

# GENOME-SCALE ALGORITHM DESIGN

by Veli Mäkinen, Djamel Belazzougui, Fabio Cunial and Alexandru I. Tomescu

Cambridge University Press, 2nd edition, 2023

www.genome-scale.info

## Exercises for Chapter 9. Burrows–Wheeler indexes

- 9.1 Consider the original definition of the Burrows–Wheeler transform with the cyclic shifts, without a unique endmarker  $\#$  added to the end of the string. Assume you have a fast algorithm to compute our default transform that assumes the endmarker is added. Show that you can feed that algorithm an input such that you can extract the original cyclic shift transform efficiently from the output.
- 9.2 Describe an algorithm that builds the succinct suffix array of a string  $T$  from the BWT index of  $T$ , within the bounds of Theorem 9.6.
- 9.3 Consider the strategy for answering a set of `occ` locate queries in batch described in Section 9.2.4. Suppose that, once all positions in  $\text{SA}_{T\#}$  have been marked, a bitvector `marked[1..n + 1]` is indexed to answer rank queries. Describe an alternative way to implement the batch extraction within the same asymptotic time bound as the one described in Section 9.2.4, but *without resorting to radix-sort*. Which strategy is more space-efficient? Which one is faster? The answer may vary depending on whether `occ` is small or large.
- 9.4 Recall that at the very end of the construction algorithm described in Section 9.3 we need to convert the BWT of the padded string  $X$  into the BWT of string  $T\#$ . Prove that  $\text{BWT}_X = T[n] \cdot \#^{k-1} \cdot W$  and  $\text{BWT}_{T\#} = T[n] \cdot W$ , where  $n = |T|$ ,  $k = |X| - n$ , and string  $W$  is a permutation of  $T[1..n - 1]$ .
- 9.5 Assume that you have a machine with  $P$  processors that share the same memory. Adapt Algorithm 9.3 to make use of all  $P$  processors, and estimate the resulting speedup with respect to Algorithm 9.3.
- 9.6 Recall that Property 1 is key for implementing Algorithm 9.3 with just one BWT. Prove this property.
- 9.7 In many applications, the operation `enumerateLeft`( $i, j$ ), where  $[i..j] = \mathbb{I}(W, T)$ , is immediately followed by operations `extendLeft`( $c, \mathbb{I}(W, T), \mathbb{I}(\underline{W}, \underline{T})$ ) applied to every distinct character  $c$  returned by `enumerateLeft`( $i, j$ ).
  - (a) Show how to use operation `rangeListExtended`( $T, i, j, l, r$ ) defined in Exercise 3.12 to efficiently support a new operation called `enumerateLeftExtended`( $\mathbb{I}(W, T), \mathbb{I}(\underline{W}, \underline{T})$ ) that returns the distinct characters that appear in substring  $\text{BWT}_{T\#}[i..j]$ , and for each such character  $c$  returns the pair of intervals computed by `extendLeft`( $c, \mathbb{I}(W, T), \mathbb{I}(\underline{W}, \underline{T})$ ). Implement also the symmetric operation `enumerateRightExtended`.
  - (b) Reimplement Algorithm 9.3 to use the operation `enumerateLeftExtended` instead of `enumerateLeft`, and `extendLeft`.
- 9.8 Recall that Property 2 is key for navigating the Burrows–Wheeler index of a labeled tree top-down (Section 9.5.1). Prove this property.

9.9 Consider the Burrows–Wheeler index of a tree  $\mathcal{T}$  described in Section 9.5.

- (a) Given a character  $c \in [1..\sigma]$ , let  $\bar{c}$  be the starting position of its interval  $\bar{c}$  in  $\text{BWT}_{\mathcal{T}}$ , and let  $C'[1..\sigma]$  be an array such that  $C'[c] = \bar{c}$ . Describe an  $O(n)$ -time algorithm to compute  $C'$  from  $C$  and  $\text{last}$ .
- (b) Recall that  $\mathbb{I}(v) = [\bar{v}..\bar{v}]$  is the interval in  $\text{BWT}_{\mathcal{T}}$  that contains all the children of node  $v$ , and that  $\mathbb{I}'(v)$  is the position of the arc  $(v, \text{parent}(v))$  in  $\text{BWT}_{\mathcal{T}}$ . Let  $D[1..n+m-1]$  be an array such that  $D[i] = \bar{v}$  where  $\mathbb{I}'(v) = i$ , or  $-1$  if  $v$  is a leaf. Describe an  $O(n+m)$ -time algorithm to compute  $D$  from  $C'$ ,  $\text{last}$  and  $\text{labels}$ . *Hint.* Use the same high-level strategy as the one in Lemmas 8.7 and 9.9.
- (c) Give the pseudocode of an  $O(n)$ -time algorithm that reconstructs the original tree  $\mathcal{T}$  from  $\text{BWT}_{\mathcal{T}}$  and  $D$ , in depth-first order.

9.10 In order to build the Burrows–Wheeler index of a tree  $\mathcal{T}$  described in Section 9.5, we can adapt the prefix-doubling approach used in Lemma 8.8 for building suffix arrays. In particular, assume that we want to assign to every node  $v$  the order  $R(v)$  of its path to the root, among all such paths in  $\mathcal{T}$ . We define the  $i$ th contraction of tree  $\mathcal{T}$  as the graph  $\tilde{T}^i = (V, E^i)$  such that  $(u, v) \in E^i$  iff  $v_1, v_2, \dots, v_{2^{i+1}}$  is a path in  $\mathcal{T}$  with  $v_0 = u$  and  $v_{2^{i+1}} = v$ . Clearly  $\tilde{T}^0 = \mathcal{T}$ , and  $E^{i+1}$  is the set of all pairs  $(u, v)$  such that  $P = u, w, v$  is a path in  $\tilde{T}^i$ . Show how to iteratively refine the value  $R(v)$  for every node, starting from the initial approximation  $R^0(v) = \overline{\ell(v)} = C'[\ell(v)]$  computed as in Exercise 9.9. Describe the time and space complexity of your algorithm.

9.11 Assume that we want to generalize the recursive,  $O(\sigma+n)$ -time algorithm for building the suffix array  $\text{SA}_{S\#}$  of a string  $S$  to a similar algorithm for building the Burrows–Wheeler index of a labeled tree  $\mathcal{T}$  described in Section 9.5. Do the results in Definition 8.9, Lemma 8.10, and Lemma 8.11 still hold? Does the overall generalization still achieve linear time for any  $\mathcal{T}$ ? If not, for which topologies or labelings of  $\mathcal{T}$  does the generalization achieve linear time?

9.12 Assume that a numerical value  $\text{id}(v)$  must be associated with every vertex  $v$  of a labeled, strongly distinguishable DAG  $G$ : for example,  $\text{id}(v)$  could be the probability of reaching vertex  $v$  from  $s$ . Assume that, whenever there is exactly one path  $P$  connecting vertex  $u$  to vertex  $v$ , we have that  $\text{id}(v) = f(\text{id}(u), |P|)$ , where  $f$  is a known function. Describe a method to augment the Burrows–Wheeler index of  $G$  described in Section 9.6 with the values of  $\text{id}$ , and its space and time complexity. *Hint.* Adapt the strategy used for strings in Section 9.2.3.

9.13 Consider the Burrows–Wheeler index of a labeled, strongly distinguishable DAG described in Section 9.6, and assume that  $\text{labels}$  is represented with  $\sigma+1$  bitvectors  $B_c[1..m]$  for all  $c \in [1..\sigma+1]$ , where  $B_c[i] = 1$  if and only if  $\text{labels}[i] = c$ . Describe a more space-efficient encoding of  $\text{labels}$ . Does it speed up the operations described in Section 9.6?

9.14 Given any directed labeled graph  $G = (V, E, \Sigma)$ , Algorithm 9.4 builds a reverse-deterministic graph  $G' = (V', E', \Sigma)$  in which a vertex  $v \in V'$  corresponds to a subset of  $V$ . Adapt the algorithm so that a vertex  $v \in G'$  corresponds to a subset of  $E$ . Show an example in which the number of vertices in the reverse-deterministic graph produced by this version of the algorithm is smaller than the number of vertices in the reverse-deterministic graph produced by Algorithm 9.4.

- 9.15 Recall that in Section 9.6.3 a labeled DAG  $G$  is transformed into a reverse-deterministic DAG  $G'$ , and then  $G'$  is itself transformed into a strongly distinguishable DAG  $G^*$ . Show that  $|G^*| \in O(2^{|G|})$  in the worst case.
- 9.16 Consider the algorithm to enforce distinguishability described in Section 9.6.3.
- Show that this algorithm does not also enforce reverse-determinism.
  - Assume that the input to this algorithm contains two arcs  $(u, w)$  and  $(v, w)$  such that  $\ell(u) = \ell(v)$ : give an upper bound on the number of vertices that result from the projection of  $u$  and  $v$  to the output.
  - Modify the algorithm to build a distinguishable, rather than a strongly distinguishable, DAG.
  - Modify the algorithm to decide whether the input graph is distinguishable or not.
  - Implement the algorithm using just sorts, scans, and joins of lists of tuples of integers, *without storing or manipulating strings* `path` *explicitly*. Give the pseudocode of your implementation.
  - Give the time and space complexity of the algorithm as a function of the size of its input and its output.
- 9.17 Describe a labeled, strongly distinguishable DAG that recognizes an infinite language. Describe an infinite language that cannot be recognized by any labeled, strongly distinguishable DAG.
- 9.18 Consider the frequency-oblivious representation of a de Bruijn graph described in Section 9.7.1. Show how to implement the function `getArc` when `labels` contains characters in  $\Sigma \cup \Sigma'$ .
- 9.19 Recall that marking the starting position of every interval in  $\text{BWT}_{T\#}$  that corresponds to a  $k$ -mer is a key step for building both the frequency-aware and the frequency-oblivious representation of a de Bruijn graph. Prove that Lemma 9.22 performs this marking correctly.
- 9.20 Some applications in high-throughput sequencing require the de Bruijn graph of a set of strings  $\mathcal{R} = \{R^1, R^2, \dots, R^r\}$ . Describe how to adapt the data structures in Sections 9.7.1 and 9.7.2 to represent the de Bruijn graph of string  $R = R^1 \# R^2 \# \dots \# R^r \#$ , where  $\# \notin [1..\sigma]$ . Give an upper bound on the space taken by the frequency-oblivious and by the frequency-aware representation.
- 9.21 Describe the frequency-oblivious and the frequency-aware representations of the de Bruijn graph  $\text{DBG}_{T,k}$  when every string of length  $k$  on a reference alphabet  $\Sigma$  occurs at most once in  $T$ .
- 9.22 Develop an alternative representation for de Bruijn graphs using Bloom filters of Section 3.5.2.
- 9.23 Some long read sequencers cannot measure the length of a run of the same character accurately, so reads are often preprocessed by compressing every such run into a single character. Consider the graph  $G = (V, E)$  where the set of vertices  $V$  coincides with the set of all distinct strings of length  $k$  on alphabet  $\Sigma$  such that no two consecutive characters are equal, and where there is an arc  $(u, v)$  in  $E$  iff the

suffix of  $u$  of length  $k - 1$  coincides with the prefix of  $v$  of length  $k - 1$ .  $G$  is called *Kautz graph* of order  $k$  on  $\Sigma$ .

- (i) Let  $G' = (V', E')$  be the de Bruijn graph of order  $k - 1$  on an alphabet  $\Sigma'$  with  $|\Sigma'| = |\Sigma| - 1$  (so  $V'$  is the set of all distinct strings of length  $k - 1$  on  $\Sigma'$ ). Show that there is a surjective map  $f : G \mapsto G'$  such that if  $(u, v) \in E$  then  $(f(u), f(v)) \in E'$ .
- (ii) Consider the undirected version of  $G$ . For a node  $u$ , let  $N^L(u)$  be the set of its left-neighbors, let  $N^R(u)$  be the set of its right-neighbors, and let  $N(u) = N^L(u) \cup N^R(u)$ . Show that: (1)  $(u, v) \in E \Rightarrow |N(u) \cap N(v)| \leq 1$ ; if such a common neighbor  $w$  exists, then  $u, v$  and  $w$  are periodic strings with periods  $abc, bca,$  and  $cab$ , respectively, for  $\{a, b, c\} \subseteq \Sigma$ . (2) Given a node  $u$ , if  $N^L(u) \cap N^R(u) \neq \emptyset$ , then such a common node  $v$  must have period  $ba$  and  $u$  must have period  $ab$ .