



1 Eulertigs: minimum plain text representation of 2 k -mer sets without repetitions in linear time

3 Sebastian Schmidt¹  

4 University of Helsinki, Finland

5 Jarno N. Alanko  

6 University of Helsinki, Finland

7 — Abstract —

8 A fundamental operation in computational genomics is to reduce the input sequences to their
9 constituent k -mers. For maximum performance of downstream applications it is important to store
10 the k -mers in small space, while keeping the representation easy and efficient to use (i.e. without
11 k -mer repetitions and in plain text). Recently, heuristics were presented to compute a near-minimum
12 such representation. We present an algorithm to compute a minimum representation in optimal
13 (linear) time and use it to evaluate the existing heuristics. For that, we present a formalisation
14 of arc-centric bidirected de Bruijn graphs and carefully prove that it accurately models the k -mer
15 spectrum of the input. Our algorithm first constructs the de Bruijn graph in linear time in the
16 length of the input strings (for a fixed-size alphabet). Then it uses a Eulerian-cycle-based algorithm
17 to compute the minimum representation, in time linear in the size of the output.

18 **2012 ACM Subject Classification** Applied computing → Computational biology; Theory of com-
19 putation → Data compression; Theory of computation → Graph algorithms analysis; Theory of
20 computation → Data structures design and analysis

21 **Keywords and phrases** Spectrum preserving string sets, Eulerian cycle, Suffix tree, Bidirected
22 arc-centric de Bruijn graph, k -mer based methods

23 **Digital Object Identifier** 10.4230/LIPIcs.WABI.2022.2

24 **Related Version** This paper is also available on biorxiv <https://www.biorxiv.org/content/10.1101/2022.05.17.492399>

26 **Funding** *Sebastian Schmidt*: Funded by the European Research Council (ERC) under the European
27 Union's Horizon 2020 research and innovation programme (grant agreement No. 851093, SAFE BIO)

28 *Jarno N. Alanko*: Funded by NIH NIAID grant No. R01HG011392 and Academy of Finland grant
29 339070.

¹ corresponding author



2:2 Eulertigs: minimum plain text k -mer sets without repetitions in linear time

30 **1** Introduction

31 Motivation

32 A k -mer is a DNA string of length k that is considered equal to itself and its reverse
33 complement. A common pattern in bioinformatics is to reduce a set of input strings to their
34 constituent k -mers. Such representations are at the core of many bioinformatics pipelines –
35 see e.g. Schmidt et al. [23] or Brinda et al. [6] for an overview of applications. The wide-spread
36 use of k -mer sets has prompted the question of what is the smallest *plain text representation*
37 for a set of k -mers. Here, a plain text representation means a set of strings that have the
38 same set of k -mers as the input strings, i.e. the *spectrum* is preserved. Such representations
39 are also called *spectrum preserving string sets* (SPSS) [22], or *simplitigs* [6]. This has the
40 following advantages over encoded representations:

- 41 ■ When storing k -mer sets to disk, plain text may remove the need of decompression before
42 usage, as some tools that usually take unitigs as input can take any other plain text
43 representation without modification (e.g. Bifrost [13]).
- 44 ■ Within an application, an encoded representation would require decoding whenever a
45 k -mer is accessed, which may slow down the application a lot compared to when each
46 k -mer is in RAM in plain text.

47 Further, in applications, it might be useful if the representation contains each k -mer exactly
48 once. This is because some applications, like e.g. SShash [21], are able to take any set of
49 k -mers as input, but cannot easily deal with duplicate k -mers in the input.

50 Related work

51 There are two heuristic approaches to the construction of a small SPSS without repeti-
52 tions, namely *prophasm* [6] and *UST* [22]. While neither of these computes a minimum
53 representation, Rahman et al. [22] also present a lower bound to the minimum size of any
54 representation without repetition, and they show that they are within 3% of this lower bound
55 in practice. They also present a counter-example showing that their lower bound is not tight.
56 Small SPSSs without repetitions are used e.g. in SShash [21] and are also computed by
57 state-of-the-art de Bruijn graph compactors like Cuttlefish 2 [15].

58 When k -mer repetitions are allowed in an SPSS, there is a known polynomially computable
59 minimum representation, namely *matchtigs* [23]. While matchtigs are expensive to compute,
60 the authors also present a more efficient greedy heuristic that is able to compute a near-
61 minimum representation on a modern server with no significant penalty in runtime (when
62 compared to computing just unitigs), but a significant increase in RAM usage.

63 In [6, 23] the authors also showed that decreasing the size of an SPSS results in signi-
64 ficantly better performance in downstream applications, i.e. when further compressing the
65 representation with general purpose compressors, or when performing k -mer-based queries.

66 The authors of both [6] and [22] consider whether computing a minimum representation
67 without repetitions may be NP-hard, as it is equivalent to computing a minimum path cover
68 in a de Bruijn graph, which is NP-hard in general graphs by reduction from Hamiltonian cycle.
69 However, computing a Hamiltonian cycle in a de Bruijn graph is actually polynomial [14].
70 The authors of [14] argue that de Bruijn graphs are a subclass of *adjoint* graphs, in which
71 solving the Hamiltonian cycle problem is equivalent to solving the Eulerian cycle problem

72 in the *original* of the adjoint graph, which can be computed in linear time². However, the
73 argument is only made for normal directed (and not bidirected) graphs, and thus is not
74 applicable to our setup, where a k -mer is also considered equal to its reverse complement.

75 Our contributions

76 Our first technical contribution is to carefully define the notion of a bidirected de Bruijn
77 graph such that the spectrum of the input is accurately modelled in the allowed walks of the
78 graph. Our definition also takes into account k -mers that are their own reverse complement.
79 This technicality is often neglected in the literature, and sidestepped by requiring that the
80 value of k is odd, in which case this special case does not occur. We give a suffix-tree-
81 based deterministic linear-time algorithm to construct such a graph, filling a theory gap
82 in the literature, as existing approaches [8, 15, 13, 1] depend on the value of k and/or
83 are probabilistic due to the use of hashing, minimizers or Bloom filters, or do not use the
84 reverse-complement-aware definition of k -mers [7].

85 Given the bidirected de Bruijn graph, we present an algorithm that computes a minimum
86 plain text representation of k -mer sets without repetitions, which runs in output sensitive
87 linear time. Steps 1 to 3 run in linear time in the number of nodes and arcs in the graph. In
88 short, it works as follows:

- 89 1. Add breaking arcs into this graph to make it Eulerian.
- 90 2. Compute a Eulerian cycle in the resulting graph.
- 91 3. Break that cycle at the breaking arcs.
- 92 4. Output the strings spelled by the resulting walks.

93 The algorithm is essentially an adaption of the `matchtigs` algorithm [23], removing the
94 possibility of joining walks by repeating k -mers. We give detailed descriptions for all these
95 steps and prove their correctness in our bidirected de Bruijn graph model. Together with our
96 linear-time de Bruijn graph construction algorithm, we obtain the main result of our paper:

97 ► **Theorem 1.** *Let k be a positive integer and let I be a set of strings of length at least k*
98 *over some alphabet Σ . Then we can compute a set of strings I' of length at least k with*
99 *minimum cumulative length and $CS_k(I) = CS_k(I')$ in $O(|I| \log |\Sigma|)$ time.*

100 where $CS_k(I) = CS_k(I')$ means that I' is an SPSS of I , and $|I|$ is the cumulative length of
101 I (see Section 2 for accurate definitions). This gives a positive answer to the open question
102 if a minimum SPSS without repetitions can be computed in polynomial time. Additionally,
103 we give an easily computable tight lower bound on the size of a minimum SPSS without
104 repetitions.

105 For our experiments, we have implemented steps 1 to 4 in Rust, taking the de Bruijn
106 graph as given. The implementation is available on github: <https://github.com/algbio/matchtigs>. Our experimental evaluation shows that our algorithm does not result in
107 significant practical improvements, but for the first time allows to benchmark the quality the
108 heuristics `prophasm` and `UST` against an optimal solution. It turns out that both produce
109 close-to-optimal results, but with a different distribution of computational resources.
110

² The original of an adjoint graph can be computed by splitting each node v into two nodes v' and v'' such that v' keeps the incoming arcs, and v'' the outgoing arcs as in [5, Figure 4]. Then, the graph is a collection of complete bipartite graphs [5]. These graphs can be contracted into single nodes, and then we add an arc between the contracted representations of each v' and v'' . This can be computed in linear time and is the original graph, since all nodes have become arcs again, and the arcs have the correct predecessors and successors.

2:4 Eulertigs: minimum plain text k -mer sets without repetitions in linear time

111 Our work also shows that using arc-centric de Bruijn graphs can aid the intuition for
112 certain problems, as in this case, the node-centric variant hides the relationship between
113 Eulerian cycles and minimum SPSS without repetition.

114 Organisation of the paper

115 In Section 2 we give preliminary definitions of well-known concepts. In Section 3 we define
116 de Bruijn graphs and prove the soundness of the definitions. In Section 4 we show how
117 to construct de Bruijn graphs by our definitions in linear time. In Section 5 we show how
118 to construct a minimum SPSS without repetitions in linear time if the de Bruijn graph is
119 given. In Section 6 we compare our algorithm and Eulertigs against strings computed with
120 prophase and UST on practical data sets.

121 **2 Preliminaries**

122 In this section we give the prerequisite knowledge required for this paper.

123 2.1 Bidirected graphs

124 In this section we define our notion of the bidirected graphs and the incidence model.

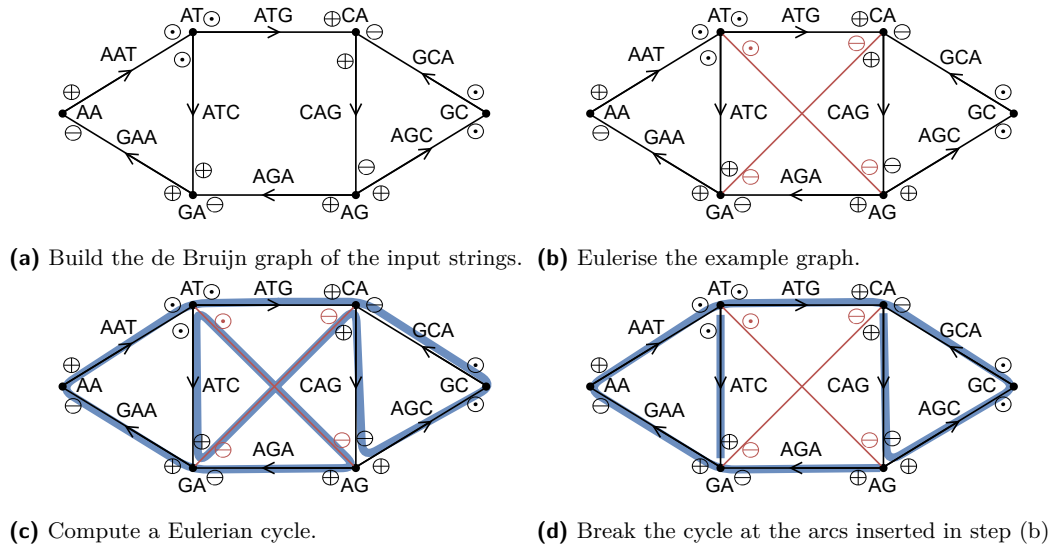
125 A multiset is defined as a set M , and an implicit function $\#_M : M \rightarrow \mathbb{Z}^+$ mapping
126 elements to their multiplicities. The cardinality is defined as $|M| := \sum_{s \in M} \#_M(s)$.

127 An *alphabet* Σ is an ordered set, and a Σ -*word* is a string of characters of that set. String
128 concatenation is written as ab for two strings a and b . The set Σ^k is the set of all Σ -words
129 of length k and the set Σ^* is the set of all Σ -words, including the empty word ϵ . Given
130 a positive integer k , the k -*suffix* $\text{suf}_k(w)$ (k -*prefix* $\text{pre}_k(w)$) of a word w is the substring
131 of its last (first) k characters. A k -*mer* is a word of length k . A *complement function*
132 over Σ is a function $\text{comp} : \Sigma \rightarrow \Sigma$ mapping characters to characters that is self-inverse
133 (i.e. $\text{comp}(\text{comp}(x)) = x$). A *reverse complement* function for alphabet Σ is a function
134 $\text{rc} : \Sigma^* \rightarrow \Sigma^*$ defined as $\text{rc}((w_1, \dots, w_\ell)) := (\text{comp}(w_\ell), \dots, \text{comp}(w_1))$, for some arbitrary
135 complement function comp . On sets, rc is defined to compute the reverse complement of
136 each element in the set. Note that rc is self-inverse. A *canonical k -mer* is a k -mer that is
137 lexicographically smaller than or equal to its reverse complement.

138 Given an integer k and an alphabet Σ , the k -*spectrum* of a set of strings $I \subseteq \bigcup_{k' \geq k} \Sigma^{k'}$
139 is a set of strings $S_k(I) := \{w \in \Sigma^k \mid \exists i \in I : w \text{ is substring of } i \text{ or } \text{rc}(i)\}$. The *canonical*
140 k -*spectrum* of I is $\text{CS}_k(I) := \{w \in S_k(I) \mid w \text{ is canonical}\}$. For simplicity, the spectrum and
141 canonical spectrum are defined for a single string w as if it were a set $\{w\}$. A *spectrum*
142 *preserving string set* of a set of strings I is a set of strings I' such that $\text{CS}_k(I) = \text{CS}_k(I')$.
143 The cumulative length of I is $\|I\| := \sum_{w \in I} |w|$.

144 Our definition of a bidirected graph is mostly standard like in e.g. [17], however we allow
145 self-complemental nodes that occur in bidirected de Bruijn graphs. A *bidirected graph* is a
146 tuple $G = (V, E, c)$ with a set of normal and *self-complemental* nodes $v \in V$, a set of arcs
147 $e \in E$, and a function $c : V \rightarrow \{1, 0\}$ marking self-complemental nodes with 1, and normal
148 nodes with 0. An *incidence* is a pair vd , where $d \in \{\oplus, \ominus, \odot\}$ is called its *sign* (e.g. $v\oplus$). The
149 negation of a sign is defined as $\neg\oplus := \ominus$, $\neg\ominus := \oplus$ and $\neg\odot := \odot$. For self-complemental nodes
150 $v \in V$, only incidences $v\odot$ are allowed, and for normal nodes only incidences $v\oplus$ and $v\ominus$ are
151 allowed. An *arc* $(v_1d_1, v'_1d'_1, \eta) \in E$ is a tuple of incidences and a unique identifier η , where η
152 can be of any type. The *reversal* of an arc is denoted by $(v_1d_1, v'_1d'_1, \eta)^{-1} := (v'_1d'_1, v_1d_1, \eta)$.

153 If not required, we may drop the identifier (i.e. just write $(v_1 \ominus, v'_1 \odot) \in E$). If a node $v \in V$
 154 is present with a \oplus (\ominus) sign in an arc, then the arc is *outgoing* (*incoming*) from (to) v .



■ **Figure 1** Overview of our algorithm executed on the input strings $\{GAATG, ATCTGCT\}$ with $k = 3$. After step (d), the resulting spelled SPSS is $\{ATC, AGAATGCTG\}$.

155 Note that, other than in standard directed graphs, in bidirected graphs arcs can be
 156 outgoing or incoming on both ends, and the order of the incidences in the arc does not affect
 157 if it is outgoing or incoming to a node. Further, our notation differs from that of standard
 158 bidirected graphs in that arcs have a direction. This is required because we will work with
 159 arc-centric de Bruijn graphs (see Section 3), which have labels on the arcs and not the nodes.
 160 Using the sign of the incidence pairs, it is possible to decide if a node is traversed forwards
 161 or backwards, but not if the arc is traversed forwards or backwards. But to decide which
 162 label (forwards or reverse complement) to use when computing the string spelled by an arc,
 163 the direction is relevant. See Figure 1 (a) for an example of a bigraph, which has labels that
 164 make it a de Bruijn graph as well.

165 A *walk* in a bigraph is a sequence of arcs $W := ((v_1 d_1, v'_1 d'_1, \eta_1), (v_2 d_2, v'_2 d'_2, \eta_2), \dots,$
 166 $(v_\ell d_\ell, v'_\ell d'_\ell, \eta_\ell))$ where for every i it holds that $(v_i d_i, v'_i d'_i, \eta_i) \in E$ or $(v'_i d'_i, v_i d_i, \eta_i) \in E$ (we
 167 can arbitrarily walk over arcs forwards and reverse), and for every $i < \ell$ it holds that $v'_i = v_{i+1}$
 168 and $d'_i = -d_{i+1}$. The *length* of a walk is $\ell = |W|$. If $v_1 = v'_\ell$ and $d_1 = -d'_\ell$, then W is a *cycle*.
 169 A bigraph is *connected*, if for each pair of nodes $v_1, v_2 \in V$ there is a walk from v_1 to v_2 .

170 For a node $v \in V$, the *multiset of incidences* is defined as $I(v) := \{vd \mid d \in \{\oplus, \ominus, \odot\}\}$,
 171 with multiplicities $\#_{I(v)}(vd) := \sum_{e \in E} \#_e(vd)$ (treating the arcs as multisets such that self-
 172 loops count as two separate incidences). For a node $v \in V$ that is not self-complemental, the
 173 *outdegree* is defined as $\delta^+(v) := \#_{I(v)}(v\oplus)$, and the *indegree* is defined as $\delta^-(v) := \#_{I(v)}(v\ominus)$.
 174 For a self-complemental node $v \in V$, the *degree* is defined as $\delta(v) := \#_{I(v)}(v\odot)$.

175 We define the *imbalance* of a node $v \in V$ that is not self-complemental as the difference
 176 of its outdegree and indegree $\text{imbalance}(v) := \delta^+(v) - \delta^-(v)$. For a self-complemental node
 177 $v \in V$ the imbalance is defined as $\text{imbalance}(v) := 1$ if $\delta(v)$ is odd, and $\text{imbalance}(v) := 0$
 178 otherwise. A node $v \in V$ is called *unbalanced*, if $\text{imbalance}(v) \neq 0$, and *balanced* otherwise.

179 A *labelled graph* is a bidirected graph $G = (V, E, c)$ where the identifiers of arcs are strings
 180 over some alphabet Σ (e.g. $(v_1 \oplus, v_2 \ominus, ACCTG) \in E$).

2:6 Eulertigs: minimum plain text k -mer sets without repetitions in linear time

181 2.2 Suffix arrays and suffix trees

182 Section 4 requires knowledge of suffix arrays and suffix trees. We assume the reader is familiar
183 with these data structures, and briefly give the relevant definitions and properties below. We
184 point the reader to Gusfield [12] and Mäkinen [18] for an in-depth treatment of the topics.

185 A suffix array SA_T for a string T is an array of length $|T|$ such that $SA_T[i]$ is the starting
186 position of the lexicographically i -th suffix of T . The suffix array interval of a string x is the
187 maximal interval $[i..j]$ such that all the suffixes pointed by $SA_T[i], \dots, SA_T[j]$ have x as a
188 prefix, or the empty interval if x is not a substring of T .

189 A suffix tree of a string T is a compacted version of the trie of all suffixes of T , such that
190 non-branching paths are merged into single arcs, with arcs pointing away from the root. The
191 compactification concatenates the labels of the arcs on the compacted path. The nodes that
192 were compacted away and are now in the middle of an arc are called implicit nodes, and the
193 rest of the nodes are explicit. A *locus* (plural *loci*) is a node that is either explicit or implicit.
194 A locus v is represented by a pair (u, d) , where u is the explicit suffix tree node at the end of
195 the arc containing v (u is equal to v if v is explicit), and d is the depth of locus v in the trie
196 of loci. The suffix array interval of a node is the interval of leaves in the subtree of the node.
197 The suffix array interval of an implicit locus (u, d) is the same as the suffix array interval of
198 u .

199 The suffix tree can be constructed in linear time in $|T|$ using e.g. Ukkonen's algorithm [24].
200 The tree comes with a function `child` that takes an explicit node and a character, and returns
201 the child at the end of the arc from that node whose label starts with the given character (if
202 such node exists). This can be implemented in $O(\log |\Sigma|)$ time by binary searching over `child`
203 pointers sorted by labels. The `child` function can also be easily implemented for implicit
204 loci. Ukkonen's algorithm also produces *suffix links* for the explicit nodes, which map from
205 the suffix tree node of a string cx to the suffix tree node of string x . It is possible to emulate
206 suffix links on the implicit loci using constant-time weighted level-ancestor queries [4] by
207 mapping $(u, d) \mapsto (f_{d-1}(SL(u)), d-1)$, where $SL(u)$ is the destination of a suffix link from
208 u , and $f_{d-1}(SL(u))$ is the furthest suffix tree ancestor from $SL(u)$ at depth at least $d-1$ in
209 the trie of loci. The inverse pointers of suffix links are called *Weiner links*, and they can also
210 be simulated on the implicit loci by mapping $(u, d) \mapsto (WL(u, c), d+1)$, where $WL(u, c)$ is
211 the destination of a Weiner link from u with character c .

212 3 De Bruijn graphs

213 The *de Bruijn graph* of order k of a set of input strings I is defined as a labelled graph
214 constructed by Algorithm 1. See Figure 1 (a) for an example. A de Bruijn graph computed
215 by this algorithm has the following property (see Appendix B for some of the proofs of this
216 section).

217 ► **Lemma 2 (Sound labels).** *Let k be a positive integer and let I be a set of strings of length*
218 *at least k . Let $G = (V, E, c)$ be the de Bruijn graph of order k constructed from I . For all*
219 *pairs of arcs $e_1 := (v_1d_1, v'_1d'_1, \eta_1), e_2 := (v_2d_2, v'_2d'_2, \eta_2) \in E$ it holds that:*

- 220 (a) $(v'_1 = v_2 \text{ and } d'_1 = -d_2)$ if and only if $\text{suf}_{k-1}(\eta_1) = \text{pre}_{k-1}(\eta_2)$,
- 221 (b) $(v'_1 = v'_2 \text{ and } d'_1 = -d'_2)$ if and only if $\text{suf}_{k-1}(\eta_1) = \text{pre}_{k-1}(\text{rc}(\eta_2))$,
- 222 (c) $(v_1 = v_2 \text{ and } d_1 = -d_2)$ if and only if $\text{suf}_{k-1}(\text{rc}(\eta_1)) = \text{pre}_{k-1}(\eta_2)$, and
- 223 (d) $(v_1 = v'_2 \text{ and } d_1 = -d'_2)$ if and only if $\text{suf}_{k-1}(\text{rc}(\eta_1)) = \text{pre}_{k-1}(\text{rc}(\eta_2))$.

224 For a walk $W := (e_1 = (v_1d_1, v'_1d'_1, \eta_1), \dots, e_\ell = (v_\ell d_\ell, v'_\ell d'_\ell, \eta_\ell))$ in a de Bruijn graph, its
225 *sequence of k -mers* is $K := (\kappa_1, \dots, \kappa_\ell)$, where for each i we define κ_i as η_i if $e_i \in E$, and

■ **Algorithm 1** DEBRUIJNGRAPH

Input: An integer k and a set of strings I where each string has length at least k .
Output: A de Bruijn graph $G = (V, E, c)$ of order k .

```

1  $V \leftarrow CS_{k-1}(I)$  /* the nodes are the canonical  $(k-1)$ -mers */
2 foreach  $v \in V$  do
3    $\lfloor$  if  $rc(v) = v$  then  $c(v) \leftarrow 1$  else  $c(v) \leftarrow 0$ 
4 foreach  $\eta \in CS_k(I)$  do
5    $w \leftarrow pre_{k-1}(\eta)$  /* compute  $v$  */
6    $v \leftarrow canonical\ w$ 
7    $w' \leftarrow suf_{k-1}(\eta)$  /* compute  $v'$  */
8    $v' \leftarrow canonical\ w'$ 
9   if  $c(v) = 1$  then  $d \leftarrow \ominus$  /* compute the direction of  $v$  */
10  else if  $v = w$  then  $d \leftarrow \oplus$ 
11  else  $d \leftarrow \ominus$ 
12  if  $c(v') = 1$  then  $d' \leftarrow \ominus$  /* compute the direction of  $v'$  */
13  else if  $v' \neq w'$  then  $d' \leftarrow \oplus$  /* note that  $\neq$  differs from  $=$  above */
14  else  $d' \leftarrow \ominus$ 
15   $e \leftarrow (vd, v'd', \eta)$  /* insert the arc into the graph */
16   $E \leftarrow E \cup \{e\}$ 

```

226 as $rc(\eta_i)$ if $e_1^{-1} \in E$. The string $spell(W)$ is the string *spelled* by W , which is defined as its
227 *collapsed* sequence of kmers, i.e. its sequence of k -mers gets concatenated while overlapping
228 consecutive k -mers by $k - 1$. This is computed by Algorithm 2. We prove the following
229 lemmas to show that our definition of the $spell(\cdot)$ function is sound for our purposes, i.e.
230 correctly spells the string belonging to a walk in a de Bruijn graph.

231 ► **Lemma 3** (Sound sequence of k -mers). *Let k be a positive integer and let I be a set of*
232 *strings of length at least k . Let $G = (V, E, c)$ be the de Bruijn graph of order k constructed*
233 *from I . Let $W := (e_1 = (v_1d_1, v'_1d'_1, \eta_1), \dots, e_\ell = (v_\ell d_\ell, v'_\ell d'_\ell, \eta_\ell))$ be a walk in G , and*
234 *$K := (\kappa_1, \dots, \kappa_\ell)$ its sequence of k -mers. Then for each consecutive pair of kmers κ_i, κ_{i+1} it*
235 *holds that $suf_{k-1}(\kappa_i) = pre_{k-1}(\kappa_{i+1})$.*

236 We define the *sequence of k -mers* $K = (\kappa_1, \dots, \kappa_\ell)$ of a string $w = (a_1, \dots, a_{\ell+k-1})$ by
237 $\kappa_i := (a_i, \dots, a_{i+k-1})$ for each i .

■ **Algorithm 2** SPELL

Input: A de Bruijn graph $G = (V, E, c)$ of order k and a walk
 $W = (e_1 := (v_1d_1, v'_1d'_1, \eta_1), \dots, e_\ell := (v_\ell d_\ell, v'_\ell d'_\ell, \eta_\ell))$.
Output: The string s spelled by W , i.e. $spell(W)$.

```

1 if  $W$  is empty then
2    $s \leftarrow \epsilon$ 
3 else /* compute the sequence of kmers from  $W$  */
4   foreach  $i \in (1, \dots, \ell)$  do /* iterate the sequence in order */
5      $\lfloor$  if  $e_i \in E$  then  $\kappa_i \leftarrow \eta_i$ 
6      $\lfloor$  else  $\kappa_i \leftarrow rc(\eta_i)$  /*  $e_i^{-1} \in E$  */
7    $s \leftarrow k - 1$  prefix of  $\kappa_1$ 
8   foreach  $i \in (1, \dots, \ell)$  do /* iterate the sequence in order */
9      $\lfloor$  append the last character from  $\kappa_i$  to  $s$ 

```

2:8 Eulertigs: minimum plain text k -mer sets without repetitions in linear time

238 ► **Lemma 4** (Sound spell). *Let k be a positive integer and let I be a set of strings of length*
 239 *at least k . Let $G = (V, E, c)$ be the de Bruijn graph of order k constructed from I . Let W be*
 240 *a walk in G , K_W its sequence of k -mers and K'_W the sequence of k -mers of $\text{spell}(W)$. Then*
 241 $K_W = K'_W$.

242 ► **Lemma 5** (Complete representation). *Let k be a positive integer and let I be a set of*
 243 *strings of length at least k . Let $G = (V, E, c)$ be the de Bruijn graph of order k constructed*
 244 *from I . Let w be a string with $\text{CS}_k(w) \subseteq \text{CS}_k(I)$. Then there exists a walk W in G with*
 245 $\text{spell}(W) = w$.

246 **Proof.** Let $K_w = (\kappa_1, \dots, \kappa_\ell)$ be the sequence of k -mers of w . We construct $W = (e_1 =$
 247 $(v_1 d_1, v'_1 d'_1, \eta_1), \dots, e_\ell = (v_\ell d_\ell, v'_\ell d'_\ell, \eta_\ell))$ as follows: for each i , let η_i be the canonical of κ_i
 248 and $f_i \in E$ be the arc whose identifier is η_i . We set $e_i = f_i$ if κ_i is canonical, and $e_i = f_i^{-1}$
 249 otherwise.

250 For W to fulfil the definition of a walk we need that $v'_i = v_{i+1}$ and $d'_i = -d'_{i+1}$ for all i .
 251 Using Lemma 2, we get:

- 252 ■ If $e_i, e_{i+1} \in E$, then $\text{suf}_{k-1}(\eta_i) = \text{suf}_{k-1}(\kappa_i) = \text{pre}_{k-1}(\kappa_{i+1}) = \text{pre}_{k-1}(\eta_{i+1})$. Therefore,
 253 by Lemma 2 (a), it holds that $v'_i = v_{i+1}$ and $d'_i = -d'_{i+1}$.
- 254 ■ If $e_i, e_{i+1}^{-1} \in E$, then $\text{suf}_{k-1}(\eta_i) = \text{suf}_{k-1}(\kappa_i) = \text{pre}_{k-1}(\kappa_{i+1}) = \text{pre}_{k-1}(\text{rc}(\eta_{i+1}))$. There-
 255 fore, by Lemma 2 (b), it holds that $v'_i = v_{i+1}$ and $d'_i = -d'_{i+1}$.
- 256 ■ If $e_i^{-1}, e_{i+1} \in E$, then $\text{suf}_{k-1}(\text{rc}(\eta_i)) = \text{suf}_{k-1}(\kappa_i) = \text{pre}_{k-1}(\kappa_{i+1}) = \text{pre}_{k-1}(\eta_{i+1})$.
 257 Therefore, by Lemma 2 (c), it holds that $v'_i = v_{i+1}$ and $d'_i = -d'_{i+1}$.
- 258 ■ If $e_i^{-1}, e_{i+1}^{-1} \in E$, then $\text{suf}_{k-1}(\text{rc}(\eta_i)) = \text{suf}_{k-1}(\kappa_i) = \text{pre}_{k-1}(\kappa_{i+1}) = \text{pre}_{k-1}(\text{rc}(\eta_{i+1}))$.
 259 Therefore, by Lemma 2 (d), it holds that $v'_i = v_{i+1}$ and $d'_i = -d'_{i+1}$.

260 To complete the proof we need to show that $\text{spell}(W) = w$. By definition, the sequence
 261 of k -mers K_W of W is equivalent to K_w . And since W is a walk, by Lemma 4 we get that
 262 the sequence of k -mers of $\text{spell}(W)$ is equivalent to K_W , and therefore $\text{spell}(W) = w$. ◀

263 A *walk cover* \mathcal{W} of a bigraph G is a set of walks such that for each arc $e \in E$ it holds that e
 264 is part of some walk $W \in \mathcal{W}$, or e^{-1} is part of some walk $W \in \mathcal{W}$.

265 ► **Theorem 6** (Dualism between SPSS and walk cover). *Let k be a positive integer and let*
 266 *I and I' be sets of strings of length at least k . Let $G = (V, E, c)$ be the de Bruijn graph of*
 267 *order k constructed from I . Then it holds that $\text{CS}_k(I) = \text{CS}_k(I')$, if and only if there is a*
 268 *walk cover \mathcal{W} in G that spells the strings in I' .*

269 **Proof.** If $\text{CS}_k(I') \subseteq \text{CS}_k(I)$, then for each string $w' \in I'$ it holds that $\text{CS}_k(w') \subseteq \text{CS}_k(I)$.
 270 Therefore, by Lemma 5, there exists a walk w in G with $\text{spell}(w) = w'$. Then, the set of all
 271 such walks \mathcal{W} spells I' . Further, because $\text{CS}_k(I) \subseteq \text{CS}_k(I')$, the identifier η of each arc $e \in E$
 272 is in $\text{CS}_k(I')$, and therefore in the sequence of kmers $K_{w'}$ of some string $w' \in I'$ (possibly as
 273 a reverse complement). By Lemma 4 it holds that $K_{w'} = K_w$, where K_w is the sequence of
 274 k -mers of walk w . By the definition of the sequence of k -mers of a walk, this implies that w
 275 visits e (possible in reverse direction). Since this holds for each $e \in E$, it holds that \mathcal{W} is a
 276 walk cover of G .

277 Assume that there is a walk cover \mathcal{W} in G that spells the strings in I' , and let $w \in \mathcal{W}$ be a
 278 walk, K_w its sequence of k -mers, $w' := \text{spell}(w)$ and $K_{w'}$ the sequence of k -mers of w' . Then,
 279 by Lemma 4, $K_w = K_{w'}$, which, by the definition of the sequence of k -mers of a walk implies
 280 that $\text{CS}_k(I) \subseteq \text{CS}_k(I')$. And since \mathcal{W} is a walk cover of G , we get $\text{CS}_k(I) = \text{CS}_k(I')$. ◀

281 ► **Corollary 7.** *By setting $I = I'$ in Theorem 6 we can confirm that our definition of a*
282 *de Bruijn graph is sound in that there is a set of walks that spells the strings used for its*
283 *construction.*

284 A compacted de Bruijn graph is constructed from a de Bruijn graph by contracting all nodes
285 $v \in V$ that are either self-complemental and have exactly two arcs that have exactly one
286 incidence to v each, or that are not self-complemental and have exactly one incoming and
287 one outgoing arc. For simplicity, we use uncompact de Bruijn graphs in our theoretical
288 sections, however all results equally apply to compacted de Bruijn graphs.

289 4 Linear-time construction of compacted bidirected de Bruijn graphs

290 In this section, we fill a gap in the literature by describing on a high level an algorithm to
291 construct the bidirectional de Bruijn graph of a set of input strings in time linear in the total
292 length of the input strings, independent of the value of k .

293 4.1 Algorithm

294 Let $I = \{w_1, \dots, w_m\}$ be the set of input strings. Consider the following concatenation:

$$T = w_1\$w_2\$ \dots \$w_m\$rc(w_1)\$rc(w_2)\$ \dots \$rc(w_m)\$,$$

295 where $\$$ is a special character outside of the alphabet Σ of the input strings. We require an
296 index on T that can answer the following queries: `extendRight`, `extendLeft`, `contractRight`
297 and `contractLeft` in constant time. The extension operations take as input a character
298 $c \in \Sigma$ and the interval of a string x in the suffix array of T , and return the suffix array
299 intervals of xc in the case of `extendRight` and cx in the case of `extendLeft`. The contraction
300 operations are the inverse operations of these, mapping the suffix array intervals of xc to
301 x in the case of `contractRight` and cx to x in the case of `contractLeft`. For efficiency,
302 we also require operations `enumerateRight` and `enumerateLeft`, which take a string x and
303 give all characters such that `extendRight` and `extendLeft` respectively return a non-empty
304 interval, in time that is linear in the number of such characters. Implementations for all the
305 six subroutines are given in Section 4.2.

306 Using these operations, we can simulate the regular non-bidirected de Bruijn graph of
307 T . Each k -mer of the input strings for a fixed k corresponds to a disjoint interval in the
308 suffix array of T . The nodes are represented by their suffix array intervals. The outgoing
309 arcs from a $(k-1)$ -mer x are those characters c where `extendRight`(x, c) returns a non-
310 empty interval. We can enumerate all the characters c with this property in constant time
311 using `enumerateRight`(x). The incoming arcs can be enumerated symmetrically with the
312 `enumerateLeft`(x). Finally, we can find the destination or origin of an arc labelled with x
313 by running a `contractLeft` or `contractRight` operation respectively on x .

314 To construct the bidirected de Bruijn graph, we merge together nodes that are the reverse
315 complement of each other. To find which nodes are complementary, we scan the input strings
316 I while maintaining the suffix array interval of the current k -mer using `extendRight` and
317 `contractLeft` operations, while at the same time maintaining the suffix array interval of
318 the reverse complement using `extendLeft` and `contractRight` operations. Whenever we
319 merge two nodes, we combine the incoming and outgoing arcs, assigning the incidences of the
320 arcs according to the incidence rules in our definition. We are able to tell in constant time
321 which k -mer of a pair of complementary k -mers is canonical by comparing the suffix array

2:10 Eulertigs: minimum plain text k -mer sets without repetitions in linear time

322 intervals of the k -mers: the k -mer whose suffix array interval has a smaller starting point is
323 the canonical k -mer. If the starting points are the same, the k -mer is self-complemental.

324 Using the `enumerateRight` and `enumerateLeft` functions, we can check if a node would
325 be contracted in a compacted de Bruijn graph. By extending k -mers over such nodes, we
326 can in linear time also output only the arcs and nodes of a compacted de Bruijn graph. For
327 storing the labels, we use one pointer into the input strings to store a single k -mer, as well
328 as a flag that is set whenever the label is not canonical. If a label has multiple k -mers, then
329 we store the remaining k -mers as explicit strings, however without their overlap with the
330 “pointer- k -mer”. This way, we can store each label in $O(\ell)$ space, where ℓ is the number of
331 k -mers in the label. We additionally store the first and last character of each label, as an
332 easy way to make the SPELL function run in output sensitive linear time.

333 4.2 Implementation of the subroutines

334 All required the subroutines `extendRight`, `extendLeft`, `contractRight`, `contractLeft`,
335 `enumerateRight` and `enumerateLeft` can be implemented with the suffix tree of T by
336 simulating the trie of the suffix tree loci as described in Section 2.2. The suffix array intervals
337 of explicit nodes can be stored with the nodes, so that we can operate on loci (u, d) and
338 retrieve the suffix array intervals on demand. The operation `extendRight` follows an arc
339 from a locus to a child, and the operation `contractRight` is implemented by going to the
340 parent of the current locus. The operation `contractLeft` follows a suffix link from the
341 current locus, and `extendLeft` follows a Weiner link. The operations `enumerateRight` and
342 `enumerateLeft` are implemented by storing the children and the Weiner links from explicit
343 suffix tree nodes as neighbor lists. The total number of these links is linear in $|T|$ [18].
344 With this implementation, the slowest operations are `extendRight` and `extendLeft`, taking
345 $O(\log |\Sigma|)$ time to binary search the neighbor lists. We therefore obtain the following result:

346 ► **Theorem 8.** *The compacted arc-centric bidirected de Bruijn graph of order k of a set of*
347 *input strings I from the alphabet Σ can be constructed in time $O(|I| \log |\Sigma|)$.*

348 We note that the same operations can also be implemented on top of the bidirectional BWT
349 index of Belazzougui and Cunial [2], using the data structures of Belazzougui et al. [3] for
350 the enumeration operations. This gives an index that supports all the required subroutines
351 in *constant time*. The drawback of the bidirectional BWT index is that only randomized
352 construction algorithms are known, but the expected time is still linear in $|T|$. We leave as
353 an open problem the construction of the compacted arc-centric bidirected de Bruijn graph in
354 deterministic linear time independent of the alphabet size.

355 5 Linear-time minimum SPSS without repetitions

356 Let I be a set of strings. To compute an SPSS without repetitions we first build a compacted
357 de Bruijn graph G from I . Because of Theorem 6, finding an SPSS is equivalent to finding
358 a walk cover in G . Further, with Lemma 4, we get that an SPSS without repetitions is
359 equivalent to a walk cover that visits each arc exactly once (either once forwards, or once
360 reverse, but not both forwards and reverse). We call such a walk cover a *unique walk cover*.

361 For minimality, observe that the cumulative length of an SPSS S relates to its equivalent
362 set of walks \mathcal{W} as follows:

$$363 \quad ||S|| = \sum_{W \in \mathcal{W}} (k - 1 + |W|) \quad (1)$$

364

365 This is because in Algorithm 2, in Line 7, $k - 1$ characters are appended to the result, and
366 then in the loop in Line 8, one additional character per arc in W is appended. We cannot
367 alter the sum $\sum_{W \in \mathcal{W}} |W|$, since we need to cover all arcs in G . However we can alter the
368 number of strings, and decreasing or increasing this number by one will decrease or increase
369 the cumulative length of S by $k - 1$. Therefore, finding a minimum SPSS of I without
370 repetitions equals finding a unique walk cover of G that has a minimum number of walks.

371 Note that computing a minimum SPSS in a bigraph that is not connected is equivalent to
372 separately computing an SPSS in each maximal connected subgraph. Therefore we restrict
373 to connected bigraphs from here on.

374 5.1 A lower bound for an SPSS without repetitions

375 Using the imbalance of the nodes of a bigraph, we can derive a lower bound for the number
376 of walks in a walk cover.

377 ► **Lemma 9.** *Let $v \in V$ be an unbalanced node in a bigraph $G = (V, E, c)$. Then in a unique*
378 *walk cover \mathcal{W} of G , either at least $|\text{imbalance}(v)|$ walks start in v , or at least $|\text{imbalance}(v)|$*
379 *walks end in v .*

380 **Proof.** If v is self-complemental, then its imbalance is 1, so by definition v has an odd number
381 of incident arcs. Each walk that does not start or end in v needs to enter and leave v via
382 two distinct arcs whenever it visits v . But since the number of incident arcs is odd, there is
383 at least one arc that cannot be covered this way, implying that a walk needs to start or end
384 in this arc.

385 If v is not self-complemental and has a positive imbalance, then it has $\text{imbalance}(v)$ more
386 outgoing arcs than incoming arcs. Since walks need to leave v with the opposite sign than
387 they entered v , at least $\text{imbalance}(v)$ arcs cannot be covered by walks that do not start or
388 end in v . If v has negative imbalance, the situation is symmetric. ◀

389 ► **Definition 10 (Imbalance of a bigraph).** *The imbalance $\text{imbalance}(G)$ of a bigraph $G =$*
390 *(V, E, c) is the sum of the absolute imbalance of all nodes $\sum_{v \in V} |\text{imbalance}(v)|$.*

391 ► **Theorem 11 (Lower bound).** *Let G be a bigraph. A walk cover \mathcal{W} of G has a minimum*
392 *string count of $\text{imbalance}(G)/2$.*

393 **Proof.** Let $v \in V$ be an unbalanced node. Then, by Lemma 9 at least $|\text{imbalance}(v)|$ walks
394 start in v or at least $|\text{imbalance}(v)|$ walks end in v . Since each walk has exactly one start
395 node and one end node, \mathcal{W} has a minimum string count of $\text{imbalance}(G)/2$. ◀

396 5.2 Eulerising a bigraph

397 A directed graph is called *Eulerian*, if all nodes have indegree equal to outdegree, i.e. are
398 balanced [10]. If the graph is strongly connected³, then this is equivalent to the graph
399 admitting a *Eulerian cycle*, i.e. a cycle that visits each arc exactly once. The same notion
400 can be used with bidirected graphs, using our definition of imbalance.

401 ► **Definition 12 (Eulerian bigraph).** *A bigraph is Eulerian, if all nodes have imbalance zero.*

402 A connected bigraph can be transformed into a Eulerian bigraph by adding arcs using
403 Algorithm 3. See Figure 1 (b) for an example.

³ Strongly connected means that there is a directed path from each node v_1 to each node v_2 .

2:12 EulerTigs: minimum plain text k -mer sets without repetitions in linear time

Algorithm 3 EULERISE

Input: Bigraph $G = (V, E, c)$.
Output: Eulerised bigraph $G' = (V, E', c)$.

```

1  $G' \leftarrow G$  /*  $G$  and  $G'$  share  $V$  and  $c$  */
2  $L \leftarrow$  empty list /* collect missing incidences to balance  $G'$  */
3 foreach  $v \in V$  do
4    $i \leftarrow$  imbalance( $v$ )
5   if  $c(v) = 1$  then
6     if  $i \neq 0$  then append  $v \ominus$  to  $L$ 
7   else
8     if  $i > 0$  then append  $i$  copies of  $v \ominus$  to  $L$ 
9     if  $i < 0$  then append  $i$  copies of  $v \oplus$  to  $L$ 
10 while  $|L| > 0$  do /* insert missing incidences as arcs */
11    $vd \leftarrow$  remove the first incidence from  $L$ 
12    $v'd' \leftarrow$  remove the first incidence from  $L$ 
13   insert 1 arc  $(vd, v'd', |L|)$  into  $E'$  /* use distinct identifiers */

```

404 ► **Lemma 13.** *The imbalance of a bigraph is even.*

405 **Proof.** Adding or removing an arc changes the imbalance of two nodes by 1, or of one node
406 by two. In both cases, the imbalance of the graph can only change by -2 , 0 , or 2 . Since the
407 imbalance of a graph without arcs is 0 , this implies that there can be no graph with odd
408 imbalance. ◀

409 ► **Lemma 14.** *Given a connected bigraph $G = (V, E, c)$, Algorithm 3 outputs a Eulerian
410 bigraph $G' = (V, E', c)$.*

411 **Proof.** Algorithm 3 is well-defined, since by Lemma 13, it holds that L has even length
412 in each iteration of the loop in Line 10, so the removal operation in Line 12 always has
413 something to remove.

414 The output of Algorithm 3 is a valid bigraph, since for self-complemental nodes $v \in V$,
415 only incidences $v \ominus$ are added to G' , and for not self-complemental nodes $v \in V$, only
416 incidences $v \oplus$ and $v \ominus$ are added to G' .

417 Further, the output is a Eulerian bigraph, because for all $v \in V$, it holds that imbalance(v)
418 is 0 , by the following argument:

- 419 ■ If $c(v) = 1$ and v has imbalance zero in G , then its imbalance stays the same in G' . If
420 it has imbalance 1 , then one incident arc is inserted, making its degree even and its
421 imbalance therefore zero.
- 422 ■ If $c(v) = 0$ and v has positive imbalance i in G , then i incoming arcs are added to v
423 (counting incoming self-loops twice), and no outgoing arcs are added. Therefore, it has
424 imbalance zero in G' . By symmetry, if v has negative imbalance in G , it has imbalance
425 zero in G' . ◀

426 ► **Lemma 15.** *Given a bigraph $G = (V, E, c)$, Algorithm 3 terminates after $O(|V| + |E|)$
427 steps.*

428 **Proof.** For the list data structure we choose a doubly linked list, and for the graph an
429 adjacency list (and array with an entry for each node containing a doubly linked list for the
430 arcs).

431 The loop in Line 3 runs $|V|$ times and each iteration runs in $O(|\text{imbalance}(v)|)$ for a node
 432 v , because a doubly linked list supports appending in constant time. The sum of absolute
 433 imbalances of all nodes cannot exceed $2|E|$, because each arc adds at most 1 to the absolute
 434 imbalance of at most two nodes, or adds at most 2 to the absolute imbalance of at most one
 435 node. Therefore, the length of list L after completing the loop is at most $2|E|$, and the loop
 436 runs in $O(|V| + |E|)$ time.

437 The loop in Line 10 runs at most $|L| \leq 2|E|$ times and performs only constant-time
 438 operations, since L is a doubly linked list and we can insert arcs into an adjacency list in
 439 constant time. Therefore, this loop also runs in $O(|V| + |E|)$ time. ◀

440 With Lemmas 14 and 15 we get the following.

441 ▶ **Theorem 16.** *Algorithm 3 is correct and runs in $O(|V| + |E|)$ time.*

442 5.3 Computing a Eulerian cycle in a bigraph

■ Algorithm 4 EULERIANCYCLE

Input: Connected Eulerian bigraph $G = (V, E, c)$.
Output: Eulerian cycle W .

```

1 while  $|E| > 0$  do
2   if  $|W| = 0$  then
3      $(vd, v'd', \eta) \leftarrow$  remove some arc from  $E$ 
4      $W' \leftarrow ((vd, v'd', \eta))$  /* doubly linked list */
5   else /* search a used arc that connects to an unused arc */
6      $(vd, v'd', \eta) \leftarrow$  dereference  $first\_unfinished$ 
7     while  $E$  has no arc with incidence  $v'-d'$  do
8       advance  $first\_unfinished$  to the next arc in  $W$ 
9        $(vd, v'd', \eta) \leftarrow$  dereference  $first\_unfinished$ 
10    // extend  $W'$  without repeating arcs until it closes a cycle
11    while  $E$  contains an arc  $e = (v_e d_e, v'_e d'_e, \eta_e)$  with incidence  $v'-d'$  do
12      remove  $e$  from  $E$ 
13      if  $v_e d_e = v'-d'$  then  $(vd, v'd', \eta) \leftarrow (v_e d_e, v'_e d'_e, \eta_e)$ 
14      else  $(vd, v'd', \eta) \leftarrow (v'_e d'_e, v_e d_e, \eta_e)$  /*  $v'_e d'_e = v'-d'$  */
15      append  $(vd, v'd', \eta)$  to  $W'$ 
16    if  $|W| = 0$  then
17       $W \leftarrow W'$ 
18       $first\_unfinished \leftarrow$  pointer to the first arc in  $W$ 
19    else
20      insert  $W'$  after  $first\_unfinished$  in  $W$ 
21       $W' \leftarrow ()$  /* empty doubly linked list */

```

443 After Eulerising the bigraph, we can compute a Eulerian cycle using Algorithm 4. We
 444 do this similarly to Hierholzer's classic algorithm for Eulerian cycles [10]. First we find an
 445 arbitrary cycle. Then, as long as there are unused arcs left, we search along the current cycle
 446 for unused arcs, and find additional cycles through such unused arcs. We integrate each of
 447 those additional cycles into the main cycle. See Figure 1 (c) for an example of a Eulerian
 448 cycle.

449 ▶ **Lemma 17.** *Given a connected Eulerian bigraph $G = (V, E, c)$, Algorithm 4 terminates*
 450 *and outputs a Eulerian cycle W .*

2:14 Eulertigs: minimum plain text k -mer sets without repetitions in linear time

451 **Proof.** For $W = (e_1 = (v_1d_1, v'_1d'_1, \eta_1), \dots, e_\ell = (v_\ell d_\ell, v'_\ell d'_\ell, \eta_\ell))$ to be a Eulerian cycle, it
452 must be a cycle that contains each arc exactly once.

453 The sequence W' constructed by the loop in Line 10 is a walk by construction, and since
454 G is Eulerian it is a cycle after the loop terminates. After finding the initial cycle in the first
455 iteration of the outer loop, each additional cycle is started from a node on the initial cycle,
456 and is a cycle again. Therefore it can be inserted into the original cycle without breaking its
457 cycle property.

458 Since each arc is deleted when being added to W' , there is no duplicate arc in W . And if
459 the algorithm terminates, then $|E| = 0$ (Line 1), so W contains all arcs.

460 For termination, consider that if W is not complete after the first iteration of the outer
461 loop, then the loop in Line 7 searches for an unused arc using the *first_unfinished* pointer.
462 Since the prefix of W up to including *first_unfinished* is never modified (Line 19), and
463 *first_unfinished* is only advanced when its pointee cannot reach any arc anymore, it holds
464 that no arc in W can reach an arc in E when *first_unfinished* gets advanced over the end
465 of W . Since G was initially Eulerian and only Eulerian cycles have been removed from G ,
466 this implies that all nodes visited by W are still balanced and therefore have no incident
467 arcs anymore. And since G was originally connected, W has visited all nodes, i.e. $|E| = 0$.
468 Therefore, *first_unfinished* cannot be advanced over the end of W , because the outer loop
469 terminates before that.

470 To complete the proof of termination, consider that in each iteration of the outer loop, at
471 least one arc gets removed from E . In the first iteration, this happens at least in Line 3, and
472 in all following iterations, this happens in Line 11. ◀

473 ▶ **Lemma 18.** *Given a connected Eulerian bigraph $G = (V, E, c)$, Algorithm 4 terminates*
474 *after $O(|V| + |E|)$ steps.*

475 **Proof.** We use a doubly linked list for W and W' , and an adjacency list for G . Then all
476 lines can be executed in constant time.

477 The loop in Line 10 removes one arc from E each iteration, so it runs at most $|E|$ times in
478 total (over all iterations of the outer loop). The loop in Line 7 advances *first_unfinished*
479 each iteration. Since the algorithm is correct by Lemma 17, $|W| \leq |E|$ and *first_unfinished*
480 never runs over the end of *first_unfinished*, so the loop runs at most $|E|$ times in total
481 (over all iterations of the outer loop).

482 The condition for the loop in Line 10 is true at least once in each iteration of the outer
483 loop, since the preceding branch sets up $(vd, v'd', \eta)$ such that it has a successor (in the first
484 iteration because of Eulerianess). So in each iteration of the outer loop, at least one arc gets
485 removed, so the outer loop runs at most $|E|$ times in total.

486 As a result, all loops individually run at most $|E|$ times, therefore Algorithm 4 terminates
487 after $O(|V| + |E|)$ steps. ◀

488 With Lemmas 17 and 18 we get the following.

489 ▶ **Theorem 19.** *Algorithm 4 is correct and runs in $O(|V| + |E|)$ time.*

490 5.4 Computing a minimum SPSS without repetitions

491 We convert the Eulerian cycle into a walk cover of the original bigraph by breaking it at all
492 arcs inserted by Algorithm 3, and removing those arcs (see Figure 1 (d) for an example).
493 This results in a walk cover with either one walk, if Algorithm 3 inserted zero or one
494 arcs, or $\text{imbalance}(G)/2$ arcs, if Algorithm 3 inserted more arcs. By Theorem 11, this is a

495 minimum number of walks, and therefore the SPSS spelled by these walks is minimum as
 496 well. Constructing the de Bruijn graph takes $O(|I| \log \Sigma)$ time, and it has $O(|I|)$ k -mers,
 497 so it holds that $|V| \in O(|I|)$ and $|E| \in O(|I|)$. Further, spelling the walk cover takes
 498 time linear to the cumulative length of the spelled strings. Since we compute a minimum
 499 representation, it holds that the output is not larger than the total length of the input strings.
 500 Therefore we get:

501 ► **Theorem 1.** *Let k be a positive integer and let I be a set of strings of length at least k*
 502 *over some alphabet Σ . Then we can compute a set of strings I' of length at least k with*
 503 *minimum cumulative length and $CS_k(I) = CS_k(I')$ in $O(|I| \log |\Sigma|)$ time.*

504 6 Experiments

genome	algorithm	CL ratio	SC ratio	time [s]	memory [GiB]
C. elegans (reads)	unitigs	1.789	2.831	1888	5.97
	UST	1.035	1.080	2738 (1.45)	15.2 (2.54)
	Eulertigs	1	1	3735 (1.98)	25.0 (4.19)
B. mori (reads)	unitigs	1.912	3.136	7737	9.36
	UST	1.050	1.118	10937 (1.41)	52.4 (5.60)
	Eulertigs	1	1	13793 (1.78)	79.4 (8.48)
H. sapiens (reads)	unitigs	1.418	2.143	56966	13.0
	UST	1.016	1.044	57736 (1.01)	16.4 (1.26)
	Eulertigs	1	1	58861 (1.03)	29.2 (2.25)
C. elegans	unitigs	1.060	3.154	54.7	1.22
	UST	1.002	1.089	58.0 (1.06)	1.22 (1.00)
	Eulertigs	1	1	65.9 (1.21)	1.22 (1.00)
B. mori	unitigs	1.262	3.310	224	3.32
	UST	1.018	1.156	258 (1.16)	3.32 (1.00)
	Eulertigs	1	1	315 (1.41)	3.32 (1.00)
H. sapiens	unitigs	1.195	3.532	3166	10.0
	UST	1.015	1.192	3369 (1.06)	10.0 (1.00)
	Eulertigs	1	1	3717 (1.17)	10.0 (1.00)

■ **Table 1** Experiments on references and read sets of single genomes with $k = 51$ and a min abundance of 10 for human and 1 for the others. The CL and SC ratios are compared to the CL-optimal Eulertigs. For time and memory, we report the total time and maximum memory required to compute the tigs from the respective data set. BCALM2 directly computes unitigs, while UST- and Eulertigs require a run of BCALM2 first before they can be computed themselves. Prophasm can only be run for $k \leq 32$, which does not make sense for large genomes. The number in parentheses behind time and memory indicates the slowdown/increase over computing just unitigs with BCALM2. BCALM2 was run with 28 threads, while all other tools support only one thread. The lengths of the genomes are 100Mbp for C. elegans, 482Mbp for B. mori and 3.21Gbp for H. sapiens and the read data sets have a coverage of 64x for C. elegans, 58x for B. mori and 300x for H. sapiens.

505 We ran our experiments on a server running Linux with two 64-core AMD EPYC 7H12
 506 processors with 2 logical cores per physical core, 1.96TiB RAM and an SSD. Our data sets

2:16 Eulertigs: minimum plain text k -mer sets without repetitions in linear time

pangenome	tigs	CL ratio	SC ratio	time [s]	memory [MiB]
1102x <i>N. gonorrhoeae</i>	unitigs	1.615	3.052	24.2	4328
	UST	1.022	1.074	26.1 (1.08)	4328 (1.00)
	prophasm	1.00004	1.00013	774 (31.9)	208 (0.05)
	Eulertigs	1	1	26.9 (1.11)	4328 (1.00)
616x <i>S. pneumoniae</i>	unitigs	1.679	3.055	21.4	3135
	UST	1.027	1.081	25.9 (1.21)	3135 (1.00)
	prophasm	1.00004	1.00012	436 (20.3)	434 (0.14)
	Eulertigs	1	1	28.5 (1.33)	3135 (1.00)
3682x <i>E. coli</i>	unitigs	1.705	3.092	334	7146
	UST	1.031	1.092	416 (1.24)	7146 (1.00)
	prophasm	1.00008	1.00023	7456 (22.3)	7221 (1.01)
	Eulertigs	1	1	471 (1.41)	7146 (1.00)

■ **Table 2** Experiments on (references of) pangenomes with $k = 31$ and a min abundance of 1. The CL and SC ratios are compared to the CL-optimal Eulertigs. For time and memory, we report the total time and maximum memory required to compute the tigs from the respective data set. BCALM2 directly computes unitigs, while UST- and Eulertigs require a run of BCALM2 first before they can be computed themselves. Prophasm is run directly on the source data. The number in parentheses behind time and memory indicates the slowdown/increase over computing just unitigs with BCALM2. BCALM2 was run with 28 threads, while all other tools support only one thread. The *N. gonorrhoeae* pangenome contains 8.36 million unique kmers, the *S. pneumoniae* pangenome contains 19.3 million unique kmers and the *E. coli* pangenome contains 341 million unique kmers.

507 are the same as in [23], and we also adapted their metrics *cumulative length* (CL), which is
 508 the total count of characters in all strings, and *string count* (SC), which is the number of
 509 strings. Our implementation does not use the formalisation of bidirected graphs introduced
 510 in this work, but instead uses the formalisation from [23]. For constructing de Bruijn graphs,
 511 we do not implement our purely theoretical linear time algorithm, since practical de Bruijn
 512 graph construction is a well-researched field [8, 13, 15, 9, 20, 19], and we want to focus
 513 more on computing the compressed representation from unitigs. UST only supports unitigs
 514 constructed by BCALM2 [8], since it needs certain additional data. BCALM2 is not a linear
 515 time algorithm, but works efficient in practice. Therefore, we use BCALM2 to construct
 516 a node-centric de Bruijn graph, and then convert it to an arc-centric variant using a hash
 517 table.

518 Our experimental pipeline is constructed with [16] and using the bioconda software
 519 repository [11]. We ran all multithreaded tools with up to 28 threads and never used more
 520 than 128 cores of our machine at once to prevent hyperthreading from affecting our timing.
 521 The code to reproduce our experiments is available at [https://doi.org/10.5281/zenodo.](https://doi.org/10.5281/zenodo.6538261)
 522 **6538261**.

523 The performance figures are all very similar, with two exceptions. Prophasm does not
 524 support parallel computation at the moment, therefore its runtime is much higher. Compared
 525 to that, all other algorithms use parallel computation to compute unitigs, but computing
 526 the final tigs from unitigs seems to be negligible compared to computing the de Bruijn
 527 topology. Moreover, running UST or Eulertigs on read data sets of larger genomes consumes
 528 significantly more memory than computing just unitigs. This is likely because BCALM2
 529 uses external memory to compute unitigs, while the other tools simply load the whole set of

530 unitigs into memory.

531 It is notable that the Eulertigs algorithm is always slower than UST. This may be because
532 of the Eulertig algorithm being more complex, but also because our loading and storing
533 routines might not be as efficient. While UST uses node-centric de Bruijn graphs and can
534 therefore directly make use of the topology output by BCALM2 (which is a fasta file with
535 arcs stored as custom annotations), we need to convert the graph into arc-centric format.
536 This is supported by e.g. the B. mori short read data set, on which the computation of
537 Eulertigs uses only 11% of the runtime for the algorithm itself, while 89% are from loading
538 the graph (including the conversion to arc-centric) and storing the result.

539 In terms of CL, we see that the SPSS computed with UST mostly remains within the
540 expected 3% of the lower bound, but they are up to 5% above the lower bound on more
541 compressible data sets. The SPSS computed by prophase is very close to the optimum in
542 all cases, and we assume that this difference in quality is because prophase extends paths
543 both forwards and backwards, while the UST heuristic merely extends them forwards.

544 Looking at SC, we see that Eulertigs are always the lowest, which is due to the string
545 count directly being connected to the cumulative length by Equation (1). This also explains
546 the correlation between CL and SC, which can be observed in all cases.

547 **7** Conclusions

548 We have presented a linear and hence optimal algorithm for computing a minimum SPSS
549 without repetitions for a fixed alphabet size. This closes the open question about its
550 complexity raised in [6, 22]. Using our optimal algorithm, we were able to accurately evaluate
551 the existing heuristics and show that they are very close to the optimum in practice. Further,
552 we have published our algorithm as a command-line tool on github, allowing it to easily be
553 used in other projects.

554 Further, we have presented how bidirected de Bruijn graphs can be formalised without
555 excluding any corner cases. We have also shown how such a graph can be constructed in
556 linear time for a fixed-size alphabet. The construction of the compacted arc-centric bidirected
557 de Bruijn graph in linear time independent of the alphabet size stays an open problem.

558 ——— References ———

- 559 **1** Anton Bankevich, Andrey V Bzikadze, Mikhail Kolmogorov, Dmitry Antipov, and Pavel A
560 Pevzner. Multiplex de bruijn graphs enable genome assembly from long, high-fidelity reads.
561 *Nature biotechnology*, pages 1–7, 2022.
- 562 **2** Djamal Belazzougui and Fabio Cunial. Fully-functional bidirectional burrows-wheeler indexes
563 and infinite-order de bruijn graphs. In *30th Annual Symposium on Combinatorial Pattern
564 Matching (CPM 2019)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019.
- 565 **3** Djamal Belazzougui, Fabio Cunial, Juha Kärkkäinen, and Veli Mäkinen. Versatile succinct
566 representations of the bidirectional burrows-wheeler transform. In *European Symposium on
567 Algorithms*, pages 133–144. Springer, 2013.
- 568 **4** Djamal Belazzougui, Dmitry Kosolobov, Simon J Puglisi, and Rajeev Raman. Weighted
569 ancestors in suffix trees revisited. In *32nd Annual Symposium on Combinatorial Pattern
570 Matching*, 2021.
- 571 **5** Jacek Blazewicz, Alain Hertz, Daniel Kobler, and Dominique de Werra. On some properties
572 of dna graphs. *Discrete Applied Mathematics*, 98(1-2):1–19, 1999.
- 573 **6** Karel Břinda, Michael Baym, and Gregory Kucherov. Simplitigs as an efficient and scalable
574 representation of de Bruijn graphs. *Genome Biology*, 22(1):1–24, 2021.

2:18 Eulertigs: minimum plain text k -mer sets without repetitions in linear time

- 575 7 Bastien Cazaux, Thierry Lecroq, and Eric Rivals. From indexing data structures to de Bruijn
576 graphs. In *Symposium on combinatorial pattern matching*, pages 89–99. Springer, 2014.
- 577 8 Rayan Chikhi, Antoine Limasset, and Paul Medvedev. Compacting de bruijn graphs from
578 sequencing data quickly and in low memory. *Bioinformatics*, 32(12):i201–i208, 2016.
- 579 9 Victoria G Crawford, Alan Kuhnle, Christina Boucher, Rayan Chikhi, and Travis Gagie.
580 Practical dynamic de bruijn graphs. *Bioinformatics*, 34(24):4189–4195, 2018.
- 581 10 Herbert Fleischner. *Eulerian graphs and related topics*. Elsevier, 1990.
- 582 11 Björn Grüning, Ryan Dale, Andreas Sjödin, Brad A Chapman, Jillian Rowe, Christopher H
583 Tomkins-Tinch, Renan Valieris, and Johannes Köster. Bioconda: sustainable and comprehens-
584 ible software distribution for the life sciences. *Nature Methods*, 15(7):475–476, 2018.
- 585 12 Dan Gusfield. *Algorithms on Strings, Trees, and Sequences - Computer Science and Computa-*
586 *tional Biology*. Cambridge University Press, 1997. doi:10.1017/cbo9780511574931.
- 587 13 Guillaume Holley and Páll Melsted. Bifrost: highly parallel construction and indexing of
588 colored and compacted de Bruijn graphs. *Genome Biology*, 21(1):1–20, 2020.
- 589 14 Marta Kasprzak. Classification of de Bruijn-based labeled digraphs. *Discrete Applied Math-*
590 *ematics*, 234:86–92, 2018. Special Issue on the Ninth International Colloquium on Graphs and
591 Optimization (GO IX), 2014. URL: <https://www.sciencedirect.com/science/article/pii/S0166218X16304826>, doi:<https://doi.org/10.1016/j.dam.2016.10.014>.
- 592 15 Jamshed Khan, Marek Kokot, Sebastian Deorowicz, and Rob Patro. Scalable, ultra-fast, and
593 low-memory construction of compacted de bruijn graphs with cuttlefish 2. *bioRxiv*, 2021.
- 594 16 Johannes Köster and Sven Rahmann. Snakemake—a scalable bioinformatics workflow engine.
595 *Bioinformatics*, 28(19):2520–2522, 2012.
- 596 17 Vamsi Kundeti, Sanguthevar Rajasekaran, and Heiu Dinh. An efficient algorithm for chinese
597 postman walk on bi-directed de bruijn graphs. In Weili Wu and Ovidiu Daescu, editors,
598 *Combinatorial Optimization and Applications*, pages 184–196, Berlin, Heidelberg, 2010. Springer
599 Berlin Heidelberg.
- 600 18 Veli Mäkinen, Djamel Belazzougui, Fabio Cunial, and Alexandru I Tomescu. *Genome-scale*
601 *algorithm design*. Cambridge University Press, 2015.
- 602 19 Martin D Muggli, Bahar Alipanahi, and Christina Boucher. Building large updatable colored
603 de bruijn graphs via merging. *Bioinformatics*, 35(14):i51–i60, 2019.
- 604 20 Martin D Muggli, Alexander Bowe, Noelle R Noyes, Paul S Morley, Keith E Belk, Robert
605 Raymond, Travis Gagie, Simon J Puglisi, and Christina Boucher. Succinct colored de bruijn
606 graphs. *Bioinformatics*, 33(20):3181–3187, 2017.
- 607 21 Giulio Ermanno Pibiri. Sparse and skew hashing of k -mers. *bioRxiv*, 2022. URL:
608 <https://www.biorxiv.org/content/early/2022/04/04/2022.01.15.476199>, arXiv:
609 <https://www.biorxiv.org/content/early/2022/04/04/2022.01.15.476199.full.pdf>,
610 doi:10.1101/2022.01.15.476199.
- 611 22 Amatur Rahman and Paul Medvedev. Representation of k -mer sets using spectrum-preserving
612 string sets. *Journal of Computational Biology*, 28(4):381–394, 2021.
- 613 23 Sebastian Schmidt, Shahbaz Khan, Jarno Alanko, and Alexandru I Tomescu. Matchtigs:
614 minimum plain text representation of k mer sets. *bioRxiv*, 2021.
- 615 24 Esko Ukkonen. On-line construction of suffix trees. *Algorithmica*, 14(3):249–260, 1995.
- 616

617 **A** Authors contributions

618 JNA and SS discovered the problem, SS solved the problem when the de Bruijn graph is
 619 given and wrote most of the manuscript, JNA designed the linear-time de Bruijn graph
 620 construction algorithm and wrote Section 4. SS implemented the algorithm and conducted
 621 and evaluated the experiments. All authors reviewed and approved the final version of the
 622 manuscript.

623 **B** Omitted proofs

624 ► **Lemma 2** (Sound labels). *Let k be a positive integer and let I be a set of strings of length
 625 at least k . Let $G = (V, E, c)$ be the de Bruijn graph of order k constructed from I . For all
 626 pairs of arcs $e_1 := (v_1d_1, v'_1d'_1, \eta_1), e_2 := (v_2d_2, v'_2d'_2, \eta_2) \in E$ it holds that:*

- 627 (a) $(v'_1 = v_2 \text{ and } d'_1 = -d_2)$ if and only if $\text{suf}_{k-1}(\eta_1) = \text{pre}_{k-1}(\eta_2)$,
 628 (b) $(v'_1 = v'_2 \text{ and } d'_1 = -d'_2)$ if and only if $\text{suf}_{k-1}(\eta_1) = \text{pre}_{k-1}(\text{rc}(\eta_2))$,
 629 (c) $(v_1 = v_2 \text{ and } d_1 = -d_2)$ if and only if $\text{suf}_{k-1}(\text{rc}(\eta_1)) = \text{pre}_{k-1}(\eta_2)$, and
 630 (d) $(v_1 = v'_2 \text{ and } d_1 = -d'_2)$ if and only if $\text{suf}_{k-1}(\text{rc}(\eta_1)) = \text{pre}_{k-1}(\text{rc}(\eta_2))$.

631 **Proof.** Observe that the values of w and w' computed in Lines 5 and 7 of Algorithm 1
 632 are equal to $\text{pre}_{k-1}(\eta_1)$ and $\text{suf}_{k-1}(\eta_1)$ for e_1 and equal to $\text{pre}_{k-1}(\eta_2)$ and $\text{suf}_{k-1}(\eta_2)$ for
 633 e_2 . Further, observe that the values of v and v' computed in Lines 6 and 8 are equal to v_1
 634 and v'_1 for e_1 and equal to v_2 and v'_2 for e_2 . This makes v_1, v'_1, v_2 and v'_2 the canonicals of
 635 $\text{pre}_{k-1}(\eta_1), \text{suf}_{k-1}(\eta_1), \text{pre}_{k-1}(\eta_2)$ and $\text{suf}_{k-1}(\eta_2)$. Finally, observe that the sign values d
 636 and d' computed in Lines 9–14 are equal to d_1 and d'_1 for e_1 and equal to d_2 and d'_2 for e_2 .

- 637 (a) If $v'_1 = v_2$ and $d'_1 = -d_2$, then $w'_1 = w_2$ for all possible values of d'_1 , and therefore
 638 $\text{suf}_{k-1}(\eta_1) = \text{pre}_{k-1}(\eta_2)$.
 639 If $\text{suf}_{k-1}(\eta_1) = \text{pre}_{k-1}(\eta_2)$, then $w'_1 = w_2$, and therefore $v'_1 = v_2$ because v'_1 and v_2 are
 640 the canonicals of w'_1 and w_2 . Additionally, $d'_1 = -d_2$ for all possible values of d'_1 .
 641 (b) If $v'_1 = v'_2$ and $d'_1 = -d'_2$, then $w'_1 = \text{rc}(w'_2)$ for all possible values of d'_1 , and therefore
 642 $\text{suf}_{k-1}(\eta_1) = \text{rc}(\text{suf}_{k-1}(\eta_2)) = \text{pre}_{k-1}(\text{rc}(\eta_2))$.
 643 If $\text{suf}_{k-1}(\eta_1) = \text{pre}_{k-1}(\text{rc}(\eta_2))$, then $w'_1 = \text{rc}(w'_2)$, and therefore $v'_1 = v'_2$ because v'_1 and
 644 v'_2 are the canonicals of w'_1 and w'_2 . Additionally, $d'_1 = -d'_2$ for all possible values of d'_1 .
 645 (c) If $v_1 = v_2$ and $d_1 = -d_2$, then $\text{rc}(w_1) = w_2$ for all possible values of d_1 , and therefore
 646 $\text{suf}_{k-1}(\text{rc}(\eta_1)) = \text{rc}(\text{pre}_{k-1}(\eta_1)) = \text{pre}_{k-1}(\eta_2)$.
 647 If $\text{suf}_{k-1}(\text{rc}(\eta_1)) = \text{pre}_{k-1}(\eta_2)$, then $w_1 = \text{rc}(w_2)$, and therefore $v_1 = v_2$ because v_1 and
 648 v_2 are the canonicals of w_1 and w_2 . Additionally, $d_1 = -d_2$ for all possible values of d_1 .
 649 (d) This case is equivalent to the first case when swapping e_1 and e_2 , because $\text{suf}_{k-1}(\eta_1) =$
 650 $\text{pre}_{k-1}(\eta_2) \iff \text{suf}_{k-1}(\text{rc}(\eta_2)) = \text{pre}_{k-1}(\text{rc}(\eta_1))$. ◀

651 ► **Lemma 3** (Sound sequence of k -mers). *Let k be a positive integer and let I be a set of
 652 strings of length at least k . Let $G = (V, E, c)$ be the de Bruijn graph of order k constructed
 653 from I . Let $W := (e_1 = (v_1d_1, v'_1d'_1, \eta_1), \dots, e_\ell = (v_\ell d_\ell, v'_\ell d'_\ell, \eta_\ell))$ be a walk in G , and
 654 $K := (\kappa_1, \dots, \kappa_\ell)$ its sequence of k -mers. Then for each consecutive pair of k -mers κ_i, κ_{i+1} it
 655 holds that $\text{suf}_{k-1}(\kappa_i) = \text{pre}_{k-1}(\kappa_{i+1})$.*

656 **Proof.** Let $i \in \{1, \dots, \ell - 1\}$. By the definition of walk it holds that $v'_i = v_{i+1}$ and $d'_i = -d_{i+1}$.
 657 We can apply Lemma 2 case by case.

- 658 (a) If $e_i, e_{i+1} \in E$, then by Lemma 2 (a), it holds that $\text{suf}_{k-1}(\eta_i)$ equals $\text{pre}_{k-1}(\eta_{i+1})$. By
 659 definition, $\kappa_i = \eta_i$ and $\kappa_{i+1} = \eta_{i+1}$, so $\text{suf}_{k-1}(\kappa_i) = \text{pre}_{k-1}(\kappa_{i+1})$.

2:20 Eulertigs: minimum plain text k -mer sets without repetitions in linear time

- 660 (b) If $e_i, e_{i+1}^{-1} \in E$, then by Lemma 2 (b) applied to e_i, e_{i+1}^{-1} , it holds that $\text{suf}_{k-1}(\eta_i)$
 661 equals $\text{pre}_{k-1}(\text{rc}(\eta_{i+1}))$. By definition, $\kappa_i = \eta_i$ and $\kappa_{i+1} = \text{rc}(\eta_{i+1})$, so $\text{suf}_{k-1}(\kappa_i) =$
 662 $\text{pre}_{k-1}(\kappa_{i+1})$.
- 663 (c) If $e_i^{-1}, e_{i+1} \in E$, then by Lemma 2 (c) applied to e_i^{-1}, e_{i+1} , it holds that $\text{suf}_{k-1}(\text{rc}(\eta_i))$
 664 equals $\text{pre}_{k-1}(\eta_{i+1})$. By definition, $\kappa_i = \text{rc}(\eta_i)$ and $\kappa_{i+1} = \eta_{i+1}$, so $\text{suf}_{k-1}(\kappa_i) =$
 665 $\text{pre}_{k-1}(\kappa_{i+1})$.
- 666 (d) If $e_i^{-1}, e_{i+1}^{-1} \in E$, then by Lemma 2 (d) applied to e_i^{-1}, e_{i+1}^{-1} , it holds that $\text{suf}_{k-1}(\text{rc}(\eta_i))$
 667 equals $\text{pre}_{k-1}(\text{rc}(\eta_{i+1}))$. By definition, $\kappa_i = \text{rc}(\eta_i)$ and $\kappa_{i+1} = \text{rc}(\eta_{i+1})$, so $\text{suf}_{k-1}(\kappa_i) =$
 668 $\text{pre}_{k-1}(\kappa_{i+1})$. ◀

669 ▶ **Lemma 4 (Sound spell).** *Let k be a positive integer and let I be a set of strings of length*
 670 *at least k . Let $G = (V, E, c)$ be the de Bruijn graph of order k constructed from I . Let W be*
 671 *a walk in G , K_W its sequence of k -mers and K'_W the sequence of k -mers of $\text{spell}(W)$. Then*
 672 $K_W = K'_W$.

673 **Proof.** Let $(\kappa_1, \dots, \kappa_\ell) := K_W$. We use induction over the length of W . For an empty W ,
 674 K is empty, $\text{spell}(W)$ is empty, and therefore K' is empty as well. For $|W| = 1$, Algorithm 2
 675 outputs $\text{spell}(W) = \kappa_1$ and it holds that $K'_W = (\kappa_1) = K_W$.

676 For $|W| \geq 2$ we consider that $K_X = K'_X$ holds for a prefix X of W with $|X| = |W| - 1$.
 677 When $i = |W|$ at the beginning of the loop in Line 8, then $s = \text{spell}(X)$. By Lemma 3 it
 678 holds that the last $k - 1$ characters of s are equal to the first $k - 1$ characters of κ_ℓ . Therefore,
 679 by appending the last character from κ_ℓ to s , κ_ℓ is appended to K'_X forming K'_W . Therefore,
 680 last k -mer of K'_W equals the last k -mer of K_W , and the first $\ell - 1$ k -mers of K'_W equal those
 681 of K_W by induction. ◀

682 **C Pseudocode for linear-time construction of compacted de Bruijn** 683 **graphs**

684 The pseudocode for computing a compacted de Bruijn graph in linear time is given by
 685 Algorithm 6 which uses Algorithm 5 as a subroutine. The data structure D used by the
 686 algorithms is that described in Section 4. Note that if we compute the arc labels as plain
 687 strings as in Algorithm 1, we need up to $O(k)$ bits to store a single- k -mer arc. And since arcs
 688 are not substrings of input strings (but potentially combinations of input strings), we would
 689 need a string set of up to $O(k||I||)$ characters to store all arc labels without referring to the
 690 input strings. This contradicts the algorithm being linear in $||I||$. However, we can store the
 691 labels as tuples (p, η, q, r) , where $p\eta q$ is the label where p and q are explicit strings while η is
 692 a pointer to a k -mer in the input. If r is true, then the label must be reverse complemented
 693 to match that defined by Algorithm 1. With this fix, the size of each label is linear in the
 694 number of k -mers it represents, and in total the de Bruijn graph represents $O(||I||)$ k -mers.

695 The comparison on Line 16 of Algorithm 6 can be done in linear time in $|\eta_1| + |\eta_2|$ by
 696 finding the suffix array intervals of $\eta_1\eta\eta_2$ and $\text{rc}(\eta_1\eta\eta_2)$ with `extendLeft` and `extendRight`
 697 from η and $\text{rc}(\eta)$ respectively, and comparing the starts of the intervals. This way, the total
 698 time taken by all those comparisons is proportional to the sum of $|\eta_1| + |\eta_2|$ over all unitigs,
 699 which is linear in $||I||$ because each character of η_1 and η_2 can be mapped to a distinct edge
 700 in the non-compacted de Bruijn graph of $||I||$. Therefore, the algorithm can be implemented
 701 to run in $O(||I||)$ time.

702 Our pseudocode does not compute the first and last character of each arc-label, but this
 703 can be easily computed in constant time using w_i , η_1 and η_2 in Algorithm 6.

■ **Algorithm 5** FINDUNITIGEND

Input: A data structure D , a pair of suffix-intervals $[a_f, b_f], [a_r, b_r]$, an array S_E mapping from suffix-space to boolean, an array S_V mapping from suffix space to nodes, a set of nodes V . Each node in V contains a parameter c .

Output: A node v at the end of the unitig and a sign d , as well as the updated S_E, S_V, V and the label η of the traversed path.

```

1  $[a_f, b_f] \leftarrow \text{contractLeft}(D, [a_f, b_f])$ 
2  $[a_r, b_r] \leftarrow \text{contractRight}(D, [a_r, b_r])$ 
3  $\eta \leftarrow \epsilon$ 
  // extend over  $(k-1)$ -mers that have indegree and outdegree of 1
4 while  $|\text{enumerateRight}(D, [a_f, b_f]) \cup \text{rc}(\text{enumerateLeft}(D, [a_r, b_r]))| =$ 
    $|\text{enumerateLeft}(D, [a_f, b_f]) \cup \text{rc}(\text{enumerateRight}(D, [a_r, b_r]))| = 1$  do
5    $\{\sigma\} \leftarrow \text{enumerateRight}(D, [a_f, b_f]) \cup \text{rc}(\text{enumerateLeft}(D, [a_r, b_r]))$ 
6    $\eta \leftarrow \eta\sigma$ 
7    $[a_f, b_f] \leftarrow \text{extendRight}(D, [a_f, b_f], \sigma)$ 
8    $[a_r, b_r] \leftarrow \text{extendLeft}(D, [a_r, b_r], \text{rc}(\sigma))$ 
9   foreach  $h \in [a_f, b_f] \cup [a_r, b_r]$  do  $S_E[h] \leftarrow \text{true}$ 
10   $[a_f, b_f] \leftarrow \text{contractLeft}(D, [a_f, b_f])$ 
11   $[a_r, b_r] \leftarrow \text{contractRight}(D, [a_r, b_r])$ 
12 if  $S_V[a_f] = \perp$  then
13   insert node  $v$  into  $V$ 
14   foreach  $h \in [a_f, b_f] \cup [a_r, b_r]$  do  $S_V[h] \leftarrow v$ 
15 else  $v \leftarrow S_V[a_f]$ 
16 if  $a_f = a_r$  then  $c(v) \leftarrow 1$ ;  $d \leftarrow \odot$  /*  $v$  self-complemental */
17 else if  $a_f < a_r$  then  $c(v) \leftarrow 0$ ;  $d \leftarrow \ominus$  /*  $v$  canonical */
18 else  $c(v) \leftarrow 0$ ;  $d \leftarrow \oplus$  /*  $v$  not canonical */
19 return  $(v, d, S_E, S_V, V, \eta)$ 

```

2:22 Eulertigs: minimum plain text k -mer sets without repetitions in linear time

■ Algorithm 6 LINEARCOMPACTEDDBG

Input: An integer k and a set of strings $I = (w_1, \dots, w_\ell)$ where each string has length at least k .

Output: A de Bruijn graph $G = (V, E, c)$ of order k .

```

1  $V \leftarrow \emptyset; E \leftarrow \emptyset$  //  $c$  is stored as parameter of each node
  //  $\$$  is a special character outside of the alphabet
2  $T \leftarrow \$w_1\$w_2\$ \dots \$w_\ell\$rc(w_1)\$rc(w_2)\$ \dots \$rc(w_\ell)\$$ 
3  $S_V \leftarrow$  array of length  $|T|$  filled with  $\perp$  mapping from suffix space to nodes in  $V$ 
4  $S_E \leftarrow$  array of length  $|T|$  filled with false marking used  $k$ -mers
5 build data structure  $D$  over  $T$  // See text in Section 4
6 foreach  $w_i \in I$  do
7    $[a_f, b_f] \leftarrow \text{find}(D, \text{pre}_k(w_i))$  // Suffix array interval of  $\text{pre}_k(w_i)$ 
8    $[a_r, b_r] \leftarrow \text{find}(D, \text{pre}_k(\text{rc}(w_i)))$  // Suffix array interval of  $\text{pre}_k(\text{rc}(w_i))$ 
9   foreach  $j \in (k + 1, \dots, |w_i|)$  do
10    if  $S_E[a_f] = \text{false}$  then // create arc from unused  $k$ -mer
11      foreach  $h \in [a_f, b_f] \cup [a_r, b_r]$  do  $S_E[h] \leftarrow \text{true}$ 
12       $\eta \leftarrow$  pointer to  $\text{pre}_k(w_i)$ 
        // find unitig start by finding the end on the rev. comp.
13       $(v_1, d_1, S_E, S_V, V, \eta_1) \leftarrow \text{FINDUNITIGEND}(D, [a_r, b_r], [a_f, b_f], S_E, S_V, V)$ 
        // find unitig end
14       $(v_2, d_2, S_E, S_V, V, \eta_2) \leftarrow \text{FINDUNITIGEND}(D, [a_f, b_f], [a_r, b_r], S_E, S_V, V)$ 
        // Reverse because finding the start was done in reverse
15       $\eta_1 \leftarrow \text{rc}(\eta_1)$ 
16      if  $\text{rc}(\eta_1\eta_2) < \eta_1\eta_2$  then // arc labels are always canonical
17        swap  $v_1$  and  $v_2$ 
18        swap  $d_1$  and  $d_2$ 
19         $d_1 \leftarrow \neg d_1$ 
20         $d_2 \leftarrow \neg d_2$ 
21         $r \leftarrow \text{true}$ 
22      else
23         $r \leftarrow \text{false}$ 
24      insert  $e = (v_1d_1, v_2d_2, (\eta_1, \eta, \eta_2, r))$  into  $E$ 
25       $[a_f, b_f] \leftarrow \text{extendRight}(D, [a_f, b_f], w_i[j])$ 
26       $[a_r, b_r] \leftarrow \text{extendLeft}(D, [a_r, b_r], \text{rc}(w_i)[j])$ 
27       $[a_f, b_f] \leftarrow \text{contractLeft}(D, [a_f, b_f])$ 
28       $[a_r, b_r] \leftarrow \text{contractRight}(D, [a_r, b_r])$ 

```
