R}^-$, $y_1, y_2 \in [1 \mathinner{\ldotp\ldotp} q]$ are (uniquely) defined.

We start by showing that $\mathrm{str}(u_1) = \mathrm{pow}(H) \cdot T[e^{\mathrm{full}}(\mathrm{SA}[i_1]) \mathinner{\ldotp\ldotp} n]$ and $\mathrm{str}(u_2) = \mathrm{pow}(H) \cdot T[e^{\mathrm{full}}(\mathrm{SA}[i_2]) \mathinner{\ldotp\ldotp} n]$. As noted in Section 6.3.2, the sequence $(T[A_{\mathsf{Z}}[i] \mathinner{\ldotp\ldotp} n])_{i \in [1 \mathinner{\ldotp\ldotp} q]}$ is lexicographically sorted. Thus, by definition of $u_1$ and $u_2$, we have $\mathrm{str}(u_1) = T[A_{\mathsf{Z}}[y_1] \mathinner{\ldotp\ldotp} n]$ and $\mathrm{str}(u_2) = T[A_{\mathsf{Z}}[y_2] \mathinner{\ldotp\ldotp} n]$. As also noted in Section 6.3.2, for every $i \in [1 \mathinner{\ldotp\ldotp} q]$, $T[A_{\mathsf{Z}}[i] \mathinner{\ldotp\ldotp} n] = \mathrm{pow}(H_i) \cdot T[e^{\mathrm{full}}(r_i^{\mathrm{lex}}) \mathinner{\ldotp\ldotp} n]$, where $H_i = \text{L-root}(r_i^{\mathrm{lex}})$. Combining that with the assumptions $e^{\mathrm{full}}(r_{y_1}^{\mathrm{lex}}) = e^{\mathrm{full}}(\mathrm{SA}[i_1])$ and $e^{\mathrm{full}}(r_{y_2}^{\mathrm{lex}}) = e^{\mathrm{full}}(\mathrm{SA}[i_2])$, we therefore obtain

$$\mathrm{str}(u_1) = \mathrm{pow}(H_{y_1}) \cdot T[e^{\mathrm{full}}(\mathrm{SA}[i_1]) \mathinner{\ldotp\ldotp} n],$$
$$\mathrm{str}(u_2) = \mathrm{pow}(H_{y_2}) \cdot T[e^{\mathrm{full}}(\mathrm{SA}[i_2]) \mathinner{\ldotp\ldotp} n].$$

To obtain $\mathrm{str}(u_1) = \mathrm{pow}(H) \cdot T[e^{\mathrm{full}}(\mathrm{SA}[i_1]) \mathinner{\ldotp\ldotp} n]$ and $\mathrm{str}(u_2) = \mathrm{pow}(H) \cdot T[e^{\mathrm{full}}(\mathrm{SA}[i_2]) \mathinner{\ldotp\ldotp} n]$ it thus remains to show $H_{y_1} = H_{y_2} = H$. To this end, we first note that by $e^{\mathrm{full}}(r_{y_1}^{\mathrm{lex}}) = e^{\mathrm{full}}(\mathrm{SA}[i_1])$, (resp. $e^{\mathrm{full}}(r_{y_2}^{\mathrm{lex}}) = e^{\mathrm{full}}(\mathrm{SA}[i_2])$), and Lemmas 5.12 and 5.14, the positions $e^{\mathrm{full}}(r_{y_1}^{\mathrm{lex}})$ and $\mathrm{SA}[i_1]$ (resp. $e^{\mathrm{full}}(r_{y_2}^{\mathrm{lex}})$ and $\mathrm{SA}[i_2]$) belong to the same contiguous block of elements from $\mathsf{R}$. Next, by $|Q| \geq 3\tau - 1$ and $\mathrm{SA}[i_1] \in \mathrm{Occ}(Q, T)$, we obtain $\mathrm{lcp}(T[\mathrm{SA}[i_1] \mathinner{\ldotp\ldotp} n], Q) \geq 3\tau - 1$. Thus, by combining Lemma 5.12 and Lemma 6.9, we have $H_{y_1} = \text{L-root}(r_{y_1}^{\mathrm{lex}}) = \text{L-root}(\mathrm{SA}[i_1]) = \text{L-root}(Q) = \text{L-root}(v) = H$. Analogously, $H_{y_2} = H$.

By the above and Observation 4.2, we thus have $\text{str}(u) = \text{str}(\text{LCA}(u_1, u_2)) = \text{pow}(H) \cdot T[e^{\text{full}}(\text{SA}[i_1]) \mathinner{.\,.} e^{\text{full}}(\text{SA}[i_1]) + \ell)$, where $\ell = \text{LCE}(e^{\text{full}}(\text{SA}[i_1]), e^{\text{full}}(\text{SA}[i_2]))$. Moreover, by $\text{SA}[i_1] \in \text{Occ}(Q, T)$, we have $e^{\text{full}}(\text{SA}[i_1]) \in \text{Occ}(Q_{\text{suf}}, T)$. Thus, it remains to show that $\ell = |Q_{\text{suf}}|$. For this, recall that by $\text{SA}[i_1], \text{SA}[i_2] \in \text{Occ}(Q, T)$ we also have $e^{\text{full}}(\text{SA}[i_1]) - \text{SA}[i_1] = e^{\text{full}}(\text{SA}[i_2]) - \text{SA}[i_2] = e^{\text{full}}(Q) - 1$. Therefore, by $e^{\text{full}}(Q) - 1 \le e(Q) - 1 < |Q|$, we have $|Q| = \text{LCE}(\text{SA}[i_1], \text{SA}[i_2]) = e^{\text{full}}(Q) - 1 + \text{LCE}(e^{\text{full}}(\text{SA}[i_1]), e^{\text{full}}(\text{SA}[i_2]))$. Equivalently, $\text{LCE}(e^{\text{full}}(\text{SA}[i_1]), e^{\text{full}}(\text{SA}[i_2])) = |Q| - e^{\text{full}}(Q) + 1 = |Q_{\text{suf}}|$. We therefore obtained $\text{str}(u) = \text{pow}(H) \cdot Q_{\text{suf}}$, i.e., $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\mathsf{Z}}}(v) = u$. $\qquad\square$

**Proposition 7.31.** *Let $v$ be a periodic explicit node $v$ of $\mathcal{T}_{\text{st}}$ satisfying $e(v) \le |\text{str}(v)|$ and $\text{type}(v) = -1$. Given the data structure from Section 7.3.1 and $\text{repr}(v)$, we can in $\mathcal{O}(\log\log n)$ time compute the pointer to the node $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\mathsf{Z}}}(v)$.*

*Proof.* Denote $(b, e) = \text{repr}(v)$, $i_1 = b + 1$, and $i_2 = e$. First, using Proposition 5.24 in $\mathcal{O}(1)$ time we compute the L-exp$(\text{SA}[i_1])$ and $\delta^{\mathsf{s}}(\text{SA}[i_1])$. Next, as explained in the proof of Proposition 5.25, given $i_1$, L-exp$(\text{SA}[i_1])$, and $\delta^{\mathsf{s}}(\text{SA}[i_1])$, in $\mathcal{O}(\log\log n)$ time we compute $y_1 \in [1 \mathinner{.\,.} q]$ satisfying $e^{\text{full}}(r_{y_1}^{\text{lex}}) = e^{\text{full}}(\text{SA}[i_1])$. Note that Proposition 5.25 requires that $\text{SA}[i_1] \in \mathsf{R}^-$, which holds by $\text{SA}[i_1] \in \text{Occ}(Q, T)$, where $Q = \text{str}(v)$ (see the proof of Lemma 7.30). Analogously we compute $y_2 \in [1 \mathinner{.\,.} q]$ satisfying $e^{\text{full}}(r_{y_2}^{\text{lex}}) = e^{\text{full}}(\text{SA}[i_2])$. Then, using Proposition 4.3 in $\mathcal{O}(1)$ time we compute the pointers to the $y_1$th and $y_2$th leftmost leaves $u_1$ and $u_2$ (respectively) of $\mathcal{T}_{\mathsf{Z}}$. Then, again using Proposition 4.3, in $\mathcal{O}(1)$ time we compute and return the pointer to $u = \text{LCA}(u_1, u_2)$. By Lemma 7.30, it holds $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\mathsf{Z}}}(v) = u$. $\qquad\square$

**Mapping from $\mathcal{T}_{\mathsf{Z}}$ to $\mathcal{T}_{\text{st}}$** Let $u$ be a node of $\mathcal{T}_{\mathsf{Z}}$ such that there exists $H \in \text{Roots}$ for which $\text{pow}(H)$ is a prefix of $\text{str}(u)$. For any $\ell \ge 0$, we define $\text{pseudoinv}_{\mathcal{T}_{\mathsf{Z}}}(\ell, u)$ as follows. If, letting $s = \ell \bmod |H|$ and $k = \lfloor \frac{\ell}{|H|} \rfloor$, it holds $\mathsf{P} := \{j \in \mathsf{R}_{s,H}^- : \text{L-exp}(j) = k\} = \emptyset$, then $\text{pseudoinv}_{\mathcal{T}_{\mathsf{Z}}}(\ell, u) := (0, 0)$. Otherwise (i.e., $\mathsf{P} \neq \emptyset$), letting $b_{\mathsf{P}}, e_{\mathsf{P}} \in [0 \mathinner{.\,.} n]$ be such that $\{\text{SA}[i]\}_{i \in (b_{\mathsf{P}} \mathinner{.\,.} e_{\mathsf{P}}]} = \mathsf{P}$, $b_H, e_H \in [0 \mathinner{.\,.} q]$ be such that $\{r_i^{\text{lex}}\}_{i \in (b_H \mathinner{.\,.} e_H]} = \mathsf{R}_H'^-$, and letting $z_1 = \text{lrank}(u)$ and $z_2 = \text{rrank}(u)$, we define $\text{pseudoinv}_{\mathcal{T}_{\mathsf{Z}}}(\ell, u) := (e_{\mathsf{P}} - c_1, e_{\mathsf{P}} - c_2)$, where

$$c_1 := \text{rcount}_{A_{\text{len}}}(\ell, e_H) - \text{rcount}_{A_{\text{len}}}(\ell, z_1),$$
$$c_2 := \text{rcount}_{A_{\text{len}}}(\ell, e_H) - \text{rcount}_{A_{\text{len}}}(\ell, z_2).$$

*Remark* 7.32. To see that $H$ is well-defined, recall (see the proof of Proposition 6.17) that $\{\text{pow}(H)\}_{H \in \text{Roots}}$ is prefix-free. Thus, at most one element of $\{\text{pow}(H)\}_{H \in \text{Roots}}$ can be a prefix of $\text{str}(u)$. Furthermore, since $X \neq Y$ implies $\text{pow}(X) \neq \text{pow}(Y)$, $\text{pow}(H)$ uniquely identifies $H$.

To see that $b_{\mathsf{P}}$ and $e_{\mathsf{P}}$ are well-defined, recall that by Lemma 5.11, if $\mathsf{P} \neq \emptyset$, then all positions in $\mathsf{P}$ occupy a contiguous block in SA (see also the proof of Proposition 5.19).

Finally, to show that $b_H$ and $e_H$ are well-defined, note that $\text{pow}(H)$ being a prefix of $\text{str}(u)$ implies, by definition of $\mathcal{T}_{\mathsf{Z}}$, that there exists $i \in [1 \mathinner{.\,.} q]$ such that $H = \text{L-root}(r_i^{\text{lex}})$. Recall (see the proof of Proposition 5.21), that for any $i, i' \in [1 \mathinner{.\,.} q]$, $i < i'$ implies $\text{L-root}(r_i^{\text{lex}}) \preceq \text{L-root}(r_{i'}^{\text{lex}})$. Thus, there exists a unique $(b_H, e_H)$ (with $0 \le b_H < e_H \le q$) such that $\{r_i^{\text{lex}}\}_{i \in (b_H \mathinner{.\,.} e_H]} = \mathsf{R}_H'^-$.

*Remark* 7.33. Note that similarly as for $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\mathsf{S}}}$ (see Section 7.2.2), the mapping $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\mathsf{Z}}}$ is not necessarily injective, and hence it may not have an inverse (see also Remark 7.17). To perform the mapping from $\mathcal{T}_{\mathsf{Z}}$ to $\mathcal{T}_{\text{st}}$, we will use the above function. Although it is well-defined for every $\ell$ and $u$ (specified as above), its value is not always meaningful. Below we show a simple but useful condition where it is, and in the following sections we show the more subtle uses.

**Lemma 7.34.** *Let $P \in [0 \mathinner{.\,.} \sigma)^m$ be a periodic pattern satisfying $e(P) \le m$ and $\text{type}(P) = -1$. Denote $H = \text{L-root}(P)$, $s = \text{L-head}(P)$, $k = \text{L-exp}(P)$, $\ell = e^{\text{full}}(P) - 1$, and $P' = \text{pow}(H) \cdot$*

$P(\ell \mathinner{\ldotp\ldotp} m]$. *Assume that* $\mathsf{P} := \{j \in \mathsf{R}^-_{s,H} : \mathrm{L\text{-}exp}(j) = k\} \neq \emptyset$ *and let* $b_\mathsf{P}, e_\mathsf{P} \in [0 \mathinner{\ldotp\ldotp} n]$ *be such that* $\{\mathrm{SA}[i]\}_{i \in (b_\mathsf{P} \mathinner{\ldotp\ldotp} e_\mathsf{P}]} = \mathsf{P}$ *and* $b_H, e_H \in [0 \mathinner{\ldotp\ldotp} q]$ *be such that* $\{r^{\mathrm{lex}}_i\}_{i \in (b_H \mathinner{\ldotp\ldotp} e_H]} = \mathsf{R}'^-_H$. *Finally, let* $(b_{\mathrm{pre}}, e_{\mathrm{pre}})$ *be such that* $b_{\mathrm{pre}} = |\{i \in [1 \mathinner{\ldotp\ldotp} q] : T[A_\mathsf{Z}[i] \mathinner{\ldotp\ldotp} n] \prec P'\}|$ *and* $(b_{\mathrm{pre}} \mathinner{\ldotp\ldotp} e_{\mathrm{pre}}] = \{i \in [1 \mathinner{\ldotp\ldotp} q] : P'$ *is a prefix of* $T[A_\mathsf{Z}[i] \mathinner{\ldotp\ldotp} n]\}$. *Then, it holds*

$$(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T)) = (e_\mathsf{P} - c_1, e_\mathsf{P} - c_2),$$

*where* $c_1 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, b_{\mathrm{pre}})$ *and* $c_2 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_{\mathrm{pre}})$.

*Proof.* The proof consists of two steps:

1. First, we prove that $|\mathrm{Occ}(P,T)| = c_1 - c_2 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_{\mathrm{pre}}) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, b_{\mathrm{pre}})$. By Lemma 6.13, $\mathrm{Occ}(P,T)$ is a disjoint union of $\mathrm{Occ}^{\mathsf{a}}(P,T)$ and $\mathrm{Occ}^{\mathsf{s}}(P,T)$ (see the beginning of Section 6.3.4 for definitions). Moreover, since $e(P) \le m$, Lemma 6.14 and its symmetric version (adapted according to Lemma 6.9) imply that $\mathrm{Occ}^{\mathsf{a}}(P,T) = \emptyset$. Thus, we need to prove $|\mathrm{Occ}^{\mathsf{s}}(P,T)| = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_{\mathrm{pre}}) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, b_{\mathrm{pre}})$. By Lemma 6.11(2), it holds $\mathrm{Occ}(P,T) \subseteq \mathsf{R}^-$. Thus, $\mathrm{Occ}^{\mathsf{s}}(P,T) = \mathrm{Occ}^{\mathsf{s}-}(P,T)$. Recall now that $A_\mathsf{Z}[i] = e^{\mathrm{full}}(r^{\mathrm{lex}}_i) - |\mathrm{pow}(\mathrm{L\text{-}root}(r^{\mathrm{lex}}_i))|$. Since the set $\{\mathrm{pow}(H) : H \in \mathrm{Roots}\}$ is prefix-free, it follows, letting $H_j = \mathrm{L\text{-}root}(j)$ (where $j \in \mathsf{R}$), that

$$
\begin{aligned}
\{r^{\mathrm{lex}}_i\}_{i \in (b_{\mathrm{pre}} \mathinner{\ldotp\ldotp} e_{\mathrm{pre}}]} &= \{j \in \mathsf{R}'^- : \mathrm{pow}(H) \cdot P(\ell \mathinner{\ldotp\ldotp} m] \text{ is a prefix of } T[e^{\mathrm{full}}(j) - |\mathrm{pow}(H_j)| \mathinner{\ldotp\ldotp} n]\} \\
&= \{j \in \mathsf{R}'^- : \mathrm{pow}(H) \cdot P(\ell \mathinner{\ldotp\ldotp} m] \text{ is a prefix of } \mathrm{pow}(H_j) \cdot T[e^{\mathrm{full}}(j) \mathinner{\ldotp\ldotp} n]\} \\
&= \{j \in \mathsf{R}'^-_H : P(\ell \mathinner{\ldotp\ldotp} m] \text{ is a prefix of } T[e^{\mathrm{full}}(j) \mathinner{\ldotp\ldotp} n]\}
\end{aligned}
$$

By Lemma 6.16, we thus have $|\mathrm{Occ}^{\mathsf{s}-}(P,T)| = |\{i \in (b_{\mathrm{pre}} \mathinner{\ldotp\ldotp} e_{\mathrm{pre}}] : e^{\mathrm{full}}(r^{\mathrm{lex}}_i) - r^{\mathrm{lex}}_i \ge e^{\mathrm{full}}(P) - 1\}| = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_{\mathrm{pre}}) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, b_{\mathrm{pre}})$ (recall, that $A_{\mathrm{len}}[i] = e^{\mathrm{full}}(r^{\mathrm{lex}}_i) - r^{\mathrm{lex}}_i$; see Section 5.3.2).

2. Second, we prove that $\mathrm{RangeBeg}(P,T) = e_\mathsf{P} - c_1$. We start by observing that since $P$ is periodic and satisfies $\mathrm{type}(P) = -1$, it follows from Lemmas 6.19 and 6.20 that $\mathrm{RangeBeg}(P,T) = \mathrm{RangeBeg}(X,T) + \delta(P,T) = \mathrm{RangeBeg}(X,T) + \delta^{\mathsf{a}}(P,T) - \delta^{\mathsf{s}}(P,T)$, where $X = P[1 \mathinner{\ldotp\ldotp} 3\tau - 1]$. On the other hand, combining the equalities $\mathrm{L\text{-}head}(P) = s$, $\mathrm{L\text{-}root}(P) = H$, $\mathrm{L\text{-}exp}(P) = k$, and $\mathrm{type}(P) = -1$ with the definition of $\mathsf{P}$ yields $\mathrm{RangeBeg}(X,T) + \delta^{\mathsf{a}}(P,T) = e_\mathsf{P}$. Consequently, we obtain $\mathrm{RangeBeg}(P,T) = e_\mathsf{P} - \delta^{\mathsf{s}}(P,T)$. It thus remains to show $\delta^{\mathsf{s}}(P,T) = c_1$. By utilizing that by definition of the sequence $(r^{\mathrm{lex}}_i)_{i \in [1 \mathinner{\ldotp\ldotp} q]}$, for every $i, i' \in (b_H \mathinner{\ldotp\ldotp} e_H]$, $i < i'$ implies $T[e^{\mathrm{full}}(r^{\mathrm{lex}}_i) \mathinner{\ldotp\ldotp} n] \preceq T[e^{\mathrm{full}}(r^{\mathrm{lex}}_{i'}) \mathinner{\ldotp\ldotp} n]$, it follows by the above formula for $\{r^{\mathrm{lex}}_i\}_{i \in (b_{\mathrm{pre}} \mathinner{\ldotp\ldotp} e_{\mathrm{pre}}]}$ that

$$\{r^{\mathrm{lex}}_i\}_{i \in (b_{\mathrm{pre}} \mathinner{\ldotp\ldotp} e_H]} = \{j \in \mathsf{R}'^-_H : P(\ell \mathinner{\ldotp\ldotp} m] \preceq T[e^{\mathrm{full}}(j) \mathinner{\ldotp\ldotp} n]\}.$$

By Lemma 6.22, we thus have $\delta^{\mathsf{s}}(P,T) = |\{i \in (b_{\mathrm{pre}} \mathinner{\ldotp\ldotp} e_H] : e^{\mathrm{full}}(r^{\mathrm{lex}}_i) - r^{\mathrm{lex}}_i \ge e^{\mathrm{full}}(P) - 1\}| = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, b_{\mathrm{pre}}) = c_1$. $\qquad\square$

*Remark 7.35.* Note that since the range $(b_{\mathrm{pre}} \mathinner{\ldotp\ldotp} e_{\mathrm{pre}}]$ is well-defined even if $e_{\mathrm{pre}} - b_{\mathrm{pre}} = 0$, the above lemma holds even if $|\mathrm{Occ}(P,T)| = 0$.

**Lemma 7.36.** *Let $v$ be an explicit periodic node of $\mathcal{T}_{\mathrm{st}}$ such that $e(v) \le |\mathrm{str}(v)|$ and $\mathrm{type}(v) = -1$. Let $u = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_\mathsf{Z}}(v)$ and $\ell = e^{\mathrm{full}}(v) - 1$. Then, it holds $\mathrm{repr}(v) = \mathrm{pseudoinv}_{\mathcal{T}_\mathsf{Z}}(\ell, u)$.*

*Proof.* Denote $H = \mathrm{L\text{-}root}(v)$, $P = \mathrm{str}(v)$, and $(b, e) = \mathrm{pseudoinv}_{\mathcal{T}_\mathsf{Z}}(\ell, u)$. Let $s = \ell \bmod |H|$, $k = \lfloor \frac{\ell}{|H|} \rfloor$, and $\mathsf{P} = \{j \in \mathsf{R}^-_{s,H} : \mathrm{L\text{-}exp}(j) = k\}$. Note that we then have $\mathrm{L\text{-}head}(P) = (e^{\mathrm{full}}(P) - 1) \bmod |H| = s$ and $\mathrm{L\text{-}exp}(P) = \lfloor \frac{e^{\mathrm{full}}(P) - 1}{|H|} \rfloor = k$. Observe that this implies $\mathsf{P} \neq \emptyset$. To see this, consider any $j \in \mathrm{Occ}(\mathrm{str}(v), T) = \mathrm{Occ}(P,T)$. By Lemma 6.9, it follows that $j \in \mathsf{R}$,

L-root$(j)$ = L-root$(P)$ = $H$, and L-head$(j)$ = L-head$(P)$ = $s$, i.e., $j \in \mathsf{R}_{s,H}$. Furthermore, by $e(P) \le |P|$ and type$(P) = -1$ we obtain from Lemma 6.11(2) that L-exp$(j)$ = L-exp$(P)$ = $k$ and type$(j)$ = type$(P)$ = $-1$. Thus, $j \in \mathsf{P}$ and consequently $\mathsf{P} \ne \emptyset$. By definition of pseudoinv$_{\mathcal{T}_Z}(\ell, u)$, we thus obtain that $(b, e) = (e_{\mathsf{P}} - c_1, e_{\mathsf{P}} - c_2)$, where $b_{\mathsf{P}}, e_{\mathsf{P}} \in [0 \mathinner{.\,.} n]$ are such that $\{\mathrm{SA}[i]\}_{i \in (b_{\mathsf{P}} \mathinner{.\,.} e_{\mathsf{P}}]} = \mathsf{P}$, $b_H, e_H \in [0 \mathinner{.\,.} q]$ are such that $\{r_i^{\mathrm{lex}}\}_{i \in (b_H \mathinner{.\,.} e_H]} = \mathsf{R}'^{-}_H$, $z_1 = \mathrm{lrank}(u)$, $z_2 = \mathrm{rrank}(u)$, and

$$c_1 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, z_1),$$
$$c_2 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, z_2).$$

By definition of $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_Z}(v)$, we have $\mathrm{str}(u) = \mathrm{pow}(H) \cdot P[e^{\mathrm{full}}(P) \mathinner{.\,.} |P|]$. Thus, denoting $P' = \mathrm{pow}(H) \cdot P[e^{\mathrm{full}}(P) \mathinner{.\,.} |P|]$, by definition of $\mathcal{T}_Z$, we have $\mathrm{lrank}(u) = |\{i \in [1 \mathinner{.\,.} q] : T[A_Z[i] \mathinner{.\,.} n] \prec P'\}|$ and $(\mathrm{lrank}(u) \mathinner{.\,.} \mathrm{rrank}(u)] = \{i \in [1 \mathinner{.\,.} q] : P' \text{ is a prefix of } T[A_Z[i] \mathinner{.\,.} n]\}$. By Lemma 7.34, this implies that $(\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T)) = (e_{\mathsf{P}} - (\mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(u))), e_{\mathsf{P}} - (\mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(u)))) = (e_{\mathsf{P}} - c_1, e_{\mathsf{P}} - c_2) = (b, e)$. This immediately implies $\mathrm{repr}(v) = \mathrm{pseudoinv}_{\mathcal{T}_Z}(\ell, u)$. $\qquad\square$

**Proposition 7.37.** *Let $H \in \mathrm{Roots}$ and let $u$ be a node of $\mathcal{T}_Z$ such that $\mathrm{pow}(H)$ is a prefix of $\mathrm{str}(u)$. Given the data structure from Section 7.3.1, a pointer to $u$, and integers $\mathrm{int}(H)$ and $\ell \ge 0$, we can in $\mathcal{O}(\log\log n)$ time compute the pair $\mathrm{pseudoinv}_{\mathcal{T}_Z}(\ell, u)$.*

*Proof.* Let $p := |H|$. We first compute $s := \ell \bmod p$ and $k = \lfloor \frac{\ell}{p} \rfloor$. Next, using the lookup tables $L_{\mathrm{pref}}$ and $L_{\mathrm{range}}$, we compute in $\mathcal{O}(1)$ time the pair $(b_{s,H}, e_{s,H}) = (\mathrm{RangeBeg}(X, T), \mathrm{RangeEnd}(X, T))$, where $X = \mathrm{Pref}_{3\tau - 1}(s, H)$. By Lemma 6.9, we then have that $b_{s,H} = e_{s,H}$ holds if and only if $\mathsf{R}_{s,H} = \emptyset$, and if $b_{s,H} \ne e_{s,H}$ then $\{\mathrm{SA}[i] : i \in (b_{s,H} \mathinner{.\,.} e_{s,H}]\} = \mathsf{R}_{s,H}$. If $b_{s,H} = e_{s,H}$, we return $\mathrm{pseudoinv}_{\mathcal{T}_Z}(\ell, u) = (0, 0)$. Let us now assume $b_{s,H} \ne e_{s,H}$.

Next, using the data structure from Section 5.3.2, as explained in the proof of Proposition 5.19, in $\mathcal{O}(1)$ time we compute the pair $(b_{\mathsf{P}}, e_{\mathsf{P}})$ satisfying $\{\mathrm{SA}[i]\}_{i \in (b_{\mathsf{P}} \mathinner{.\,.} e_{\mathsf{P}}]} = \mathsf{P}$, where $\mathsf{P} = \{j \in \mathsf{R}^{-}_{s,H} : \text{L-exp}(j) = k\}$. More precisely, first, in $\mathcal{O}(1)$ time we compute $d = \mathsf{rank}_{B_{\mathrm{exp}},1}(e_{s,H}) - \mathsf{rank}_{B_{\mathrm{exp}},1}(b_{s,H})$. If $d = 0$, then $\mathsf{R}^{-}_{s,H} = \emptyset$, and hence we return $\mathrm{pseudoinv}_{\mathcal{T}_Z}(\ell, u) = (0, 0)$. Otherwise, in $\mathcal{O}(1)$ time we retrieve $k_{\min} = L_{\mathrm{minexp}}[\mathrm{int}(X)]$. Then, letting $k_{\max} = k_{\min} + d - 1$, we have $[k_{\min} \mathinner{.\,.} k_{\max}] = \{\text{L-exp}(j) : j \in \mathsf{R}^{-}_{s,H}\}$. If $k \notin [k_{\min} \mathinner{.\,.} k_{\max}]$, then $\mathsf{P} = \emptyset$, and thus we return $\mathrm{pseudoinv}_{\mathcal{T}_Z}(\ell, u) = (0, 0)$. Otherwise, we have two cases. Let $p = \mathsf{rank}_{B_{\mathrm{exp}},1}(b_{s,H})$. If $k = k_{\min}$, then in $\mathcal{O}(1)$ time we compute $(b_{\mathsf{P}}, e_{\mathsf{P}}) = (b_{s,H}, \mathsf{select}_{B_{\mathrm{exp}},1}(p + 1))$. If $k > k_{\min}$, in $\mathcal{O}(1)$ time we compute $(b_{\mathsf{P}}, e_{\mathsf{P}}) = (\mathsf{select}_{B_{\mathrm{exp}},1}(p + k - k_{\min}), \mathsf{select}_{B_{\mathrm{exp}},1}(p + k + 1 - k_{\min}))$.

For the final step, we first in $\mathcal{O}(1)$ time compute $e_H = \sum_{H' \preceq H} |\mathsf{R}'^{-}_{H'}|$ using the lookup table $L_{\mathrm{runs}}$ stored as part of the structure from Section 5.3.2. Then, it holds that there exists $b_H < e_H$ such that $\{r_i^{\mathrm{lex}}\}_{i \in (b_H \mathinner{.\,.} e_H]} = \mathsf{R}'^{-}_H$. Then, in $\mathcal{O}(1)$ time we obtain $z_1 = \mathrm{lrank}(u)$ and $z_2 = \mathrm{rrank}(u)$ (Proposition 4.3). Finally, in $\mathcal{O}(\log\log n)$ time we compute $c_1 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, z_1)$ and $c_2 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, z_2)$ and return $\mathrm{pseudoinv}_{\mathcal{T}_Z}(\ell, u) = (e_{\mathsf{P}} - c_1, e_{\mathsf{P}} - c_2)$. The range counting queries are implemented using the structure from Proposition 2.3 for the array $A$, which is stored as part of the structure from Section 5.3.2. $\qquad\square$

**Handling Nodes Satisfying $e(v) > |\mathrm{str}(v)|$** Next, we present a combinatorial result describing how to compute the value $e(v)$, and to check if it holds $e(v) > |\mathrm{str}(v)|$. We then show how to compute $(\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T))$ and $(\mathrm{RangeBeg}(Pc, T), \mathrm{RangeEnd}(Pc, T))$ for any periodic pattern $P \in [0 \mathinner{.\,.} \sigma)^{+}$ satisfying $e(P) > |P|$. We will use it to efficiently perform queries on periodic nodes $v$ of $\mathcal{T}_{\mathrm{st}}$ satisfying $e(v) > |\mathrm{str}(v)|$.

**Lemma 7.38.** *Let $v$ be an explicit periodic node of $\mathcal{T}_{\mathrm{st}}$. Let $i_1 = \mathrm{lrank}(v) + 1$ and $i_2 = \mathrm{rrank}(v)$.*

1. *It holds* $\mathrm{SA}[i_1], \mathrm{SA}[i_2] \in \mathsf{R}$ *and* $e(v) = 1 + \min(e(\mathrm{SA}[i_1]) - \mathrm{SA}[i_1], e(\mathrm{SA}[i_2]) - \mathrm{SA}[i_2])$.
2. $e(v) \leq |\mathrm{str}(v)|$ *holds if and only if* $T[\mathrm{SA}[i_1] + e(v) - 1] = T[\mathrm{SA}[i_2] + e(v) - 1]$.

*Proof.* Denote $\ell = \mathrm{sdepth}(v)$, $b = \mathrm{lrank}(v)$, and $e = \mathrm{rrank}(v)$.

1. Let $s = \mathrm{L\text{-}head}(v)$, $H = \mathrm{L\text{-}root}(v)$, $p = |H|$, and $Q = \mathrm{str}(v)$. By definition, we have $b < e$ and $\{\mathrm{SA}[i]\}_{i \in (b\,..\,e]} = \mathrm{Occ}(Q, T)$. On the other hand, by $|Q| \geq 3\tau - 1$ and Lemma 6.9, for every $j \in \mathrm{Occ}(Q, T)$ it holds $j \in \mathsf{R}_{s,H}$. In particular, we thus obtain $\mathrm{SA}[i_1], \mathrm{SA}[i_2] \in \mathsf{R}$.

Next, we prove the following two facts.

- First, we show that there exists $t \in \{1, 2\}$ satisfying $e(v) - 1 = e(\mathrm{SA}[i_t]) - \mathrm{SA}[i_t]$. By definition, it holds $\ell = \mathrm{LCE}(\mathrm{SA}[i_1], \mathrm{SA}[i_2])$ and $T[\mathrm{SA}[i_1]\,..\,\mathrm{SA}[i_1] + \ell) = T[\mathrm{SA}[i_2]\,..\,\mathrm{SA}[i_2] + \ell)$. If $e - b = 1$, then any $t \in \{1, 2\}$ satisfies the claim, since then $\ell = n - \mathrm{SA}[i_t] + 1$, and thus it follows from $\mathrm{SA}[i_t] \in \mathrm{Occ}(Q, T)$ that

$$\begin{aligned} e(v) - 1 &= p + \mathrm{lcp}(Q(0\,..\,\ell - p], Q(p\,..\,\ell]) \\ &= p + \mathrm{lcp}(T[\mathrm{SA}[i_t]\,..\,\mathrm{SA}[i_t] + \ell - p), T[\mathrm{SA}[i_t] + p\,..\,\mathrm{SA}[i_t] + \ell)) \\ &= p + \mathrm{lcp}(T[\mathrm{SA}[i_t]\,..\,n - p], T[\mathrm{SA}[i_t] + p\,..\,n]) \\ &= p + \mathrm{LCE}(\mathrm{SA}[i_t], \mathrm{SA}[i_t] + p) \\ &= e(\mathrm{SA}[i_t]) - \mathrm{SA}[i_t]. \end{aligned}$$

  Assume now $e - b > 1$. Then, $T[\mathrm{SA}[i_1] + \ell] \neq T[\mathrm{SA}[i_2] + \ell]$. [11] Thus, by $T[\mathrm{SA}[i_1] + \ell - p] = T[\mathrm{SA}[i_2] + \ell - p]$ there exists $t \in \{1, 2\}$ such that $T[\mathrm{SA}[i_t] + \ell] \neq T[\mathrm{SA}[i_t] + \ell - p]$. For such $t$, we have $\mathrm{LCE}(\mathrm{SA}[i_t], p + \mathrm{SA}[i_t]) \leq \ell - p$ and hence $e(\mathrm{SA}[i_t]) - \mathrm{SA}[i_t] = p + \mathrm{LCE}(\mathrm{SA}[i_t], \mathrm{SA}[i_t] + p) \leq \ell$. We therefore obtain $e(\mathrm{SA}[i_t]) - \mathrm{SA}[i_t] = p + \mathrm{lcp}(T[\mathrm{SA}[i_t]\,..\,\mathrm{SA}[i_t] + \ell - p), T[\mathrm{SA}[i_t] + p\,..\,\mathrm{SA}[i_t] + \ell))$. On the other hand, by $Q = T[\mathrm{SA}[i_t]\,..\,\mathrm{SA}[i_t] + \ell)$ we have $e(v) - 1 = p + \mathrm{lcp}(Q(0\,..\,\ell - p], Q(p\,..\,\ell]) = p + \mathrm{lcp}(T[\mathrm{SA}[i_t]\,..\,\mathrm{SA}[i_t] + \ell - p), T[\mathrm{SA}[i_t] + p\,..\,\mathrm{SA}[i_t] + \ell))$. Therefore, $e(v) - 1 = e(\mathrm{SA}[i_t]) - \mathrm{SA}[i_t]$.
- Second, we show that for every $i \in (b\,..\,e]$, it holds $e(v) - 1 \leq e(\mathrm{SA}[i]) - \mathrm{SA}[i]$. For this, recall that $e(\mathrm{SA}[i]) - \mathrm{SA}[i] = p + \mathrm{LCE}(\mathrm{SA}[i], \mathrm{SA}[i] + p)$. Therefore, by $\mathrm{SA}[i] \in \mathrm{Occ}(Q, T)$, we obtain $e(v) - 1 = e(Q) - 1 = p + \mathrm{lcp}(Q(0\,..\,\ell - p], Q(p\,..\,\ell]) = p + \mathrm{lcp}(T[\mathrm{SA}[i]\,..\,\mathrm{SA}[i] + \ell - p), T[\mathrm{SA}[i] + p\,..\,\mathrm{SA}[i] + \ell)) \leq p + \mathrm{LCE}(\mathrm{SA}[i], \mathrm{SA}[i] + p) = e(\mathrm{SA}[i]) - \mathrm{SA}[i]$.

By the above two facts, we obtain $\min(e(\mathrm{SA}[i_1]) - \mathrm{SA}[i_1], e(\mathrm{SA}[i_2]) - \mathrm{SA}[i_2]) = \min(e(\mathrm{SA}[i_t]) - \mathrm{SA}[i_t], e(\mathrm{SA}[i_{3-t}]) - \mathrm{SA}[i_{3-t}]) = e(v) - 1$.

2. We start by showing that $\mathrm{SA}[i_1] + e(v) - 1, \mathrm{SA}[i_2] + e(v) - 1 \leq n$. Observe that for every $j \in \mathsf{R}$, by the uniqueness of $T[n]$, it holds $e(j) \leq n$. Consider any $i \in (b\,..\,e]$. Above, we proved $e(v) - 1 \leq e(\mathrm{SA}[i]) - \mathrm{SA}[i]$. Thus, we obtain $\mathrm{SA}[i] + e(v) - 1 \leq e(\mathrm{SA}[i]) \leq n$. In particular, $\mathrm{SA}[i_1] + e(v) - 1, \mathrm{SA}[i_2] + e(v) - 1 \leq n$. We now prove the equivalence. Recall, that $|\mathrm{str}(v)| = \ell = \mathrm{LCE}(\mathrm{SA}[i_1], \mathrm{SA}[i_2])$ holds by definition. Let us first assume $e(v) \leq \ell$. By the assumption $T[\mathrm{SA}[i_1]\,..\,\mathrm{SA}[i_1] + \ell) = T[\mathrm{SA}[i_2]\,..\,\mathrm{SA}[i_2] + \ell)$, this immediately implies $T[\mathrm{SA}[i_1] + e(v) - 1] = T[\mathrm{SA}[i_2] + e(v) - 1]$. To show the opposite implication, assume by contraposition that $e(v) > \ell$. Since by definition we have $e(v) \leq |\mathrm{str}(v)| + 1$, we must have $e(v) = \ell + 1$. Then, by definition of LCE, we have $T[\mathrm{SA}[i_1] + e(v) - 1] = T[\mathrm{SA}[i_1] + \ell] \neq T[\mathrm{SA}[i_2] + \ell] = T[\mathrm{SA}[i_2] + e(v) - 1]$. $\qquad\square$

---

[11] To see that symbols $T[\mathrm{SA}[i_1] + \ell]$ and $T[\mathrm{SA}[i_2] + \ell]$ are well-defined, observe that by $\ell > 0$ and $b + 1 < e$, it follows that $\mathrm{SA}[i_1] + \ell - 1 \neq \mathrm{SA}[i_2] + \ell - 1$. On the other hand, we have $T[\mathrm{SA}[i_1] + \ell - 1] = T[\mathrm{SA}[i_2] + \ell - 1]$. Thus, by the uniqueness of $T[n]$ we must have $\mathrm{SA}[i_1] + \ell - 1 < n$ and $\mathrm{SA}[i_2] + \ell - 1 < n$

**Proposition 7.39.** *Let $P \in [0 \mathinner{..} \sigma)^+$ be a periodic pattern satisfying $e(P) > |P|$. Given the structure from Section 7.3.1, and the values $\mathrm{L\text{-}head}(P)$, $\mathrm{L\text{-}root}(P)$, and $|P|$, we can in $\mathcal{O}(\log \log n)$ time compute the pair $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$.*

*Proof.* Denote $s = \mathrm{L\text{-}head}(P)$, $H = \mathrm{L\text{-}root}(P)$, and $m = |P|$. First, in $\mathcal{O}(1)$ time we compute $k := \mathrm{L\text{-}exp}(P) = \lfloor \frac{m-s}{|H|} \rfloor$ and $t := \mathrm{L\text{-}tail}(P) = m - s - k|H|$. Next, using the lookup table $L_{\mathrm{pref}}$, in $\mathcal{O}(1)$ time we compute $X := \mathrm{Pref}_{3\tau-1}(s, H) = P[1 \mathinner{..} 3\tau-1]$.

Next, we compute $|\mathrm{Occ}(P,T)|$. Recall that by Lemma 6.13, $|\mathrm{Occ}(P,T)| = |\mathrm{Occ}^{\mathsf{a}}(P,T)| + |\mathrm{Occ}^{\mathsf{s}}(P,T)| = |\mathrm{Occ}^{\mathsf{a}-}(P,T)| + |\mathrm{Occ}^{\mathsf{a}+}(P,T)| + |\mathrm{Occ}^{\mathsf{s}-}(P,T)| + |\mathrm{Occ}^{\mathsf{s}+}(P,T)|$ (see Section 6.3.4).

- To compute $|\mathrm{Occ}^{\mathsf{a}-}(P,T)|$, we proceed as in the proof of Proposition 6.15, except for one modification: Since we already have $\mathrm{L\text{-}head}(P)$, $\mathrm{L\text{-}root}(P)$, and $\mathrm{L\text{-}exp}(P)$ (note that we do not need $\mathrm{L\text{-}tail}(P)$ here since we assumed $e(P) = |P| + 1$), we can skip the first step which takes $\mathcal{O}(1 + m/\log_\sigma n)$ time. Note that after such modification, we no longer need the packed representation of the whole pattern $P$, but only $P[1 \mathinner{..} 3\tau - 1]$, which we computed above. The rest of the algorithm in Proposition 6.15 takes $\mathcal{O}(1)$ time. The structures from Proposition 6.15 that we used (augmented bitvector $B_{\mathrm{exp}}$ and lookup tables $L_{\mathrm{minexp}}$ and $L_{\mathrm{range}}$) are components of the structure from Section 7.3.1.
- To compute $|\mathrm{Occ}^{\mathsf{s}-}(P,T)|$, we proceed as in Proposition 6.17, except for two modifications. First, we again already have $\mathrm{L\text{-}head}(P)$, $\mathrm{L\text{-}root}(P)$, and $\mathrm{L\text{-}exp}(P)$, which lets us skip the first step taking $\mathcal{O}(1 + m/\log_\sigma n)$ time. Second, rather than computing $b_{\mathrm{pre}}$ and $e_{\mathrm{pre}}$ in $\mathcal{O}(m/\log_\sigma n + \log \log n)$ time, we use the lookup table $L_{\mathrm{runs}}$ stored in the structure from Section 7.3.1. More precisely, $b_{\mathrm{pre}}$ and $e_{\mathrm{pre}}$ are obtained in $\mathcal{O}(1)$ time by looking up in $L_{\mathrm{runs}}$ the pair associated with the key $(H, H')$, where $H'$ is a length-$t$ prefix of $H$ (note that $P[e^{\mathrm{full}}(P) \mathinner{..} m] = H'$). The rest of the algorithm in Proposition 6.17 takes $\mathcal{O}(\log \log n)$ time. Again, the components used in Proposition 6.17 are present in structure from Section 7.3.1.

The values $|\mathrm{Occ}^{\mathsf{a}+}(P,T)|$ and $|\mathrm{Occ}^{\mathsf{s}+}(P,T)|$ are computed analogously (see the proof of Proposition 6.18) using the symmetric components of the structure from Section 7.3.1. We can thus compute $|\mathrm{Occ}(P,T)|$ in $\mathcal{O}(\log \log n)$ time.

The next step of the algorithm is to compute $\delta(P,T)$ (Section 6.3.4). Observe (see Section 6.3.1) that $e(P) > |P|$ implies $\mathrm{type}(P) = -1$. Recall that for such $P$, by Lemma 6.20, $\delta(P,T) = \delta^{\mathsf{a}}(P,T) - \delta^{\mathsf{s}}(P,T)$.

- To compute $\delta^{\mathsf{a}}(P,T)$, we proceed as in the proof of Proposition 6.21, employing the same modification as when computing $|\mathrm{Occ}^{\mathsf{a}-}(P,T)|$ above. Thus, the computation takes $\mathcal{O}(1)$ time. Proposition 6.21 uses the structure from Proposition 6.15 and, as above, the used components are already present in the structure from Section 7.3.1.
- To compute $\delta^{\mathsf{s}}(P,T)$, we observe that for a periodic pattern $P$ satisfying $e(P) > |P|$, it holds by Lemma 6.9(2) that $\mathrm{Pos}^{\mathsf{s}}(P,T) = \mathrm{Occ}^{\mathsf{s}-}(P,T)$. Consequently, we can compute $\delta^{\mathsf{s}}(P,T) = |\mathrm{Occ}^{\mathsf{s}-}(P,T)|$ as above in $\mathcal{O}(\log \log n)$ time.

Combining the above two steps, the computation of $\delta(P,T)$ takes $\mathcal{O}(\log \log n)$ time.

We use the above values to obtain $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$ as follows. By Lemma 6.19, $\mathrm{RangeBeg}(P,T) = \mathrm{RangeBeg}(X,T) + \delta(P,T)$, where $X = P[1 \mathinner{..} 3\tau-1]$. The value $\mathrm{RangeBeg}(X,T)$ is obtained in $\mathcal{O}(1)$ time using the lookup table $L_{\mathrm{range}}$. We thus obtain $\mathrm{RangeBeg}(P,T)$. By definition, we then compute $\mathrm{RangeEnd}(P,T) = \mathrm{RangeEnd}(P,T) + |\mathrm{Occ}(P,T)|$. In total, the query takes $\mathcal{O}(\log \log n)$ time. $\qquad\square$

*Remark 7.40.* Note that the above result holds even if $\mathrm{Occ}(P,T) = \emptyset$. Thus, it is more general than the result needed to support efficient processing of periodic nodes $v$ of $\mathcal{T}_{\mathrm{st}}$ satisfying $e(v) > |\mathrm{str}(v)|$, since for such nodes we have $\mathrm{Occ}(\mathrm{str}(v), T) \neq \emptyset$.

**Proposition 7.41.** *Let $P \in [0 \mathinner{.\,.} \sigma)^+$ be a periodic pattern satisfying $e(P) > |P|$. Given any $c \in [0 \mathinner{.\,.} \sigma)$, the structure from Section 7.3.1, and the values* L-head$(P)$, L-root$(P)$, *and* $|P|$, *we can in* $\mathcal{O}(\log \log n)$ *time compute the pair* $(\text{RangeBeg}(Pc, T), \text{RangeEnd}(Pc, T))$.

*Proof.* Denote $P' = Pc$ and $m = |P| + 1 = |P'|$. Observe, that since $P$ is periodic and it is a prefix of $P'$, by Lemma 6.10, $P'$ is also periodic and it holds L-head$(P) = $ L-head$(P')$ and L-root$(P) = $ L-root$(P')$. Let us denote $s = $ L-head$(P) = $ L-head$(P')$ and $H = $ L-root$(P) = $ L-root$(P')$. By the assumption, we have $e(P) = m$. First, in $\mathcal{O}(1)$ time we compute $t := $ L-tail$(P) = (m - 1 - s) \bmod |H|$. We then check if $e(P') \le |P'|$ by comparing $c$ to $H[t + 1]$. If $c = H[t + 1]$, then we have $e(P') > |P'|$. Since we have L-head$(P') = s$, L-root$(P') = H$, and $|P'| = m$, in $\mathcal{O}(\log \log n)$ time we thus compute and return $(\text{RangeBeg}(P', T), \text{RangeEnd}(P', T))$ using Proposition 7.39. Let us thus assume $c \ne H[t + 1]$, i.e., $e(P') \le |P'|$. We then compute type$(P')$ by comparing $c$ with $H[t + 1]$. Let us assume that $c \prec H[t + 1]$, i.e., type$(P) = -1$ (the case type$(P) = +1$ is handled symmetrically). We now execute the modified algorithm from Proposition 6.24 for $P'$. The modification is to replace implementation of operations taking $\Theta(m/\log_\sigma n)$ time with faster alternatives, exploiting the fact that by $e(P') = e(P)$, L-head$(P') = $ L-head$(P)$, and L-root$(P') = $ L-root$(P)$ it follows that $e^{\text{full}}(P') = e^{\text{full}}(P) = e(P) - $ L-tail$(P) = m - t$ and thus $P'[e^{\text{full}}(P') \mathinner{.\,.} m]$ is of length $t + 1 \le \tau$ (importantly, the modified algorithm will not use the components of the data structures in Section 6.3 which are not part of the structure from Section 7.3.1). More precisely:

- First, using $L_{\text{pref}}$, in $\mathcal{O}(1)$ time we compute $X = \text{Pref}_{3\tau - 1}(s, H) = P'[1 \mathinner{.\,.} 3\tau - 1]$.
- We then compute $|\text{Occ}(P', T)|$. First, note that since $e(P') \le |P'|$ and type$(P') = -1$, it follows by Lemma 6.11(2) that $\text{Occ}(P', T) \subseteq \mathsf{R}^-$, and that for every $j \in \text{Occ}(P', T)$ it holds L-exp$(j) = $ L-exp$(P')$. Thus, the sets $\text{Occ}^{\mathsf{a}}(P', T)$ and $\text{Occ}^{\mathsf{s}+}(P', T)$ are empty, and hence it remains to explain the computation of $|\text{Occ}^{\mathsf{s}-}(P', T)|$ (Proposition 6.17). Observe, that the expensive operations are the computation of L-head$(P')$, L-root$(P')$, L-exp$(P')$, and the pair $(b_{\text{pre}}, e_{\text{pre}})$. Observe, however, that here we already have $s = $ L-head$(P')$, $H = $ L-root$(P')$, and $e(P') = m$. This lets us deduce $k := $ L-exp$(P') = \lfloor \frac{m - 1 - s}{|H|} \rfloor$ in $\mathcal{O}(1)$ time. As for the computation of $(b_{\text{pre}}, e_{\text{pre}})$, we first in $\mathcal{O}(1)$ time compute $H' := P'[e^{\text{full}}(P') \mathinner{.\,.} m] = H[1 \mathinner{.\,.} t + 1]$, and then obtain $(b_{\text{pre}}, e_{\text{pre}})$ by looking up the pair associated with the key $(H, H')$ in the lookup table $L_{\text{runs}}$. The rest of the algorithm in Proposition 6.17 takes $\mathcal{O}(\log \log n)$ time.
- Finally, we compute $\delta(P', T)$. By type$(P') = -1$ and Lemma 6.20, it holds $\delta(P', T) = \delta^{\mathsf{a}}(P', T) - \delta^{\mathsf{s}}(P', T)$. To compute $\delta^{\mathsf{a}}(P', T)$, we proceed as in the proof of Proposition 6.21. The string $X$ was already obtained above. The expensive step in Proposition 6.21 is the computation of L-root$(P')$ and L-exp$(P')$. As noted above, here we already have L-root$(P') = H$, and in $\mathcal{O}(1)$ time we can compute L-exp$(P') = \lfloor \frac{e(P') - 1 - \text{L-head}(P')}{|\text{L-root}(P')|} \rfloor = \lfloor \frac{m - 1 - s}{|H|} \rfloor$. The rest of the algorithm in Proposition 6.21 takes $\mathcal{O}(1)$ time. We then compute $\delta^{\mathsf{s}}(P', T)$ using a modified Proposition 6.23. The expensive part is the computation of $x$ and $x'$. After those are computed, the rest takes $\mathcal{O}(\log \log n)$ time. Here, we obtain $x$ by observing that it is equal to $b_{\text{pre}}$ (which was computed above), and then obtain $x'$ using $L_{\text{runs}}$ (this only requires knowing L-root$(P')$, which we already have).

Note that all components of the structure from Propositions 6.17, 6.21, and 6.23 that we used are also components of the structure from Section 7.3.1. Using the above values, we now obtain $(\text{RangeBeg}(P', T), \text{RangeEnd}(P', T))$ as follows. By Lemma 6.19, it holds $\text{RangeBeg}(P', T) = \text{RangeBeg}(X, T) + \delta(P', T)$, where $X = P'[1 \mathinner{.\,.} 3\tau - 1]$. The value $\text{RangeBeg}(X, T)$ is obtained in $\mathcal{O}(1)$ time using the lookup table $L_{\text{range}}$. We thus obtain $\text{RangeBeg}(P', T)$. By definition, we then compute $\text{RangeEnd}(P', T) = \text{RangeEnd}(P', T) + |\text{Occ}(P', T)|$. $\qquad \square$

*Remark* 7.42. Note that, analogously to Proposition 7.39 (see Remark 7.40), the above result holds even if $\mathrm{Occ}(Pc, T) = \emptyset$.

### 7.3.3 Implementation of $\mathrm{LCA}(u, v)$

**Lemma 7.43.** *Let $v_1$ and $v_2$ be explicit nodes of $\mathcal{T}_{\mathrm{st}}$ such that $\mathrm{LCA}(v_1, v_2)$ is periodic and it holds $e(\mathrm{LCA}(v_1, v_2)) \leq |\mathrm{str}(\mathrm{LCA}(v_1, v_2))|$ and $\mathrm{type}(\mathrm{LCA}(v_1, v_2)) = -1$. Then, $v_1$ and $v_2$ are periodic and it holds $e(v_1) \leq |\mathrm{str}(v_1)|$, $e(v_2) \leq |\mathrm{str}(v_2)|$, and $\mathrm{type}(v_1) = \mathrm{type}(v_2) = -1$. Moreover,*

$$\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(\mathrm{LCA}(v_1, v_2)) = \mathrm{LCA}(\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v_1), \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v_2)).$$

*Proof.* Denote $v = \mathrm{LCA}(v_1, v_2)$, $Q = \mathrm{str}(v)$, $H = \text{L-root}(v)$, and $s = \text{L-head}(v)$. Let also $Q_1 = \mathrm{str}(v_1)$. By $|Q| \geq 3\tau - 1$ and since $v$ is an ancestor of $v_1$, we have $\mathrm{lcp}(Q, Q_1) \geq 3\tau - 1$. Consequently, by Lemma 6.10, the node $v_1$ is periodic and it holds $\text{L-root}(v_1) = \text{L-root}(Q_1) = \text{L-root}(Q) = \text{L-root}(v) = H$ and $\text{L-head}(v_1) = \text{L-head}(Q_1) = \text{L-head}(Q) = \text{L-head}(v) = s$. Furthermore, by $e(Q) \leq |Q|$ and $\mathrm{type}(Q) = -1$, it holds $Q[e(Q)] \prec Q[e(Q) - p]$. Since $Q$ is a prefix of $Q_1$, this immediately implies $e(Q_1) = e(Q) \leq |Q| \leq |Q_1|$ and $Q_1[e(Q_1)] = Q[e(Q)] \prec Q[e(Q) - p] = Q_1[e(Q_1) - p]$, i.e., $\mathrm{type}(Q_1) = -1$. We have thus shown $e(v_1) \leq |\mathrm{str}(v_1)|$ and $\mathrm{type}(v_1) = -1$. Analogously, we obtain that $v_2$ is periodic and it holds $e(v_2) = e(v)$, $\text{L-root}(v_2) = H$, $\text{L-head}(v_2) = s$, $e(v_2) \leq |\mathrm{str}(v_2)|$, and $\mathrm{type}(v_2) = -1$. We have thus shown that $u_1 = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v_1)$ and $u_2 = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v_2)$ are well-defined (see Section 7.3.2).

Let $u = \mathrm{LCA}(u_1, u_2)$, $\ell' = \mathrm{sdepth}(u)$, and $\ell = \mathrm{sdepth}(v)$. By Observation 4.2, we have $\ell = \mathrm{lcp}(\mathrm{str}(v_1), \mathrm{str}(v_2))$, $\ell' = \mathrm{lcp}(\mathrm{str}(u_1), \mathrm{str}(u_2))$, $\mathrm{str}(v) = \mathrm{str}(v_1)[1 \mathinner{.\,.} \ell]$, and $\mathrm{str}(u) = \mathrm{str}(u_1)[1 \mathinner{.\,.} \ell']$. Denote $\delta = e^{\mathrm{full}}(v)$. As observed above, $e(v_1) = e(v)$, $\text{L-head}(v_1) = \text{L-head}(v)$, and $\text{L-root}(v_1) = \text{L-root}(v)$. Thus, $e^{\mathrm{full}}(v_1) = 1 + \text{L-head}(v_1) + |\text{L-root}(v_1)| \cdot \lfloor \frac{e(v_1) - 1 - \text{L-head}(v_1)}{|\text{L-root}(v_1)|} \rfloor = 1 + \text{L-head}(v) + |\text{L-root}(v)| \cdot \lfloor \frac{e(v) - 1 - \text{L-head}(v)}{|\text{L-root}(v)|} \rfloor = e^{\mathrm{full}}(v) = \delta$. Analogously, $e^{\mathrm{full}}(v_2) = \delta$. By definition of $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v_1)$ and $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v_2)$, we have $\mathrm{str}(u_1) = \mathrm{pow}(H) \cdot \mathrm{str}(v_1)[e^{\mathrm{full}}(v_1) \mathinner{.\,.} |\mathrm{str}(v_1)|] = \mathrm{pow}(H) \cdot \mathrm{str}(v_1)[\delta \mathinner{.\,.} |\mathrm{str}(v_1)|]$ and $\mathrm{str}(u_2) = \mathrm{pow}(H) \cdot \mathrm{str}(v_2)[e^{\mathrm{full}}(v_2) \mathinner{.\,.} |\mathrm{str}(v_2)|] = \mathrm{pow}(H) \cdot \mathrm{str}(v_2)[\delta \mathinner{.\,.} |\mathrm{str}(v_2)|]$. Thus, $\ell' = \mathrm{lcp}(\mathrm{str}(u_1), \mathrm{str}(u_2)) = |\mathrm{pow}(H)| + (\mathrm{lcp}(\mathrm{str}(v_1), \mathrm{str}(v_2)) - \delta + 1) = |\mathrm{pow}(H)| + \ell - \delta + 1$. Consequently, $\mathrm{str}(u) = \mathrm{str}(u_1)[1 \mathinner{.\,.} \ell'] = \mathrm{pow}(H) \cdot \mathrm{str}(v_1)[\delta \mathinner{.\,.} \delta + \ell' - |\mathrm{pow}(H)| - 1] = \mathrm{pow}(H) \cdot \mathrm{str}(v_1)[\delta \mathinner{.\,.} \ell] = \mathrm{pow}(H) \cdot \mathrm{str}(v)[\delta \mathinner{.\,.} \ell] = \mathrm{pow}(H) \cdot \mathrm{str}(v)[e^{\mathrm{full}}(v) \mathinner{.\,.} |\mathrm{str}(v)|]$. Thus, by definition of $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(u)$, and since no two nodes of $\mathcal{T}_{\mathsf{Z}}$ have the same value of $\mathrm{str}$, we therefore obtain $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v) = u$. $\qquad\square$

**Lemma 7.44.** *Let $v_1$ and $v_2$ be explicit nodes of $\mathcal{T}_{\mathrm{st}}$ such that $\mathrm{LCA}(v_1, v_2)$ is periodic. Denote $v = \mathrm{LCA}(v_1, v_2)$, $i_{\min} = \min(\mathrm{lrank}(v_1), \mathrm{lrank}(v_2)) + 1$, and $i_{\max} = \max(\mathrm{rrank}(v_1), \mathrm{rrank}(v_2))$. Then, it holds:*

1. *$\mathrm{SA}[i_{\min}], \mathrm{SA}[i_{\max}] \in \mathsf{R}$ and $e(v) = 1 + \min(e(\mathrm{SA}[i_{\min}]) - \mathrm{SA}[i_{\min}], e(\mathrm{SA}[i_{\max}]) - \mathrm{SA}[i_{\max}])$.*
2. *$e(v) \leq |\mathrm{str}(v)|$ holds if and only if $T[\mathrm{SA}[i_{\min}] + e(v) - 1] = T[\mathrm{SA}[i_{\max}] + e(v) - 1]$.*

*Proof.* Denote $i_1 = \mathrm{lrank}(v_1) + 1$, $i_2 = \mathrm{rrank}(v_1)$, $i_3 = \mathrm{lrank}(v_2) + 1$, and $i_4 = \mathrm{rrank}(v_2)$.

1. Let $H = \text{L-root}(v)$, and $p = |H|$. We start by noting that $\mathrm{SA}[i_{\min}], \mathrm{SA}[i_{\max}] \in \mathsf{R}$ follows by Lemma 7.38(1). Next, we prove the formula for $e(v)$.

First, we show that $e(v) = \min(e(v_1), e(v_2))$. Observe that if $P$ is a prefix of $S$, both $P$ and $S$ and periodic, and $\text{L-root}(S) = H$, then, $e(S) = 1 + p + \mathrm{lcp}(S(0 \mathinner{.\,.} |S| - p], S(p \mathinner{.\,.} |S|]) \geq 1 + p + \mathrm{lcp}(S(0 \mathinner{.\,.} |P| - p], S(p \mathinner{.\,.} |P|]) = 1 + p + \mathrm{lcp}(P(0 \mathinner{.\,.} |P| - p], P(p \mathinner{.\,.} |P|]) = e(P)$. Since $\mathrm{str}(v)$ is a prefix of $\mathrm{str}(v_1)$ and $\mathrm{str}(v_2)$, we thus obtain $e(v_1) \geq e(v)$ and $e(v_2) \geq e(v)$. It remains to show that there exists $t \in \{1, 2\}$ such that $e(v_t) = e(v)$. Consider two cases:

- If $e(v) = |\mathrm{str}(v)| + 1$, then there are two possibilities. Either for some $t \in \{1, 2\}$, we have $|\mathrm{str}(v_t)| = |\mathrm{str}(v)|$, in which case $\mathrm{str}(v_t) = \mathrm{str}(v)$ and thus $e(v_t) = e(v)$ follows. The other

possibility is that $|\mathrm{str}(v_1)| > |\mathrm{str}(v)|$ and $|\mathrm{str}(v_2)| > |\mathrm{str}(v)|$. Since $\mathrm{str}(v)$ is the longest common prefix of $\mathrm{str}(v_1)$ and $\mathrm{str}(v_2)$, we then have $\mathrm{str}(v_1)[e(v)] = \mathrm{str}(v_1)[|\mathrm{str}(v)| + 1] \neq \mathrm{str}(v_2)[|\mathrm{str}(v) + 1] = \mathrm{str}(v_2)[e(v)]$. Thus, there exists $t \in \{1, 2\}$ such that $\mathrm{str}(v_t)[e(v)] \neq \mathrm{str}(v)[e(v) - p]$. By definition, for such $t$ we have $e(v_t) = e(v)$.

- Let us now assume $e(v) \leq |\mathrm{str}(v)|$. This implies that $\mathrm{str}(v)[e(v)] = \mathrm{str}(v_1)[e(v)] = \mathrm{str}(v_2)[e(v)]$ and $\mathrm{str}(v)[e(v)] \neq \mathrm{str}(v)[e(v) - p]$. Thus, by $\mathrm{str}(v)[1 \mathinner{..} e(v)] = \mathrm{str}(v_1)[1 \mathinner{..} e(v)] = \mathrm{str}(v_2)[1 \mathinner{..} e(v)]$ we obtain $e(v_1) = e(v_2) = e(v)$.

We have thus shown that there exists $t \in \{1, 2\}$ such that $e(v) = e(v_t)$. Combined with $e(v_1) \geq e(v)$ and $e(v_2) \geq e(v)$, this yields $\min(e(v_1), e(v_2)) = \min(e(v_t), e(v_{3-t})) = e(v)$.

By the above and Lemma 7.38(1) for $v_1$ and $v_2$, it holds $e(v) = 1 + \min_{t \in [1 \mathinner{..} 4]}\{e(\mathrm{SA}[i_t]) - \mathrm{SA}[i_t]\}$. To show that this is equal to the expression for $e(v)$ from the claim, we first observe that letting $X = \mathrm{str}(v)[1 \mathinner{..} 3\tau - 1]$ and $(b, e) = (\mathrm{RangeBeg}(X, T), \mathrm{RangeEnd}(X, T))$, we have $i_t \in (b \mathinner{..} e]$ for all $t \in [1 \mathinner{..} 4]$. Observe that by Lemma 5.11, the sequence $(e(\mathrm{SA}[i]) - \mathrm{SA}[i])_{i=b+1}^{e}$ is bitonic, i.e., there exists $m \in (b \mathinner{..} e]$ such that $e(\mathrm{SA}[b+1]) - \mathrm{SA}[b+1] \leq e(\mathrm{SA}[b+2]) - \mathrm{SA}[b+2] \leq \cdots \leq e(\mathrm{SA}[m]) - \mathrm{SA}[m]$ and $e(\mathrm{SA}[m]) - \mathrm{SA}[m] \geq e(\mathrm{SA}[m+1]) - \mathrm{SA}[m+1] \geq \cdots \geq e(\mathrm{SA}[e]) - \mathrm{SA}[e]$. This implies that for every triple $k_1, k_2, k_3 \in (b \mathinner{..} e]$, the inequalities $k_1 \leq k_2 \leq k_3$ imply $\min(e(\mathrm{SA}[k_1]) - \mathrm{SA}[k_1], e(\mathrm{SA}[k_3]) - \mathrm{SA}[k_3]) = \min_{t \in [1 \mathinner{..} 3]}\{e(\mathrm{SA}[k_t]) - \mathrm{SA}[k_t]\}$. For a proof, consider two cases:

- If $k_2 < m$, then by the bitonic property, we have $e(\mathrm{SA}[k_2]) - \mathrm{SA}[k_2] \geq e(\mathrm{SA}[k_1]) - \mathrm{SA}[k_1]$. Thus, the expression $e(\mathrm{SA}[k_2]) - \mathrm{SA}[k_2]$ has no effect on the minimum.
- If $k_2 \geq m$, then by the bitonic property, we have $e(\mathrm{SA}[k_2]) - \mathrm{SA}[k_2] \geq e(\mathrm{SA}[k_3]) - \mathrm{SA}[k_3]$. Thus, the expression $e(\mathrm{SA}[k_2]) - \mathrm{SA}[k_2]$ can again be excluded in the minimum.

By the above, letting $i'_{\min} = \min_{t \in [1 \mathinner{..} 4]}\{i_t\}$ and $i'_{\max} = \max_{i \in [1 \mathinner{..} 4]}\{i_t\}$, we thus have $e(v) = 1 + \min_{t \in [1 \mathinner{..} 4]}\{e(\mathrm{SA}[i_t]) - \mathrm{SA}[i_t]\} = 1 + \min(e(\mathrm{SA}[i'_{\min}]) - \mathrm{SA}[i'_{\min}], e(\mathrm{SA}[i'_{\max}]) - \mathrm{SA}[i'_{\max}])$.

It remains to show that $i'_{\min} = i_{\min}$ and $i'_{\max} = i_{\max}$. For this, it suffices to note that by definition, we have $i_1 \leq i_2$ and $i_3 \leq i_4$, thus, $i'_{\min} = \min_{t \in [1 \mathinner{..} 4]}\{i_t\} = \min(i_1, i_3) = i_{\min}$ and analogously, $i'_{\max} = \max_{t \in [1 \mathinner{..} 4]}\{i_t\} = \max(i_2, i_4) = i_{\max}$.

2. As observed in the proof of Lemma 7.38(2), it holds $\mathrm{SA}[i_1]+e(v_1)-1 \leq n$, $\mathrm{SA}[i_2]+e(v_1)-1 \leq n$, $\mathrm{SA}[i_3] + e(v_2) - 1 \leq n$, and $\mathrm{SA}[i_4] + e(v_2) - 1 \leq n$. Thus, by $e(v) = \min(e(v_1), e(v_2))$, for every $t \in [1 \mathinner{..} 4]$, we have $\mathrm{SA}[i_t] + e(v) - 1 \leq n$. In particular, $\mathrm{SA}[i_{\min}] + e(v) - 1 \leq n$ and $\mathrm{SA}[i_{\max}] + e(v) - 1 \leq n$. We now prove the equivalence. Let us first assume $e(v) \leq |\mathrm{str}(v)|$. Then, since $\mathrm{str}(v)$ is a prefix of both $\mathrm{str}(v_1)$ and $\mathrm{str}(v_2)$ and $\mathrm{SA}[i_{\min}], \mathrm{SA}[i_{\max}] \in \mathrm{Occ}(\mathrm{str}(v_1), T) \cup \mathrm{Occ}(\mathrm{str}(v_2), T)$, it follows that $\mathrm{SA}[i_{\min}], \mathrm{SA}[i_{\max}] \in \mathrm{Occ}(\mathrm{str}(v), T)$. Therefore, $T[\mathrm{SA}[i_{\min}] + e(v) - 1] = T[\mathrm{SA}[i_{\max}] + e(v) - 1]$ follows immediately. To show the opposite implication, assume by contraposition that $e(v) = |\mathrm{str}(v)| + 1$. Then, there are two possibilities. Either for some $t \in \{1, 2\}$ we have $\mathrm{str}(v_t) = \mathrm{str}(v)$, in which case $\mathrm{str}(v_t)$ is a prefix of $\mathrm{str}(v_{3-t})$, which in turn implies $i_{\min} = \mathrm{lrank}(v_t) + 1$ and $i_{\max} = \mathrm{rrank}(v_t)$. Then, $e(v) = e(v_t)$ and by applying Lemma 7.38(2) to $v_t$, we obtain $T[\mathrm{SA}[i_{\min}] + e(v) - 1] = T[\mathrm{SA}[\mathrm{lrank}(v_t) + 1] + e(v_t) - 1] \neq T[\mathrm{SA}[\mathrm{rrank}(v_t)] + e(v_t) - 1] = T[\mathrm{SA}[i_{\max}] + e(v) - 1]$. The other possibility is that $|\mathrm{str}(v_1)| > |\mathrm{str}(v)|$ and $|\mathrm{str}(v_2)| > |\mathrm{str}(v)|$. Then, since $\mathrm{str}(v)$ is the longest common prefix of $\mathrm{str}(v_1)$ and $\mathrm{str}(v_2)$, neither of $v_1$ or $v_2$ is an ancestor of the other, and hence either it holds $i_{\min} = i_1 \leq i_2 < i_3 \leq i_4 = i_{\max}$ or $i_{\min} = i_3 \leq i_4 < i_1 \leq i_2 = i_{\max}$. In the first case $\mathrm{SA}[i_{\min}] \in \mathrm{Occ}(\mathrm{str}(v_1), T)$ and $\mathrm{SA}[i_{\max}] \in \mathrm{Occ}(\mathrm{str}(v_2), T)$, and in the second case $\mathrm{SA}[i_{\min}] \in \mathrm{Occ}(\mathrm{str}(v_2), T)$ and $\mathrm{SA}[i_{\max}] \in \mathrm{Occ}(\mathrm{str}(v_1), T)$. Therefore, in both cases we have $T[\mathrm{SA}[i_{\min}]+e(v)-1] = T[\mathrm{SA}[i_{\min}]+|\mathrm{str}(v)|] \neq T[\mathrm{SA}[i_{\max}]+|\mathrm{str}(v)|] = T[\mathrm{SA}[i_{\max}]+e(v)-1]$. □

*Remark* 7.45. Observe that in Lemma 7.44, it does not necessarily hold that $\mathrm{lrank}(\mathrm{LCA}(v_1, v_2)) = i_{\min}$ or $\mathrm{rrank}(\mathrm{LCA}(v_1, v_2)) = i_{\max}$. Thus, the lemma does not immediately follow as a corollary

from Lemma 7.38.

**Proposition 7.46.** *Let $v_1$ and $v_2$ be explicit nodes of $\mathcal{T}_{st}$ such that $\text{LCA}(v_1, v_2)$ is periodic. Given the data structure from Section 7.3.1 and the pairs $\text{repr}(v_1)$ and $\text{repr}(v_2)$, we can in $\mathcal{O}(\log \log n)$ time compute $\text{repr}(\text{LCA}(v_1, v_2))$.*

*Proof.* Denote $v = \text{LCA}(v_1, v_2)$, $\text{repr}(v_1) = (b_1, e_1)$ and $\text{repr}(v_2) = (b_2, e_2)$ (recall that for $i \in \{1, 2\}$, we have $b_i = \text{lrank}(v_i)$ and $e_i = \text{rrank}(v_i)$).

First, in $\mathcal{O}(1)$ time we compute $i_{\min} = \min(b_1, b_2) + 1$ and $i_{\max} = \max(e_1, e_2)$. By Lemma 7.44(1), we have $\text{SA}[i_{\min}], \text{SA}[i_{\max}] \in \mathsf{R}$. Using Proposition 5.26, in $\mathcal{O}(\log \log n)$ time we compute $j_{\min} = \text{SA}[i_{\min}]$ and $j_{\max} = \text{SA}[i_{\max}]$. Next, using Proposition 5.15 in $\mathcal{O}(1)$ time we compute $H := \text{L-root}(j_{\min})$, $s := \text{L-head}(j_{\min})$, $k_{\min} = \text{L-exp}(j_{\min})$, $k_{\max} = \text{L-exp}(j_{\max})$, $t_{\min} = \text{L-tail}(j_{\min})$, and $t_{\max} = \text{L-tail}(j_{\max})$. Observe that since $v$ is periodic, and $j_{\min}, j_{\max} \in \text{Occ}(\text{str}(v_1), T) \cup \text{Occ}(\text{str}(v_2), T) \subseteq \text{Occ}(\text{str}(v), T)$, it follows by Lemmas 6.9 and 5.11, that $\text{L-root}(v) = \text{L-root}(j_{\max}) = H$ and $\text{L-head}(v) = \text{L-head}(j_{\max}) = s$. In $\mathcal{O}(1)$ time we thus compute $e_{\min} := e(j_{\min}) = j_{\min} + s + k_{\min}|H| + t_{\min}$ and $e_{\max} := e(j_{\max}) = j_{\max} + s + k_{\max}|H| + t_{\max}$. Next, in $\mathcal{O}(1)$ time we compute $e_v := e(v) = 1 + \min(e_{\min} - j_{\min}, e_{\max} - j_{\max})$ (see Lemma 7.44(1)). Using Lemma 7.44(2), we then in $\mathcal{O}(1)$ time check if it holds $e(v) \leq |\text{str}(v)|$ by comparing $T[j_{\min} + e_v - 1]$ with $T[j_{\max} + e_v - 1]$. Consider two cases:

- Let $T[j_{\min} + e_v - 1] = T[j_{\max} + e_v - 1]$, i.e., $e(v) \leq |\text{str}(v)|$. Recall now that $j_{\min} \in \text{Occ}(\text{str}(v), T)$. In $\mathcal{O}(1)$ time we thus compute $\text{type}(v)$ by comparing $T[j_{\min} + e_v - 1]$ with $T[j_{\min} + e_v - 1 - |H|]$. Let us assume that $T[j_{\min} + e_v - 1] \prec T[j_{\min} + e_v - 1 - |H|]$, i.e., $\text{type}(v) = -1$ (the case $\text{type}(v) = +1$ is handled symmetrically, using the part of the structure from Section 7.3.1 adapted according to Lemma 5.11). By Lemma 7.43, we now have that $v_1$ and $v_2$ are periodic and it holds $e(v_1) \leq |\text{str}(v_1)|$, $e(v_2) \leq |\text{str}(v_2)|$, and $\text{type}(v_1) = \text{type}(v_2) = -1$. Using Proposition 7.31, in $\mathcal{O}(\log \log n)$ time we compute pointers to $u_1 = \text{map}_{\mathcal{T}_{st}, \mathcal{T}_Z}(v_1)$ and $u_2 = \text{map}_{\mathcal{T}_{st}, \mathcal{T}_Z}(v_2)$. Using the representation of $\mathcal{T}_Z$ stored as part of the structure in Section 7.3.1, and Proposition 4.3, in $\mathcal{O}(1)$ time we compute a pointer to $u = \text{LCA}(u_1, u_2)$. By Lemma 7.43, it holds $\text{map}_{\mathcal{T}_{st}, \mathcal{T}_Z}(v) = u$. Our goal is to exploit this connection to compute $\text{repr}(v)$. In $\mathcal{O}(1)$ time we compute $k := \text{L-exp}(v) = \lfloor \frac{e_v - 1 - s}{|H|} \rfloor$ and $\ell := e^{\text{full}}(v) - 1 = s + k|H|$. Using Proposition 7.37, in $\mathcal{O}(\log \log n)$ time we then compute the pair $(b, e) = \text{pseudoinv}_{\mathcal{T}_Z}(\ell, u)$. As noted above, it holds $\text{map}_{\mathcal{T}_{st}, \mathcal{T}_Z}(v) = u$. Thus, by Lemma 7.36, we have $\text{repr}(v) = (b, e)$.
- Let $T[j_{\min} + e_v - 1] \neq T[j_{\max} + e_v - 1]$, i.e., $e(v) > |\text{str}(v)|$. Letting $P = \text{str}(v)$, we then have $e(P) > |P|$, $\text{L-head}(P) = s$, $\text{L-root}(P) = H$, and $|P| = e_v - 1$. Using Proposition 7.41, we thus compute $(b, e) = (\text{RangeBeg}(P, T), \text{RangeEnd}(P, T))$ in $\mathcal{O}(\log \log n)$ time, and return $\text{repr}(v) = (b, e)$. $\qquad \square$

### 7.3.4 Implementation of $\text{child}(v, c)$

**Lemma 7.47.** *Let $c \in [0 .. \sigma)$ and $v$ be an explicit periodic internal node of $\mathcal{T}_{st}$ satisfying $e(v) \leq |\text{str}(v)|$ and $\text{type}(v) = -1$. Let $u = \text{map}_{\mathcal{T}_{st}, \mathcal{T}_Z}(v)$. If $\text{child}(u, c) = \bot$ then $\text{child}(v, c) = \bot$. Otherwise, letting $u' = \text{child}(u, c)$, it holds*

$$\text{repr}(\text{child}(v, c)) = \begin{cases} (b, e) & \text{if } b \neq e, \\ (0, 0) & \text{otherwise,} \end{cases}$$

*where $(b, e) = \text{pseudoinv}_{\mathcal{T}_Z}(\ell, u')$ and $\ell = e^{\text{full}}(v) - 1$.*

*Proof.* Let $H = \text{L-root}(v)$, $s = \text{L-head}(v)$, $k = \text{L-exp}(v)$, and $\mathsf{P} = \{j \in \mathsf{R}^-_{s,H} : \text{L-exp}(j) = k\}$. We first show that $\mathsf{P} \neq \emptyset$. Consider any $j \in \text{Occ}(\text{str}(v), T)$. By Lemma 6.9, $j \in \mathsf{R}$,

L-root$(j)$ = L-root$(v)$ = $H$, and L-head$(j)$ = L-head$(v)$ = $s$, i.e., $j \in \mathsf{R}_{s,H}$. Furthermore, by $e(v) \leq |\mathrm{str}(v)|$ and type$(v) = -1$ we obtain from Lemma 6.11(2) that L-exp$(j)$ = L-exp$(v)$ = $k$ and type$(j)$ = type$(v)$ = $-1$. Thus, $j \in \mathsf{P}$, and hence $\mathsf{P} \neq \emptyset$. Let $b_{\mathsf{P}}, e_{\mathsf{P}} \in [0\mathinner{.\,.}n]$ be such that $\{\mathrm{SA}[i]\}_{i \in (b_{\mathsf{P}}\mathinner{.\,.}e_{\mathsf{P}}]} = \mathsf{P}$, and $b_H, e_H \in [0\mathinner{.\,.}q]$ be such that $\{r_i^{\mathrm{lex}}\}_{i \in (b_H\mathinner{.\,.}e_H]} = \mathsf{R}'^{-}_H$.

Denote $P = \mathrm{str}(v)c$ and $P' = \mathrm{pow}(H) \cdot P[e^{\mathrm{full}}(P)\mathinner{.\,.}|P|]$. Using the above notation, we now establish the characterization of $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$ with the help of Lemma 7.34. First, we observe that since $\mathrm{str}(v)$ is periodic, it follows by Lemma 6.10 that $P$ is periodic and it holds L-root$(P)$ = L-root$(v)$ = $H$ and L-head$(P)$ = L-head$(v)$ = $s$. Moreover, since $e(v) \leq |\mathrm{str}(v)|$ and type$(v) = -1$, it follows by Lemma 6.11(1), that $e(P) = e(v)$, $e^{\mathrm{full}}(P) - 1 = e^{\mathrm{full}}(v) - 1 = \ell$, L-exp$(P)$ = L-exp$(v)$ = $k$, and type$(P)$ = type$(v)$ = $-1$. In particular, this implies that the assumptions of Lemma 7.34 as satisfied. More precisely, $e(P) = e(v) \leq |\mathrm{str}(v)| \leq |P|$. On the other hand, as shown above, $\{j \in \mathsf{R}^{-}_{s,H} : \text{L-exp}(j) = k\} \neq \emptyset$. Observe also that by $e^{\mathrm{full}}(P) - 1 = \ell$, we have $P' = \mathrm{pow}(H) \cdot P[\ell\mathinner{.\,.}|P|]$. Putting all this together, by Lemma 7.34 we obtain that $(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T)) = (e_{\mathsf{P}} - c_1, e_{\mathsf{P}} - c_2)$, where $c_1 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, b_{\mathrm{pre}})$, $c_2 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_{\mathrm{pre}})$, $b_{\mathrm{pre}} = |\{i \in [1\mathinner{.\,.}q] : T[A_{\mathsf{Z}}[i]\mathinner{.\,.}n] \prec P'\}|$, and $(b_{\mathrm{pre}}\mathinner{.\,.}e_{\mathrm{pre}}] = \{i \in [1\mathinner{.\,.}q] : P' \text{ is a prefix of } T[A_{\mathsf{Z}}[i]\mathinner{.\,.}n]\}$.

We are now ready to show the first claim. Recall, that by definition of $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v)$, we have $\mathrm{str}(u) = \mathrm{pow}(H) \cdot \mathrm{str}(v)(\ell\mathinner{.\,.}|\mathrm{str}(v)|]$. Thus, it holds $\mathrm{str}(u)c = P'$. By definition of $\mathcal{T}_{\mathsf{Z}}$ and child$(u, c)$, we thus obtain that child$(u, c) = \bot$ implies $e_{\mathrm{pre}} - b_{\mathrm{pre}} = 0$. Consequently, by the above characterization, it holds

$$\begin{aligned}
|\mathrm{Occ}(P,T)| &= \mathrm{RangeEnd}(P,T) - \mathrm{RangeBeg}(P,T) \\
&= (e_{\mathsf{P}} - c_2) - (e_{\mathsf{P}} - c_1) \\
&= \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_{\mathrm{pre}}) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, b_{\mathrm{pre}}) \\
&= 0.
\end{aligned}$$

Thus, child$(v, c) = \bot$.

Let us now assume child$(u, c) = u' \neq \bot$. Using the above notation, we first show the characterization of $\mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{Z}}}(\ell, u')$. First, note that $\mathrm{pow}(H)$ is a prefix of $\mathrm{str}(u')$ (since it is a prefix of $\mathrm{str}(u)$). Next, note that $\ell \bmod |H| = (e^{\mathrm{full}}(v) - 1) \bmod |H| = \text{L-head}(v) = s$ and $\lfloor \frac{\ell}{|H|} \rfloor = \lfloor \frac{e^{\mathrm{full}}(v) - 1}{|H|} \rfloor = \text{L-exp}(v) = k$. As shown above, the set $\{j \in \mathsf{R}^{-}_{s,H} : \text{L-exp}(j) = k\}$ is nonempty. This implies that $\mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{Z}}}(\ell, u) = (e_{\mathsf{P}} - c'_1, e_{\mathsf{P}} - c'_2)$, where $c'_1 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(u'))$ and $c'_2 = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(u'))$. It remains to observe that by definition of $\mathcal{T}_{\mathsf{Z}}$ and the facts that child$(u, c) = u'$ and $\mathrm{str}(u)c = P'$, we have $\mathrm{lrank}(u') = b_{\mathrm{pre}}$ and $\mathrm{rrank}(u') = e_{\mathrm{pre}}$. Thus, we have $c'_1 = c_1$ and $c'_2 = c_2$, and consequently

$$\begin{aligned}
(\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T)) &= (e_{\mathsf{P}} - c_1, e_{\mathsf{P}} - c_2) \\
&= (e_{\mathsf{P}} - c'_1, e_{\mathsf{P}} - c'_2) \\
&= \mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{Z}}}(\ell, u') \\
&= (b, e).
\end{aligned}$$

By the above, if $b \neq e$, then $\mathrm{Occ}(P,T) \neq \emptyset$. This implies child$(v, c) \neq \bot$ and $\mathrm{repr}(\mathrm{child}(v, c)) = (\mathrm{RangeBeg}(P,T), \mathrm{RangeEnd}(P,T))$. We thus indeed have $\mathrm{repr}(\mathrm{child}(v, c)) = (b, e)$. Otherwise (i.e., if $b = e$), by the above we have $\mathrm{Occ}(P,T) = \emptyset$. This implies child$(v, c) = \bot$ and hence indeed we also have $\mathrm{repr}(\mathrm{child}(v, c)) = (0, 0)$. $\square$

*Remark* 7.48. Note that, similarly as in Lemma 7.22 (see Remark 7.23), even though in the above result we have $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v) = u$ and child$(u, c)$ contains information used to determine child$(v, c)$, it does not necessarily hold that $\mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(\mathrm{child}(v, c)) = \mathrm{child}(u, c)$.

**Proposition 7.49.** *Let $v$ be an explicit periodic internal node of $\mathcal{T}_{st}$. Given the data structure from Section 7.3.1, $\mathrm{repr}(v)$, and $c \in [0\mathinner{.\,.}\sigma)$, in $\mathcal{O}(\log\log n)$ time we can compute $\mathrm{repr}(\mathrm{child}(v,c))$.*

*Proof.* Denote $i_1 = \mathrm{lrank}(v)+1$ and $i_2 = \mathrm{rrank}(v)$. By Lemma 7.38(1), it holds $\mathrm{SA}[i_1], \mathrm{SA}[i_2] \in \mathsf{R}$. Using Proposition 5.26, in $\mathcal{O}(\log\log n)$ time we compute $j_1 = \mathrm{SA}[i_1]$ and $j_2 = \mathrm{SA}[i_2]$. Next, using Proposition 5.15 in $\mathcal{O}(1)$ time we compute $H = \mathrm{L\text{-}root}(j_1)$, $s = \mathrm{L\text{-}head}(j_1)$, $k_1 = \mathrm{L\text{-}exp}(j_1)$, $k_2 = \mathrm{L\text{-}exp}(j_2)$, $t_1 = \mathrm{L\text{-}tail}(j_1)$, and $t_2 = \mathrm{L\text{-}tail}(j_2)$. Observe that since $v$ is periodic, and $j_1, j_2 \in \mathrm{Occ}(\mathrm{str}(v), T)$, it follows by Lemmas 6.9 and 5.11 that $\mathrm{L\text{-}root}(v) = \mathrm{L\text{-}root}(j_2) = H$ and $\mathrm{L\text{-}head}(v) = \mathrm{L\text{-}head}(j_2) = s$. In $\mathcal{O}(1)$ time we thus compute $e_1 := e(j_1) = j_1 + s + k_1|H| + t_1$ and $e_2 := e(j_2) = j_2 + s + k_2|H| + t_2$. Next, in $\mathcal{O}(1)$ time we compute $e_v := e(v) = 1 + \min(e_1 - j_1, e_2 - j_2)$ (see Lemma 7.38(1)). Using Lemma 7.38(2), we then in $\mathcal{O}(1)$ time check if it holds $e(v) \le |\mathrm{str}(v)|$ by comparing $T[j_1 + e_v - 1]$ with $T[j_2 + e_v - 1]$. Consider two cases:

- Let $T[j_1 + e_v - 1] = T[j_2 + e_v - 1]$, i.e., $e(v) \le |\mathrm{str}(v)|$. In $\mathcal{O}(1)$ time we compute $\mathrm{type}(v)$ by comparing $T[j_1 + e_v - 1]$ with $T[j_1 + e_v - 1 - |H|]$. Let us assume that $T[j_1 + e_v - 1] \prec T[j_1 + e_v - 1 - |H|]$, i.e., $\mathrm{type}(v) = -1$ (the case $\mathrm{type}(v) = +1$ it handled symmetrically, using the part of the structure from Section 7.3.1 adapted according to Lemma 5.11). Using Proposition 7.31, in $\mathcal{O}(\log\log n)$ time we compute a pointer to $u = \mathrm{map}_{\mathcal{T}_{st}, \mathcal{T}_{Z}}(v)$. Using the representation of $\mathcal{T}_{Z}$ stored as part of the structure in Section 7.3.1, and Proposition 4.3, in $\mathcal{O}(\log\log n)$ time we check if $\mathrm{child}(u,c) = \bot$. If so, by Lemma 7.47 we have $\mathrm{child}(v,c) = \bot$, and thus we return $\mathrm{repr}(\mathrm{child}(v,c)) = (0,0)$. Otherwise ($\mathrm{child}(u,c) \ne \bot$), we obtain a pointer to $u' = \mathrm{child}(u,c)$. In $\mathcal{O}(1)$ time we now compute $k := \mathrm{L\text{-}exp}(v) = \lfloor \frac{e_v - 1 - s}{|H|} \rfloor$ and $\ell := e^{\mathrm{full}}(v) - 1 = s + k|H|$. Using Proposition 7.37, in $\mathcal{O}(\log\log n)$ time we then compute the pair $(b,e) = \mathrm{pseudoinv}_{\mathcal{T}_{Z}}(\ell, u')$. If $b = e$ then by Lemma 7.47 it holds $\mathrm{child}(v,c) = \bot$ and hence we return $\mathrm{repr}(\mathrm{child}(v,c)) = (0,0)$. Otherwise, by Lemma 7.47, it holds $\mathrm{repr}(\mathrm{child}(v,c)) = (b,e)$. We thus return $(b,e)$.
- Let $T[j_1 + e_v - 1] \ne T[j_2 + e_v - 1]$, i.e., $e(v) > |\mathrm{str}(v)|$. Denote $P = \mathrm{str}(v)$. We then have $e(P) > |P|$, $\mathrm{L\text{-}head}(P) = s$, $\mathrm{L\text{-}root}(P) = H$, and $|P| = e_v - 1$. Using Proposition 7.41, in $\mathcal{O}(\log\log n)$ time we compute $(b,e) = (\mathrm{RangeBeg}(Pc, T), \mathrm{RangeEnd}(Pc, T))$. If $b = e$, then $\mathrm{Occ}(P, T) = \mathrm{Occ}(\mathrm{str}(v)c, T) = \emptyset$, and hence $\mathrm{child}(v,c) = \bot$. We thus return $\mathrm{repr}(\mathrm{child}(v,c)) = (0,0)$. Otherwise, we return that $\mathrm{repr}(\mathrm{child}(v,c)) = (b,e)$. $\qquad\square$

### 7.3.5 Implementation of $\mathrm{pred}(v,c)$

**Lemma 7.50.** *Let $c \in [0\mathinner{.\,.}\sigma)$ and $v$ be an explicit periodic internal node of $\mathcal{T}_{st}$ satisfying $e(v) \le |\mathrm{str}(v)|$ and $\mathrm{type}(v) = -1$. Let $u = \mathrm{map}_{\mathcal{T}_{st}, \mathcal{T}_{Z}}(v)$. If $\mathrm{pred}(u,c) = \bot$ then $\mathrm{RangeBeg}(\mathrm{str}(v)c, T) = \mathrm{RangeBeg}(\mathrm{str}(v), T)$. Otherwise, letting $u' = \mathrm{pred}(u,c)$, it holds*

$$\mathrm{RangeBeg}(\mathrm{str}(v)c, T) = e,$$

*where $(b,e) = \mathrm{pseudoinv}_{\mathcal{T}_{Z}}(\ell, u')$ and $\ell = e^{\mathrm{full}}(v) - 1$.*

*Proof.* We start by characterizing $\mathrm{RangeBeg}(\mathrm{str}(v), T)$. Let $H = \mathrm{L\text{-}root}(v)$, $s = \mathrm{L\text{-}head}(v)$, $k = \mathrm{L\text{-}exp}(v)$, and $\mathsf{P} = \{j \in \mathsf{R}_{s,H}^- : \mathrm{L\text{-}exp}(j) = k\}$. In the proof of Lemma 7.47, we showed that $\mathsf{P} \ne \emptyset$. Let $b_{\mathsf{P}}, e_{\mathsf{P}} \in [0\mathinner{.\,.}n]$ be such that $\{\mathrm{SA}[i]\}_{i \in (b_{\mathsf{P}}\mathinner{.\,.}e_{\mathsf{P}}]} = \mathsf{P}$, and $b_H, e_H \in [0\mathinner{.\,.}q]$ be such that $\{r_i^{\mathrm{lex}}\}_{i \in (b_H\mathinner{.\,.}e_H]} = \mathsf{R}_H'^-$. We now additionally note that by definition of $\mathrm{map}_{\mathcal{T}_{st}, \mathcal{T}_{Z}}(v)$, we have $\mathrm{str}(u) = \mathrm{pow}(H) \cdot \mathrm{str}(v)(\ell\mathinner{.\,.}|\mathrm{str}(v)|]$. Thus, by definition of $\mathcal{T}_{Z}$, it holds $|\{i \in [1\mathinner{.\,.}q] : T[A_{Z}[i]\mathinner{.\,.}n] \prec \mathrm{pow}(H) \cdot \mathrm{str}(v)(\ell\mathinner{.\,.}|\mathrm{str}(v)|]\}| = \mathrm{lrank}(u)$. By Lemma 7.34 for pattern $\mathrm{str}(v)$, we thus obtain $\mathrm{RangeBeg}(\mathrm{str}(v), T) = e_{\mathsf{P}} - (\mathrm{rcount}_{A_{\mathrm{len}}}(\ell, e_H) - \mathrm{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(u)))$.

Next, we characterize $\mathrm{RangeBeg}(\mathrm{str}(v)c, T)$. Denote $P = \mathrm{str}(v)c$ and $P' = \mathrm{pow}(H) \cdot P[e^{\mathrm{full}}(P)\mathinner{.\,.}|P|]$. In the proof of Lemma 7.47, we observed that $P$ is periodic and it holds

L-root$(P)$ = L-root$(v)$ = $H$, L-head$(P)$ = L-head$(v)$ = $s$, $e(P) = e(v)$, $e^{\text{full}}(P) - 1 = e^{\text{full}}(v) - 1 = \ell$, L-exp$(P)$ = L-exp$(v)$ = $k$, and type$(P)$ = type$(v)$ = $-1$. Moreover, we noted that $e(P) \le |P|$ and $P' = \text{pow}(H) \cdot P(\ell \mathinner{.\,.} |P|]$. Finally, putting all this together, we observed that $(\text{RangeBeg}(P,T), \text{RangeEnd}(P,T)) = (e_{\mathsf{P}} - c_1, e_{\mathsf{P}} - c_2)$, where $c_1 = \mathsf{rcount}_{A_{\text{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\text{len}}}(\ell, b_{\text{pre}})$, $c_2 = \mathsf{rcount}_{A_{\text{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\text{len}}}(\ell, e_{\text{pre}})$, $b_{\text{pre}} = |\{i \in [1 \mathinner{.\,.} q] : T[A_{\mathsf{Z}}[i] \mathinner{.\,.} n] \prec P'\}|$, and $(b_{\text{pre}} \mathinner{.\,.} e_{\text{pre}}] = \{i \in [1 \mathinner{.\,.} q] : P' \text{ is a prefix of } T[A_{\mathsf{Z}}[i] \mathinner{.\,.} n]\}$.

Let us first assume $\text{pred}(u,c) = \bot$. By definition, this implies $|\{i \in [1 \mathinner{.\,.} q] : T[A_{\mathsf{Z}}[i] \mathinner{.\,.} n] \prec \text{str}(u)c\}| = \text{lrank}(u)$. Recall, however, that $\text{str}(u) = \text{pow}(H) \cdot \text{str}(v)(\ell \mathinner{.\,.} |\text{str}(v)|]$. Thus, $\text{str}(u)c = P'$ and consequently $b_{\text{pre}} = \text{lrank}(u)$. Using the above characterization, we thus have $\text{RangeBeg}(\text{str}(v)c, T) = e_{\mathsf{P}} - (\mathsf{rcount}_{A_{\text{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\text{len}}}(\ell, \text{lrank}(u)))$. Since above we also established that $\text{RangeBeg}(\text{str}(v), T) = e_{\mathsf{P}} - (\mathsf{rcount}_{A_{\text{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\text{len}}}(\ell, \text{lrank}(u)))$, we have thus proved that $\text{pred}(u,c) = \bot$ implies $\text{RangeBeg}(\text{str}(v)c, T) = \text{RangeBeg}(\text{str}(v), T)$.

Let us now assume $\text{pred}(u,c) = u' \ne \bot$. By definition of $\text{pred}(u,c)$, this implies $|\{i \in [1 \mathinner{.\,.} q] : T[A_{\mathsf{Z}}[i] \mathinner{.\,.} n] \prec \text{str}(u)c\}| = \text{rrank}(u')$. By recalling again that $\text{str}(u)c = P'$, we thus have $b_{\text{pre}} = \text{rrank}(u')$. By the above characterization, we thus have $\text{RangeBeg}(\text{str}(v)c, T) = e_{\mathsf{P}} - (\mathsf{rcount}_{A_{\text{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\text{len}}}(\ell, \text{rrank}(u')))$. On the other hand, by definition of $(b,e) = \text{pseudoinv}_{\mathcal{T}_{\mathsf{Z}}}(\ell, u')$, we have $e = e_{\mathsf{P}} - (\mathsf{rcount}_{A_{\text{len}}}(\ell, e_H) - \mathsf{rcount}_{A_{\text{len}}}(\ell, \text{rrank}(u')))$. We thus obtain $\text{RangeBeg}(\text{str}(v)c, T) = e$. $\qquad\square$

**Proposition 7.51.** *Let $v$ be an explicit periodic internal node of $\mathcal{T}_{\text{st}}$. Given the data structure from Section 7.3.1, repr$(v)$, and $c \in [0 \mathinner{.\,.} \sigma)$, in $\mathcal{O}(\log \log n)$ time we can compute $\text{RangeBeg}(\text{str}(v)c, T)$.*

*Proof.* Denote $i_1 = \text{lrank}(v) + 1$ and $i_2 = \text{rrank}(v)$. By Lemma 7.38(1), it holds $\text{SA}[i_1], \text{SA}[i_2] \in \mathsf{R}$. Using Proposition 5.26, in $\mathcal{O}(\log \log n)$ time we compute $j_1 = \text{SA}[i_1]$ and $j_2 = \text{SA}[i_2]$. Next, using Proposition 5.15 in $\mathcal{O}(1)$ time we compute $H = \text{L-root}(j_1)$, $s = \text{L-head}(j_1)$, $k_1 = \text{L-exp}(j_1)$, $k_2 = \text{L-exp}(j_2)$, $t_1 = \text{L-tail}(j_1)$, and $t_2 = \text{L-tail}(j_2)$. Observe that since $v$ is periodic, and $j_1, j_2 \in \text{Occ}(\text{str}(v), T)$, it follows by Lemmas 6.9 and 5.11 that $\text{L-root}(v) = \text{L-root}(j_2) = H$ and $\text{L-head}(v) = \text{L-head}(j_2) = s$. In $\mathcal{O}(1)$ time we thus compute $e_1 := e(j_1) = j_1 + s + k_1|H| + t_1$ and $e_2 := e(j_2) = j_2 + s + k_2|H| + t_2$. Next, in $\mathcal{O}(1)$ time we compute $e_v := e(v) = 1 + \min(e_1 - j_1, e_2 - j_2)$ (see Lemma 7.38(1)). Using Lemma 7.38(2), we then in $\mathcal{O}(1)$ time check if it holds $e(v) \le |\text{str}(v)|$ by comparing $T[j_1 + e_v - 1]$ with $T[j_2 + e_v - 1]$. Consider two cases:

- Let $T[j_1 + e_v - 1] = T[j_2 + e_v - 1]$, i.e., $e(v) \le |\text{str}(v)|$. In $\mathcal{O}(1)$ time we compute type$(v)$ by comparing $T[j_1 + e_v - 1]$ with $T[j_1 + e_v - 1 - |H|]$. Let us assume that $T[j_1 + e_v - 1] \prec T[j_1 + e_v - 1 - |H|]$, i.e., type$(v) = -1$ (the case type$(v) = +1$ it handled symmetrically, using the part of the structure from Section 7.3.1 adapted according to Lemma 5.11). Using Proposition 7.31, in $\mathcal{O}(\log \log n)$ time we compute a pointer to $u = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\mathsf{Z}}}(v)$. Using the representation of $\mathcal{T}_{\mathsf{Z}}$ stored as part of the structure in Section 7.3.1, and Proposition 4.3, in $\mathcal{O}(\log \log n)$ time we check if $\text{pred}(u,c) = \bot$. If so, by Lemma 7.50 we have $\text{RangeBeg}(\text{str}(v)c, T) = \text{RangeBeg}(\text{str}(v), T)$, and thus we return $\text{rrank}(v)$ as the answer. Otherwise ($\text{pred}(u,c) \ne \bot$), we obtain a pointer to $u' = \text{pred}(u,c)$. In $\mathcal{O}(1)$ time we now compute $k := \text{L-exp}(v) = \lfloor \frac{e_v - 1 - s}{|H|} \rfloor$ and $\ell := e^{\text{full}}(v) - 1 = s + k|H|$. Using Proposition 7.37, in $\mathcal{O}(\log \log n)$ time we then compute the pair $(b,e) = \text{pseudoinv}_{\mathcal{T}_{\mathsf{Z}}}(\ell, u')$. By Lemma 7.50, we then have $\text{RangeBeg}(\text{str}(v)c, T) = e$. Thus, we return $e$ as the answer.
- Let $T[j_1 + e_v - 1] \ne T[j_2 + e_v - 1]$, i.e., $e(v) > |\text{str}(v)|$. Denote $P = \text{str}(v)$. We then have $e(P) > |P|$, $\text{L-head}(P) = s$, $\text{L-root}(P) = H$, and $|P| = e_v - 1$. Using Proposition 7.41, in $\mathcal{O}(\log \log n)$ time we compute $(b,e) = (\text{RangeBeg}(Pc, T), \text{RangeEnd}(Pc, T))$. We then return $b$ as the answer. $\qquad\square$

### 7.3.6 Implementation of $\text{WA}(v, d)$

**Lemma 7.52.** *Let $v$ be an explicit periodic node of $\mathcal{T}_{\text{st}}$ satisfying $\text{type}(v) = -1$ and $d$ be an integer satisfying $e(v) \leq d \leq |\text{str}(v)|$. Then, letting $u = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\text{Z}}}(v)$, it holds*

$$\text{repr}(\text{WA}(v, d)) = \text{pseudoinv}_{\mathcal{T}_{\text{Z}}}(\ell, \widehat{u}),$$

*where $\ell = e^{\text{full}}(v) - 1$, $H = \text{L-root}(v)$, and $\widehat{u} = \text{WA}(u, d - \ell + |\text{pow}(H)|)$.*

*Proof.* As in the proof of Lemma 7.27, let us denote $f^{(0)}(x) = x$ and $f^{(i)}(x) = f(f^{(i-1)}(x))$ for $i \in \mathbb{Z}_+$. Let

$$\mathcal{V} := \{\text{parent}^{(i)}(v) : i \in \mathbb{Z}_{\geq 0} \text{ and } \text{sdepth}(\text{parent}^{(i)}(v)) \geq e(v)\} \text{ and}$$
$$\mathcal{U} := \{\text{parent}^{(i)}(u) : i \in \mathbb{Z}_{\geq 0} \text{ and } \text{sdepth}(\text{parent}^{(i)}(u)) \geq |\text{pow}(H)| + \text{L-tail}(v) + 1\}$$

By $e(v) \leq |\text{str}(v)|$, $\text{type}(v) = -1$, and Lemma 6.11(1), for every $v' \in \mathcal{V}$ it holds that $\text{str}(v')$ is periodic, and we have $e(v') = e(v) \leq |\text{str}(v')|$, $\text{type}(v') = \text{type}(v) = -1$, and $e^{\text{full}}(v') = e^{\text{full}}(v) = \ell + 1$. Thus, for every $v' \in \mathcal{V}$, the node $u' = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\text{Z}}}(v')$ is well-defined and satisfies $\text{str}(u') = \text{pow}(H) \cdot \text{str}(v')[e^{\text{full}}(v') \mathinner{..} |\text{str}(v')|] = \text{pow}(H) \cdot \text{str}(v')(\ell \mathinner{..} |\text{str}(v')|]$. In particular, $\text{str}(u) = \text{pow}(H) \cdot \text{str}(v)(\ell \mathinner{..} |\text{str}(v)|]$. Since for any $v' \in \mathcal{V}$, $\text{str}(v') = \text{str}(v)[1 \mathinner{..} |\text{str}(v')|]$, we thus obtain that for $u' = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\text{Z}}}(v')$ it holds $\text{str}(u') = \text{pow}(H) \cdot \text{str}(v')(\ell \mathinner{..} |\text{str}(v')|] = \text{pow}(H) \cdot \text{str}(v)(\ell \mathinner{..} |\text{str}(v')|] = \text{str}(u)[1 \mathinner{..} |\text{str}(u')|]$. i.e., $u'$ is an ancestor of $u$. Moreover, $\text{sdepth}(u') = |\text{pow}(H)| + |\text{str}(v')| - e^{\text{full}}(v') + 1 = |\text{pow}(H)| + |\text{str}(v')| - e^{\text{full}}(v) + 1 \geq |\text{pow}(H)| + e(v) - e^{\text{full}}(v) + 1 = |\text{pow}(H)| + \text{L-tail}(v) + 1$. Consequently, $\mathcal{U}' := \{\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\text{Z}}}(v') : v' \in \mathcal{V}\}$ satisfies $\mathcal{U}' \subseteq \mathcal{U}$. Note also, that $v' \neq v''$ implies $\text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\text{Z}}}(v') \neq \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_{\text{Z}}}(v'')$.

For any $u' \in \mathcal{U}$, denote $(s(u'), t(u')) = \text{pseudoinv}_{\mathcal{T}_{\text{Z}}}(\ell, u')$. We prove the following property of $\mathcal{U}'$. Let $w, w' \in \mathcal{U}$ be such that $w = \text{parent}(w')$. We claim, that $(s(w), t(w)) \neq (s(w'), t(w'))$ implies $w \in \mathcal{U}'$. The proof consists of five steps:

1. For any node $y$ of $\mathcal{T}_{\text{Z}}$ such that $\text{pow}(H)$ is a prefix of $\text{str}(y)$, by $S_y$ we denote a string such that $\text{str}(y) = \text{pow}(H) \cdot S_y$. Let $P$ and $P'$ be such that $PP'$ is a prefix of $\text{str}(v)$, and it holds $|P| = \ell$ and $|P'| = \text{L-tail}(v) + 1$ (which is defined by $e(v) \leq |\text{str}(v)|$). Consider any node $y$ of $\mathcal{T}_{\text{Z}}$ such that $\text{pow}(H) \cdot P'$ is a prefix of $\text{str}(y)$ (note that although this includes all nodes in $\mathcal{U}$, it is possible that $y \notin \mathcal{U}$). We prove that for any such $y$, it holds $|\text{Occ}(PS_y, T)| = \text{rcount}_{A_{\text{len}}}(\ell, \text{rrank}(y)) - \text{rcount}_{A_{\text{len}}}(\ell, \text{lrank}(y))$. First, observe that since $e(\text{str}(v)) = |PP'|$, we obtain by $\text{lcp}(PS_y, \text{str}(v)) \geq |PP'|$ and Lemma 6.11(1), that $PS_y$ is periodic, $e(PS_y) = e(\text{str}(v)) = |PP'| \leq |PS_y|$, $e^{\text{full}}(PS_y) = e^{\text{full}}(\text{str}(v)) = \ell + 1$, and $\text{type}(PS_y) = \text{type}(\text{str}(v)) = -1$. By Lemma 6.13, $\text{Occ}(PS_y, T)$ is thus a disjoint union of $\text{Occ}^{\text{a}}(PS_y, T)$ and $\text{Occ}^{\text{s}}(PS_y, T)$ (see the beginning of Section 6.3.4 for definitions). By $e(PS_y) \leq |PS_y|$, Lemma 6.14 and its symmetric version (adapted according to Lemma 6.9) moreover imply that $\text{Occ}^{\text{a}}(PS_y, T) = \emptyset$. Finally, by $\text{type}(PS_y) = -1$ and Lemma 6.11(2), it follows that $\text{Occ}(PS_y, T) \subseteq \mathsf{R}^-$. Thus, $\text{Occ}^{\text{s}}(PS_y, T) = \text{Occ}^{\text{s}-}(PS_y, T)$ and consequently $\text{Occ}(PS_y, T) = \text{Occ}^{\text{s}-}(PS_y, T)$. It thus remains to prove $|\text{Occ}^{\text{s}-}(PS_y, T)| = \text{rcount}_{A_{\text{len}}}(\ell, \text{rrank}(y)) - \text{rcount}_{A_{\text{len}}}(\ell, \text{lrank}(y))$. Recall that the set $\{\text{pow}(H) : H \in \text{Roots}\}$ is prefix-free. Letting $H_j = \text{L-root}(j)$ (where $j \in \mathsf{R}$), it follows by definition of $\mathcal{T}_{\text{Z}}$ that:

$$\{r_i^{\text{lex}}\}_{i \in (\text{lrank}(y) \mathinner{..} \text{rrank}(y)]} = \{j \in \mathsf{R}'^- : \text{pow}(H) \cdot S_y \text{ is a prefix of } T[e^{\text{full}}(j) - |\text{pow}(H_j)| \mathinner{..} n]\}$$
$$= \{j \in \mathsf{R}'^- : \text{pow}(H) \cdot S_y \text{ is a prefix of } \text{pow}(H_j) \cdot T[e^{\text{full}}(j) \mathinner{..} n]\}$$
$$= \{j \in \mathsf{R}'^-_H : S_y \text{ is a prefix of } T[e^{\text{full}}(j) \mathinner{..} n]\}.$$

Finally, note that by $e^{\mathrm{full}}(PS_y) = \ell + 1$, we have $(PS_y)[e^{\mathrm{full}}(PS_y) \mathinner{.\,.} |PS_y|]] = S_y$. Thus, by the above and Lemma 6.16, $|\mathrm{Occ}^{\mathsf{s}-}(PS_y, T)| = |\{i \in (\mathrm{lrank}(y) \mathinner{.\,.} \mathrm{rrank}(y)] : A_{\mathrm{len}}[i] \geq e^{\mathrm{full}}(PS_y) - 1\}| = |\{i \in (\mathrm{lrank}(y) \mathinner{.\,.} \mathrm{rrank}(y)] : A_{\mathrm{len}}[i] \geq \ell\}| = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(y)) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(y))$.

2. We prove that there exists $c' \in [0 \mathinner{.\,.} \sigma)$ such that $|\mathrm{Occ}(PS_w c', T)| > 0$ (where $S_w$ is defined as above). First, note that by $u = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v)$, it holds $\mathrm{str}(u) = \mathrm{pow}(H) \cdot \mathrm{str}(v)[e^{\mathrm{full}}(v) \mathinner{.\,.} |\mathrm{str}(v)|] = \mathrm{pow}(H) \cdot \mathrm{str}(v)(\ell \mathinner{.\,.} |\mathrm{str}(v)|]$. On the other hand, by definition of $S_u$, we have $\mathrm{str}(u) = \mathrm{pow}(H) \cdot S_u$. Thus, $S_u = \mathrm{str}(v)(\ell \mathinner{.\,.} |\mathrm{str}(v)|]$. By $P = \mathrm{str}(v)[1 \mathinner{.\,.} \ell]$, we thus obtain $\mathrm{str}(v) = PS_u$. Consequently, since $v$ is a node of $\mathcal{T}_{\mathrm{st}}$, we have $|\mathrm{Occ}(PS_u, T)| = |\mathrm{Occ}(\mathrm{str}(v), T)| > 0$. Observe now that since $w'$ is an ancestor of $u$, the string $\mathrm{str}(w') = \mathrm{pow}(H) \cdot S_{w'}$ is a prefix of $\mathrm{str}(u) = \mathrm{pow}(H) \cdot S_u$. This implies that $S_{w'}$ is a prefix of $S_u$, and hence $PS_{w'}$ is a prefix of $PS_u$. Consequently, $|\mathrm{Occ}(PS_{w'}, T)| \geq |\mathrm{Occ}(PS_u, T)| > 0$. In particular, since $PS_w$ is a prefix of $PS_{w'}$, letting $c' \in [0 \mathinner{.\,.} \sigma)$ be such that $\mathrm{child}(w, c') = w'$, we have $|\mathrm{Occ}(PS_w c', T)| > 0$.

3. Let $s = \ell \bmod |H|$, $k = \lfloor \frac{\ell}{|H|} \rfloor$. Let also $\mathsf{P} := \{j \in \mathsf{R}^-_{s,H} : \mathrm{L\text{-}exp}(j) = k\}$, $b_{\mathsf{P}}, e_{\mathsf{P}} \in [0 \mathinner{.\,.} n]$ be such that $\{\mathrm{SA}[i]\}_{i \in (b_{\mathsf{P}} \mathinner{.\,.} e_{\mathsf{P}}]} = \mathsf{P}$, and $b_H, e_H \in [0 \mathinner{.\,.} q]$ be such that $\{r_i^{\mathrm{lex}}\}_{i \in (b_H \mathinner{.\,.} e_H]} = \mathsf{R}'^-_H$. Note that $\mathrm{L\text{-}head}(v) = s$, $\mathrm{L\text{-}root}(v) = H$, $\mathrm{L\text{-}exp}(v) = k$, $e(v) \leq |\mathrm{str}(v)|$, and $\mathrm{type}(v) = -1$ imply that $\mathsf{P} \neq \emptyset$ (it suffices to take $j = \mathrm{SA}[i]$ for any $i \in (\mathrm{lrank}(v) \mathinner{.\,.} \mathrm{rrank}(v)]$, and apply Lemma 6.9 and Lemma 6.11(2)). Therefore, $(b_{\mathsf{P}}, e_{\mathsf{P}})$ and $(b_H, e_H)$ are well-defined. Denote $\delta = e_{\mathsf{P}} - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, e_H)$. By definition of $\mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{Z}}}(\ell, w)$ and $\mathrm{pseudoinv}_{\mathcal{T}_{\mathsf{Z}}}(\ell, w')$, we then have $(s(w), t(w)) = (\delta + \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(w)), \delta + \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(w)))$ and $(s(w'), t(w')) = (\delta + \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(w')), \delta + \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(w')))$. Thus, the assumption $(s(w), t(w)) \neq (s(w'), t(w'))$, or equivalently, $s(w) \neq s(w')$ or $t(w) \neq t(w')$, implies

$$\mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(w)) \neq \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(w')) \text{ or}$$
$$\mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(w)) \neq \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(w')).$$

4. By definition, the values $\mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(\widehat{w})) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(\widehat{w}))$ over all children $\widehat{w}$ of $w$ sum up to $\mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(w)) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(w))$. Thus, it follows by Step 3 that there exists a child $w'' \neq w'$ of $w$ such that $\mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(w'')) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(w'')) > 0$. By Step 1, for such $w''$, we thus have $|\mathrm{Occ}(PS_{w''}, T)| = \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{rrank}(w'')) - \mathsf{rcount}_{A_{\mathrm{len}}}(\ell, \mathrm{lrank}(w'')) > 0$. In particular, letting $c'' \in [0 \mathinner{.\,.} \sigma)$ be such that $\mathrm{child}(w, c'') = w''$, it holds $|\mathrm{Occ}(PS_w c'', T)| > 0$. Note that $w'' \neq w'$ implies $c'' \neq c'$.

5. We have thus proved (Steps 2 and 4) that there exist $c', c'' \in [0 \mathinner{.\,.} \sigma)$ such that $c' \neq c''$, $|\mathrm{Occ}(PS_w c', T)| > 0$, and $|\mathrm{Occ}(PS_w c'', T)| > 0$. This implies that there exists a node $v'$ in $\mathcal{T}_{\mathrm{st}}$ such that $\mathrm{str}(v') = PS_w$. As observed in Step 1, $PS_w$ is periodic, and it holds $e(PS_w) \leq |PS_w|$, $\mathrm{type}(PS_w) = -1$, and $e^{\mathrm{full}}(PS_w) = |P| + 1$. Thus, the node $u' = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v')$ is defined and satisfies $\mathrm{str}(u') = \mathrm{pow}(H) \cdot S_w$. This implies $u' = w$, and consequently, $w \in \mathcal{U}'$.

We are now ready to prove the main claim. Let $v' = \mathrm{WA}(v, d)$ and $v'' = \mathrm{parent}(v')$. We then have $\mathrm{sdepth}(v'') < d \leq \mathrm{sdepth}(v')$. Moreover, by $e(v) \leq d$, we have $v' \in \mathcal{V}$. Let $u' = \mathrm{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v')$. As observed earlier, we then have $u' \in \mathcal{U}'$, and hence $d - \ell + |\mathrm{pow}(H)| \leq \mathrm{sdepth}(u')$. This implies that $\widehat{u} = \mathrm{WA}(u', d - \ell + |\mathrm{pow}(H)|)$ satisfies $d - \ell + |\mathrm{pow}(H)| \leq \mathrm{sdepth}(\widehat{u}) \leq \mathrm{sdepth}(u')$. By $d - \ell + |\mathrm{pow}(H)| \geq e(v) - \ell + |\mathrm{pow}(H)| = e(v) - (e^{\mathrm{full}}(v) - 1) + |\mathrm{pow}(H)| = \mathrm{L\text{-}tail}(v) + 1 + |\mathrm{pow}(H)|$, this implies that $\widehat{u} \in \mathcal{U}$. Let $k \in \mathbb{Z}_{\geq 0}$ be such that $\widehat{u} = \mathrm{parent}^{(k)}(u')$. This implies that $\mathrm{parent}^{(i)}(u') \notin \mathcal{U}'$ holds for $i \in [1 \mathinner{.\,.} k]$, since otherwise it would contradict $v' = \mathrm{WA}(v, d)$. If $k = 0$ then we trivially have $(s(u'), t(u')) = (s(\widehat{u}), t(\widehat{u}))$. Otherwise, by (the contraposition of)

81

the above property of $\mathcal{U}'$ we have

$$
\begin{aligned}
(s(u'), t(u')) &= (s(\text{parent}(u')), t(\text{parent}(u'))) \\
&= \ldots \\
&= (s(\text{parent}^{(k)}(u')), t(\text{parent}^{(k)}(u'))) \\
&= (s(\widehat{u}), t(\widehat{u})).
\end{aligned}
$$

Recall now that $e(v') \leq |\text{str}(v')|$, $\text{type}(v') = -1$, and $e^{\text{full}}(v') = e^{\text{full}}(v) = \ell + 1$. Thus, by Lemma 7.36, we have $\text{repr}(v') = \text{pseudoinv}_{\mathcal{T}_Z}(e^{\text{full}}(v') - 1, u') = \text{pseudoinv}_{\mathcal{T}_Z}(\ell, u')$. Consequently, $\text{repr}(\text{WA}(v, d)) = \text{pseudoinv}_{\mathcal{T}_Z}(\ell, u') = (s(u'), t(u')) = (s(\widehat{u}), t(\widehat{u})) = \text{pseudoinv}_{\mathcal{T}_Z}(\ell, \widehat{u})$. $\qquad\square$

**Proposition 7.53.** *Let $v$ be an explicit periodic node of $\mathcal{T}_{\text{st}}$. Given the data structure from Section 7.3.1, $\text{repr}(v)$, and an integer $d$ satisfying $3\tau - 1 \leq d \leq |\text{str}(v)|$, in $\mathcal{O}(\log \log n)$ time we can compute $\text{repr}(\text{WA}(v, d))$.*

*Proof.* Denote $i_1 = \text{lrank}(v) + 1$ and $i_2 = \text{rrank}(v)$. By Lemma 7.38(1), it holds $\text{SA}[i_1], \text{SA}[i_2] \in \mathsf{R}$. Using Proposition 5.26, in $\mathcal{O}(\log \log n)$ time we first compute $j_1 = \text{SA}[i_1]$ and $j_2 = \text{SA}[i_2]$. Next, using Proposition 5.15, in $\mathcal{O}(1)$ time we compute $H = \text{L-root}(j_1)$, $s = \text{L-head}(j_1)$, $k_1 = \text{L-exp}(j_1)$, $k_2 = \text{L-exp}(j_2)$, $t_1 = \text{L-tail}(j_1)$, and $t_2 = \text{L-tail}(j_2)$. Since $v$ is periodic, and $j_1, j_2 \in \text{Occ}(\text{str}(v), T)$, it follows by Lemmas 6.9 and 5.11 that $\text{L-root}(v) = \text{L-root}(j_2) = H$ and $\text{L-head}(v) = \text{L-head}(j_2) = s$. In $\mathcal{O}(1)$ time we thus compute $e_1 := e(j_1) = j_1 + s + k_1|H| + t_1$ and $e_2 := e(j_2) = j_2 + s + k_2|H| + t_2$. Next, in $\mathcal{O}(1)$ time we compute $e_v := e(v) = 1 + \min(e_1 - j_1, e_2 - j_2)$ (see Lemma 7.38(1)). We then consider two cases:

- Assume $e_v \leq d$. Then, to obtain $\text{repr}(\text{WA}(v, d))$ we follow Lemma 7.52. First, in $\mathcal{O}(1)$ time we compute $\text{type}(v)$ by comparing $T[j_1 + e_v - 1]$ with $T[j_1 + e_v - 1 - |H|]$. Let us assume that $T[j_1 + e_v - 1] \prec T[j_1 + e_v - 1 - |H|]$, i.e., $\text{type}(v) = -1$ (the case $\text{type}(v) = +1$ it handled symmetrically, using the part of the structure from Section 7.3.1 adapted according to Lemma 5.11). Using Proposition 7.31, in $\mathcal{O}(\log \log n)$ time we compute a pointer to $u = \text{map}_{\mathcal{T}_{\text{st}}, \mathcal{T}_Z}(v)$. In $\mathcal{O}(1)$ time we also calculate $|\text{pow}(H)| = |H|\lceil \frac{\tau}{|H|} \rceil$. Using the representation of $\mathcal{T}_Z$ stored as part of the structure in Section 7.3.1, and Proposition 4.3, in $\mathcal{O}(\log \log n)$ time we compute a pointer to $\widehat{u} = \text{WA}(u, d - \ell + |\text{pow}(H)|)$. In $\mathcal{O}(1)$ time we now compute $k := \text{L-exp}(v) = \lfloor \frac{e_v - 1 - s}{|H|} \rfloor$ and $\ell := e^{\text{full}}(v) - 1 = s + k|H|$. Using Proposition 7.37, in $\mathcal{O}(\log \log n)$ time we then compute the pair $(b, e) = \text{pseudoinv}_{\mathcal{T}_Z}(\ell, \widehat{u})$. By Lemma 7.52, it holds $\text{repr}(\text{WA}(v, d)) = (b, e)$. We thus return $(b, e)$.

- Assume $e_v > d$. Let $v' = \text{WA}(v, d)$, $S = \text{str}(v')$, and $S' = S[1 .. d]$. Since, by definition, $v'$ does not have an ancestor $v''$ in $\mathcal{T}_{\text{st}}$ satisfying $\text{sdepth}(v'') \geq d$, it holds $\text{repr}(v') = (\text{RangeBeg}(S, T), \text{RangeEnd}(S, T)) = (\text{RangeBeg}(S', T), \text{RangeEnd}(S', T))$. We thus focus on computing the latter pair. First, we observe that since $v'$ is an ancestor $v$, we have $S' = \text{str}(v)[1 .. d]$. Therefore, since $\text{str}(v)$ is periodic, and it holds $3\tau - 1 \leq d$, we obtain by Lemma 6.10 that $S'$ is periodic, and it holds $\text{L-root}(S') = \text{L-root}(v) = H$ and $\text{L-head}(S') = \text{L-head}(v) = s$. To show $e(S') > |S'|$, let us denote $Q = \text{str}(v)[1 .. e(v)]$. By definition, we have $e(Q) = 1 + p + \text{lcp}(Q, Q(p .. |Q|]) = |Q| + 1$. Thus, we must have $\text{lcp}(Q, Q(p .. |Q|]) = |Q| - p$. Consequently, since by $e_v = e(v) > d$ the string $S'$ is a prefix of $Q$, we have $\text{lcp}(S', S'(p .. |S'|]) = |S'| - p$, and hence $e(S') = 1 + p + \text{lcp}(S', S'(p .. |S'|]) = |S'| + 1$. Considering all the above properties of $S'$, the next step of the algorithm is therefore to compute and return the pair $(b, e) = (\text{RangeBeg}(S', T), \text{RangeEnd}(S', T))$ in $\mathcal{O}(\log \log n)$ time using Proposition 7.39. As observed above, it holds $\text{repr}(\text{WA}(v, d)) = (b, e)$. $\qquad\square$

### 7.3.7 Construction Algorithm

**Proposition 7.54.** *Given* $\mathrm{C}_{\mathrm{ST}}(T)$, *we can in* $\mathcal{O}(n/\log_\sigma n)$ *time we can augment it into a data structure from Section 7.3.1.*

*Proof.* First, we combine Propositions 5.4 and 5.27 (recall that the packed representation of $T$ is a component of $\mathrm{C}_{\mathrm{ST}}(T)$) to construct the data structure from Section 5.3.2 in $\mathcal{O}(n/\log_\sigma n)$ time. In particular, this constructs $(r_i^{\mathrm{lex}})_{i\in[1..q]}$. Using Proposition 5.15, we can now compute $A_{\mathsf{Z}}[i]$ for any $i \in [1..q]$ in $\mathcal{O}(1)$ time. Then, in $\mathcal{O}(n/\log_\sigma n)$ time we construct the data structure $\mathcal{T}_{\mathsf{Z}}$ using Proposition 4.3.

After the above components are constructed, we then analogously construct their symmetric counterparts (adapted according to Lemma 5.11). $\qquad\square$

## 7.4 The Final Data Structure

In this section, we put together Sections 7.1 to 7.3 to obtain a data structure that performs suffix tree operations in $\mathcal{O}(\log^\epsilon n)$ time.

The section is organized as follows. First, we introduce the components of the data structure (Section 7.4.1). We then describe the query algorithms for all operations in Table 1 (Sections 7.4.2 to 7.4.20). Finally, we show the construction algorithm (Section 7.4.21).

### 7.4.1 The Data Structure

**Definitions**  Recall (Section 2), that we assumed $T[n] = 0$, and that 0 that not appear anywhere else in $T$. We define $T^{\mathrm{rev}}$ as a text obtained by first reversing $T$, and then moving the symbol 0 from the beginning to the end. Formally, for every $i \in [1..n]$:

$$T^{\mathrm{rev}}[i] = \begin{cases} T[n-i] & \text{if } i \neq n, \\ T[n] & \text{if } i = n. \end{cases}$$

Observe that for every $P$ not containing the symbol 0, $j \in \mathrm{Occ}(P,T)$ holds if and only if $j' \in \mathrm{Occ}(\overline{P}, T^{\mathrm{rev}})$, where $j' = n - (j + |P| - 1)$.

*Remark* 7.55. The motivation for defining $T^{\mathrm{rev}}$ is that the standard reverse operation on $T$ (denoted $\overline{T}$) does not preserve a unique sentinel at the end.

**Components**  The data structure consists of two parts. The first part is constructed for $T$ and consists of the following two components:

1. The structure from Section 7.2.1 (used to handle nonperiodic nodes).
2. The structure from Section 7.3.1 (used to handle periodic nodes). Note that similarly as the first component it also includes $\mathrm{C}_{\mathrm{ST}}(T)$. It suffices, however, to only store one copy.

The second part contains the analogous two components for the text $T^{\mathrm{rev}}$. In this section, unless specified otherwise, we refer to the part of the structure for text $T$.

In total, the data structure takes $\mathcal{O}(n/\log_\sigma n)$ space.

### 7.4.2 Implementation of $\mathrm{sdepth}(v)$

**Proposition 7.56.** *Let $v$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$. Given the data structure from Section 7.4.1 and* $\mathrm{repr}(v)$, *we can in* $\mathcal{O}(\log^\epsilon n)$ *time compute* $\mathrm{sdepth}(v)$.

*Proof.* Denote $i_1 = \mathrm{lrank}(v) + 1$ and $i_2 = \mathrm{rrank}(v)$. Let $v_1$ and $v_2$ be the $i_1$th and $i_2$th (respectively) leftmost leaf of $\mathcal{T}_{\mathrm{st}}$. Then, $v = \mathrm{LCA}(v_1, v_2)$. By Observation 4.2, we thus have

$\text{sdepth}(v) = \text{lcp}(\text{str}(v_1), \text{str}(v_2)) = \text{LCE}(\text{SA}[i_1], \text{SA}[i_2])$. Consequently, to compute $\text{sdepth}(v)$ we proceed as follows. First, in $\mathcal{O}(\log^\epsilon n)$ time we compute $j_1 = \text{SA}[i_1]$ and $j_2 = \text{SA}[i_2]$ using Proposition 5.29. Then, using the structure to answer LCE queries (stored as part of the structure in Section 5.3.2), in $\mathcal{O}(1)$ time we compute and return $\text{sdepth}(v) = \text{LCE}(j_1, j_2)$. $\qquad\square$

### 7.4.3  Implementation of $\text{LCA}(u, v)$

**Proposition 7.57.** *Let $v_1$ and $v_2$ be explicit nodes of $\mathcal{T}_{\text{st}}$. Given the data structure from Section 7.4.1 and the pairs $\text{repr}(v_1)$ and $\text{repr}(v_2)$, we can in $\mathcal{O}(\log^\epsilon n)$ time compute the pair $\text{repr}(\text{LCA}(v_1, v_2))$.*

*Proof.* First, using Proposition 7.8, in $\mathcal{O}(1)$ time we check if $\text{sdepth}(\text{LCA}(v_1, v_2)) < 3\tau - 1$. If so, in $\mathcal{O}(1)$ time we additionally obtain $\text{repr}(\text{LCA}(v_1, v_2))$. Let us thus assume $\text{sdepth}(\text{LCA}(v_1, v_2)) \geq 3\tau - 1$. Then, Proposition 7.8 additionally indicates whether $\text{LCA}(v_1, v_2)$ is periodic. If not, we use Proposition 7.21 to compute $\text{repr}(\text{LCA}(v_1, v_2))$ in $\mathcal{O}(\log^\epsilon n)$ time. Otherwise, we obtain $\text{repr}(\text{LCA}(v_1, v_2))$ in $\mathcal{O}(\log \log n)$ time using Proposition 7.46. $\qquad\square$

### 7.4.4  Implementation of $\text{child}(v, c)$

**Proposition 7.58.** *Let $v$ be an explicit internal node of $\mathcal{T}_{\text{st}}$. Given the data structure from Section 7.4.1, $\text{repr}(v)$, and $c \in [0 \mathinner{.\,.} \sigma)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\text{repr}(\text{child}(v, c))$.*

*Proof.* First, using Proposition 7.6, in $\mathcal{O}(1)$ time we check if $v$ is periodic. If so, we obtain $\text{repr}(\text{child}(v, c))$ in $\mathcal{O}(\log \log n)$ time using Proposition 7.49. Otherwise (i.e., if $v$ is not periodic), Proposition 7.6 additionally return the information on whether it holds $\text{sdepth}(v) < 3\tau - 1$. If so, then we obtain $\text{repr}(\text{child}(v, c))$ in $\mathcal{O}(1)$ time using Proposition 7.10. Otherwise, we obtain $\text{repr}(\text{child}(v, c))$ in $\mathcal{O}(\log^\epsilon n)$ time using Proposition 7.24. $\qquad\square$

### 7.4.5  Implementation of $\text{pred}(v, c)$

**Proposition 7.59.** *Let $v$ be an explicit internal node of $\mathcal{T}_{\text{st}}$. Given the data structure from Section 7.4.1, $\text{repr}(v)$, and $c \in [0 \mathinner{.\,.} \sigma)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\text{RangeBeg}(\text{str}(v)c, T)$.*

*Proof.* First, using Proposition 7.6, in $\mathcal{O}(1)$ time we check if $v$ is periodic. If so, we obtain $\text{RangeBeg}(\text{str}(v)c, T)$ in $\mathcal{O}(\log \log n)$ time using Proposition 7.51. Otherwise (i.e., if $v$ is not periodic), Proposition 7.6 additionally return the information on whether it holds $\text{sdepth}(v) < 3\tau - 1$. If so, then we obtain $\text{RangeBeg}(\text{str}(v)c, T)$ in $\mathcal{O}(1)$ time using Proposition 7.11. Otherwise, we obtain $\text{RangeBeg}(\text{str}(v)c, T)$ in $\mathcal{O}(\log^\epsilon n)$ time using Proposition 7.26. $\qquad\square$

**Proposition 7.60.** *Let $v$ be an explicit internal node of $\mathcal{T}_{\text{st}}$. Given the data structure from Section 7.4.1, $\text{repr}(v)$, and $c \in [0 \mathinner{.\,.} \sigma)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\text{repr}(\text{pred}(v, c))$.*

*Proof.* Denote $(b, e) = \text{repr}(v)$. First, using Proposition 7.59, in $\mathcal{O}(\log^\epsilon n)$ time we compute $i = \text{RangeBeg}(\text{str}(v)c, T)$. Observe that by definition of $\text{pred}(v, c)$ we then have $\text{pred}(v, c) = \perp$ if and only if $i = b$. If $i = b$, we thus return $\text{repr}(\text{pred}(v, c)) = (0, 0)$. Let us thus assume $i \neq b$. Observe that we then have $\text{SA}[i] \in \text{Occ}(\text{str}(\text{pred}(v, c)), T)$, and moreover, $\text{pred}(v, c) = \text{child}(v, c')$, where $c' = T[\text{SA}[i] + \text{sdepth}(v)]$. We thus proceed as follows. First, using Proposition 5.29, in $\mathcal{O}(\log^\epsilon n)$ time we compute $j = \text{SA}[i]$. Next, using Proposition 7.56, in $\mathcal{O}(\log^\epsilon n)$ time we compute $\ell = \text{sdepth}(v)$. In $\mathcal{O}(1)$ time we then obtain $c' = T[j + \ell]$. Finally, using Proposition 7.58, in $\mathcal{O}(\log^\epsilon n)$ time we compute and return $\text{repr}(\text{child}(v, c')) = \text{repr}(\text{pred}(v, c))$. $\qquad\square$

### 7.4.6 Implementation of $\mathrm{WA}(v, d)$

**Proposition 7.61.** *Let $v$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$. Given the data structure from Section 7.4.1, repr$(v)$, and an integer $d$ satisfying $0 \leq d \leq |\mathrm{str}(v)|$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute repr$(\mathrm{WA}(v, d))$.*

*Proof.* If $d < 3\tau - 1$, we obtain repr$(\mathrm{WA}(v, d))$ in $\mathcal{O}(1)$ time using Proposition 7.13. Let us thus assume $d \geq 3\tau - 1$. This implies sdepth$(v) \geq 3\tau - 1$. First, using Proposition 7.6, in $\mathcal{O}(1)$ time we determine whether $v$ is periodic. If not, then in $\mathcal{O}(\log^\epsilon n)$ time we compute repr$(\mathrm{WA}(v, d))$ using Proposition 7.28. Otherwise, we obtain repr$(\mathrm{WA}(v, d))$ using Proposition 7.53 in $\mathcal{O}(\log \log n)$ time. $\square$

### 7.4.7 Implementation of $\mathrm{wlink}(v, c)$

**Proposition 7.62.** *Let $P \in [0 \mathinner{.\,.} \sigma)^m$. Given the data structure from Section 7.4.1, the value $|P|$, any $j \in \mathrm{Occ}(P, T)$, and any $c \in [0 \mathinner{.\,.} \sigma)$, in $\mathcal{O}(\log^\epsilon n)$ time we can check if $\mathrm{Occ}(Pc, T) \neq \emptyset$, and if so, return some position $j' \in \mathrm{Occ}(Pc, T)$.*

*Proof.* We start by checking if $P$ contains the symbol 0. For this, we simply check if $j + |P| = n + 1$. If so, we return that $\mathrm{Occ}(Pc, T) = \emptyset$. Let us thus assume $j + |P| \leq n$.

Using Proposition 5.28, in $\mathcal{O}(\log^\epsilon n)$ time we compute $i = \mathrm{ISA}[j]$. Let $(b, e) = (i - 1, i)$, and observe that we then have $(b, e) = \mathrm{repr}(v)$, where $v$ is a leaf of $\mathcal{T}_{\mathrm{st}}$ satisfying $\mathrm{str}(v) = T[j \mathinner{.\,.} n]$. Next, using Proposition 7.61, in $\mathcal{O}(\log^\epsilon n)$ time we compute the pair $(b', e') = \mathrm{repr}(\mathrm{WA}(v, |P|))$ (we can use it, since $|P| \leq n - j + 1 = \mathrm{sdepth}(v)$). We then have:

$$
\begin{aligned}
(b', e') &= (\mathrm{RangeBeg}(\mathrm{str}(v)[1 \mathinner{.\,.} |P|], T), \mathrm{RangeEnd}(\mathrm{str}(v)[1 \mathinner{.\,.} |P|], T)) \\
&= (\mathrm{RangeBeg}(T[j \mathinner{.\,.} j + |P|), T), \mathrm{RangeEnd}(T[j \mathinner{.\,.} j + |P|), T)) \\
&= (\mathrm{RangeBeg}(P, T), \mathrm{RangeEnd}(P, T)).
\end{aligned}
$$

Next, note that it holds $(b', e') = \mathrm{repr}(v')$ for some node $v'$ such that sdepth$(v') \geq |P|$. To check if sdepth$(v') = |P|$, in $\mathcal{O}(\log^\epsilon n)$ time we compute $j_1 = \mathrm{SA}[b' + 1]$ and $j_2 = \mathrm{SA}[e']$ using Proposition 5.29. As explained in the proof of Proposition 7.56, we then have sdepth$(v') = |P|$ if and only if $T[j_1 + |P|] \neq T[j_2 + |P|]$, which we can check in $\mathcal{O}(1)$ time (note that $T[j_1 + |P|]$ and $T[j_2 + |P|]$ are well-defined, since $j_1, j_2 \in \mathrm{Occ}(P, T)$ and we assumed that $P$ does not contain symbol $T[n] = 0$). Consider two cases:

- If $T[j_1 + |P|] = T[j_2 + |P|]$, then $v'$ satisfies sdepth$(v') > |P|$. In that case we check if $c = T[j_1 + |P|]$. If so, we have $j_1 \in \mathrm{Occ}(Pc, T)$ and hence we return $j_1$. Otherwise, we return that $\mathrm{Occ}(Pc, T) = \emptyset$.
- Otherwise (i.e., if $T[j_1 + |P|] \neq T[j_2 + |P|]$), we have sdepth$(v') = |P|$. Using Proposition 7.58, we then compute the pair $(b'', e'') = \mathrm{repr}(\mathrm{child}(v', c))$ in $\mathcal{O}(\log^\epsilon n)$ time. If $(b'', e'') = (0, 0)$, then we return that $\mathrm{Occ}(Pc, T) = \emptyset$. Otherwise, we have $\mathrm{Occ}(Pc, T) \neq \emptyset$. We then use Proposition 5.29 to compute $j' = \mathrm{SA}[e''] \in \mathrm{Occ}(Pc, T)$ in $\mathcal{O}(\log^\epsilon n)$ time. $\square$

**Proposition 7.63.** *Let $v$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$. Given the data structure from Section 7.4.1, repr$(v)$, and $c \in [0 \mathinner{.\,.} \sigma)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute repr$(\mathrm{wlink}'(v, c))$.*

*Proof.* Denote $(b, e) = \mathrm{repr}(v)$ and $P = \mathrm{str}(v)$. The algorithm consists of two steps:

1. The first step is to determine if $\mathrm{Occ}(cP, T) \neq \emptyset$, and if so, to compute some $j' \in \mathrm{Occ}(cP, T)$. First, using Proposition 7.56, in $\mathcal{O}(\log^\epsilon n)$ time we compute $\ell := \mathrm{sdepth}(v) = |P|$. Using Proposition 5.29, in $\mathcal{O}(\log^\epsilon n)$ time we also compute $j = \mathrm{SA}[e]$. We then have $j \in \mathrm{Occ}(P, T)$.

We now check if $j+\ell-1 = n$. If so, then by the uniqueness of $T[n]$, we have $\mathrm{Occ}(P,T) = \{j\}$. In that case, we have $\mathrm{Occ}(cP,T) \neq \emptyset$ if and only if $T[j-1] = c$, which we can check in $\mathcal{O}(1)$ time. If $T[j-1] = c$, in $\mathcal{O}(1)$ time we then obtain $j' \in \mathrm{Occ}(cP,T)$, where $j' = j-1$. Let us now assume that $j + \ell - 1 \neq n$. We now check if $c = 0$. If so, then $\mathrm{Occ}(cP,T) \neq \emptyset$ holds if and only if $\ell = 0$. We can again check this condition in $\mathcal{O}(1)$ time. Moreover, if $\ell = 0$, then we have $j' \in \mathrm{Occ}(cP,T)$, where $j' = n$. Let us thus assume that $c \neq 0$. Observe that then, letting $j^{\mathrm{rev}} := n - (j+\ell-1)$, it holds $j^{\mathrm{rev}} \in \mathrm{Occ}(\overline{P}, T^{\mathrm{rev}})$. Denote $\ell' = \ell + 1$. Using Proposition 7.62 for the text $T^{\mathrm{rev}}$, in $\mathcal{O}(\log^\epsilon n)$ time we check if $\mathrm{Occ}(\overline{P}c, T^{\mathrm{rev}}) = \emptyset$ (note that we have $|\overline{P}c| = \ell'$). If so, we have $\mathrm{Occ}(cP,T) = \emptyset$, since $\overline{cP} = \overline{P}c$ and hence $\mathrm{Occ}(\overline{P}c, T^{\mathrm{rev}}) = \emptyset$ holds if and only if $\mathrm{Occ}(cP,T) = \emptyset$. Otherwise (i.e., if $\mathrm{Occ}(\overline{P}c, T^{\mathrm{rev}}) \neq \emptyset$), Proposition 7.62 returns some position $j_c^{\mathrm{rev}} \in \mathrm{Occ}(\overline{P}c, T^{\mathrm{rev}})$. Letting $j' := n - (j_c^{\mathrm{rev}} + \ell' - 1)$, we then have $j' \in \mathrm{Occ}(cP,T)$.

2. If in the first step we found that $\mathrm{Occ}(cP,T) = \emptyset$, then by definition it holds $\mathrm{wlink}'(v,c) = \bot$, and hence we return $\mathrm{repr}(\mathrm{wlink}'(v,c)) = (0,0)$. Let us thus assume that $\mathrm{Occ}(cP,T) \neq \emptyset$ and $j' \in \mathrm{Occ}(cP,T)$. We now compute the SA range containing all elements of $\mathrm{Occ}(cP,T)$. For this, we first compute $i = \mathrm{ISA}[j']$ using Proposition 5.28 in $\mathcal{O}(\log^\epsilon n)$ time. Letting $(b',e') = (i-1,i)$, we then have $(b',e') = \mathrm{repr}(v')$, where $v'$ is a leaf of $\mathcal{T}_{\mathrm{st}}$ satisfying $\mathrm{str}(v') = T[j' \mathinner{.\,.} n]$. Using Proposition 7.61, in $\mathcal{O}(\log^\epsilon n)$ time, we compute $(b'',e'') = \mathrm{repr}(\mathrm{WA}(v',\ell'))$. We then have

$$\begin{aligned}(b'',e'') &= (\mathrm{RangeBeg}(\mathrm{str}(v')[1 \mathinner{.\,.} \ell'], T), \mathrm{RangeEnd}(\mathrm{str}(v')[1 \mathinner{.\,.} \ell'], T)) \\ &= (\mathrm{RangeBeg}(cP,T), \mathrm{RangeEnd}(cP,T)) \\ &= \mathrm{repr}(\mathrm{wlink}'(v,c)).\end{aligned}$$

In total, the query takes $\mathcal{O}(\log^\epsilon n)$ time. $\qquad\square$

**Proposition 7.64.** *Let $v$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$. Given the data structure from Section 7.4.1, $\mathrm{repr}(v)$, and $c \in [0 \mathinner{.\,.} \sigma)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\mathrm{repr}(\mathrm{wlink}(v,c))$.*

*Proof.* As observed at the beginning of Section 7, $\mathrm{wlink}(v,c) \neq \bot$ holds if and only if $\mathrm{wlink}'(v,c) \neq \bot$ and $\mathrm{sdepth}(\mathrm{wlink}'(v,c)) = \mathrm{sdepth}(v) + 1$. Therefore, we can use $\mathrm{wlink}'(v,c)$ to compute $\mathrm{wlink}(v,c)$. First, using Proposition 7.63, in $\mathcal{O}(\log^\epsilon n)$ time we compute $(b,e) = \mathrm{repr}(\mathrm{wlink}'(v,c))$. If $(b,e) = (0,0)$, then by the above we have $\mathrm{wlink}(v,c) = \bot$, and hence return $\mathrm{repr}(\mathrm{wlink}(v,c)) = (0,0)$. Otherwise, using Proposition 7.56, in $\mathcal{O}(\log^\epsilon n)$ time we compute $\ell = \mathrm{sdepth}(v)$ and $\ell' = \mathrm{sdepth}(\mathrm{wlink}'(v,c))$. If $\ell' = \ell + 1$, then we return that $\mathrm{repr}(\mathrm{wlink}(v,c)) = (b,e)$. Otherwise, we have $\mathrm{wlink}(v,c) = \bot$ and we return $\mathrm{repr}(\mathrm{wlink}(v,c)) = (0,0)$. $\qquad\square$

### 7.4.8 Implementation of $\mathrm{slink}(v)$

**Proposition 7.65.** *Let $v \neq \mathrm{root}(\mathcal{T}_{\mathrm{st}})$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$. Given the data structure from Section 7.4.1 and $\mathrm{repr}(v)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\mathrm{repr}(\mathrm{slink}(v))$.*

*Proof.* Denote $(b,e) = \mathrm{repr}(v)$, $P = \mathrm{str}(v)$, and $P' = P[2 \mathinner{.\,.} |P|]$. Recall, that for every $v \neq \mathrm{root}(\mathcal{T}_{\mathrm{st}})$, $\mathrm{slink}(v)$ is an explicit node of $\mathcal{T}_{\mathrm{st}}$. Thus, to compute $\mathrm{repr}(\mathrm{slink}(v))$, we need to determine $(\mathrm{RangeBeg}(P',T), \mathrm{RangeEnd}(P',T))$.

First, using Proposition 7.56, in $\mathcal{O}(\log^\epsilon n)$ time we compute $\ell := \mathrm{sdepth}(v) = |P|$. Next, using Proposition 5.29, in $\mathcal{O}(\log^\epsilon n)$ time we compute $j = \mathrm{SA}[e]$. We then have $j \in \mathrm{Occ}(P,T)$. Then, $j' := j+1$ satisfies $j' \in \mathrm{Occ}(P',T)$. Using Proposition 5.28, in $\mathcal{O}(\log^\epsilon n)$ time we compute $i = \mathrm{ISA}[j']$. Letting $(b',e') = (i-1,i)$, we then have $(b',e') = \mathrm{repr}(v')$, where $v'$ is a leaf

of $\mathcal{T}_{\mathrm{st}}$ satisfying $\mathrm{str}(v') = T[j' \mathinner{.\,.} n]$. Using Proposition 7.61, in $\mathcal{O}(\log^\epsilon n)$ time, we compute $(b'', e'') = \mathrm{repr}(\mathrm{WA}(v', \ell - 1))$. We then have

$$
\begin{aligned}
(b'', e'') &= (\mathrm{RangeBeg}(\mathrm{str}(v')[1 \mathinner{.\,.} \ell - 1], T), \mathrm{RangeEnd}(\mathrm{str}(v')[1 \mathinner{.\,.} \ell - 1], T)) \\
&= (\mathrm{RangeBeg}(T[j' \mathinner{.\,.} j' + \ell - 1), T), \mathrm{RangeEnd}(T[j' \mathinner{.\,.} j' + \ell - 1), T)) \\
&= (\mathrm{RangeBeg}(P', T), \mathrm{RangeEnd}(P', T)) \\
&= \mathrm{repr}(\mathrm{slink}(v)).
\end{aligned}
$$

In total, the query takes $\mathcal{O}(\log^\epsilon n)$ time. $\qquad\square$

### 7.4.9  Implementation of $\mathrm{slink}(v, i)$

**Proposition 7.66.** *Let $i \in \mathbb{Z}_+$ and let $v$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$ satisfying $\mathrm{sdepth}(v) \geq i$. Given the data structure from Section 7.4.1, $\mathrm{repr}(v)$, and the value $i$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\mathrm{repr}(\mathrm{slink}(v, i))$.*

*Proof.* Denote $(b, e) = \mathrm{repr}(v)$, $P = \mathrm{str}(v)$, and $P' = P[i+1 \mathinner{.\,.} |P|]$. Note that since for every $v \neq \mathrm{root}(\mathcal{T}_{\mathrm{st}})$, $\mathrm{slink}(v)$ is an explicit node of $\mathcal{T}_{\mathrm{st}}$ (Proposition 7.65), it follows that for every explicit node $v$ of $\mathcal{T}_{\mathrm{st}}$ that satisfies $\mathrm{sdepth}(v) \geq i$, $\mathrm{slink}(v, i)$ is an explicit node of $\mathcal{T}_{\mathrm{st}}$. Thus, to compute $\mathrm{repr}(\mathrm{slink}(v))$, we need to determine $(\mathrm{RangeBeg}(P', T), \mathrm{RangeEnd}(P', T))$.

The procedure is a generalization of the one explained in the proof of Proposition 7.66. First, using Proposition 7.56, in $\mathcal{O}(\log^\epsilon n)$ time we compute $\ell := \mathrm{sdepth}(v) = |P|$. Next, using Proposition 5.29, in $\mathcal{O}(\log^\epsilon n)$ time we compute $j = \mathrm{SA}[e]$. We then have $j \in \mathrm{Occ}(P, T)$. Then, $j' := j + i$ satisfies $j' \in \mathrm{Occ}(P', T)$. Using Proposition 5.28, in $\mathcal{O}(\log^\epsilon n)$ time we compute $i' = \mathrm{ISA}[j']$. Letting $(b', e') = (i' - 1, i')$, we then have $(b', e') = \mathrm{repr}(v')$, where $v'$ is a leaf of $\mathcal{T}_{\mathrm{st}}$ satisfying $\mathrm{str}(v') = T[j' \mathinner{.\,.} n]$. Using Proposition 7.61, in $\mathcal{O}(\log^\epsilon n)$ time, we compute $(b'', e'') = \mathrm{repr}(\mathrm{WA}(v', \ell - i))$. We then have $(b'', e'') = (\mathrm{RangeBeg}(\mathrm{str}(v')[1 \mathinner{.\,.} \ell - i], T)$, $\mathrm{RangeEnd}(\mathrm{str}(v')[1 \mathinner{.\,.} \ell - i], T)) = (\mathrm{RangeBeg}(T[j' \mathinner{.\,.} j' + \ell - i), T), \mathrm{RangeEnd}(T[j' \mathinner{.\,.} j' + \ell - i), T)) = (\mathrm{RangeBeg}(P', T), \mathrm{RangeEnd}(P', T)) = \mathrm{repr}(\mathrm{slink}(v, i))$. $\qquad\square$

### 7.4.10  Implementation of $\mathrm{parent}(v)$

**Lemma 7.67.** *Let $v \neq \mathrm{root}(\mathcal{T}_{\mathrm{st}})$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$. Let $\mathrm{repr}(v) = (b, e)$. If $b \neq 0$ (resp. $e \neq n$) then, letting $v_1$ and $v_2$ be the $b$th and $(b + 1)$st (resp. $e$th and $(e + 1)$st) leftmost leaves of $\mathcal{T}_{\mathrm{st}}$, the following conditions are equivalent:*

1. *$\mathrm{leftsibling}(v) \neq \bot$ (resp. $\mathrm{rightsibling}(v) \neq \bot$),*
2. *$\mathrm{parent}(v) = \mathrm{LCA}(v_1, v_2)$.*

*Proof.* Assume $\mathrm{leftsibling}(v) = v_s \neq \bot$ (resp. $\mathrm{rightsibling}(v) = v_s \neq \bot$). By $\mathrm{repr}(v) = (b, e)$, we have $\mathrm{repr}(v_1) = (b - 1, b)$ (resp. $\mathrm{repr}(v_1) = (e - 1, e)$), and $\mathrm{repr}(v_2) = (b, b + 1)$ (resp. $\mathrm{repr}(v_2) = (e, e + 1)$). This implies that $v_1$ is in the subtree rooted in $v_s$ (resp. $v$) and $v_2$ in the subtree rooted in $v$ (resp. $v_s$). Consequently, $\mathrm{LCA}(v_1, v_2) = \mathrm{LCA}(v, v_s)$. On the other hand, since $v_s$ is a sibling of $v$, we have $\mathrm{LCA}(v, v_s) = \mathrm{parent}(v)$. Thus, $\mathrm{parent}(v) = \mathrm{LCA}(v_1, v_2)$.

We show that $\mathrm{parent}(v) = \mathrm{LCA}(v_1, v_2)$ implies $\mathrm{leftsibling}(v) \neq \bot$ (resp. $\mathrm{rightsibling}(v) \neq \bot$) by contraposition. Assume $\mathrm{leftsibling}(v) = \bot$ (resp. $\mathrm{rightsibling}(v) = \bot$) and denote $v_p = \mathrm{parent}(v)$. Observe that then $\mathrm{repr}(v_p) = (b, e_p)$ for some $e_p > b$ (resp. $\mathrm{repr}(v_p) = (b_p, e)$ for some $b_p < e$). By $\mathrm{repr}(v_1) = (b - 1, b)$ (resp. $\mathrm{repr}(v_2) = (e, e + 1)$), the node $v_1$ (resp. $v_2$) is thus not in the subtree rooted in $v_p$. On the other hand, $\mathrm{repr}(v_2) = (b, b + 1)$ (resp. $\mathrm{repr}(v_1) = (e - 1, e)$) implies that $v_p$ is an ancestor of $v_2$ (resp. $v_1$). Therefore, the node $\mathrm{LCA}(v_1, v_2) = \mathrm{LCA}(v_1, v_p)$ (resp. $\mathrm{LCA}(v_1, v_2) = \mathrm{LCA}(v_p, v_2)$) is a proper ancestor of $v_p$. In particular, $\mathrm{LCA}(v_1, v_2) \neq v_p$. $\qquad\square$

**Proposition 7.68.** *Let $v \neq \operatorname{root}(\mathcal{T}_{\mathrm{st}})$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$. Given the data structure from Section 7.4.1 and* $\operatorname{repr}(v)$, *in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\operatorname{repr}(\operatorname{parent}(v))$.*

*Proof.* Let $\operatorname{repr}(v) = (b, e)$. We first construct a set of pairs $\mathcal{P}$ as follows.

- If $b = 0$, we skip this step. Otherwise, let $v_1$ and $v_2$ be the leftmost $b$th and $(b+1)$st leaves of $\mathcal{T}_{\mathrm{st}}$, $(b_1, e_1) = (b-1, b)$, and $(b_2, e_2) = (b, b+1)$. We then have $\operatorname{repr}(v_1) = (b_1, e_1)$ and $\operatorname{repr}(v_2) = (b_2, e_2)$. Using Proposition 7.57, in $\mathcal{O}(\log^\epsilon n)$ we obtain $\operatorname{repr}(v')$, where $v' = \operatorname{LCA}(v_1, v_2)$. Note that $v'$ is an ancestor of $v$. We add $\operatorname{repr}(v')$ to $\mathcal{P}$.
- If $e = n$, we skip this step. Otherwise, let $v_1'$ and $v_2'$ be the leftmost $e$th and $(e+1)$st leaves of $\mathcal{T}_{\mathrm{st}}$, $(b_1', e_1') = (e-1, e)$, and $(b_2', e_2') = (e, e+1)$. We repeat the same procedure as above, again adding $\operatorname{repr}(v'')$ (where $v'' = \operatorname{LCA}(v_1', v_2')$) to $\mathcal{P}$.

Recall now that we assumed $|T| \geq 2$ and that $T[n]$ is unique in $T$. This implies that the root of $\mathcal{T}_{\mathrm{st}}$ has at least two children. On the other hand, any other non-leaf node has at least two children by definition. This implies that for every explicit node $v \neq \operatorname{root}(\mathcal{T}_{\mathrm{st}})$, it holds that either $\operatorname{leftsibling}(v) \neq \bot$ or $\operatorname{rightsibling}(v) \neq \bot$. Therefore, by Lemma 7.67, there exists $(b_p, e_p) \in \mathcal{P}$ such that $(b_p, e_p) = \operatorname{repr}(\operatorname{parent}(p))$. Since each of the nodes $u$ corresponding to an element in $\mathcal{P}$ is an ancestor of $v$, to compute $\operatorname{parent}(v)$, it suffices to compute $\operatorname{sdepth}(u)$ for all candidates $u$ and return the pair $\operatorname{repr}(u)$ corresponding to $u$ with the largest value. We obtain $\operatorname{sdepth}(u)$ using Proposition 7.56 in $\mathcal{O}(\log^\epsilon n)$ time. By $|\mathcal{P}| \leq 2$, the whole procedure takes $\mathcal{O}(\log^\epsilon n)$ time. $\square$

*Remark* 7.69. It might appear that the computation of $\operatorname{parent}(v)$ could be implemented by modifying the definition of the $\operatorname{WA}(v, d)$ to instead return the deepest ancestor $v'$ of $v$ satisfying $\operatorname{sdepth}(v') \leq d$ (rather than the most shallow ancestor $v'$ of $v$ satisfying $\operatorname{sdepth}(v') \geq d$). Observe, however that as shown in the proof of Lemma 7.27 (resp. Lemma 7.52), $\operatorname{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{S}}}(v)$ (resp. $\operatorname{map}_{\mathcal{T}_{\mathrm{st}}, \mathcal{T}_{\mathsf{Z}}}(v)$) always returns the *lowest* of all nodes $u'$ of $\mathcal{T}_{\mathsf{S}}$ (resp. $\mathcal{T}_{\mathsf{Z}}$) satisfying $(s(u'), t(u')) = \operatorname{repr}(v)$. This enforces the current definition of $\operatorname{WA}(v, d)$ and implies that the implementation of $\operatorname{parent}(v)$ with $\operatorname{WA}(v, d)$ would require a binary search. Thus, to achieve faster time, $\operatorname{parent}(v)$ is implemented as above.

### 7.4.11 Implementation of $\operatorname{firstchild}(v)$

**Proposition 7.70.** *Let $v$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$. Given the data structure from Section 7.4.1 and* $\operatorname{repr}(v)$, *in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\operatorname{repr}(\operatorname{firstchild}(v))$.*

*Proof.* Denote $(b, e) = \operatorname{repr}(v)$ and $P = \operatorname{str}(v)$. First, we check if $b + 1 = e$. If so, then $v$ is a leaf and hence we return $\operatorname{repr}(\operatorname{firstchild}(v)) = (0, 0)$ (note that here we used that $|T| \geq 2$ and that $T[n]$ is unique in $T$, since this implies that every non-leaf node of $\mathcal{T}_{\mathrm{st}}$, including the root, has at least two children).

Let us thus assume $b + 1 \neq e$. Denote $v' = \operatorname{firstchild}(v)$ and $P' = \operatorname{str}(v')$. We then have $v' \neq \bot$. Observe that letting $(b', e') = (b, b+1)$, it holds $(b', e') = \operatorname{repr}(v'')$, where $v''$ is a leaf of $\mathcal{T}_{\mathrm{st}}$ such that $\operatorname{str}(v')$ is a prefix of $\operatorname{str}(v'')$. On the other hand, by definition, we have $\operatorname{sdepth}(v') \geq \operatorname{sdepth}(v) + 1$, and there is no ancestor of $v'$ at depth $d \in (\operatorname{sdepth}(v) \mathinner{.\,.} \operatorname{sdepth}(v'))$. Therefore, we must have $v' = \operatorname{WA}(v'', \operatorname{sdepth}(v) + 1)$. We thus proceed as follows. First, using Proposition 7.56, in $\mathcal{O}(\log^\epsilon n)$ time we compute $\ell := \operatorname{sdepth}(v) = |P|$. Next, using Proposition 7.61, in $\mathcal{O}(\log^\epsilon n)$ time we compute $(b'', e'') = \operatorname{repr}(\operatorname{WA}(v'', \ell + 1))$. We then have $\operatorname{repr}(\operatorname{firstchild}(v)) = (b'', e'')$. In total, the query takes $\mathcal{O}(\log^\epsilon n)$ time. $\square$

### 7.4.12 Implementation of lastchild($v$)

**Proposition 7.71.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$. Given the data structure from Section 7.4.1 and $\text{repr}(v)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\text{repr}(\text{lastchild}(v))$.*

*Proof.* Denote $(b, e) = \text{repr}(v)$ and $P = \text{str}(v)$. The algorithm is symmetrical to the one presented in the proof of Proposition 7.70, i.e., rather than setting $(b', e') = (b, b + 1)$, we set $(b', e') = (e - 1, e)$. For such pair, it holds $(b', e') = \text{repr}(v'')$, where $v''$ is a leaf of $\mathcal{T}_{\text{st}}$ such that, letting $v' = \text{lastchild}(v)$, the string $\text{str}(v')$ is a prefix of $\text{str}(v'')$. $\qquad\square$

### 7.4.13 Implementation of rightsibling($v$)

**Proposition 7.72.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$. Given the data structure from Section 7.4.1 and $\text{repr}(v)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\text{repr}(\text{rightsibling}(v))$.*

*Proof.* Denote $(b, e) = \text{repr}(v)$. We start by checking if $(b, e) = (0, n)$. If so, then $v = \text{root}(\mathcal{T}_{\text{st}})$. In that case, we have $\text{rightsibling}(v) = \bot$ and hence we return $\text{repr}(\text{rightsibling}(v)) = (0, 0)$.

Let us thus assume that $(b, e) \neq (0, n)$, i.e., $v \neq \text{root}(\mathcal{T}_{\text{st}})$. Using Proposition 7.68, in $\mathcal{O}(\log^\epsilon n)$ time we compute $(b', e') = \text{repr}(\text{parent}(v))$. We then have $b' \leq b < e \leq e'$. Next, we compare $e$ and $e'$. If $e = e'$ then, by definition, $v$ is the rightmost child of its parent and hence we return $\text{repr}(\text{rightsibling}(v)) = (0, 0)$. Let us thus assume $e < e'$. We then have $\text{rightsibling}(v) \neq \bot$. Moreover, letting $v' = \text{rightsibling}(v)$ and $P' = \text{str}(v')$, it then holds $\text{SA}[e + 1] \in \text{Occ}(P', T)$. This implies that, letting $(b'', e'') = (e, e + 1)$, we have $(b'', e'') = \text{repr}(v'')$, where $v''$ is a leaf of $\mathcal{T}_{\text{st}}$ such that $P'$ is a prefix of $\text{str}(v'')$. Moreover, it holds $\text{sdepth}(v') \geq \text{sdepth}(\text{parent}(v)) + 1$, and the node $v'$ does not have any ancestors at depth $d \in (\text{sdepth}(\text{parent}(v)) .. \text{sdepth}(v'))$. Consequently, $v' = \text{WA}(v'', \text{sdepth}(\text{parent}(v)) + 1)$. We thus proceed as follows. First, using Proposition 7.56, in $\mathcal{O}(\log^\epsilon n)$ time we compute $\ell := \text{sdepth}(\text{parent}(v))$. Then, using Proposition 7.61, in $\mathcal{O}(\log^\epsilon n)$ time we compute $(b''', e''') = \text{repr}(\text{WA}(v'', \ell + 1))$. By the above discussion, we have $(b''', e''') = \text{repr}(\text{rightsibling}(v))$. $\qquad\square$

### 7.4.14 Implementation of leftsibling($v$)

**Proposition 7.73.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$. Given the data structure from Section 7.4.1 and $\text{repr}(v)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\text{repr}(\text{leftsibling}(v))$.*

*Proof.* Denote $(b, e) = \text{repr}(v)$. The algorithm is symmetrical to the one presented in the proof of Proposition 7.72. More precisely, we replace the check $e = e'$ with $b = b'$. We also set $(b'', e'') = (b - 1, b)$ instead of $(b'', e'') = (e, e + 1)$. $\qquad\square$

### 7.4.15 Implementation of isleaf($v$)

**Proposition 7.74.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$. Given $\text{repr}(v)$, we can check if $v$ is a leaf in $\mathcal{O}(1)$ time.*

*Proof.* Recall, that $|T| \geq 2$ and that $T[n]$ is unique in $T$. This implies that the root of $\mathcal{T}_{\text{st}}$ has at least two children. On the other hand, any other non-leaf node has at least two children by definition. Thus, letting $(b, e) = \text{repr}(v)$, and recalling that $\text{repr}(v) = (\text{lrank}(v), \text{rrank}(v))$, the node $v$ is a leaf if and only if $b + 1 = e$, which we can check in $\mathcal{O}(1)$ time. $\qquad\square$

### 7.4.16 Implementation of index($v$)

**Proposition 7.75.** *Let $v$ be an explicit node of $\mathcal{T}_{\text{st}}$. Given the data structure from Section 7.4.1 and $\text{repr}(v)$, in $\mathcal{O}(\log^\epsilon n)$ time we can compute $\text{index}(v)$.*

*Proof.* Denote $(b, e) = \operatorname{repr}(v)$ and $P = \operatorname{str}(v)$. Recall, that it holds $\operatorname{repr}(v) = (\operatorname{RangeBeg}(P, T),$ $\operatorname{RangeEnd}(P, T))$, and hence $\operatorname{Occ}(P, T) = \{\operatorname{SA}[i]\}_{i \in (b..e]}$. Thus, to obtain $\operatorname{index}(v)$, it suffices to compute $j = \operatorname{SA}[i]$ for any $i \in (b..e]$. Using Proposition 5.29, this takes $\mathcal{O}(\log^\epsilon n)$ time. $\qquad\square$

### 7.4.17 Implementation of $\operatorname{count}(v)$

**Proposition 7.76.** *Let $v$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$. Given $\operatorname{repr}(v)$, we can compute $\operatorname{count}(v)$ in $\mathcal{O}(1)$ time.*

*Proof.* Denote $(b, e) = \operatorname{repr}(v)$. Since by definition we have $\operatorname{Occ}(\operatorname{str}(v), T) = \{\operatorname{SA}[i]\}_{i \in (b..e]}$, in $\mathcal{O}(1)$ time we return $\operatorname{count}(v) = |\operatorname{Occ}(\operatorname{str}(v), T)| = e - b$. $\qquad\square$

### 7.4.18 Implementation of $\operatorname{letter}(v, i)$

**Proposition 7.77.** *Let $v$ be an explicit node of $\mathcal{T}_{\mathrm{st}}$ and $i \in [1..|\operatorname{str}(v)|]$. Given the data structure from Section 7.4.1, $\operatorname{repr}(v)$, and the value $i$, we can compute $\operatorname{letter}(v, i)$ in $\mathcal{O}(\log^\epsilon n)$ time.*

*Proof.* It suffices to find any $j \in \operatorname{Occ}(\operatorname{str}(v), T)$ and return $T[j + i - 1]$. Using Proposition 7.75, we find $j$ in $\mathcal{O}(\log^\epsilon n)$ time. We then return the output symbol in $\mathcal{O}(1)$ time (recall, that the packed representation of $T$ is stored as part of $\mathrm{C}_{\mathrm{ST}}(T)$). $\qquad\square$

### 7.4.19 Implementation of $\operatorname{isancestor}(u, v)$

**Proposition 7.78.** *Let $u$ and $v$ be explicit nodes of $\mathcal{T}_{\mathrm{st}}$. Given $\operatorname{repr}(u)$ and $\operatorname{repr}(v)$, we can check if $u$ is an ancestor of $v$ in $\mathcal{O}(1)$ time.*

*Proof.* Denote $(b, e) = \operatorname{repr}(u)$, $(b', e') = \operatorname{repr}(v)$. The node $u$ is an ancestor of $v$ if and only if $\operatorname{Occ}(\operatorname{str}(v), T) \subseteq \operatorname{Occ}(\operatorname{str}(u), T)$, which (by definition) holds if and only if $b \leq b' < e' \leq e$. $\qquad\square$

### 7.4.20 Implementation of $\operatorname{findleaf}(j)$

**Proposition 7.79.** *Let $j \in [1..n]$. Given the data structure from Section 7.4.1 and the position $j$, in $\mathcal{O}(\log^\epsilon n)$ time we can return $\operatorname{repr}(v)$, where $v$ is a leaf of $\mathcal{T}_{\mathrm{st}}$ satisfying $\operatorname{str}(v) = T[j..n]$.*

*Proof.* Let $(b, e) = \operatorname{repr}(v)$. Observe that by definition of $v$, we have $\operatorname{Occ}(\operatorname{str}(v), T) = \{j\}$. Thus, it must hold $\{\operatorname{SA}[i]\}_{i \in (b..e]} = \{j\}$. This implies, that it suffices to compute $i = \operatorname{ISA}[j]$ and then return that $\operatorname{repr}(v) = (i - 1, i)$. Using Proposition 5.28, the computation of $\operatorname{ISA}[j]$ takes $\mathcal{O}(\log^\epsilon n)$ time. $\qquad\square$

### 7.4.21 Construction Algorithm

**Proposition 7.80.** *Given the packed representation of a text $T \in [0..\sigma)^n$, we can construct the data structure from Section 7.4.1 in $\mathcal{O}(n \min(1, \log \sigma / \sqrt{\log n}))$ time and $\mathcal{O}(n / \log_\sigma n)$ working space.*

*Proof.* The first part of the structure is constructed as follows. First, from a packed representation of $T$, we construct $\mathrm{C}_{\mathrm{ST}}(T)$ in $\mathcal{O}(n / \log_\sigma n)$ time using Proposition 7.14. Then, using Propositions 7.29 and 7.54, we augment $\mathrm{C}_{\mathrm{ST}}(T)$ in $\mathcal{O}(n \min(1, \log \sigma / \sqrt{\log n}))$ and $\mathcal{O}(n / \log_\sigma n)$ time (respectively) and using $\mathcal{O}(n / \log_\sigma n)$ working space into the two components of the structure from Section 7.4.1. The overall runtime is thus $\mathcal{O}(n \min(1, \log \sigma / \sqrt{\log n}))$.

Next, we compute $T^{\mathrm{rev}}$. With the help of the lookup table $L_{\mathrm{rev}}$, we first compute $T^{\mathrm{rev}}[1..n) = \overline{T[1..n)}$ in $\mathcal{O}(n / \log_\sigma n)$ time. In $\mathcal{O}(1)$ time we then append the sentinel $T^{\mathrm{rev}}[n] := 0$. After that, analogously as above, we construct the structures from Sections 7.2.1 and 7.3.1 for $T^{\mathrm{rev}}$, i.e., the second part of the structure from Section 7.4.1. $\qquad\square$

## 7.5 Summary

By combining Propositions 7.56 to 7.58, 7.60, 7.61, 7.64 to 7.66, 7.68, and 7.70 to 7.80 we obtain the following main result of this section.

**Theorem 7.81.** *Given any constant $\epsilon \in (0, 1)$ and the packed representation of a text $T \in [0\mathinner{.\,.}\sigma)^n$ with $2 \le \sigma < n^{1/7}$, in $\mathcal{O}(n \min(1, \log \sigma / \sqrt{\log n}))$ time and $\mathcal{O}(n / \log_\sigma n)$ working space we can construct a representation of the suffix tree of $T$ occupying $\mathcal{O}(n / \log_\sigma n)$ space and supporting all standard operations (see Table 1) in $\mathcal{O}(\log^\epsilon n)$ time.*

We also immediately obtain the following general reduction.

**Theorem 7.82.** *Consider a data structure answering prefix rank and selection queries that, for any string of length $m$ over alphabet $[0\mathinner{.\,.}\sigma)^\ell$, achieves the following complexities:*

1. *Space usage $S(m, \ell, \sigma)$,*
2. *Preprocessing time $P_t(m, \ell, \sigma)$,*
3. *Preprocessing space $P_s(m, \ell, \sigma)$,*
4. *Query time $Q(m, \ell, \sigma)$.*

*For every $T \in [0\mathinner{.\,.}\sigma)^n$ with $2 \le \sigma < n^{1/7}$, there exist $m = \mathcal{O}(n / \log_\sigma n)$ and $\ell = \mathcal{O}(\log_\sigma n)$ such that, given the packed representation of $T$, we can in $\mathcal{O}(n / \log_\sigma n + P_t(m, \ell, \sigma))$ time and $\mathcal{O}(n / \log_\sigma n + P_s(m, \ell, \sigma))$ working space construct a representation of the suffix tree of $T$ occupying $\mathcal{O}(n / \log_\sigma n + S(m, \ell, \sigma))$ space and supporting all standard operations (see Table 1) in $\mathcal{O}(\log \log n + Q(m, \ell, \sigma))$ time.*

# References

[1] Donald Adjeroh, Tim Bell, and Amar Mukherjee. *The Burrows-Wheeler Transform: Data Compression, Suffix Arrays, and Pattern Matching.* Springer, Boston, MA, USA, 2008. `doi:10.1007/978-0-387-78909-5`.

[2] Amihood Amir, Gad M. Landau, Moshe Lewenstein, and Dina Sokol. Dynamic text and static pattern matching. *ACM Trans. Algorithms*, 3(2):19, 2007. `doi:10.1145/1240233.1240242`.

[3] Maxim Babenko, Paweł Gawrychowski, Tomasz Kociumaka, and Tatiana Starikovskaya. Wavelet trees meet suffix trees. In *26th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 572–591, 2015. `doi:10.1137/1.9781611973730.39`.

[4] Jérémy Barbay, Francisco Claude, Travis Gagie, Gonzalo Navarro, and Yakov Nekrich. Efficient fully-compressed sequence representations. *Algorithmica*, 69(1):232–268, 2014. `doi:10.1007/s00453-012-9726-3`.

[5] Djamal Belazzougui. Linear time construction of compressed text indices in compact space. In *46th Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 148–193, 2014. `doi:10.1145/2591796.2591885`.

[6] Djamal Belazzougui, Travis Gagie, Paweł Gawrychowski, Juha Kärkkäinen, Alberto Ordóñez Pereira, Simon J. Puglisi, and Yasuo Tabei. Queries on LZ-bounded encodings. In *Data Compression Conference (DCC)*, pages 83–92, 2015. `doi:10.1109/DCC.2015.69`.

[7] Djamal Belazzougui and Gonzalo Navarro. Alphabet-independent compressed text indexing. *ACM Trans. Algorithms*, 10(4):23:1–23:19, 2014. `doi:10.1145/2635816`.

[8] Djamal Belazzougui and Gonzalo Navarro. Optimal lower and upper bounds for representing sequences. *ACM Trans. Algorithms*, 11(4):31:1–31:21, 2015. `doi:10.1145/2629339`.

[9] Djamal Belazzougui and Simon J. Puglisi. Range predecessor and Lempel-Ziv parsing. In *27th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 2053–2071, 2016. `doi:10.1137/1.9781611974331.ch143`.

[10] Michael A. Bender and Martin Farach-Colton. The LCA problem revisited. In *4th Latin American Symposium on Theoretical Informatics (LATIN)*, pages 88–94, 2000. `doi:10.1007/10719839_9`.

[11] Philip Bille, Mikko Berggren Ettienne, Inge Li Gørtz, and Hjalte Wedel Vildhøj. Time-space trade-offs for Lempel-Ziv compressed indexing. *Theor. Comput. Sci.*, 713:66–77, 2018. `doi:10.1016/j.tcs.2017.12.021`.

[12] Philip Bille, Inge Li Gørtz, and Frederik Rye Skjoldjensen. Deterministic indexing for packed strings. In *28th Annual Symposium on Combinatorial Pattern Matching (CPM)*, pages 6:1–6:11, 2017. `doi:10.4230/LIPIcs.CPM.2017.6`.

[13] Philip Bille, Gad M. Landau, Rajeev Raman, Kunihiko Sadakane, Srinivasa Rao Satti, and Oren Weimann. Random access to grammar-compressed strings and trees. *SIAM J. Comput.*, 44(3):513–539, 2015. `doi:10.1137/130936889`.

[14] Christina Boucher, Ondrej Cvacho, Travis Gagie, Jan Holub, Giovanni Manzini, Gonzalo Navarro, and Massimiliano Rossi. PFP compressed suffix trees. In *24th Symposium on Algorithm Engineering and Experiments (ALENEX)*, pages 60–72, 2021. `doi:10.1137/1.9781611976472.5`.

[15] Michael Burrows and David J. Wheeler. A block-sorting lossless data compression algorithm. Technical Report 124, Digital Equipment Corporation, 1994. URL: `https://www.hpl.hp.com/techreports/Compaq-DEC/SRC-RR-124.pdf`.

[16] Manuel Cáceres and Gonzalo Navarro. Faster repetition-aware compressed suffix trees based on block trees. *Inf. Comput.*, 285(Part):104749, 2022. `doi:10.1016/j.ic.2021.104749`.

[17] Ho-Leung Chan, Wing-Kai Hon, Tak Wah Lam, and Kunihiko Sadakane. Compressed indexes for dynamic text collections. *ACM Trans. Algorithms*, 3(2):21, 2007. `doi:10.1145/1240233.1240244`.

[18] Timothy M. Chan and Mihai Pătraşcu. Counting inversions, offline orthogonal range counting, and related problems. In *21st Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 161–173, 2010. `doi:10.1137/1.9781611973075.15`.

[19] Bernard Chazelle. A functional approach to data structures and its use in multidimensional searching. *SIAM J. Comput.*, 17(3):427–462, 1988. `doi:10.1137/0217026`.

[20] Anders Roy Christiansen, Mikko Berggren Ettienne, Tomasz Kociumaka, Gonzalo Navarro, and Nicola Prezza. Optimal-time dictionary-compressed indexes. *ACM Trans. Algorithms*, 17(1):8:1–8:39, 2021. `doi:10.1145/3426473`.

[21] David R. Clark. *Compact Pat Trees*. PhD thesis, University of Waterloo, 1998. URL: `https://uwspace.uwaterloo.ca/bitstream/handle/10012/64/nq21335.pdf`.

[22] Richard Cole, Tsvi Kopelowitz, and Moshe Lewenstein. Suffix trays and suffix trists: Structures for faster text indexing. *Algorithmica*, 72(2):450–466, 2015. `doi:10.1007/s00453-013-9860-6`.

[23] T. M. Cover and J. A. Thomas. *Elements of information theory*. Wiley, 2nd edition, 2006. `doi:10.1002/047174882X`.

[24] Maxime Crochemore, Christophe Hancart, and Thierry Lecroq. *Algorithms on strings*. Cambridge University Press, Cambridge, UK, 2007. `doi:10.1017/cbo9780511546853`.

[25] Martin Farach and S. Muthukrishnan. Perfect hashing for strings: Formalization and algorithms. In *7th Annual Symposium on Combinatorial Pattern Matching (CPM)*, pages 130–140, 1996. `doi:10.1007/3-540-61258-0\_11`.

[26] Martin Farach-Colton, Paolo Ferragina, and S. Muthukrishnan. On the sorting-complexity of suffix tree construction. *J. ACM*, 47(6):987–1011, 2000. `doi:10.1145/355541.355547`.

[27] Paolo Ferragina and Giovanni Manzini. Opportunistic data structures with applications. In *41st IEEE Annual Symposium on Foundations of Computer Science (FOCS)*, pages 390–398, 2000. `doi:10.1109/SFCS.2000.892127`.

[28] Paolo Ferragina and Giovanni Manzini. Indexing compressed text. *J. ACM*, 52(4):552–581, 2005. `doi:10.1145/1082036.1082039`.

[29] Nathan J. Fine and Herbert S. Wilf. Uniqueness theorems for periodic functions. *Proc. Am. Math. Soc.*, 16(1):109–114, 1965. `doi:10.2307/2034009`.

[30] Johannes Fischer and Paweł Gawrychowski. Alphabet-dependent string searching with Wexponential search trees. In *26th Annual Symposium on Combinatorial Pattern Matching (CPM)*, pages 160–171, 2015. Full version: `https://arxiv.org/abs/1302.3347`. `doi:10.1007/978-3-319-19929-0_14`.

[31] Johannes Fischer, Veli Mäkinen, and Gonzalo Navarro. Faster entropy-bounded compressed suffix trees. *Theor. Comput. Sci.*, 410(51):5354–5364, 2009. `doi:10.1016/j.tcs.2009.09.012`.

[32] Travis Gagie, Paweł Gawrychowski, Juha Kärkkäinen, Yakov Nekrich, and Simon J. Puglisi. A faster grammar-based self-index. In *6th International Conference on Language and Automata Theory and Applications (LATA)*, pages 240–251, 2012. `doi:10.1007/978-3-642-28332-1_21`.

[33] Travis Gagie, Paweł Gawrychowski, Juha Kärkkäinen, Yakov Nekrich, and Simon J. Puglisi. LZ77-based self-indexing with faster pattern matching. In *11th Latin American Symposium on Theoretical Informatics (LATIN)*, pages 731–742, 2014. `doi:10.1007/978-3-642-54423-1_63`.

[34] Travis Gagie, Gonzalo Navarro, and Nicola Prezza. Fully functional suffix trees and optimal text searching in BWT-runs bounded space. *J. ACM*, 67(1):1–54, apr 2020. `doi:10.1145/3375890`.

[35] Younan Gao, Meng He, and Yakov Nekrich. Fast preprocessing for optimal orthogonal range reporting and range successor with applications to text indexing. In *28th Annual European Symposium on Algorithms (ESA)*, pages 54:1–54:18, 2020. `doi:10.4230/LIPIcs.ESA.2020.54`.

[36] Pawel Gawrychowski, Adam Karczmarz, Tomasz Kociumaka, Jakub Lacki, and Piotr Sankowski. Optimal dynamic strings. In *29th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1509–1528, 2018. Full version: `https://arxiv.org/abs/1511.02612`. `doi:10.1137/1.9781611975031.99`.

[37] Simon Gog. *Compressed suffix trees: design, construction, and applications.* PhD thesis, University of Ulm, 2011. URL: `http://vts.uni-ulm.de/docs/2011/7786/vts_7786_11228.pdf`.

[38] Simon Gog, Timo Beller, Alistair Moffat, and Matthias Petri. From theory to practice: Plug and play with succinct data structures. In *13th International Symposium on Experimental Algorithms (SEA)*, pages 326–337, 2014. `doi:10.1007/978-3-319-07959-2_28`.

[39] Simon Gog, Juha Kärkkäinen, Dominik Kempa, Matthias Petri, and Simon J. Puglisi. Fixed block compression boosting in FM-indexes: Theory and practice. *Algorithmica*, 81(4):1370–1391, 2019. `doi:10.1007/s00453-018-0475-9`.

[40] Simon Gog, Alistair Moffat, and Matthias Petri. CSA++: Fast pattern search for large alphabets. In *19th Workshop on Algorithm Engineering and Experiments (ALENEX)*, pages 73–82, 2017. `doi:10.1137/1.9781611974768.6`.

[41] Roberto Grossi, Ankur Gupta, and Jeffrey Scott Vitter. High-order entropy-compressed text indexes. In *14th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 841–850, 2003. URL: `https://dl.acm.org/doi/10.5555/644108.644250`.

[42] Roberto Grossi and Jeffrey Scott Vitter. Compressed suffix arrays and suffix trees with applications to text indexing and string matching (extended abstract). In *32nd Annual ACM Symposium on Theory of Computing (STOC)*, pages 397–406, 2000. `doi:10.1145/335305.335351`.

[43] Roberto Grossi and Jeffrey Scott Vitter. Compressed suffix arrays and suffix trees with applications to text indexing and string matching. *SIAM J. Comput.*, 35(2):378–407, 2005. `doi:10.1137/S0097539702402354`.

[44] Dan Gusfield. *Algorithms on Strings, Trees, and Sequences - Computer Science and Computational Biology.* Cambridge University Press, 1997. `doi:10.1017/cbo9780511574931`.

[45] Torben Hagerup. Sorting and searching on the word RAM. In *15th Annual Symposium on Theoretical Aspects of Computer Science (STACS)*, pages 366–398, 1998. `doi:10.1007/BFb0028575`.

[46] Wing-Kai Hon, Kunihiko Sadakane, and Wing-Kin Sung. Breaking a time-and-space barrier in constructing full-text indices. In *44th IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 251–260, 2003. `doi:10.1109/SFCS.2003.1238199`.

[47] Guy Jacobson. Space-efficient static trees and graphs. In *30th IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 549–554, 1989. `doi:10.1109/SFCS.1989.63533`.

[48] Juha Kärkkäinen, Dominik Kempa, and Simon J. Puglisi. Hybrid compression of bitvectors for the FM-index. In *Data Compression Conference (DCC)*, pages 302–311, 2014. `doi:10.1109/DCC.2014.87`.

[49] Juha Kärkkäinen, Peter Sanders, and Stefan Burkhardt. Linear work suffix array construction. *J. ACM*, 53(6):918–936, 2006. `doi:10.1145/1217856.1217858`.

[50] Dominik Kempa and Tomasz Kociumaka. String synchronizing sets: Sublinear-time BWT construction and optimal LCE data structure. In *51st Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 756–767, 2019. `doi:10.1145/3313276.3316368`.

[51] Dominik Kempa and Tomasz Kociumaka. Resolution of the Burrows-Wheeler Transform conjecture. In *61st IEEE Annual Symposium on Foundations of Computer Science (FOCS)*, pages 1002–1013, 2020. `doi:10.1109/FOCS46700.2020.00097`.

[52] Dominik Kempa and Tomasz Kociumaka. Dynamic suffix array with polylogarithmic queries and updates. In *54th Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 1657–1670, 2022. `doi:10.1145/3519935.3520061`.

[53] Dominik Kempa and Nicola Prezza. At the roots of dictionary compression: String attractors. In *50th Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 827–840, 2018. `doi:10.1145/3188745.3188814`.

[54] Tomasz Kociumaka. *Efficient Data Structures for Internal Queries in Texts*. PhD thesis, University of Warsaw, 2018. URL: `https://depotuw.ceon.pl/bitstream/handle/item/3614/1000-DR-INF-170341.pdf`.

[55] Ben Langmead, Cole Trapnell, Mihai Pop, and Steven L Salzberg. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.*, 10(3):R25, 2009. `doi:10.1186/gb-2009-10-3-r25`.

[56] Heng Li and Richard Durbin. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinform.*, 25(14):1754–1760, 2009. `doi:10.1093/bioinformatics/btp324`.

[57] Ruiqiang Li, Chang Yu, Yingrui Li, Tak Wah Lam, Siu-Ming Yiu, Karsten Kristiansen, and Jun Wang. SOAP2: An improved ultrafast tool for short read alignment. *Bioinform.*, 25(15):1966–1967, 2009. `doi:10.1093/bioinformatics/btp336`.

[58] Veli Mäkinen, Djamal Belazzougui, Fabio Cunial, and Alexandru I. Tomescu. *Genome-scale algorithm design: Biological sequence analysis in the era of high-throughput sequencing.* Cambridge University Press, Cambridge, UK, 2015. `doi:10.1017/cbo9781139940023`.

[59] Veli Mäkinen and Gonzalo Navarro. Dynamic entropy-compressed sequences and full-text indexes. *ACM Trans. Algorithms*, 4(3):32:1–32:38, 2008. `doi:10.1145/1367064.1367072`.

[60] Udi Manber and Eugene W. Myers. Suffix arrays: A new method for on-line string searches. *SIAM J. Comput.*, 22(5):935–948, 1993. `doi:10.1137/0222058`.

[61] J. Ian Munro, Gonzalo Navarro, and Yakov Nekrich. Space-efficient construction of compressed indexes in deterministic linear time. In *28th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 408–424, 2017. `doi:10.1137/1.9781611974782.26`.

[62] J. Ian Munro, Gonzalo Navarro, and Yakov Nekrich. Fast compressed self-indexes with deterministic linear-time construction. *Algorithmica*, 82(2):316–337, 2020. `doi:10.1007/s00453-019-00637-x`.

[63] J. Ian Munro, Gonzalo Navarro, and Yakov Nekrich. Text indexing and searching in sublinear time. In *31st Annual Symposium on Combinatorial Pattern Matching (CPM)*, pages 24:1–24:15, 2020. `doi:10.4230/LIPIcs.CPM.2020.24`.

[64] J. Ian Munro, Yakov Nekrich, and Jeffrey Scott Vitter. Fast construction of wavelet trees. *Theor. Comput. Sci.*, 638:91–97, 2016. `doi:10.1016/j.tcs.2015.11.011`.

[65] Gonzalo Navarro. *Compact data structures: A practical approach.* Cambridge University Press, Cambridge, UK, 2016. `doi:10.1017/cbo9781316588284`.

[66] Gonzalo Navarro. Indexing highly repetitive string collections, part I: Repetitiveness measures. *ACM Comput. Surv.*, 54(2), 2021. `doi:10.1145/3434399`.

[67] Gonzalo Navarro. Indexing highly repetitive string collections, part II: Compressed indexes. *ACM Comput. Surv.*, 54(2), 2021. `doi:10.1145/3432999`.

[68] Gonzalo Navarro and Veli Mäkinen. Compressed full-text indexes. *ACM Comput. Surv.*, 39(1):2, 2007. `doi:10.1145/1216370.1216372`.

[69] Gonzalo Navarro and Yakov Nekrich. Time-optimal top-k document retrieval. *SIAM J. Comput.*, 46(1):80–113, 2017. `doi:10.1137/140998949`.

[70] Gonzalo Navarro and Nicola Prezza. Universal compressed text indexing. *Theor. Comput. Sci.*, 762:41–50, 2019. `doi:10.1016/j.tcs.2018.09.007`.

[71] Takaaki Nishimoto, Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Dynamic index and LZ factorization in compressed space. *Discret. Appl. Math.*, 274:116–129, 2020. `doi:10.1016/j.dam.2019.01.014`.

[72] Enno Ohlebusch. *Bioinformatics algorithms: Sequence analysis, genome rearrangements, and phylogenetic reconstruction.* Oldenbusch Verlag, Ulm, Germany, 2013.

[73] Enno Ohlebusch, Johannes Fischer, and Simon Gog. CST++. In *17th International Symposium on String Processing and Information Retrieval (SPIRE)*, pages 322–333, 2010. `doi:10.1007/978-3-642-16321-0_34`.

[74] Nicola Prezza. A framework of dynamic data structures for string processing. In *16th International Symposium on Experimental Algorithms (SEA)*, pages 11:1–11:15, 2017. `doi:10.4230/LIPIcs.SEA.2017.11`.

[75] Nicola Prezza and Giovanna Rosone. Space-efficient construction of compressed suffix trees. *Theor. Comput. Sci.*, 852:138–156, 2021. `doi:10.1016/j.tcs.2020.11.024`.

[76] Mihai Pătraşcu. Lower bounds for 2-dimensional range counting. In *39th Annual ACM Symposium on Theory of Computing (STOC)*, pages 40–46, 2007. `doi:10.1145/1250790.1250797`.

[77] Luís M. S. Russo, Gonzalo Navarro, and Arlindo L. Oliveira. Fully compressed suffix trees. *ACM Trans. Algorithms*, 7(4):53:1–53:34, 2011. `doi:10.1145/2000807.2000821`.

[78] Milan Růžić. Constructing efficient dictionaries in close to sorting time. In *35th International Colloquium on Automata, Languages, and Programming (ICALP)*, pages 84–95, 2008. `doi:10.1007/978-3-540-70575-8_8`.

[79] Kunihiko Sadakane. Succinct representations of lcp information and improvements in the compressed suffix arrays. In *13th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 225–232, 2002. URL: `http://dl.acm.org/citation.cfm?id=545381.545410`.

[80] Kunihiko Sadakane. Compressed suffix trees with full functionality. *Theory Comput. Syst.*, 41(4):589–607, 2007. `doi:10.1007/s00224-006-1198-x`.

[81] Peter Weiner. Linear pattern matching algorithms. In *14th Annual Symposium on Switching and Automata Theory (SWAT/FOCS)*, pages 1–11, 1973. `doi:10.1109/SWAT.1973.13`.

[82] Jacob Ziv and Abraham Lempel. A universal algorithm for sequential data compression. *Trans. Inf. Theory*, 23(3):337–343, 1977. `doi:10.1109/TIT.1977.1055714`.