

The Complexity of Approximate Pattern Matching on De Bruijn Graphs

Daniel Gibney¹, Sharma V. Thankachan², and Srinivas Aluru¹

¹ School of Computational Science and Engineering, Georgia Institute of Technology, Atlanta, USA

² Department of Computer Science, University of Central Florida, Orlando, USA

Email: daniel.j.gibney@gmail.com, sharma.thankachan@ucf.edu, aluru@cc.gatech.edu

Abstract. Aligning a sequence to a walk in a labeled graph is a problem of fundamental importance to Computational Biology. For finding a walk in an arbitrary graph with $|E|$ edges that exactly matches a pattern of length m , a lower bound based on the Strong Exponential Time Hypothesis (SETH) implies an algorithm significantly faster than $\mathcal{O}(|E|m)$ time is unlikely [Equi *et al.*, ICALP 2019]. However, for many special graphs, such as de Bruijn graphs, the problem can be solved in linear time [Bowe *et al.*, WABI 2012]. For approximate matching, the picture is more complex. When edits (substitutions, insertions, and deletions) are only allowed to the pattern, or when the graph is acyclic, the problem is again solvable in $\mathcal{O}(|E|m)$ time. When edits are allowed to arbitrary cyclic graphs, the problem becomes NP-complete, even on binary alphabets [Jain *et al.*, RECOMB 2019]. These results hold even when edits are restricted to only substitutions. Despite the popularity of de Bruijn graphs in Computational Biology, the complexity of approximate pattern matching on de Bruijn graphs remained open. We investigate this problem and show that the properties that make de Bruijn graphs amenable to efficient exact pattern matching do not extend to approximate matching, even when restricted to the substitutions only case. Specifically, we prove that determining the existence of a matching walk in a de Bruijn graph is NP-complete when substitutions are allowed to the graph. In addition, we demonstrate that an algorithm significantly faster than $\mathcal{O}(|E|m)$ is unlikely for de Bruijn graphs in the case where only substitutions are allowed to the pattern. This stands in contrast to pattern-to-text matching where exact matching is solvable in linear time, like on de Bruijn graphs, but approximate matching under substitutions is solvable in subquadratic $\tilde{O}(n\sqrt{m})$ time, where n is the text's length [Abrahamson, SIAM J. Computing 1987].

27 **1 Introduction**

28 De Bruijn graphs are an essential tool in Computational Biology. Their role in de novo assembly
 29 spans back to the 1980s [39], and their application in assembly has been researched extensively since
 30 then [10,11,17,32,37,38,42,45]. More recently, de Bruijn graphs have been applied in metagenomics
 31 and in the representation of large collections of genomes [14,26,29,36,44] and for solving other
 32 problems such as read-error correction [31,34] and compression [9,23]. Due to the popularity of de
 33 Bruijn graphs in the modeling of sequencing data, an algorithm to efficiently find *walks* in a de
 34 Bruijn graph matching (or approximately matching) a given query pattern would be a significant
 35 advancement. For example, in metagenomics, such an algorithm could quickly detect the presence of
 36 a particular species within genetic material obtained from an environmental sample. Or, in the case
 37 of read-error correction, such an algorithm could be used to efficiently find the best mapping of reads
 38 onto a ‘cleaned’ reference de Bruijn graph with low-frequency k-mers removed [31]. To facilitate
 39 such tasks, several algorithms (often seed-and-extend type heuristics) and software tools have been
 40 developed that perform pattern matching on de Bruijn (and general) graphs [5,21,22,28,30,33,35,41].

41 The importance of pattern matching on labeled graphs in Computational Biology and other
 42 fields has caused a recent surge of interest in the theoretical aspects of this problem. In turn,
 43 this has led to many new fascinating algorithmic and computational complexity results. However,
 44 even with this improved understanding of the theory of pattern matching on labeled graphs, our
 45 knowledge is still lacking in many respects concerning specific, yet extremely relevant, graph classes.
 46 An overview of the current state of knowledge is provided in Table 1.

	Exact Matching	Approximate Matching
Easy	<u>Solvable in Linear Time</u> <ul style="list-style-type: none"> • Wheeler Graphs [16] (e.g. de Bruijn graphs, NFAs for multiple strings)	<u>Solvable in $\mathcal{O}(E m)$ time</u> <ul style="list-style-type: none"> • DAGs: Substitutions/Edits to graph [28] • General graphs: Substitutions/Edits to pattern [6] • de Bruijn Graphs: Substitutions to pattern -No strongly Sub-$\mathcal{O}(E m)$ algorithm (this paper)
	<u>NO Strongly Sub-$\mathcal{O}(E m)$ Algorithm</u> <ul style="list-style-type: none"> • General graphs [13,19] (including DAGs with total degree ≤ 3)	<u>NP-Complete</u> <ul style="list-style-type: none"> • General graphs: Substitutions/Edits to vertex labels [6,25] • de Bruijn Graphs: Substitutions to vertex labels (this paper)

Table 1. Overview of the computational complexity of pattern matching on labeled graphs

47 For general graphs, we can consider exact and approximate matching. For exact matching,
 48 conditional lower-bounds based on the Strong Exponential Time Hypothesis (SETH), and other
 49 conjectures in circuit complexity, indicate that an $\mathcal{O}(|E|m^{1-\varepsilon} + |E|^{1-\varepsilon}m)$ time algorithm with any
 50 constant $\varepsilon > 0$, for a graph with $|E|$ edges and a pattern of length m , is highly unlikely (as is
 51 the ability to shave more than a constant number of logarithmic factors from the $\mathcal{O}(|E|m)$ time
 52 complexity) [13,19]. These results hold for even very restricted types of graphs, for example, DAGs
 53 with maximum total degree three and binary alphabets. For approximate matching, when edits are
 54 only allowed in the pattern, the problem is solvable in $\mathcal{O}(|E|m)$ time [6]. If edits are also permitted
 55 in the graph, but the graph is a DAG, matching can be done in the same time complexity [28].
 56 However, the problem becomes NP-complete when edits are allowed in arbitrary cyclic graphs. This
 57 was originally proven in [6] for large alphabets and more recently proven for binary alphabets in
 58 [25]. These results hold even when edits are restricted to only substitutions. The distinction between

59 modifications to the graph and modifications to the pattern is important as these two problems
60 are fundamentally different. When changes are made to cyclic graphs the same modification can
61 be encountered multiple times while matching a pattern with no additional cost. Furthermore,
62 algorithmic solutions appearing in [28,35,41] are for the case where modifications are performed
63 only to the pattern.

64 De Bruijn graphs are an interesting case from a theoretical perspective. Many graphs allow for
65 extending Burrows-Wheeler Transformation (BWT) based techniques for efficient pattern matching.
66 Sufficient conditions for doing this are captured by the definition of Wheeler graphs, introduced in
67 [16], and further studied in [3,4,12,15,20]. De Bruijn graphs are themselves Wheeler graphs, hence
68 on a de Bruijn graph the exact pattern matching problem is solvable in linear time. However, the
69 complexity of finding an approximate match in a de Bruijn graph when permitting modifications
70 to the graph or modifications to the pattern remained open [25].

71 Our work makes two important contributions (see Table 1). First, we prove that for de Bruijn
72 graphs, despite exact matching being solvable in linear time, the approximate matching problem
73 with vertex label substitutions is NP-complete. Second, we prove that a strongly subquadratic time
74 algorithm for the approximate pattern matching problem on de Bruijn graphs, where substitutions
75 are only allowed in the pattern, is not possible under SETH. Such a proof confirms the optimality
76 of the known quadratic time algorithms when considering polynomial factors. Note that pattern-
77 to-text matching (under substitutions) can be solved in sub-quadratic $\tilde{O}(n\sqrt{m})$ time, where n is
78 the length of the text [2,7]. Our results demonstrate that *the properties that make de Bruijn graphs*
79 *amenable to efficient exact pattern matching do not extend to approximate matching.* To the best
80 of our knowledge, these are the first such results for any type of Wheeler graph.

81 1.1 Technical Background and Our Results

82 Some necessary terminology is presented next.

83 **Notation for edges:** For a directed edge from a vertex u to a vertex v we will use the notation
84 (u, v) . Additionally, we will refer to u as the *tail* of (u, v) , and v as the *head* of (u, v) .

85 **Walks versus paths:** A distinction must be made between the concept of a *walk* and a *path* in
86 a graph. A walk is a sequence of vertices v_1, v_2, \dots, v_t such that for each $i \in [1, t - 1]$, $(v_i, v_{i+1}) \in E$.
87 Vertices can be repeated in a walk. A path is a walk where vertices are not repeated. The length
88 of a walk is defined as the number of edges in the walk, $t - 1$, or equivalently one less than the
89 number of vertices in the sequence (counted with multiplicity). This work will be concerning the
90 existence of walks.

91 **Induced subgraphs:** An induced subgraph of a graph $G = (V, E)$ consists of a subset of
92 vertices $V' \subseteq V$, and all edges $(u, v) \in E$ such that $u, v \in V'$. This is in contrast to an arbitrary
93 subgraph of G , where an edge can be omitted from the subgraph, even if both of its incident vertices
94 are included.

95 **De Bruijn graphs:** A full de Bruijn graph is a directed graph whose vertices consists of all
96 k -mers (strings of length k) from an alphabet $\Sigma = \{0, 1, \dots, \sigma - 1\}$ in which there is a directed edge
97 from each vertex $s_1s_2\dots s_k \in \Sigma^k$ to the σ vertices $s_2s_3\dots s_k\alpha$, $\alpha \in \Sigma$. We will work with induced
98 subgraphs of full de Bruijn graphs in this paper. We consider a vertex v as being labeled with a
99 single symbol $L(v) \in \Sigma$. The k -mer corresponding to the vertex is referred to as its *implicit label*.
100 Formally, the implicit label of a vertex v is $L(u_1)L(u_2)\dots L(u_{k-1})L(v)$ where $u_1, u_2, \dots, u_{k-1}, v$ is
101 any length $k - 1$ walk ending at v . When no such walk exists, we assume v has an implicit label

102 compatible with the other implicit labels in the graph. This is equivalent to the notion of a de
103 Bruijn graph constructed from k -mers commonly used in Computational Biology.

104 **Strings and Matching:** For a string S of length n indexed from 1 to n , we use $S[i]$ to denote
105 the i^{th} symbol in S . We use $S[i, j]$ to denote the substring $S[i]S[i + 1] \dots S[j]$. If $j < i$, then we take
106 $S[i, j]$ as the empty string. As mentioned above, we will consider every vertex v as labeled with a
107 single symbol $L(v) \in \Sigma = \{0, 1, \dots, \sigma - 1\}$. We say that a pattern P matches a walk v_1, v_2, \dots, v_m
108 iff $P[i] = L(v_i)$ for every $i \in [1, m]$.

109 With these definitions in hand, we can formally define our first problem.

110 *Problem 1 (Approximate matching with vertex label substitutions).* Given a vertex labeled graph
111 $D = (V, E)$ with alphabet Σ , pattern $P[1, m]$, and integer $\delta \geq 0$, determine if there exists a walk
112 in D matching P after at most δ substitutions to the vertex labels.

113 **Theorem 1.** *Problem 1 is NP-complete on de Bruijn graphs having an alphabet of size $\sigma \geq 4$.*

114 Theorem 1 is proven in Section 2. Intuitively, our reduction transforms a general directed graph
115 into a de Bruijn that maintains key topological properties related to the existence of walks. In our
116 proof, the k -mer size is dependent on the alphabet size σ . For example, when $\sigma = 4$, we require
117 $k = \Theta(\log^2 |V|)$, but when $\sigma = \Theta(\log |V|)$, we can reduce the value of k to $\Theta(\log |V|)$. An extension
118 of the hardness result for large σ on connected graphs is given at the end of Section 2.1.

119 The distinct problem of approximately matching a pattern to a *path* in a de Bruijn graph was
120 shown to be NP-complete in [30]. As mentioned by the authors of that work, the techniques used
121 there do not appear to be easily adaptable to the problem for walks. Our approach uses edge
122 transformations more closely inspired by those used in [27] for proving hardness on the paired de
123 Bruijn sound cycle problem.

124 *Problem 2 (Approximate matching with substitutions within the pattern).* Given a vertex labeled
125 graph $D = (V, E)$ with alphabet Σ , pattern $P[1, m]$, and integer $\delta \geq 0$, determine if there exists a
126 walk in D matching P after at most δ substitutions to the symbols in P .

127 For Problem 2, in Section 3 we provide a SETH-based hardness result. The Strong Exponential
128 Time Hypothesis, or SETH, is frequently utilized in establishing the conditional optimality of
129 polynomial time algorithms, often for problems on strings or graphs [1,8,13,18,19,24]. We refer the
130 reader to [43] for the definition of SETH and for the reduction to the Orthogonal Vectors problem
131 (OV), which is utilized to prove Theorem 2.

132 **Theorem 2.** *Conditioned on SETH, for all constants $\varepsilon > 0$, there does not exist an algorithm for*
133 *Problem 2 on de Bruijn graphs running in time $\mathcal{O}(|E|m^{1-\varepsilon} + |E|^{1-\varepsilon}m)$. This holds for alphabets of*
134 *size four (or greater) and k -mer size that is $\Theta(\log |E|)$.*

135 2 NP-Completeness of Problem 1 on De Bruijn Graphs

136 Our proof of NP-completeness uses a reduction from the Hamiltonian Cycle Problem on directed
137 graphs. To present the reduction, we introduce the concept of *merging* two vertices. To merge
138 vertices u and v , we create a new vertex w . We then take all edges with either u or v as their head
139 and make w their new head. Next, we take all edges with either u or v as their tail and make w
140 their new tail. This makes the edges (u, v) and (v, u) (if they existed) into self-loops for w . If two
141 self-loops are formed, we delete one of them. Finally, we delete the original vertices u and v .

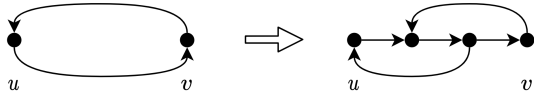


Fig. 1. Gadget to remove cycles of length 2 from the initial input graph.

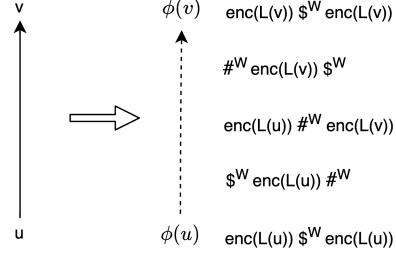


Fig. 2. The transformation from edges to paths used in our reduction.

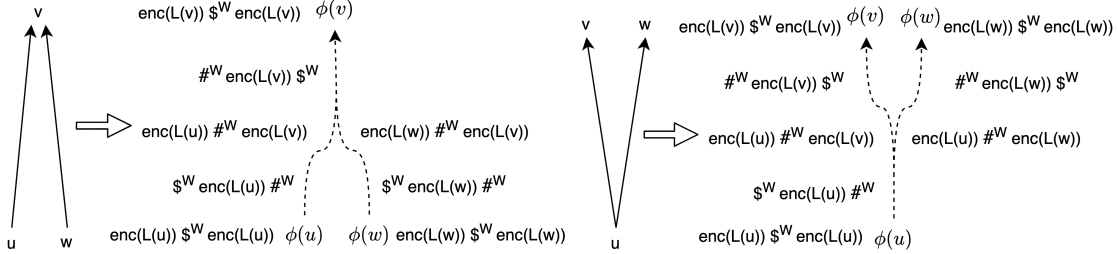


Fig. 3. Vertices with the same implicit label are merged while transforming D to D' , causing edges with shared vertices to become paths with shared vertices.

142 2.1 Reduction

143 We start with an instance of the Hamiltonian cycle problem on a directed graph where the number of
 144 edges is linear in the number of vertices. The Hamiltonian cycle problem remains NP-complete even
 145 when restricted to such graphs [40]. We can assume there are no self-loops or vertices with in-degree
 146 or out-degree zero. To simplify the proof, we first eliminate any cycles of length 2 using the gadget
 147 in Figure 1. We denote the resulting graph as $D = (V, E)$ and let $n = |V|$. We assign each vertex
 148 $v \in V$ a unique integer $L(v) \in [0, n - 1]$. Let $\ell = \lceil \log n \rceil$, $\text{bin}(i)$ be the standard binary encoding of
 149 i using ℓ bits, $\sigma \in [4, 2\ell]$, and $\Sigma = \{\$, \#, 0, 1, \dots, \sigma - 3\}$. Define $\text{enc}(i) = (0^{2\ell}1)^\alpha 23\dots(\sigma - 3) \text{bin}(i)$,
 150 where if $\sigma = 4$ then the substring $23\dots(\sigma - 3)$ is the empty string, and $\alpha = 2\ell - \sigma + 4$. Let
 151 $W = |\text{enc}(i)| = (2\ell + 1)\alpha + \sigma - 4 + \ell$ and $k = 3W$. Note that if $\sigma = 2\ell$, $k = \mathcal{O}(\log n)$.

152 We construct a new (de Bruijn) graph $D' = (V', E')$ as follows: Initially D' is the empty
 153 graph. For $i = 0, 1, \dots, n - 1$, for each edge $(u, v) \in E$ where $L(v) = i$, create a new path whose
 154 concatenation of vertex labels is $\#^W \text{enc}(i) \$^W \text{enc}(i)$. The vertex u will correspond with a new
 155 vertex $\phi(u)$ at the start of this path, and the vertex v will correspond with a new vertex $\phi(v)$ at the
 156 end of this path. The vertex $\phi(v)$ has the implicit label $\text{enc}(L(v)) \$^W \text{enc}(L(v))$. The vertex $\phi(u)$ is
 157 temporarily assigned the implicit label $\text{enc}(L(u)) \$^W \text{enc}(L(u))$. See Figure 2. We call vertices with
 158 implicit labels of the form $\text{enc}(L(\cdot)) \$^W \text{enc}(L(\cdot))$ *marked vertices*. We use the notation $\phi((u, v))$ to
 159 denote the path created when applying this transformation to $(u, v) \in E$. After the path $\phi((u, v))$
 160 is created, vertices in V' having the same implicit label are merged, and parallel edges are deleted
 161 (Figure 3). See Figure 4 for a complete example. Finally, let $\delta = 2\ell(n - 1)$ and

$$P = \#^W \text{enc}(0) \$^W \text{enc}(0) \#^W \text{enc}(1) \$^W \text{enc}(1) \#^W \dots \#^W \text{enc}(n - 1) \$^W \text{enc}(n - 1) \#^W \text{enc}(0) \$^W \text{enc}(0).$$

162 We will show that there exists a walk in D' matching P with at most δ vertex label substitutions
 163 iff D contains a Hamiltonian cycle.

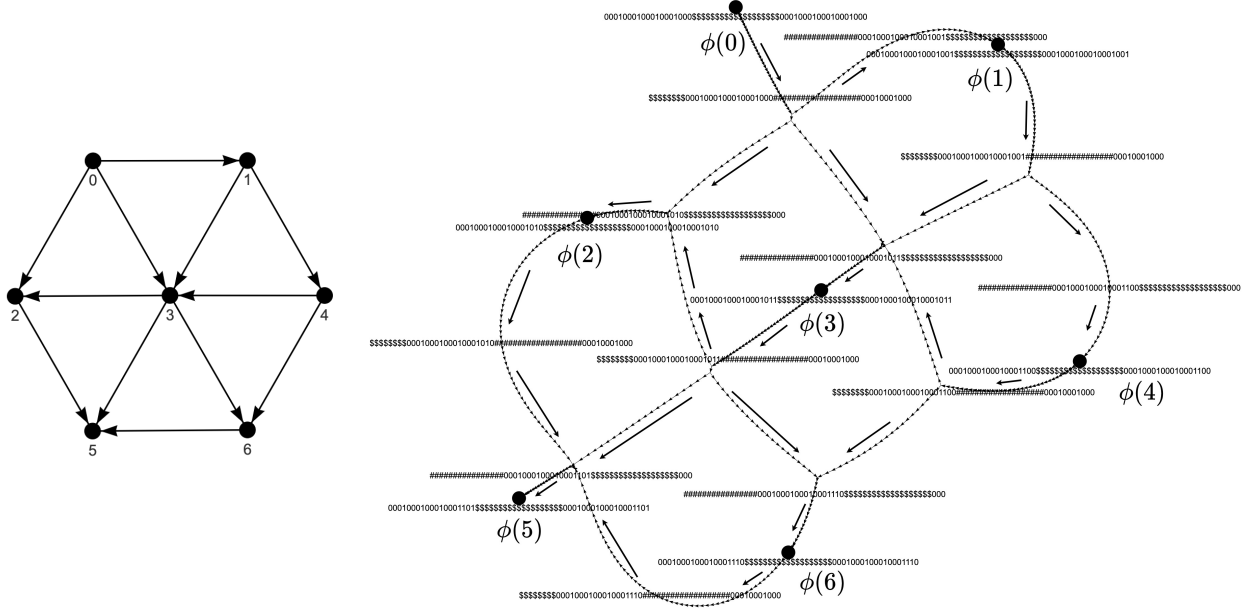


Fig. 4. (Left) A graph before the reduction is applied to it. (Right) The transformed graph. A subset of the implicit labels are shown, and the path directions are annotated by arrows beside each path. Note that $\text{enc}(\cdot)$ has been modified to have the prefix $(0^\ell 1)^{\ell+1}$ so that it fits in the figure. Also, unlike in the figure, we assume in our reduction that there are no vertices with in-degree or out-degree zero.

164 **Proof of Correctness**

165 **Lemma 1.** *The graph D' constructed as above is a de Bruijn graph.*

166 *Proof. (Overview)* Three properties must be proven: (i) Implicit labels are unique, meaning for
 167 every implicit label at most one vertex is assigned that label; (ii) No edges are missing, i.e., if the
 168 implicit label of $y \in V'$ is $S\alpha$ for some string $S[1, k - 1]$ and symbol $\alpha \in \Sigma$, and there exists a vertex
 169 $x \in V'$ with implicit label $\beta S[1, k - 1]$ for some symbol $\beta \in \Sigma$, then $(x, y) \in E'$; (iii) Implicit labels
 170 are well-defined, in that every walk of length $k - 1$ ending at a vertex $x \in V'$ matches the same
 171 string (the implicit label of x); The most involved of these is proving property (ii), which requires
 172 analyzing several cases. The full proof is given in Appendix 1. \square

173 The correctness of the reduction remains to be shown. Lemmas 2-4 establish useful structural
 174 properties of D' , Lemma 5 proves that the existence of a Hamiltonian Cycle in D implies an
 175 approximate matching in D' , and Lemmas 6-9 demonstrate the converse.

176 **Lemma 2.** *Any walk between two marked vertices $\phi(u)$ and $\phi(v)$ containing no additional marked*
 177 *vertices has length $4W$. Hence, we can conclude any such walk is a path.*

178 *Proof. (Overview)* This is proven using induction on the number of edges transformed. It is shown
 179 that for every vertex, a key property regarding the distances to its closest marked vertices continues
 180 to hold after vertices on any newly created path are merged. See Appendix 1 for the full proof. \square

181 **Lemma 3.** *For $(u_1, v_1), (u_2, v_2) \in E$, unless $u_1 = u_2$ or $v_1 = v_2$, $\phi((u_1, v_1))$ and $\phi((u_2, v_2))$ share*
 182 *no vertices.*

183 *Proof.* In the case where $\{u_1, v_1\} \cap \{u_2, v_2\} = \emptyset$ (Figure 5 left), every implicit vertex label in
184 $\phi((u_1, v_1))$ contains $\text{enc}(L(u_1))$ or $\text{enc}(L(v_1))$ (or both), and contains neither $\text{enc}(L(u_2))$ nor $\text{enc}(L(v_2))$.
185 Similarly, every implicit vertex label in $\phi((u_2, v_2))$ contains $\text{enc}(L(u_2))$ or $\text{enc}(L(v_2))$ (or both) and
186 contains neither $\text{enc}(L(u_1))$ nor $\text{enc}(L(v_1))$. This implies that none of the implicit labels match
187 between the two paths, thus no vertices are merged. In the case where $v_1 = u_2$ and $u_1 \neq v_2$ (Figure
188 5, right), the implicit labels of vertices $\phi((u_1, v_1))$ not containing $\text{enc}(L(u_1))$ have $\#$ symbols in
189 different positions than implicit labels of vertices in $\phi((u_2, v_2))$ not containing $\text{enc}(L(v_2))$, and, since
190 $v_1 \neq v_2$, cannot match the implicit labels of vertices in $\phi((u_2, v_2))$ containing $\text{enc}(L(v_2))$. Vertices
191 in $\phi((u_1, v_1))$ with implicit labels containing $\text{enc}(L(u_1))$ have $\#$ symbols in different positions than
192 implicit labels of vertices in $\phi((u_2, v_2))$ not containing $\text{enc}(L(u_2))$, and, since $u_1 \neq u_2$, cannot match
193 the implicit labels of vertices in $\phi((u_2, v_2))$ containing $\text{enc}(L(u_2))$. The case $u_1 = v_2$ and $u_2 \neq v_1$ is
194 symmetric. The case $u_1 = v_2$ and $v_1 = u_2$ cannot happen since, by the use of our gadget in Figure
195 1, D cannot contain the edges (u_1, v_1) and (v_1, u_1) . \square

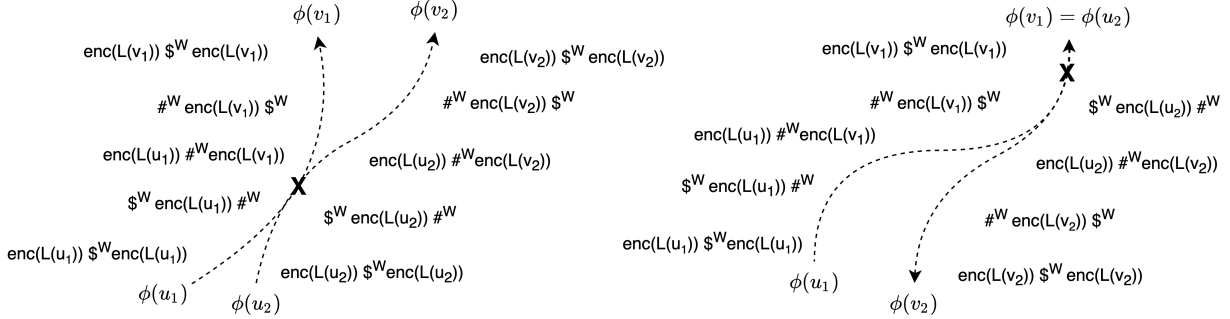


Fig. 5. Examples where paths between marked vertex cannot share any vertex: (Left) The case where $\{u_1, v_1\} \cap \{u_2, v_2\} = \emptyset$. (Right) The case where $v_1 = u_2$ and $u_1 \neq v_2$.

196 **Lemma 4.** *There exists a path from a marked vertex $\phi(u) \in V'$ to a marked vertex $\phi(v) \in V'$*
197 *containing no other marked vertices iff there is an edge $(u, v) \in E$.*

198 *Proof. (Overview)* It is clear from construction that if $(u, v) \in D$, then such a path exists in D' .
199 In the other direction, we utilize Lemmas 2 and 3 to show that such a path existing without a
200 corresponding edge would create a contradiction. The full proof is provided in Appendix 1. \square

201 **Lemma 5.** *If D has a Hamiltonian cycle, then P can be matched in D' with at most δ substitutions*
202 *to vertex labels of D' .*

203 *Proof.* To obtain a matching walk, follow the cycle corresponding to a solution in D starting with
204 the marked vertex in V' corresponding to the vertex in V with label 0. By Lemma 4, each edge
205 traversed in D corresponds to a path in D' . While traversing these paths, modify the vertex labels
206 in D' corresponding to the substrings $\text{bin}(i)$ to match P . Assuming no conflicting substitutions are
207 needed, this requires at most $2\ell(n - 1)$ substitutions.

208 It remains to be shown that no conflicting label substitutions will be necessary. Consider the
209 edges $(u_1, v_1), (u_2, v_2) \in E$ used in the Hamiltonian cycle in D . We will never have $u_1 = u_2$ or
210 $v_1 = v_2$. Hence, by Lemma 3, the sets of vertices on the paths $\phi((u_1, v_1))$ and $\phi((u_2, v_2))$ are
211 disjoint. \square

212 **Lemma 6.** *If P can be matched in D' with at most δ substitutions to vertex labels of D' , then*
 213 *all $\$$'s in P are matched with non-substituted $\$$'s in D' and all $\#$'s in P are matched with non-*
 214 *substituted $\#$'s in D' . Consequently, we can assume the only substitutions are to the vertex labels*
 215 *corresponding to $\text{bin}(i)$'s within $\text{enc}(i)$'s.*

216 *Proof. (Overview)* We establish the existence of a long, non-branching path for every marked vertex
 217 that can be traversed at most once when matching P . This, combined with maximal paths of, $\$, \#,$
 218 and $0/1$ -symbols, all being of length W , makes it so that 'shifting' P to match a portion of D forces
 219 the shift to occur throughout the walk traversed while matching P . Utilizing the large Hamming
 220 distance between shifted instances of two encodings, we can then show that not matching all non- $0/1$
 221 symbols requires more than δ substitutions. The full proof is provided in Appendix 1. \square

222 Post-substitution to vertex labels, we will refer to a vertex as marked if there exists a walk
 223 ending at it that matches a string of the form $\text{enc}(L(u))\$^W \text{enc}(L(u))$, $u \in V$. Note that this
 224 definition does not require all length $k - 1$ walks ending at such a vertex to match the same string.

225 **Lemma 7.** *If P can be matched in D' with at most δ substitutions to vertex labels of D' , then no*
 226 *additional marked vertices are created due to vertex substitutions.*

227 *Proof.* Pre-substitution, only marked vertices have implicit labels of the form $S_1\$^W S_2$ where S_1
 228 and S_2 contain no $\$$ symbols. Hence, the only way that a vertex could have a walk ending at it
 229 that matches a pattern of that form post-substitution is if either it was originally a marked vertex,
 230 or some non- $0/1$ -symbols were substituted in D' . However, by Lemma 6 the latter case cannot
 231 happen, and only originally marked vertices have walks ending at them matching strings of the
 232 form $S_1\$^W S_2$ post-substitution. \square

233 **Lemma 8.** *If P can be matched in D' with at most δ substitutions to vertex labels of D' , then each*
 234 *originally marked vertex in D' is visited exactly once, except for an originally marked vertex at the*
 235 *end of a path matching $\text{enc}(0)\$^W \text{enc}(0)$ that is visited twice.*

236 *Proof.* First, we show that all marked vertices, except the one with implicit label $\text{enc}(0)\$^W \text{enc}(0)$,
 237 are visited at most once. Pre-substitution, a marked vertex with implicit label $\text{enc}(i)\$^W \text{enc}(i)$ is
 238 at the end of a unique, branchless path of length W matching $\text{enc}(i)$. By Lemma 6, the only
 239 substitutions to this path made while matching P are substitutions making it match $\text{enc}(i')$, $i' \neq i$.
 240 If this path were modified to match $\text{enc}(i')$, $i' > 0$, then the only way the marked vertex could
 241 be visited twice while matching P is if after traversing the path, another path matching $\W is
 242 taken back to the start of this $\text{enc}(i')$ path. However, any edges leaving this marked vertex are
 243 labeled with $\#$, making this impossible. By similar reasoning, the path matching $\text{enc}(0)$ ending at
 244 a marked vertex is visited at most twice. We now show that each marked vertex is visited at least
 245 once. Suppose some marked vertex is not visited. By Lemma 7, no additional marked vertices are
 246 created. Hence, a marked vertex ending a path matching $\text{enc}(i)$, $i > 0$ is visited at least twice, or
 247 a marked vertex ending a path matching $\text{enc}(0)$ is visited at least three times, a contradiction. \square

248 **Lemma 9.** *If P can be matched in D' with at most δ substitutions to vertex labels of D' , then D*
 249 *has a Hamiltonian cycle.*

250 *Proof.* By Lemma 4, the paths between marked vertices traversed while matching with P correspond
 251 to edges between vertices in D . Combined with marked vertices being visited exactly once from
 252 Lemma 8 (except the marked vertex ending a path matching $\text{enc}(0)$), the walk matched by P in D'
 253 corresponds to a Hamiltonian cycle through D beginning and ending at the vertex labeled 0. \square

254 Finally, we note that the proof of Lemma 6 establishes that there is a unique set of at least
 255 $\Theta(k)$ vertices for every marked vertex. Additionally, at most $4W|E| = \mathcal{O}(k|V|)$ vertices are created
 256 in the reduction. Hence, $|V'| = \Theta(k|V|)$. For constant alphabet sizes, $k = \Theta(\log^2 |V|)$, implying
 257 $k = \Theta(\log^2 |V'|)$, and for $\sigma = 2\lceil \log n \rceil$, $k = \Theta(\log |V|)$, implying $k = \Theta(\log |V'|)$. This completes
 258 the proof of Theorem 1 and subsequent statements regarding k and σ .

259 *Hardness for large alphabets:* One can increase the alphabet size to $\sigma' > \sigma$ for $\sigma' \leq \frac{|V''|}{2}$ ($|V''|$
 260 being the final number of vertices in the instance of Problem 1). This can be done with only slight
 261 modifications to the existing reduction. Additionally, with the modifications, the implicit label
 262 length k remains unaltered. For every symbol $\alpha \in \{\sigma + 1, \sigma + 2, \dots, \sigma'\}$, add a new vertex y labeled
 263 with α to V' . Pick an arbitrary original vertex $x \in V'$ and add the edge (x, y) . Then y can only be
 264 used when matching P if it is matched with the end symbols in P . However, in an optimal solution,
 265 we can assume this is not done, as α matches no symbol in P .

266 3 SETH-Based Hardness for Problem 2 on De Bruijn Graphs

267 **Reduction** The Orthogonal Vectors Problem is defined as follows: given two sets of binary vectors
 268 $A, B \subseteq \{0, 1\}^d$ where $|A| = |B| = N$, determine whether there exists vectors $a \in A$ and $b \in B$
 269 such that their inner product is zero. Conditioned on SETH, a standard reduction shows that this
 270 cannot be solved in time $d^{\Theta(1)}N^{2-\varepsilon}$ for any constant $\varepsilon > 0$ [43].

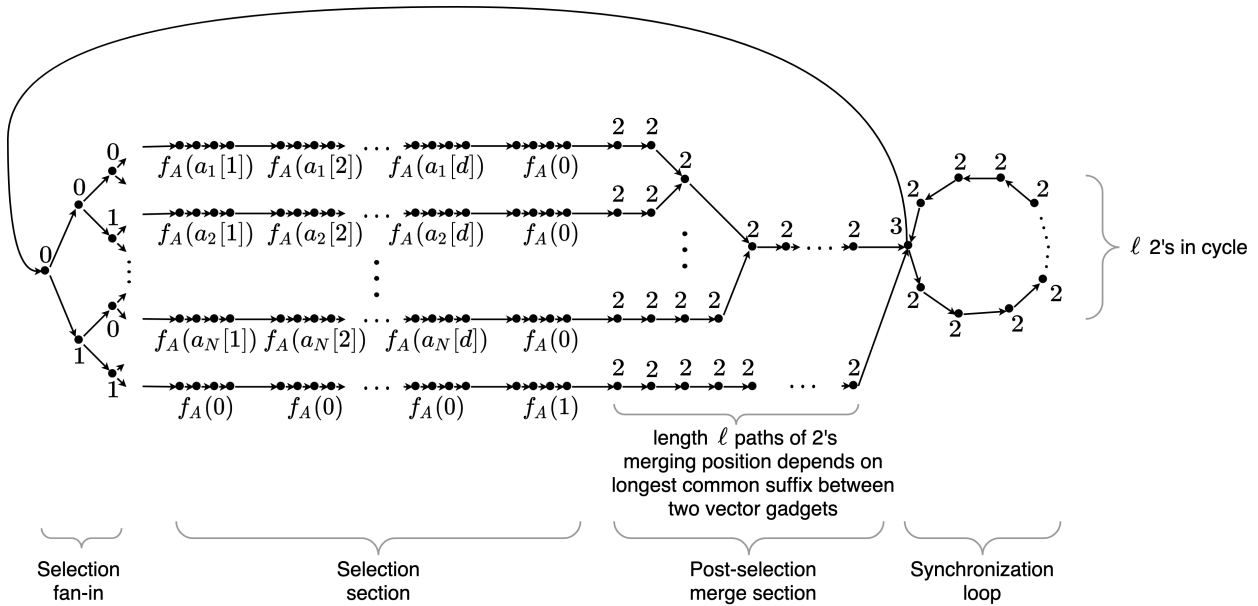


Fig. 6. An illustration of the reduction from OV to Problem 2.

271 Let the given instance of OV consist of $A, B \subseteq \{0, 1\}^d$ where $|A| = |B| = N = 2^m$ for some
 272 natural number m . Hence, we have $\lceil \log(N + 1) \rceil = \log N + 1$. This will ease computation later. We
 273 also assume that $d > \log N$. This is reasonable, as if $d \leq \log N$, then $|A|$ and $|B|$ would contain
 274 either all vectors of length d or repetitions.

275 We will next provide a formal description of the graph D our reduction creates from the set
 276 $A = \{a_1, a_2, \dots, a_N\}$ and the pattern P it creates from the set $B = \{b_1, b_2, \dots, b_N\}$. The reader
 277 may find Figure 6 helpful. The graph will consist of four sections. We name these according to
 278 their function in the reduction: the Selection fan-in, the Selection section, the Post-Selection merge
 279 section, and the Synchronization loop.

280 We start with the Selection fan-in. Let 2^c be the smallest power of 2 such that $2^c \geq N + 1$. The
 281 Selection fan-in consists of a complete binary tree with 2^c leaves where all paths are directed away
 282 from the root. The root is labeled 0 and the children of every node are labeled 0 and 1, respectively.

283 The Selection section consists of $N + 1$ paths. We first define the mappings f_A and f_B from $\{0, 1\}$
 284 to sequences of length four as $f_A(0) = 1100$, $f_A(1) = 1111$, $f_B(0) = 0110$, $f_B(1) = 0000$. These
 285 mappings have the property that $d_H(f_A(0), f_B(0)) = d_H(f_A(0), f_B(1)) = d_H(f_A(1), f_B(0)) = 2$
 286 and $d_H(f_A(1), f_B(1)) = 4$. We make the i^{th} path for $1 \leq i \leq N$ a path of $4(d + 1)$ vertices with
 287 labels matching the string $f_A(a_i[1])f_A(a_i[2])\dots f_A(a_i[d])f_A(0)$. We make the $(N + 1)^{\text{th}}$ path have
 288 $4(d + 1)$ vertices and match the string $f_A(0)^d f_A(1)$. Let s_i denoted the start vertex of path i . We
 289 arbitrarily choose $N + 1$ leaves, l_1, l_2, \dots, l_{N+1} , from the Selection fan-in and add the edges (l_i, s_i)
 290 for $1 \leq i \leq N + 1$.

291 We define the implicit label size as $k = \lceil \log(N + 1) \rceil + 4(d + 1)$ and $\ell = k - 1$. To construct
 292 the Post-selection merge section, we start with $N + 1$ length $\ell - 1$ paths, each matching the string
 293 2^ℓ . For every path in the Selection section, we add an edge from the last vertex in the path to one
 294 of the paths matching 2^ℓ . This is done so that every path matching 2^ℓ in the Post-selection merge
 295 section is connected to exactly one path from the Selection section. Next, we merge two vertices
 296 if they have the same implicit label. This is repeated until all vertices in the Post-selection merge
 297 section have a unique implicit label.

298 To construct the Synchronization loop we create a directed cycle with $\ell + 1 = k$ vertices. One of
 299 these is labeled with the symbol 3, and the rest with the symbol 2. Edges from each ending vertex
 300 in the Post-selection Merge section to the vertex labeled 3 are then added. A final edge from the
 301 vertex labeled 3 to the root of the binary tree in the Selection fan-in completes the graph, which
 302 we denote as D .

Let $t = 5d + \lceil \log(N + 1) \rceil$. To complete the reduction, we make the pattern

$$\begin{aligned}
 P = & (2^\ell 3)^t 2^{\lceil \log(N+1) \rceil} f_B(b_1[1])f_B(b_1[2]) \dots f_B(b_1[d])f_B(1) \\
 & (2^\ell 3)^t 2^{\lceil \log(N+1) \rceil} f_B(b_2[1])f_B(b_2[2]) \dots f_B(b_2[d])f_B(1) \\
 & \dots \\
 & (2^\ell 3)^t 2^{\lceil \log(N+1) \rceil} f_B(b_N[1])f_B(b_N[2]) \dots f_B(b_N[d])f_B(1)
 \end{aligned}$$

303 and the maximum number of allowed substitutions $\delta = N \lceil \log_2(N + 1) \rceil + 2(d + 1) + (2d + 4)(N - 1)$.

304 We call substrings in P of the form $f_B(b_i[1])f_B(b_i[2]) \dots f_B(b_i[d])f_B(1)$ and paths in D matching
 305 strings of the form $f_A(a_i[1])f_A(a_i[2])\dots f_A(a_i[d])f_A(0)$ vector gadgets. Note that $|E| = \mathcal{O}(dN)$ and
 306 $m = |P| = \mathcal{O}(d^2N)$. Hence, an algorithm for approximate matching running in time $\mathcal{O}(m|E|^{1-\varepsilon} +$
 307 $m^{1-\varepsilon}|E|)$ for some $\varepsilon > 0$ would imply an algorithm for OV running in time $d^{\mathcal{O}(1)}N^{2-\varepsilon}$. This implies
 308 that once the correctness of the reduction has been established, Theorem 2 follows.

309 3.1 Proof of Correctness

310 Proofs of Lemma 10 and Lemma 11 are given in Appendix 2.

311 **Lemma 10.** *The graph D is a de Bruijn graph.*

312 **Lemma 11.** *In an optimal solution, 3's in P are matched with 3's in D .*

313 **Lemma 12.** *In an optimal solution, vector gadgets in P are matched with vector gadgets in D .*

314 *Proof.* Suppose otherwise. By Lemma 11, this can only occur if some vector gadget in P is matched
315 against the Synchronization loop. This requires at least $4(d + 1)$ substitutions. We can instead
316 match the $\lceil \log(N + 1) \rceil$ 2's preceding the vector gadget in P with the Selection fan-in and the
317 vector gadget in P with the $(N + 1)^{th}$ path in the Selection section. Due to $d_H(f_A(0), f_B(0)) =$
318 $d_H(f_A(0), f_B(1)) = 2$ and $d(f_A(1), f_B(1)) = 4$, this requires $\lceil \log(N + 1) \rceil + 2d + 4$ substitutions in
319 P . Since, $\log N < d < 2d$ we have $\log N < 2d - 1$. Using that N is some power of 2, $\lceil \log(N + 1) \rceil +$
320 $2d + 4 = \log N + 1 + 2d + 4 < 4d + 4$. Hence, the cost decreases by matching the vector gadget in
321 P to a vector gadget in D instead. \square

322 **Lemma 13.** *If there exists a vector $a \in A$ and $b \in B$ such that $a \cdot b = 0$, then P can be matched
323 to D with at most δ substitutions.*

324 *Proof.* Match the vector gadget for b in P with the vector gadget for a in the Selection section
325 of D . This costs $2(d + 1)$ substitutions. Match the remaining $N - 1$ vector gadgets in P with the
326 $(N + 1)^{th}$ path in the Selection section, requiring $(2d + 4)(N - 1)$ substitutions in total. The total
327 number of substitutions of 2's in P to match the Selection fan-in is $N \lceil \log(N + 1) \rceil$. Adding these,
328 the total number of substitutions is exactly δ . The synchronization loop can be used for matching
329 all additional symbols in P without any further substitutions. \square

330 **Lemma 14.** *If P can be matched in D with at most δ substitutions, then there exists vectors $a \in A$
331 and $b \in B$ such $a \cdot b = 0$.*

332 *Proof.* By Lemma 12, we can assume vector gadgets in P are only matched against vector gadgets
333 in D . Suppose that there does not exist a pair of orthogonal vectors $a \in A$ and $b \in B$. Then, which
334 ever vector gadget in D we choose to match a vector gadget in P to, matching the vector gadget
335 requires at least $2d + 4$ substitutions. Hence, the total cost is at least $(2d + 4)N + N \lceil \log(N + 1) \rceil > \delta$,
336 proving the contrapositive of Lemma 14. \square

337 4 Discussion

338 We leave open several interesting problems. An NP-completeness proof for Problem 1 on de Bruijn
339 graphs when $k = \mathcal{O}(\log n)$ and the alphabet size is constant is still needed. Additionally, we need to
340 extend these hardness results to when substitutions are allowed in both the graph and the pattern,
341 and when insertions and deletions in some form are allowed in the graph and (or) the pattern. It
342 seems unlikely that adding more types of edit operations would make the problems computationally
343 easier, and we conjecture these variants are NP-complete on de Bruijn graphs as well. It also needs
344 to be determined whether Problem 1 is NP-complete on de Bruijn graphs with binary alphabets,
345 or whether the SETH-based hardness results hold for Problem 2 on binary alphabets. A practical
346 question is whether these problems (in general, not just on de Bruijn graphs) are hard for small
347 δ values. In applications, the allowed error thresholds are quite small. Clearly, the problems are
348 slice-wise-polynomial with respect to δ , i.e., for a constant δ it is solvable in polynomial time via
349 brute force, but are they fixed-parameter-tractable in δ ? The reduction presented here (as well as
350 the reductions in [6,25]) is based on the Hamiltonian cycle problem, where a large δ value is used.
351 This makes the existence of such a fixed-parameter-tractable algorithm a distinct possibility.

352 References

- 353 1. Abboud, A., Backurs, A., Hansen, T.D., Williams, V.V., Zamir, O.: Subtree isomorphism revisited. *ACM Trans.*
354 *Algorithms* 14(3), 27:1–27:23 (2018), <https://doi.org/10.1145/3093239>
- 355 2. Abrahamson, K.R.: Generalized string matching. *SIAM J. Comput.* 16(6), 1039–1051 (1987), [https://doi.org/](https://doi.org/10.1137/0216067)
356 [10.1137/0216067](https://doi.org/10.1137/0216067)
- 357 3. Alanko, J., D’Agostino, G., Policriti, A., Prezza, N.: Wheeler languages. *CoRR* abs/2002.10303 (2020), [https:](https://arxiv.org/abs/2002.10303)
358 [//arxiv.org/abs/2002.10303](https://arxiv.org/abs/2002.10303)
- 359 4. Alanko, J.N., Gagie, T., Navarro, G., Benkner, L.S.: Tunneling on wheeler graphs. In: *Data Compression Confer-*
360 *ence, DCC 2019, Snowbird, UT, USA, March 26-29, 2019.* pp. 122–131 (2019), [https://doi.org/10.1109/DCC.](https://doi.org/10.1109/DCC.2019.00020)
361 [2019.00020](https://doi.org/10.1109/DCC.2019.00020)
- 362 5. Almodaresi, F., Sarkar, H., Srivastava, A., Patro, R.: A space and time-efficient index for the compacted colored
363 de bruijn graph. *Bioinform.* 34(13), i169–i177 (2018), <https://doi.org/10.1093/bioinformatics/bty292>
- 364 6. Amir, A., Lewenstein, M., Lewenstein, N.: Pattern matching in hypertext. *J. Algorithms* 35(1), 82–99 (2000),
365 <https://doi.org/10.1006/jagm.1999.1063>
- 366 7. Amir, A., Lewenstein, M., Porat, E.: Faster algorithms for string matching with k mismatches. *J. Algorithms*
367 50(2), 257–275 (2004), [https://doi.org/10.1016/S0196-6774\(03\)00097-X](https://doi.org/10.1016/S0196-6774(03)00097-X)
- 368 8. Backurs, A., Indyk, P.: Which regular expression patterns are hard to match? In: *IEEE 57th Annual Symposium*
369 *on Foundations of Computer Science, FOCS 2016, 9-11 October 2016, Hyatt Regency, New Brunswick, New*
370 *Jersey, USA.* pp. 457–466 (2016), <https://doi.org/10.1109/FOCS.2016.56>
- 371 9. Benoit, G., Lemaître, C., Lavenier, D., Drezen, E., Dayris, T., Uricaru, R., Rizk, G.: Reference-free compression of
372 high throughput sequencing data with a probabilistic de bruijn graph. *BMC Bioinform.* 16, 288:1–288:14 (2015),
373 <https://doi.org/10.1186/s12859-015-0709-7>
- 374 10. Chikhi, R., Limasset, A., Jackman, S., Simpson, J.T., Medvedev, P.: On the representation of de bruijn graphs.
375 *J. Comput. Biol.* 22(5), 336–352 (2015), <https://doi.org/10.1089/cmb.2014.0160>
- 376 11. Chikhi, R., Rizk, G.: Space-efficient and exact de bruijn graph representation based on a bloom filter. *Algorithms*
377 *Mol. Biol.* 8, 22 (2013), <https://doi.org/10.1186/1748-7188-8-22>
- 378 12. Egidi, L., Louza, F.A., Manzini, G.: Space efficient merging of de bruijn graphs and wheeler graphs. *CoRR*
379 *abs/2009.03675* (2020), <https://arxiv.org/abs/2009.03675>
- 380 13. Equi, M., Grossi, R., Mäkinen, V., Tomescu, A.I.: On the complexity of string matching for graphs. In: *46th*
381 *International Colloquium on Automata, Languages, and Programming, ICALP 2019, July 9-12, 2019, Patras,*
382 *Greece.* pp. 55:1–55:15 (2019), <https://doi.org/10.4230/LIPIcs.ICALP.2019.55>
- 383 14. Flick, P., Jain, C., Pan, T., Aluru, S.: Reprint of ”a parallel connectivity algorithm for de bruijn graphs in
384 metagenomic applications”. *Parallel Comput.* 70, 54–65 (2017), [https://doi.org/10.1016/j.parco.2017.09.](https://doi.org/10.1016/j.parco.2017.09.002)
385 [002](https://doi.org/10.1016/j.parco.2017.09.002)
- 386 15. Gagie, T.: $\$r\$$ -indexing wheeler graphs. *CoRR* abs/2101.12341 (2021), <https://arxiv.org/abs/2101.12341>
- 387 16. Gagie, T., Manzini, G., Sirén, J.: Wheeler graphs: A framework for bwt-based data structures. *Theor. Comput.*
388 *Sci.* 698, 67–78 (2017), <https://doi.org/10.1016/j.tcs.2017.06.016>
- 389 17. Georganas, E., Buluç, A., Chapman, J., Olikek, L., Rokhsar, D., Yelick, K.A.: Parallel de bruijn graph construction
390 and traversal for de novo genome assembly. In: *International Conference for High Performance Computing,*
391 *Networking, Storage and Analysis, SC 2014, New Orleans, LA, USA, November 16-21, 2014.* pp. 437–448 (2014),
392 <https://doi.org/10.1109/SC.2014.41>
- 393 18. Gibney, D.: An efficient elastic-degenerate text index? not likely. In: *String Processing and Information Retrieval*
394 *- 27th International Symposium, SPIRE 2020, Orlando, FL, USA, October 13-15, 2020, Proceedings.* pp. 76–88
395 (2020), https://doi.org/10.1007/978-3-030-59212-7_6
- 396 19. Gibney, D., Hoppenworth, G., Thankachan, S.V.: Simple reductions from formula-sat to pattern matching on
397 labeled graphs and subtree isomorphism. In: *4th Symposium on Simplicity in Algorithms, SOSA 2021, Virtual*
398 *Conference, January 11-12, 2021.* pp. 232–242 (2021), <https://doi.org/10.1137/1.9781611976496.26>
- 399 20. Gibney, D., Thankachan, S.V.: On the hardness and inapproximability of recognizing wheeler graphs. In: *27th*
400 *Annual European Symposium on Algorithms, ESA 2019, September 9-11, 2019, Munich/Garching, Germany.* pp.
401 51:1–51:16 (2019), <https://doi.org/10.4230/LIPIcs.ESA.2019.51>
- 402 21. Heydari, M., Miclotte, G., de Peer, Y.V., Fostier, J.: Browniealigner: accurate alignment of illumina se-
403 quencing data to de bruijn graphs. *BMC Bioinform.* 19(1), 311:1–311:10 (2018), <https://doi.org/10.1186/s12859-018-2319-7>
- 404 22. Holley, G., Peterlongo, P.: Blastgraph: intensive approximate pattern matching in string graphs and de-bruijn
405 graphs. In: *PSC 2012* (2012)

- 407 23. Holley, G., Wittler, R., Stoye, J., Hach, F.: Dynamic alignment-free and reference-free read compression. *J.*
408 *Comput. Biol.* 25(7), 825–836 (2018), <https://doi.org/10.1089/cmb.2018.0068>
- 409 24. Hoppenworth, G., Bentley, J.W., Gibney, D., Thankachan, S.V.: The fine-grained complexity of median and center
410 string problems under edit distance. In: 28th Annual European Symposium on Algorithms, ESA 2020, September
411 7–9, 2020, Pisa, Italy (Virtual Conference). pp. 61:1–61:19 (2020), [https://doi.org/10.4230/LIPIcs.ESA.2020.](https://doi.org/10.4230/LIPIcs.ESA.2020.61)
412 61
- 413 25. Jain, C., Zhang, H., Gao, Y., Aluru, S.: On the complexity of sequence to graph alignment. In: Research in
414 Computational Molecular Biology - 23rd Annual International Conference, RECOMB 2019, Washington, DC,
415 USA, May 5–8, 2019, Proceedings. pp. 85–100 (2019), https://doi.org/10.1007/978-3-030-17083-7_6
- 416 26. Kamal, M.S., Parvin, S., Ashour, A.S., Shi, F., Dey, N.: De-bruijn graph with mapreduce framework towards
417 metagenomic data classification. *International Journal of Information Technology* 9(1), 59–75 (2017)
- 418 27. Kapun, E., Tsarev, F.: On np-hardness of the paired de bruijn sound cycle problem. In: Algorithms in Bioinform-
419 matics - 13th International Workshop, WABI 2013, Sophia Antipolis, France, September 2–4, 2013. Proceedings.
420 pp. 59–69 (2013), https://doi.org/10.1007/978-3-642-40453-5_6
- 421 28. Kavya, V.N.S., Tayal, K., Srinivasan, R., Sivadasan, N.: Sequence alignment on directed graphs. *J. Comput. Biol.*
422 26(1), 53–67 (2019), <https://doi.org/10.1089/cmb.2017.0264>
- 423 29. Li, D., Liu, C., Luo, R., Sadakane, K., Lam, T.W.: MEGAHIT: an ultra-fast single-node solution for large and
424 complex metagenomics assembly via succinct *de Bruijn* graph. *Bioinform.* 31(10), 1674–1676 (2015), <https://doi.org/10.1093/bioinformatics/btv033>
- 425 30. Limasset, A., Cazaux, B., Rivals, E., Peterlongo, P.: Read mapping on de bruijn graphs. *BMC Bioinform.* 17,
426 237 (2016), <https://doi.org/10.1186/s12859-016-1103-9>
- 427 31. Limasset, A., Flot, J., Peterlongo, P.: Toward perfect reads: self-correction of short reads via mapping on de
428 bruijn graphs. *Bioinform.* 36(2), 651 (2020), <https://doi.org/10.1093/bioinformatics/btz548>
- 429 32. Lin, Y., Shen, M.W., Yuan, J., Chaisson, M., Pevzner, P.A.: Assembly of long error-prone reads using de bruijn
430 graphs. In: Research in Computational Molecular Biology - 20th Annual Conference, RECOMB 2016, Santa
431 Monica, CA, USA, April 17–21, 2016, Proceedings. p. 265 (2016), [https://link.springer.com/content/pdf/](https://link.springer.com/content/pdf/bbm%3A978-3-319-31957-5%2F1.pdf)
432 [bbm%3A978-3-319-31957-5%2F1.pdf](https://link.springer.com/content/pdf/bbm%3A978-3-319-31957-5%2F1.pdf)
- 433 33. Liu, B., Guo, H., Brudno, M., Wang, Y.: debga: read alignment with de bruijn graph-based seed and extension.
434 *Bioinform.* 32(21), 3224–3232 (2016), <https://doi.org/10.1093/bioinformatics/btw371>
- 435 34. Morisse, P., Lecroq, T., Lefebvre, A.: Hybrid correction of highly noisy long reads using a variable-order de bruijn
436 graph. *Bioinform.* 34(24), 4213–4222 (2018), <https://doi.org/10.1093/bioinformatics/bty521>
- 437 35. Navarro, G.: Improved approximate pattern matching on hypertext. *Theor. Comput. Sci.* 237(1–2), 455–463
438 (2000), [https://doi.org/10.1016/S0304-3975\(99\)00333-3](https://doi.org/10.1016/S0304-3975(99)00333-3)
- 439 36. Pell, J., Hintze, A., Canino-Koning, R., Howe, A., Tiedje, J.M., Brown, C.T.: Scaling metagenome sequence
440 assembly with probabilistic de bruijn graphs. *Proc. Natl. Acad. Sci. USA* 109(33), 13272–13277 (2012), <https://doi.org/10.1073/pnas.1121464109>
- 441 37. Peng, Y., Leung, H.C.M., Yiu, S., Chin, F.Y.L.: IDBA - A practical iterative de bruijn graph de novo assembler.
442 In: Research in Computational Molecular Biology, 14th Annual International Conference, RECOMB 2010, Lisbon,
443 Portugal, April 25–28, 2010. Proceedings. pp. 426–440 (2010), [https://doi.org/10.1007/978-3-642-12683-3_](https://doi.org/10.1007/978-3-642-12683-3_28)
444 28
- 445 38. Peng, Y., Leung, H.C.M., Yiu, S., Lv, M., Zhu, X., Chin, F.Y.L.: Idba-tran: a more robust de novo de bruijn
446 graph assembler for transcriptomes with uneven expression levels. *Bioinform.* 29(13), 326–334 (2013), <https://doi.org/10.1093/bioinformatics/btt219>
- 447 39. Pevzner, P.A.: 1-tuple dna sequencing: computer analysis. *Journal of Biomolecular structure and dynamics* 7(1),
448 63–73 (1989)
- 449 40. Plesník, J.: The np-completeness of the hamiltonian cycle problem in planar digraphs with degree bound two.
450 *Inf. Process. Lett.* 8(4), 199–201 (1979), [https://doi.org/10.1016/0020-0190\(79\)90023-1](https://doi.org/10.1016/0020-0190(79)90023-1)
- 451 41. Rautiainen, M., Marschall, T.: Aligning sequences to general graphs in $o(v + me)$ time. *bioRxiv* p. 216127 (2017)
- 452 42. Ren, X., Liu, T., Dong, J., Sun, L., Yang, J., Zhu, Y., Jin, Q.: Evaluating de bruijn graph assemblers on 454
453 transcriptomic data. *PLoS one* 7(12), e51188 (2012)
- 454 43. Williams, V.V.: Hardness of easy problems: Basing hardness on popular conjectures such as the strong exponential
455 time hypothesis (invited talk). In: 10th International Symposium on Parameterized and Exact Computation,
456 IPEC 2015, September 16–18, 2015, Patras, Greece. pp. 17–29 (2015), [https://doi.org/10.4230/LIPIcs.IPEC.](https://doi.org/10.4230/LIPIcs.IPEC.2015.17)
457 2015.17
- 458 44. Ye, Y., Tang, H.: Utilizing de bruijn graph of metagenome assembly for metatranscriptome analysis. *Bioinform.*
459 32(7), 1001–1008 (2016), <https://doi.org/10.1093/bioinformatics/btv510>
- 460 45. Zerbino, D.R., Birney, E.: Velvet: algorithms for de novo short read assembly using de bruijn graphs. *Genome*
461 *research* 18(5), 821–829 (2008)

Appendix

1 Missing Proofs in Section 2.1

468 **Lemma 1.** *The graph D' constructed as above is a de Bruijn graph.*

469 *Proof.* There are three properties that must be proven: (i) Implicit labels are unique, meaning for
 470 every implicit label at most one vertex is assigned that label; (ii) There are no edges missing, i.e.,
 471 if the implicit label of $y \in V'$ is $S\alpha$ for some string $S[1, k-1]$ and symbol $\alpha \in \Sigma$, and there
 472 exists a vertex $x \in V'$ with implicit label $\beta S[1, k-1]$ for some symbol $\beta \in \Sigma$, then $(x, y) \in E'$; (iii)
 473 Implicit labels are well-defined, in that every walk of length $k-1$ ending at a vertex $x \in V'$ matches
 474 the same string (the implicit label of x). Property (i) holds since after every edge transformation,
 475 vertices with the same implicit label are merged, making every implicit label occur at most once. For
 476 property (ii), consider the completed D' and an arbitrary vertex y on an arbitrary path $\phi((u, v))$.
 477 Regarding a possible edge $(x, y) \in E'$, we have the following cases:

- 478 – Case: the implicit label of y is $S\alpha = \text{enc}(L(u))\$^W \text{enc}(L(u))$. Then, any potential $x \in V'$ must
 479 have an implicit label $\beta S = \beta \text{enc}(L(u))\$^W \text{enc}(L(u))[1, W-1]$. However, the only implicit labels
 480 created that have a suffix of the form $\text{enc}(L(u))\$^W \text{enc}(L(u))[1, W-i]$ have a prefix $\#^{W-i}$. This
 481 implies that $\beta = \#$, and the edge (x, y) already exists in E' (under the assumption that there
 482 are no vertices with in-degree zero in V).
- 483 – Case: the implicit label of y is $S\alpha = \text{enc}(L(u))[i, W]\$^W \text{enc}(L(u))\#^{i-1}$, $1 < i \leq W+1$. Then,
 484 any potential x must have an implicit label $\beta S = \beta \text{enc}(L(u))[i, W]\$^W \text{enc}(L(u))\#^{i-2}$. Because
 485 the only implicit labels with the substring $\$^W \text{enc}(L(u))$ have a prefix consisting of some suffix
 486 of $\text{enc}(L(u))$, this implies $\beta = \text{enc}(L(u))[i-1]$, and (x, y) already exists in E' .
- 487 – Case: the implicit label of y is $S\alpha = \$^{W-i} \text{enc}(L(u))\#^W \text{enc}(L(v))[1, i]$, $1 \leq i \leq W$. Then, any
 488 potential x must have an implicit label $\beta S = \beta \$^{W-i} \text{enc}(L(u))\#^W \text{enc}(L(v))[1, i-1]$. In the
 489 case $i < W$, since there are no runs of $\$$ in any implicit label shorter than W , this implies that
 490 $\beta = \$$, and the edge (x, y) already exists in E' . In the case where $i = W$, the only implicit label
 491 with a suffix of the form $\text{enc}(L(u))\#^W \text{enc}(L(v))[1, W-1]$, has a prefix $\$$, and the edge (x, y)
 492 already exists in E' .
- 493 – Case: the implicit label of y is $S\alpha = \text{enc}(L(u))[i, W]\#^W \text{enc}(L(v))\$^{i-1}$, $1 < i \leq W+1$. Then,
 494 any potential x must have an implicit label $\beta S = \beta \text{enc}(L(u))[i, W]\#^W \text{enc}(L(v))\$^{i-2}$. Because
 495 the only implicit labels with the substring $\#^W \text{enc}(L(v))$ have a prefix consisting of some suffix
 496 of $\text{enc}(L(u'))$ where the edge (u', v) is in D , the edge (x, y) already exists in E' . This is an
 497 interesting case, as merges can happen, i.e., $\beta \text{enc}(L(u))[i, W] = \text{enc}(L(u'))[i-1, W]$, $u' \neq u$.
- 498 – Case: the implicit label of y is $S\alpha = \#^{W-i} \text{enc}(L(v))\$^W \text{enc}(L(v))[1, i]$, $1 \leq i \leq W$. Then, any
 499 potential x must have an implicit label $\beta S = \beta \#^{W-i} \text{enc}(L(v))\$^W \text{enc}(L(v))[1, i-1]$. For $i < W$,
 500 since there are no runs of $\#$ in any implicit label shorter than W , $\beta = \#$, and the edge (x, y)
 501 already exists in E' . For $i = W$, this is equivalent to the first case.

502 We prove (iii) using induction on the number of edges transformed into paths. Our inductive
 503 hypothesis (IH) is that prior to an edge being replaced by a path, property (iii) holds for every
 504 vertex added to V' thus far. Let i denote the number of edges transformed. For $i = 1$, all vertices

505 where there exists such a walk ending at them are on the newly created path, and implicit labels
 506 are well-defined.

507 For $i > 1$, we assume the IH holds for all vertices created in the previous $i - 1$ steps of
 508 transforming edges and merging. First consider a new vertex x that is created by transforming the
 509 i^{th} edge (u_i, v_i) . Starting with $x = \phi(u_i)$. If x is merged with another transformed vertex x' having
 510 the same implicit label, then all length $k - 1$ walks ending at x' match this implicit label, and
 511 thus the IH holds for x after merging. Using a secondary induction step, we assume the IH holds
 512 post-merging for all vertices between $\phi(u_i)$ and x (not including x) on $\phi((u_i, v_i))$. Let x_{prev} be the
 513 vertex on $\phi((u_i, v_i))$ before x . Since all length $k - 1$ walks ending at x_{prev} match x_{prev} 's implicit
 514 label, the length $k - 1$ walks obtained by disregarding the vertex at the start of these walks, and
 515 adding the vertex x at the end, all match the implicit label of x . At the same time, any vertices
 516 merged with x , by the IH also have the same implicit label and hence the walks ending at them
 517 match the implicit label of x . Hence, the IH holds for x after merging it with all vertices having
 518 the same implicit label. After processing all vertices on $\phi((u_i, v_i))$, we next consider a previously
 519 created vertex $x'' \in V'$ not in $\phi((u_i, v_i))$. Consider a newly created walk W ending at x'' that is due
 520 to a vertex merging with vertices in $\phi((u_i, v_i))$. Since all length $k - 1$ walks ending at a vertex in
 521 $\phi((u_i, v_i))$ match the same implicit label, when disregarding some number of vertices at the start
 522 of a walk that ends in $\phi((u_i, v_i))$ and adding new vertices at the end to obtain W , the resulting
 523 walk W matches the implicit label for x'' , and the IH continues to hold for x'' as well. \square

524 **Lemma 2.** *Any walk between two marked vertices $\phi(u)$ and $\phi(v)$ containing no additional marked*
 525 *vertices has length $4W$. Hence, we can conclude any such walk is a path.*

526 *Proof.* We first define forward distance and backward distance. Let $x, y \in V'$. The forward distance
 527 from x to y is defined as the minimum number of edges on any path from x to y (the usual distance
 528 in a directed graph). The backward distance from x to y is defined as the minimum number of
 529 edges on any path from y to x . We say a marked vertex $\phi(u)$ is backward adjacent to x if there
 530 exists a walk from $\phi(u)$ to x not containing any other marked vertices, and x is forward adjacent
 531 $\phi(u)$ if there exists a walk from x to $\phi(u)$ not containing any other marked vertices.

532 We use induction on the number of edges transformed. Our inductive hypothesis (IH) will be
 533 that the length of all walks that end at and contain only two marked vertices is $4W$. We add to
 534 our IH that a vertex x created from an edge transformation having an implicit label of the form:

- 535 1. $\text{enc}(L(u))[j, W]\$^W \text{enc}(L(u))\#^{j-1}$, $1 \leq j \leq W$, has backward distance $j - 1$ from $\phi(u)$, which
 536 is its only backward adjacent marked vertex, and forward distance $4W - j + 1$ from all of its
 537 forward adjacent marked vertices;
- 538 2. $\$^{W-j} \text{enc}(L(u))\#^W \text{enc}(L(v))[1, j]$, $0 \leq j \leq W$, has backward distance $W + j$ from $\phi(u)$, which
 539 is its only backward adjacent marked vertex, and forward distance $3W - j$ from all of its forward
 540 adjacent marked vertices;
- 541 3. $\text{enc}(L(u))[j, W]\#^W \text{enc}(L(v))\$^{j-1}$, $1 \leq j \leq W$, has backward distance $2W + j - 1$ from all of
 542 its backward adjacent marked vertices, and forward distance $2W - j + 1$ from $\phi(v)$, which is its
 543 only forward adjacent marked vertex;
- 544 4. $\#^{W-j} \text{enc}(L(v))\$^W \text{enc}(L(v))[1, j]$, $0 \leq j \leq W$, has backward distance $3W + j$ from all of its
 545 backward adjacent marked vertices, and forward distance $W - j$ from $\phi(v)$, which is its only
 546 forward adjacent marked vertex.

547 The base case, $i = 1$, is satisfied since there exists only one such path and all stated properties hold.
 548 Now, for $i > 1$, let (u_i, v_i) be the i^{th} edge transformed. We assume the IH holds for all vertices and

549 walks created in the first $i - 1$ edge transformations. First, observe that for any walk ending at, and
 550 containing only two previously created marked vertices, for all vertices on this walk the distances
 551 from their forward adjacent marked vertices and backward adjacent marked vertices will not be
 552 altered unless one of the vertices on this walk is merged with a vertex on $\phi((u_i, v_i))$. Also, all of the
 553 stated properties in the IH also hold for $\phi((u_i, v_i))$ prior to merging any vertices. Now, let y be a
 554 vertex on $\phi((u_i, v_i))$. Starting with $y = \phi(u_i)$, and continuing from $\phi(u_i)$ to $\phi(v_i)$, we merge y with
 555 existing vertices when their implicit labels match. Because the stated distance properties hold for
 556 x and y prior to merging, they continue to hold for the vertex created from merging x and y as
 557 well. Moreover, for all of the vertices on any walk containing this now merged vertex the distances
 558 from its forward adjacent and backward adjacent marked vertices are unaltered. Because for every
 559 vertex in the new graph, these distances are unaltered, the IH regarding the length of $4W$ for walks
 560 containing only two marked vertices continues to hold as well. \square

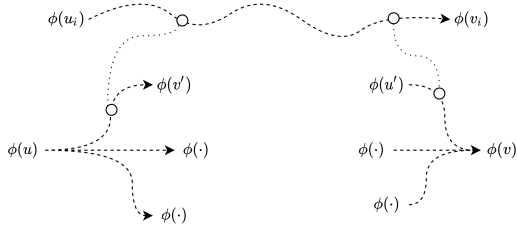


Fig. 7. In the proof of Lemma 4, we consider whether the path $\phi((u_i, v_i))$ being added could potentially cause a path between $\phi(u)$ and $\phi(v)$. The white circles connected by the thin dashed curve represent merged vertices.

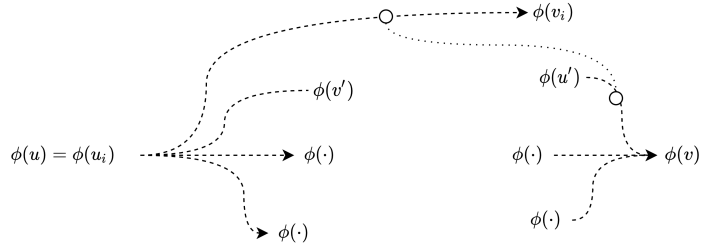


Fig. 8. In the proof of Lemma 4, the case where $u = u_i$ and $v' \neq v_i$.

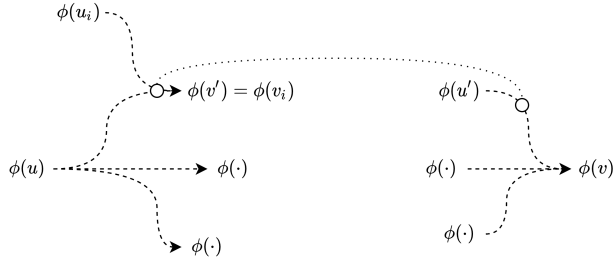


Fig. 9. In the proof of Lemma 4, the case where $u \neq u_i$, $v' = v_i$, and $u_i \neq u'$.

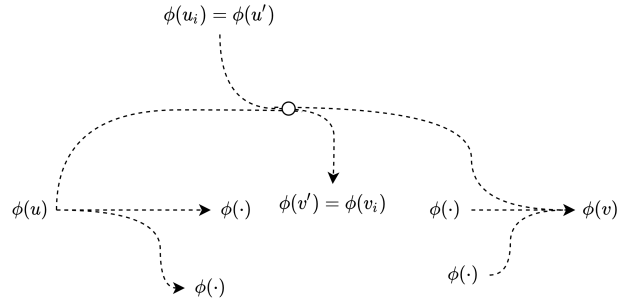


Fig. 10. In the proof of Lemma 4, the case where $u \neq u_i$, $v' = v_i$, and $u_i = u'$.

561 **Lemma 4.** *There exists a path from a marked vertex $\phi(u) \in V'$ to a marked vertex $\phi(v) \in V'$*
 562 *containing no other marked vertices iff there is an edge $(u, v) \in E$.*

563 *Proof.* It is clear from construction that if there is an edge $(u, v) \in D$, then such a walk is in D' .

564 In the other direction, suppose for the sake of contradiction that there exists such a walk starting
 565 at $\phi(u)$ and ending at $\phi(v)$ with no other marked vertices between $\phi(u)$ and $\phi(v)$ on the walk, and

566 $(u, v) \notin E$. Let the first such walk be created when transforming the i^{th} edge (u_i, v_i) . The only way
567 such a walk could exist is if some vertex in $\phi((u_i, v_i))$ is merged with a vertex on a walk $\phi((u, v'))$
568 for some $v' \neq v$, and some vertex in $\phi((u_i, v_i))$ merged with a vertex in a walk $\phi((u', v))$ for some
569 $u' \neq u$. This is since, prior to creating $\phi((u_i, v_i))$ all walks starting at $\phi(u)$ encountered some other
570 marked vertex, say $\phi(v')$, before $\phi(v)$. Similarly, there existed some set of marked vertices not
571 including $\phi(u)$ such that every walk containing a marked vertex and ending at $\phi(v)$ must include
572 at least one vertex in this set, say $\phi(u')$. See Figure 7. Consider cases:

- 573 – $u = u_i$ and $v' = v_i$: This contradicts the assumption that (u_i, v_i) is transformed on the i^{th} step.
- 574 – $u = u_i$ and $v' \neq v_i$ (Figure 8): By Lemma 3, since $u_i = u \neq u'$, $\phi((u_i, v_i))$ and $\phi((u', v))$ can only
575 share a vertex if $v_i = v$. However, this implies the edge $(u_i, v_i) = (u, v) \in E$, a contradiction.
- 576 – $u \neq u_i$ and $v' \neq v_i$: We can directly use Lemma 3 to say no such merged vertices exists between
577 $\phi((u, v'))$ and $\phi((u_i, v_i))$.
- 578 – $u \neq u_i$ and $v' = v_i$ (Figure 9): By Lemma 3, if $u_i \neq u'$, then $\phi((u_i, v_i))$ and $\phi((u', v))$ can only
579 share a vertex $v = v_i$. However, this would imply $v = v'$, a contradiction.

580 The more interesting case is if $u_i = u'$ (Figure 10). Any vertex y having an implicit label con-
581 taining $\text{enc}(L(u'))$ and occurring in $\phi((u_i, v_i))$ and $\phi((u', v))$ occurs before (has smaller backward
582 distance to $\phi(u')$) any vertex with implicit label containing $\text{enc}(L(v'))$. At the same time, any
583 vertex x occurring in $\phi((u, v'))$ and $\phi((u_i, v_i))$ has an implicit label containing $\text{enc}(L(v'))$. Since
584 the vertex x occurs later in $\phi((u_i, v_i))$ than any shared vertex y in $\phi((u_i, v_i))$ and $\phi((u', v))$, the
585 only way any vertices in $\phi((u_i, v_i))$ are in a walk from $\phi(u)$ to $\phi(v)$ not containing any other
586 marked vertices is if there is walk from x to y not containing marked vertices, however, the
587 cycle this creates contradicts Lemma 2. \square

588 **Lemma 6.** *If P can be matched in D' with at most δ substitutions to vertex labels of D' , then*
589 *all $\$$'s in P are matched with non-substituted $\$$'s in D' and all $\#$'s in P are matched with non-*
590 *substituted $\#$'s in D' . Consequently, we can assume the only substitutions are to the vertex labels*
591 *corresponding to $\text{bin}(i)$'s within $\text{enc}(i)$'s.*

592 *Proof.* We first make the following observations: pre-substitution of any of the vertex labels in D' ,

- 593 – (1) For all $u \in V$, there is exactly one path in D' that matches $\text{enc}(L(u))\#^W \text{enc}(L(u))[1, W - \ell]$,
594 and all vertices on this path have in-degree and out-degree one. This follows from the only
595 vertices with in-degree greater than one having implicit labels $\text{enc}(L(u))[i, W]\#^W \text{enc}(L(v))\$\^{i-1}$
596 where $W - \ell < i \leq W + 1$ (these vertices have vertex label $\$$). And the vertices with out-degree
597 greater than one having implicit labels of the form $\$\^{W-i} \text{enc}(L(u))\#^W \text{enc}(L(v))[1, i]$ where
598 $W - \ell \leq i \leq W$ (the last ℓ symbols in $\#^W \text{enc}(L(v))$). This path contains the marked vertex
599 $\phi(u)$. Furthermore, all marked vertices are included on exactly one such path.
- 600 – (2) Every maximal walk containing only $\$$ or $\#$ symbols is of length W , and the distance
601 from the end of any maximal walk consisting of only $\$$ symbols (or $\#$ symbols) to the start of a
602 maximal walk consisting of only $\#$ (or $\$$ symbols resp.) is W . This follows from the construction
603 of D' : every vertex added in the construction has an implicit label where all maximal substrings
604 consisting of non- $\$$ or non- $\#$ are of length W , and maximal substrings consisting of $\$$ or $\#$ are
605 of length W .

606 To see the ‘local’ number of substitutions caused by matching a $\#/\$$ -symbol in D' to a $0/1$
607 symbol in P , suppose the matching of $\text{enc}(L(u))$ in P is ‘shifted left’ by $1 \leq s < W$ so that

608 the first s symbols of some $\text{enc}(L(u))$ in P are matched against the last s symbols in some walk
609 of $\$/\#$ -symbols in D' , requiring s substitutions. Now, assuming $s < 2\ell$ and $\sigma > 4$, due to the
610 prefix $(0^{2\ell}1)^\alpha 23\dots(\sigma - 3)$ of $\text{enc}(L(u))$, matching the shifted $\text{enc}(L(u))$ in D' will require at least
611 $\alpha - 1 + (\sigma - 3)$ substitutions that do not involve a $\#$ or $\$$ symbol. When $\sigma = 4$, $\alpha = 2\ell$ the prefix
612 is $(0^{2\ell}1)^{2\ell}$ and at least $2\ell - 1$ substitutions that do not involve a $\#$ or $\$$ symbol are needed.

613 We now look at the number of substitutions needed on a ‘global’ level. Using Lemma 2, it can be
614 inferred that every walk of length $4W$ contains an originally marked vertex. Hence, while matching
615 P' at least $\lfloor |P'|/4W \rfloor = 4Wn/4W = n$ times an originally marked vertex is visited. Because every
616 substring of $P' = P[1, |P| - 4W]$ of length $3W - \ell$ is distinct, every path described in Observation 1
617 is traversed at most once while matching P' . Since each originally marked vertex is on a unique path
618 that can be traversed at most once, and we traverse at least n such paths, we traverse n distinct
619 paths described in Observation 1. We can now use Observation 2 to infer that the substitutions
620 needed to match the shifted patterns in P' must be repeated n times. Hence, to match P' the total
621 number of substitutions involving $\$/\#$ symbols is at least sn . For $s < 2\ell$ and $\sigma > 4$, the total number
622 of substitutions involving only non- $\$/\#$ symbols is at least $(\alpha + \sigma - 4)n$. Using that $\alpha = 2\ell - \sigma + 4$,
623 we have that the total number of substitutions is at least $(s + \alpha + \sigma - 4)n = (s + 2\ell)n > 2\ell(n - 1) = \delta$.
624 When $\sigma = 4$ and $\alpha = 2\ell$ the inequality above becomes $(s + 2\ell - 1)n > 2\ell(n - 1) = \delta$. Independent
625 of σ , when $2\ell \leq s < W$, then α substitutions to match the substring $(0^{2\ell}1)^\alpha$ in P may not be
626 needed, but the total number of substitutions required is still greater than δ since $sn \geq 2\ell n > \delta$.
627 A symmetric argument can be used for when the matching of P to D' is ‘shifted right’ by s so
628 that the last s symbols in $\text{enc}(L(u))$ in P are matched against the first s symbols in some walk of
629 $\$/\#$ -symbols in D' .

630 For $W < s < 4W$, it still holds that all paths described in Observation 1 are traversed exactly
631 once. Combined with Observation 2, we can infer that the substitution cost incurred when making
632 one path of length W originally matching $\#^W$'s match a substring of P without $\#$'s is incurred at
633 least n times. This results in the total number of needed substitutions being at least $nW > \delta$. \square

634 2 Missing Proofs in Section 3.1

635 **Lemma 10.** *The graph D is a de Bruijn graph.*

636 *Proof.* For each of the four graph sections discussed above, we will prove for each vertex in that
637 section that Conditions (i)-(iii) from the proof of Lemma 1 hold. That is, every vertex v , v 's implicit
638 label well-defined, unique, and there are no additional edges that should have v as their head.

639 – Selection fan-in:

- 640 • (well-defined) For any vertex v in the Selection fan-in, there are two paths of length $k - 1$
641 leading to v (one containing vertices labeled with 2's from the Post-selection merge section
642 and one containing vertices labeled with 2's from the Synchronization loop). Both match
643 the same string $2^{\ell'}3B$ where $\ell' < \ell$ and B is a binary string of length at most $\lceil \log N + 1 \rceil$.
- 644 • (unique) The binary string B could only possibly occur again as a suffix the Selection section.
645 However, all implicit labels occurring in that section contain longer binary strings. Hence
646 the implicit label occurs only once in D .
- 647 • (no missing inbound edges) Let u be any vertex such that (u, v) is in D . A vertex v in the
648 Selection fan-in has an implicit label of the form $S\alpha = 2^{\ell'}3B_{i'}[1, i]$, $\ell' < \ell$, $1 \leq i < \lceil \log N \rceil$,
649 $0 \leq i' \leq N + 1$. This implies that u has the implicit label $\beta S = \beta 2^{\ell'}3B_{i'}[1, i - 1]$. Based on

650 the limited number of implicit labels present in D , it must be that $\beta = 2$, and there exists
651 only one such u . Hence, the edge (u, v) already exists.

652 – Selection section:

- 653 • (well-defined) For a vertex v in the Selection section, there are two length $k - 1$ paths leading
654 to v (one with 2's from the Post-selection merge section and one with 2's from the Syn-
655 chronization loop). Both match a string of the form $2^{\ell'} 3B_{i'} f_A(a_i[1])f_A(a_i[2])\dots f_A(a_i[j])[1, h]$
656 where $0 \leq \ell' < \ell$ and $1 \leq h \leq 4$.
- 657 • (unique) If v has a path of length $k - 1$ matching $2^{\ell'} 3B_{i'} f_A(a_i[1])f_A(a_i[2])\dots(f_A(a_i[j]))[1, i]$,
658 then it must be in the Selection section. The substring $B_{i'}$ following the prefix $2^{\ell'} 3$ is distinct,
659 hence this implicit label only occurs once in the Selection section.
- 660 • (no missing inbound edges) Taking u and v as above, if the vertex v has an implicit label of
661 the form $S\alpha = 2^{\ell'} 3B_{i'} f_A(a_i[1])f_A(a_i[2])\dots f_A(a_i[j])[1, h]$, $1 \leq h \leq 4$, this implies that the any
662 potential u has an implicit label $\beta S = \beta 2^{\ell'} 3B_{i'}[1, h - 1]f_A(a_i[1])f_A(a_i[2])\dots f_A(a_i[j])[1, h - 1]$
663 or $\beta S = \beta 2^{\ell'} 3B_{i'}[1, h - 1]f_A(a_i[1])f_A(a_i[2])\dots f_A(a_i[j - 1])$. In either case, $\beta = 2$, and
664 the edge (u, v) already exists. If the vertex v has an implicit label of the form $S\alpha =$
665 $B_{i'} f_A(a_i[1])f_A(a_i[2])\dots f_A(a_i[d])$, then any potential vertex u has an implicit label $\beta S =$
666 $\beta B_{i'} f_A(a_i[1])f_A(a_i[2])\dots f_A(a_i[d])[1, 3]$ where β must be 3, and the edge (u, v) already exists.

667 – Post-selection merge section:

- 668 • (well-defined) For a vertex v in this section, all length $k - 1$ paths ending at v match a string
669 of the form $B2^{\ell'}$ where B is a binary string. By construction, the paths ending at v match
670 the same string (they were merged based on this condition).
- 671 • (unique) Again by construction, if another vertex v' in the Post-selection merging section
672 has a length $k - 1$ path ending at it that matches v 's implicit label v' will be merged with v .
673 At the same time, vertices in the other sections of D will not have an implicit label of the
674 form $B2^{\ell'}$.
- 675 • (no missing inbound edges) Taking u and v as above, vertex v has an implicit label of the
676 form $S\alpha = B2^{\ell'}$, $\ell' \geq 1$, this implies that any potential vertex u has an implicit label
677 $\beta S = \beta B2^{\ell'-1}$. Such a vertex u is already in the Post-selection merge section or is a vertex
678 at the end of a path in the Selection section (if $\ell' = 1$). Since appending a 2 and removing
679 β will make the implicit label of u equal to the implicit label of v , the vertex at the head of
680 the edge with tail u must have been merged with v . Hence, the edge (u, v) already exists.

681 – Synchronization loop:

- 682 • (well-defined) There are two length $k - 1$ paths to a vertex v in the synchronization loop.
683 Both match the string $2^{\ell'} 32^{\ell''}$ where $\ell' + \ell'' = k - 1 = \ell$, and ℓ' depends on v 's position
684 within the Synchronization loop.
- 685 • (unique) An implicit label for a vertex in any other section contains a symbol that is not a
686 2 or a 3. Within the synchronization loop, each implicit label clearly occurs exactly once.
- 687 • (no missing inbound edges) Taking u and v as above, vertex v has an implicit label of the form
688 $S\alpha = 2^{\ell'} 32^{\ell''}$. This implies that any potential vertex u has an implicit label $\beta S = \beta 2^{\ell'} 32^{\ell''-1}$.
689 If $\ell' < \ell$ it must be that $\beta = 2$ and the edge (u, v) already exists. If instead $\ell' = \ell$, then for
690 both $\beta S = 02^{\ell}3$ and $\beta S = 12^{\ell}3$ there already exists an edge (u, v) as well. \square

691 **Lemma 11.** *In an optimal solution, 3's in P are matched with 3's in D .*

692 *Proof.* Suppose that some 3 in P is not matched with 3 in D or with the final vertex in a path in
693 the Selection section. Since any walk between 3's in D has a length that is a multiple of k and 3 in
694 P are $k - 1$ symbols apart, all 3's must then not be matched with 3 in D . This requires at least tN
695 substitutions within P . On the other hand, when 3's in P are matched with 3's in D , there exists
696 a solution requiring at most $4d(N + 1) + N\lceil\log(N + 1)\rceil$. Specifically, this is obtained by matching
697 each vector gadget in P , $f_B(b_i[1])\dots f_B(b_i[d])$ to the $N + 1^{th}$ path in the Selection section. Since
698 $t = 5d + \lceil\log(N + 1)\rceil > 4d + \frac{4d}{N} + \lceil\log(N + 1)\rceil$ for $d = o(N)$ and N large enough, we can assume
699 that $tN > 4d(N + 1) + N\lceil\log(N + 1)\rceil$. Hence, all 3's in P are matched with the 3 in D or with
700 some final vertex in a path in the Selection section

701 Next, suppose some 3 in P is matched with the last vertex in a path in the Selection section. We
702 consider the first such occurrence. In the case where this occurrence of 3 in P is followed in P by
703 a substring $2^{\lceil\log(N+1)\rceil} f_B(a_i[1])\dots f_B(a_i[d]) f_B(1)$, a cost of at least $8(d + 1)$ is incurred, first at least
704 $4(d + 1)$ from matching the substring $2^\ell 3$ in P to a path through Selection fan-in and the Selection
705 section, then an additional $4(d + 1)$ from matching a vector gadget in P to a path of 2's in the
706 Post-selection merge section. We could have instead matched the Synchronization loop twice with
707 a cost of only $4(d + 1)$ substitutions, and started and ended at the same vertex while still matching
708 $2^\ell 3 2^{\lceil\log(N+1)\rceil} f_B(a_i[1])\dots f_B(a_i[d]) f_B(1)$. Hence, in this case, matching 3 in P with the last vertex in
709 a path in the Selection section is suboptimal. In the case where the occurrence of 3 in P is followed
710 in P by $2^\ell 3$, then the cost incurred is only $4(d + 1)$. However, we could have instead matched $2^\ell 3 2^\ell 3$
711 with the Synchronization loop twice with a substitution cost of 0, and again started and ended at
712 the same vertex. Hence, matching 3 in P with the last vertex in a path in the Selection section is
713 again suboptimal. \square