

Approximation Guarantees for Shortest Superstrings: Simpler and Better

Matthias Englert  

University of Warwick, United Kingdom

Nicolaos Matsakis  

Charles University, Czech Republic

Pavel Veselý  

Charles University, Czech Republic

Abstract

The Shortest Superstring problem is an NP-hard problem, in which given as input a set of strings, we are looking for a string of minimum length that contains all input strings as substrings. The Greedy Conjecture (Tarhio and Ukkonen, 1988) states that the GREEDY algorithm, which repeatedly merges the two strings of maximum overlap, is 2-approximate. We have recently shown (STOC 2022) that the approximation guarantee of GREEDY is at most $\frac{13+\sqrt{57}}{6} \approx 3.425$. Before that, the best established upper bound for this was 3.5 by Kaplan and Shafrir (IPL 2005), which improved upon the upper bound of 4 by Blum et al. (STOC 1991). To derive our previous result, we established two incomparable upper bounds on the overlap sum of all cycle-closing edges in an optimal cycle cover and utilized lemmas of Blum et al.

We improve the more involved one of the two bounds and, at the same time, make its proof more straightforward. This results in an improved approximation guarantee of $\frac{\sqrt{67}+2}{3} \approx 3.396$ for GREEDY. Additionally, our result implies an algorithm for the Shortest Superstring problem having an approximation guarantee of $\frac{\sqrt{67}+14}{9} \approx 2.466$, improving slightly upon the previously best guarantee of $\frac{\sqrt{57}+37}{18} \approx 2.475$ (STOC 2022).

2012 ACM Subject Classification Theory of computation → Approximation algorithms analysis

Keywords and phrases Shortest Superstring problem, Approximation Algorithms

Digital Object Identifier 10.4230/LIPIcs.ISAAC.2023.26

Funding Nicolaos Matsakis: Supported by GA ČR project 22-22997S.

Pavel Veselý: Partially supported by GA ČR project 22-22997S and by Center for Foundations of Modern Computer Science (Charles University project UNCE/SCI/004).

1 Introduction

The shortest superstring problem naturally models a scenario when we have a set of overlapping strings which we need to represent in a compressed form. However, unlike in typical lossless data compression such as Lempel-Ziv schemes, we would like the input strings to be human-readable in the result. That is, the compressed representation of input strings should be a string over the same alphabet that contains all of the strings as substrings. This viewpoint of superstrings as compressed representations has been the crux of their very recent application for representing k -mers, which are k -long substrings of a genomic sequence [19]. These k -mers are typically highly overlapping and in such cases, the shortest superstring of k -mers has length close to the theoretical minimum of the number of distinct k -mers.

Formally, we define the Shortest Superstring problem (SSP) as follows: For a given set of strings S (over a fixed alphabet), compute a minimum-length common *superstring* for the input strings, i.e., a string that contains any $s \in S$ as a substring. SSP is a classical and well-studied problem mentioned in several algorithmic textbooks, e.g., [25, 18, 9, 5]. SSP is APX-hard (i.e., it is NP-hard to obtain a $(1 + \varepsilon)$ -approximation for some $\varepsilon > 0$) and



© Matthias Englert, Nicolaos Matsakis, and Pavel Veselý;

licensed under Creative Commons License CC-BY 4.0

34th International Symposium on Algorithms and Computation (ISAAC 2023).

Editors: Satoru Iwata and Naonori Kakimura; Article No. 26; pp. 26:1–26:18

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

remains so even when restricted to binary alphabets or input strings having the same length $r \geq 3$ [24].

Therefore, assuming $P \neq NP$, the best we can hope for are constant-guarantee approximation algorithms. However, determining the best possible constant guarantee is a long-standing open problem, studied for more than three decades. First, Blum et al. [3] designed an algorithm for which they proved an upper bound of 3 on its approximation ratio. Several papers subsequently obtained better approximations using various algorithms [22, 6, 13, 1, 2, 4, 20, 16] and the currently best approximation guarantee is $\frac{37+\sqrt{57}}{18} \approx 2.475$ [7]. In contrast, the hardness result only rules out a 1.003-approximation [12].

Perhaps the most well-known approximation algorithm for SSP is GREEDY which iteratively merges two strings of maximum overlap until only one string remains (if there are more pairs of strings with maximum overlap, we choose arbitrarily). GREEDY is an appealing choice to implement in practice due to its simplicity and close-to-optimal results in experiments [8, 14, 19]. However, the worst-case behavior of GREEDY is far from understood. Blum et al. [3] showed that GREEDY is 4-approximate, an upper bound which was improved to 3.5 by Kaplan and Shafrir [11] and recently, in our previous work, to $\frac{13+\sqrt{57}}{6} \approx 3.425$ [7]. It is easy to see that GREEDY is at least 2-approximate by considering the input $\{c(ab)^k, (ba)^k, (ab)^k c\}$ for $k \rightarrow \infty$ [21]. The *Greedy Conjecture* states that this lower bound is tight [21]. Despite an extensive effort to prove or disprove this, the three works [3, 11, 7] comprise the only improvements to the approximation guarantee of GREEDY since the conjecture was first made.

Our results. We make progress on determining the optimal approximation guarantees of GREEDY and of another, more involved algorithm; the latter one improves the best proven approximation guarantee for SSP. In particular, we show the following theorems.

► **Theorem 1.** *The approximation guarantee of GREEDY is at most $\frac{\sqrt{67}+2}{3} \approx 3.396$.*

► **Theorem 2.** *An algorithm from the literature that combines GREEDY and a Max-ATSP approximation algorithm (outlined in Appendix A.2) computes a superstring of length at most $\frac{\sqrt{67}+14}{9} \approx 2.466$ times the optimal.*

Furthermore, our result implies improved approximation guarantees for two algorithms which are variants of GREEDY established in [3], namely TGREEDY and MGREEDY (outlined in Appendix A.2).

As in previous work, all our improved approximation bounds follow from a better inequality that relates certain overlaps between strings to the cost of the optimal solution.

2 The General Setting and Our Technical Contribution

Preliminaries. The set of input strings is denoted by $S = \{s_1, \dots, s_{|S|}\}$. Without loss of generality, it is assumed that no string of S is a substring of another string of S . The *length* of a string s is the number of its characters and we denote it by $|s| \in \mathbb{Z}^+$. The concatenation of two strings s and t is denoted by st . A substring of s starting at character i and ending at character $j \geq i$ of s is denoted by $s[i, j]$.

By $\text{ov}(s, t)$ we denote the maximum overlap to merge a string s to the left of a string $t \neq s$, i.e., the longest suffix of s that is a prefix of t . By $\text{ov}(s, s)$ we denote the maximum self-overlap of string s with itself, which is smaller than $|s|$. By $\text{pref}(s, t)$ we denote the prefix of s that remains after removing the overlap with t ; thus, $s = \text{pref}(s, t)\text{ov}(s, t)$ and $|\text{pref}(s, t)| = |s| - |\text{ov}(s, t)|$.

2.1 Overlap Graph, Cycle-Closing Edges, and Overlap Inequalities

The overlap graph G_{ov} plays a central role in SSP approximation, including the analysis of GREEDY. It is a complete directed graph with self-loops in which vertices correspond to the input strings, and the weight of each edge (s, t) equals the overlap length $|\text{ov}(s, t)|$.

Note that the optimal solution OPT for a fixed input corresponds to an optimal (maximum overlap) Hamiltonian path in G_{ov} ; however, finding such a path is in general a hard problem. On the other hand, finding an optimal cycle cover CC in G_{ov} can be done efficiently. In particular, in a variant of GREEDY, called MGREEDY, such a cycle cover is produced as a by-product. Observe that the total overlap of edges in CC is only larger than that of the optimal Hamiltonian path OPT; indeed, by adding the edge between the endpoints of OPT, we obtain a Hamiltonian cycle, which is a particular cycle cover (not necessarily optimal).

The GREEDY algorithm can be stated as a heuristic for a Hamiltonian path in G_{ov} : Sort the edges of G_{ov} by their overlap lengths non-increasingly, then go over the sorted list and add the i -th edge e_i to the path unless:

- (i) there would be a vertex of indegree or outdegree more than one after adding e_i (that is, edge e_i shares a head node or a tail node with an edge picked in a previous step), or
- (ii) e_i closes a cycle.

The crucial difference between GREEDY for computing an approximate superstring and MGREEDY for the optimal cycle cover CC is the condition (ii), not present in the latter, i.e., MGREEDY is defined just by condition (i). Call an edge of CC *cycle-closing* if it is the last edge of its cycle added by MGREEDY to CC (i.e., it has the smallest overlap on the cycle, breaking ties arbitrarily).

To obtain a bound on the approximation guarantee of GREEDY, we intuitively need a suitable upper bound on the total overlap of cycle-closing edges, denoted o (strictly speaking, when analyzing GREEDY we consider only the optimal cycle cover of a certain subset of nodes in G_{ov} , but this does not make a difference for our technical contribution; we explain these details in Appendix A.1). Furthermore, the overlap bound should be in terms of the *length* (and not overlap) of OPT.

This intuition was formalized in [3], who proved that $o \leq 3 \cdot n$, where n is the length of the optimal solution OPT. Moreover, they show that such a bound is sufficient for a constant upper bound on the approximation ratio of GREEDY. Later works improved the inequality to $o \leq 2.5 \cdot n$ [11] and to $o < 2.425 \cdot n$ [7]. Our technical contribution is to show that $o < 2.396 \cdot n$.

In fact, these overlap inequalities are proven and applied in a stronger form of $o < n + \beta \cdot w$, where w is a lower bound on n . To define w , we associate each edge (s, t) of the overlap graph G_{ov} also with a *length* which equals the prefix length $|\text{pref}(s, t)| = |s| - |\text{ov}(s, t)|$. Then w is the total length of all edges in the optimal cycle cover CC.

2.2 Main technical result

We now state our main technical contribution.

► **Theorem 3.** *Let S be any input set of strings, and consider an optimal superstring of length n and an optimal cycle cover CC of length w , computed using MGREEDY. Let o be the sum of overlaps of all cycle-closing edges of CC. Then it holds that*

$$o \leq n + \beta \cdot w \quad \text{for } \beta = (\sqrt{67} - 4)/3 \approx 1.396$$

132 The proofs of Theorems 1 and 2 using Theorem 3 are the same as in previous work, but we
 133 provide an outline for completeness. In Appendix A.1 we describe how Theorem 3 implies the
 134 improved upper bound on the approximation guarantees of GREEDY, using another inequality
 135 from Blum et al. [3]. Then, in Appendix A.2, we show how to derive better approximation
 136 guarantees for a family of SSP algorithms that are based on a Max-ATSP approximation
 137 algorithm; the argument is the same as in previous work (e.g., see [4, 15, 16, 7]).

138 2.3 Overview of the proof of Theorem 3

139 We build on our previous work [7], where one of the conceptual contributions was in classifying
 140 the cycles of CC into three main types. To define them, for a cycle c of CC we let

- 141 ■ $o(c)$ = the overlap of the cycle-closing edge of c , i.e., the smallest overlap on cycle c , and
- 142 ■ $w(c)$ = the total length of edges on c , i.e., the sum of prefixes of the edges of c .

143 The classification is done according to the $o(c)/w(c)$ ratio.

144 ► **Definition 4.** For parameter β defined in Theorem 3, a cycle c of CC is

- 145 ■ extra large, if $o(c) \leq \beta \cdot w(c)$,
- 146 ■ large, if $\beta \cdot w(c) < o(c) \leq 2w(c)$, and
- 147 ■ small, if $2w(c) < o(c)$.

148 The intuition behind the names is that short cycles contain highly periodic strings (e.g.,
 149 *abcabcabca*), whereas strings in large cycles are not so periodic (e.g., *abcdeabcd*)

150 In order to prove that $o \leq n + \beta \cdot w$ for $\beta = (\sqrt{67} - 4)/3$, we will assume, without loss
 151 of generality, that CC contains no extra large cycle. This follows by the argument in [7,
 152 Section 5.1], though for a different overlap to length ratio threshold between large and extra
 153 large cycles (which was suitably chosen to match the upper bound $o \leq n + 1.425w$). For
 154 completeness, we repeat the proof in Appendix B.

155 Our analysis in [7] proceeds by showing two incomparable bounds: one better if large
 156 cycles have much larger total length than small cycles, and another one for the other case.
 157 Namely, letting w_s be the sum of lengths of all small cycles and w_ℓ be the sum of lengths of
 158 large cycles, the first upper bound is

$$159 \quad o \leq n + w_s + 1.5w_\ell \tag{1}$$

160 and the second upper bound is

$$161 \quad o \leq n + w_\ell + \frac{31 + 3 \cdot \sqrt{57}}{14} w_s \approx n + w_\ell + 3.832w_s . \tag{2}$$

162 Using the better of (1) and (2) together with $w = w_s + w_\ell$, it follows that $o \leq n + 1.425w$
 163 (recall that the extra large cycles are not taken into account here).

164 Our improvement and simplification comes from a better version of the second upper
 165 bound. Specifically, we show

$$166 \quad o \leq n + w_\ell + (\gamma - 1) \cdot w_s \approx n + w_\ell + 2.884w_s , \tag{3}$$

167 where $\gamma = (\sqrt{67} + 19)/7 \approx 3.884$. In [7], the bound was shown by first modifying the
 168 input in such a way that the overlap graph G_{ov} has the property that all short cycles in the
 169 optimal cycle cover only consist of a single edge that is a self-loop. The analysis is then
 170 done utilizing this somewhat simpler cycle cover. However, the modification of the input
 171 introduces an additional loss that has to be accounted for in the bound. Our analysis is more

172 direct and works with the original optimal cycle cover, which eliminates the need for the
 173 input modification and therefore the additional loss. This brings new technical complications
 174 because certain properties no longer hold in these more general cycle covers. Nevertheless,
 175 we are able to provide a slightly simpler and more straightforward analysis.

176 **Choice of parameters.** To combine the two incomparable bounds, $o \leq n + w_s + 1.5 \cdot w_\ell$
 177 and $o \leq n + (\gamma - 1) \cdot w_s + w_\ell$, we set $\lambda = \frac{1}{2\gamma-3}$. As long as $\gamma \geq 2$, this means $\lambda \in [0, 1]$. We
 178 then multiply the first bound by $(1 - \lambda)$ and the second bound by λ and add them together.
 179 Using $w_s + w_\ell = w$ we get $o \leq n + (\frac{3}{2} - \frac{1}{4\gamma-6}) \cdot w$. In Theorem 3, we want to show that
 180 $o \leq n + \beta \cdot w$ and so if

$$181 \quad \frac{3}{2} - \frac{1}{4\gamma-6} \leq \beta \tag{4}$$

182 we are done. We will also need

$$183 \quad 3 \cdot (\beta - \frac{2}{\gamma-2}) \geq 1 \text{ (for Lemma 6)} \tag{5}$$

184 or equivalently

$$185 \quad \gamma \geq 2 + \frac{6}{3\beta-1} \text{ (for Lemma 12(b))}. \tag{6}$$

186 The maximum of these two lower bounds (4) and (5) on β is minimized for $\gamma = (\sqrt{67}+19)/7$
 187 and at this point both bounds are equal to $(\sqrt{67}-4)/3$, which is our choice for β . Apart
 188 from this, we will use a number of further inequalities that hold for this choice of parameters
 189 (but are not tight). Namely,

$$190 \quad \frac{5}{2} + \frac{1}{2(\beta-1)} \leq \gamma \text{ (for Lemma 12(c))}, \tag{7}$$

$$191 \quad \beta \geq \frac{\gamma}{\gamma-1} \text{ (for Lemma 12(d))}, \text{ and} \tag{8}$$

$$192 \quad \gamma \geq 2 \text{ (for Lemma 12(d))}. \tag{9}$$

194 **3 Analysis**

195 In this section we show our improved second bound $o \leq n + w_\ell + (\gamma - 1) \cdot w_s$, following a
 196 similar general strategy as in [7].

197 **3.1 Proof Outline**

198 Consider a directed Hamiltonian cycle CC_0 of maximum total overlap in G_{ov} . This cycle is
 199 in particular also a (not necessarily maximum) cycle cover. Therefore, the total overlap of
 200 CC_0 must be bounded from above by the total overlap of CC . Our goal is to show something
 201 stronger than this: that there is a gap between the total overlap of CC_0 and the total overlap
 202 of CC that depends in a specific way on the properties of the cycles in CC . Specifically, let \mathcal{L}
 203 and \mathcal{S} denote the sets of large and small cycles in CC , respectively, and let $|CC_i|$ denote the
 204 total overlap of a cycle cover CC_i . Then we want to show that the total overlap $|CC|$ of CC
 205 is by at least

$$206 \quad \sum_{c \in \mathcal{S}} (o(c) - \gamma \cdot w(c)) + \sum_{c \in \mathcal{L}} (o(c) - 2 \cdot w(c)) \tag{10}$$

26:6 Approximation Guarantees for Shortest Superstrings: Simpler and Better

207 larger than the total overlap $|\text{CC}_0|$ of CC_0 . Showing this is sufficient to establish $o \leq$
 208 $n + w_\ell + (\gamma - 1) \cdot w_s$ because

$$\begin{aligned}
 209 \quad n &\geq \sum_{\ell=1}^{|S|} |s_\ell| - |\text{CC}_0| \geq \sum_{\ell=1}^{|S|} |s_\ell| - |\text{CC}| + \sum_{c \in \mathcal{S}} (o(c) - \gamma \cdot w(c)) + \sum_{c \in \mathcal{L}} (o(c) - 2 \cdot w(c)) \\
 210 &\geq \sum_{c \in \mathcal{S}} w(c) + \sum_{c \in \mathcal{L}} w(c) + \sum_{c \in \mathcal{S}} (o(c) - \gamma \cdot w(c)) + \sum_{c \in \mathcal{L}} (o(c) - 2 \cdot w(c)) \\
 211 &= \sum_{c \in \mathcal{S}} (o(c) - (\gamma - 1) \cdot w(c)) + \sum_{c \in \mathcal{L}} (o(c) - w(c)) \\
 212 &= o - (\gamma - 1) \cdot \sum_{c \in \mathcal{S}} w(c) - \sum_{c \in \mathcal{L}} w(c) = o - (\gamma - 1) \cdot w_s - w_\ell. \\
 213
 \end{aligned}$$

214 **Related cycles.** Before proceeding to describe how we show (10), we borrow the following
 215 definition of *related cycles* from [7] that is useful to improve our final bounds slightly. We
 216 note that a simpler version of our proof could still be carried out without this additional
 217 concept, but at the cost of a slightly weaker bound.

218 **► Definition 5.** We define a relation R between cycles as follows. A small cycle c of CC is
 219 related to a large cycle c' of CC if $w(c) \leq (\beta/2 - 1/6) \cdot w(c')$ and there exists an edge e in G_{ov}
 220 that has one endpoint in cycle c , the other endpoint in cycle c' and satisfies $|\text{ov}(e)| \geq \beta \cdot w(c')$.
 221 In this case, we write $(c, c') \in R$.

222 In [7], the following lemma is shown. We use different values for β and γ , but the proof
 223 of the lemma only requires that $3 \cdot (\beta - 2/(\gamma - 2)) \geq 1$ and this is still satisfied for our new
 224 choice of $\beta = (\sqrt{67} - 4)/3$ and $\gamma = (5 - 3\beta)/(3 - 2\beta)$.

225 **► Lemma 6** (Lemma 7.3 in [7]). For every large cycle c' of CC , at most two different small
 226 cycles of CC are related to c' .

227 **Transforming cycle cover CC_0 into CC in small steps.** We analyze the difference of the
 228 total overlap between CC_0 and CC in small steps, gradually changing the Hamiltonian cycle
 229 CC_0 into a sequence of cycle covers $\text{CC}_0, \text{CC}_1, \text{CC}_2, \dots$ until we obtain CC . We modify a cycle
 230 cover CC_i by removing two edges $f = (v', v)$ and $f' = (u, u')$ from $\text{CC}_i \setminus \text{CC}$ and replace
 231 them with the new edges $e = (u, v)$ and $e' = (v', u')$. The resulting set of edges forms a (not
 232 necessarily optimal) cycle cover again. Furthermore, if the edges are chosen such that $e \in \text{CC}$
 233 or $e' \in \text{CC}$ (or both), then the resulting cycle cover is closer to the cycle cover CC in the sense
 234 that the cardinality of the symmetric difference of the corresponding edge sets decreases.

235 For a cycle cover CC_i , let $\mathcal{M}(\text{CC}_i)$ be the set of *small* cycles c in CC for which CC_i
 236 contains no edge with one endpoint in c and the other endpoint being a string not in c . We
 237 define

$$\begin{aligned}
 238 \quad \phi(i) &= \sum_{c \in \mathcal{M}(\text{CC}_i)} \left(\min\{|\text{ov}(\hat{e})| \mid \hat{e} \in \text{CC}_i \text{ connects two strings of } c\} - \gamma \cdot w(c) \right. \\
 239 &\quad \left. - \sum_{c': (c, c') \in R} \left(w(c') - \frac{o(c')}{2} \right) \right). \\
 240
 \end{aligned}$$

241 The idea is to perform such edge swaps to obtain a sequence $\text{CC}_0, \text{CC}_1, \text{CC}_2, \dots, \text{CC}_k = \text{CC}$
 242 of cycle covers, such that each cycle cover CC_i is closer to CC than the previous one CC_{i-1}

243 and such that $|\text{CC}_i| \geq |\text{CC}_0| + \phi(i)$. Then this implies (10) since

$$\begin{aligned}
244 \quad |\text{CC}| - |\text{CC}_0| &= |\text{CC}_k| - |\text{CC}_0| \geq \phi(k) \\
245 \quad &= \sum_{c \in \mathcal{M}(\text{CC})} \left(\min\{|\text{ov}(\hat{e})| \mid \hat{e} \in \text{CC} \text{ connects two strings of } c\} - \gamma \cdot w(c) \right) \\
246 \quad &\quad - \sum_{c':(c,c') \in R} \left(w(c') - \frac{o(c')}{2} \right) \\
247 \quad &= \sum_{c \in \mathcal{S}} \left(\min\{|\text{ov}(\hat{e})| \mid \hat{e} \in \text{CC} \text{ connects two strings of } c\} - \gamma \cdot w(c) \right) \\
248 \quad &\quad - \sum_{c':(c,c') \in R} \left(w(c') - \frac{o(c')}{2} \right) \\
249 \quad &= \sum_{c \in \mathcal{S}} \left(o(c) - \gamma \cdot w(c) - \sum_{c':(c,c') \in R} \left(w(c') - \frac{o(c')}{2} \right) \right) \\
250 \quad &= \sum_{c \in \mathcal{S}} (o(c) - \gamma \cdot w(c)) - \sum_{c \in \mathcal{S}} \sum_{c':(c,c') \in R} \left(w(c') - \frac{o(c')}{2} \right) \\
251 \quad &\geq \sum_{c \in \mathcal{S}} (o(c) - \gamma \cdot w(c)) - \sum_{c \in \mathcal{L}} (2 \cdot w(c) - o(c)), \\
252 \quad &
\end{aligned}$$

253 where the last step follows from Lemma 6 and the fact that for large cycles c' , by definition,
254 $2w(c') \geq o(c')$.

255 We use induction to show that it is possible to construct the desired sequence of cycle
256 covers that satisfies $|\text{CC}_i| \geq |\text{CC}_0| + \phi(i)$. The base case is $i = 0$ and we have $\phi(0) = 0$ because
257 $\mathcal{M}(\text{CC}_0) = \emptyset$. (Strictly speaking, it may happen that $\mathcal{M}(\text{CC}_0) \neq \emptyset$; however, in such a case,
258 the optimal Hamiltonian cycle CC_0 is a small cycle of CC , thus $\text{CC}_0 = \text{CC}$. Moreover, in such
259 a case, (1) implies $o < n + w$.)

260 In the following, we assume that we have a cycle cover CC_i with $|\text{CC}_i| \geq |\text{CC}_0| + \phi(i)$
261 and we show how to construct CC_{i+1} such that $|\text{CC}_{i+1}| \geq |\text{CC}_0| + \phi(i+1)$ and such that
262 the symmetric difference between CC_{i+1} and CC is smaller than the symmetric difference
263 between CC_i and CC . Specifically, we will identify a swap of four edges as described above to
264 obtain CC_{i+1} from CC_i such that:

- 265 ■ one of the edges that are swapped in belongs to CC , which implies that the symmetric
- 266 difference between CC_{i+1} and CC will decrease, and
- 267 ■ $|\text{CC}_{i+1}| - |\text{CC}_i| \geq \phi(i+1) - \phi(i)$.

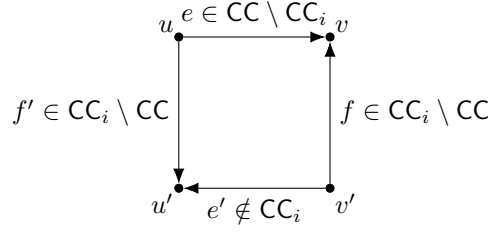
268 This proves the claim due to the induction hypothesis.

269 3.2 Important Lemmas

270 We begin with the following bound on the overlap between two strings from different cycles
271 of CC .

272 ► **Lemma 7** (Lemma 9 in [3]). *Let c and $c' \neq c$ be two cycles in CC . It holds that $|\text{ov}(s, s')| <$
273 $w(c) + w(c')$ for any two strings $s \in c$ and $s' \in c'$.*

274 When changing cycle cover CC_i into CC_{i+1} , we identify an edge $e = (u, v) \in \text{CC} \setminus \text{CC}_i$
275 that we add into CC_{i+1} . This triggers removal of edges $f = (v', v)$ and $f' = (u, u')$ from
276 CC_i and addition of one more edge $e' = (v', u')$ that does not belong to CC_i but may or
277 may not be in CC ; see Figure 1. In the following, we provide several lower bounds on



■ **Figure 1** Illustration of the notation used in lemmas in Section 3.2.

278 $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')|$, which is the total overlap length difference between CC_i
 279 and CC_{i+1} . The first lemma is the well-known *Monge Condition*.

280 ► **Lemma 8** (Lemma 7 in [3]). *Let $e = (u, v)$, $f = (v', v)$, $f' = (u, u')$, $e' = (v', u')$ be edges*
 281 *in G_{ov} , such that $\max\{|\text{ov}(e)|, |\text{ov}(e')|\} \geq \max\{|\text{ov}(f)|, |\text{ov}(f')|\}$. Then $|\text{ov}(e)| + |\text{ov}(e')| -$*
 282 *$|\text{ov}(f)| - |\text{ov}(f')| \geq 0$.*

283 The following lemma is shown in [7, Lemma 7.5] for the special case of inputs where each
 284 small cycle of CC consists of one string. Below, we generalize it for any input and cycle.

285 ► **Lemma 9.** *Let $e = (u, v)$, $f = (v', v)$, $f' = (u, u')$, and $e' = (v', u')$ be edges in G_{ov} such*
 286 *that e is an edge in cycle c in CC . Then,*

$$287 \quad |\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')| > |\text{ov}(e)| - \max\{|\text{ov}(f)|, |\text{ov}(f')|\} - w(c).$$

288 Before proving Lemma 9, we recall a few definitions from the literature. Consider a cycle
 289 c of CC having k nodes s_1, s_2, \dots, s_k . Assuming that the cycle-closing edge of c is (s_k, s_1) ,
 290 we define $s(c)$ as the string $\text{pref}(s_1, s_2)\text{pref}(s_2, s_3) \dots \text{pref}(s_k, s_1)$.

291 A *semi-infinite* string is a string obtained by concatenating an infinite number of finite
 292 strings. A semi-infinite string s is *periodic* if $s = ts$ for a non-empty string t , that is, $s = t^\infty$.

293 A string t is a *factor* of a string s if $s = t^i y$ for an integer $i > 0$, where y is a (possibly
 294 empty) prefix y of t . By $\text{factor}(s)$ of s , we denote the shortest factor of s and we define
 295 $\text{period}(s) = |\text{factor}(s)|$. Finally, we say that a string s has a *periodicity* of length q for $q \leq |s|$
 296 if s is a prefix of the semi-infinite string x^∞ for some string x of length q .

297 Next, we need a basic observation.

298 ► **Observation 10.** *Let s and t be two strings that are substrings of some string z . Then,*
 299 *$|\text{ov}(s, t)| > \min\{|s|, |t|\} - \text{period}(z)$.*

300 **Proof.** We can assume without loss of generality (w.l.o.g.) that $|s| \leq |t|$. This is because,
 301 otherwise, let s_R, t_R , and z_R be the reverse of the strings s, t , and z , respectively. We
 302 observe that $\text{ov}(t_R, s_R) = \text{ov}(s, t)$ and $\text{period}(z_R) = \text{period}(z)$. Clearly also $|s_R| = |s|$,
 303 $|t_R| = |t|$. Therefore, the inequality in the statement of the observation is equivalent to
 304 $|\text{ov}(t_R, s_R)| > \min\{|s_R|, |t_R|\} - \text{period}(z_R)$. Hence, if $|s| > |t|$ then $|t_R| \leq |s_R|$ and we can
 305 apply the arguments below to the strings t_R, s_R , and z_R instead of s, t , and z (in this order).

306 Since s and t are substrings of z we can write them as $s = z[i, i + |s| - 1]$ and $t =$
 307 $z[j, j + |t| - 1]$ for some i and j . Because of the period of z , we can assume that $i \in [1, \text{period}(z)]$
 308 and $j \in [1, \text{period}(z)]$.

309 ■ If $j \geq i$, we have $\text{ov}(s, t) = z[j, i + |s| - 1]$ and hence $|\text{ov}(s, t)| = i - j + |s| > |s| - \text{period}(z)$.

- 310 ■ If $j < i$ and $j + \text{period}(z) > |z|$, then $j + \text{period}(z) > |z| \geq i + |s| - 1$ and hence,
 311 $|\text{ov}(s, t)| \geq 0 > j - i \geq |s| - \text{period}(z)$.
- 312 ■ If $j < i$ and $j + \text{period}(z) \leq |z|$, we observe that $t = z[j, j + |t| - 1]$ also has $z[j +$
 313 $\text{period}(z), \min\{j + |t| - 1 + \text{period}(z), |z|\}]$ as a prefix (indeed, if $j + |t| - 1 + \text{period}(z) \leq |z|$
 314 this is not just a prefix of t , but exactly t). Since $i \leq j + \text{period}(z)$ and $|s| \leq |t|$, we
 315 have $\text{ov}(s, t) = z[j + \text{period}(z), i + |s| - 1]$ and hence, $|\text{ov}(s, t)| = i - j + |s| - \text{period}(z) >$
 316 $|s| - \text{period}(z)$. ◀

317 **Proof of Lemma 9.** Since $\text{ov}(f)$ and $\text{ov}(f')$ are substrings of $s(c)^\infty$, we use Observation 10
 318 to get

$$319 \quad |\text{ov}(e')| \geq |\text{ov}(\text{ov}(f), \text{ov}(f'))|$$

$$320 \quad > \min\{|\text{ov}(f)|, |\text{ov}(f')|\} - \text{period}(s(c)^\infty) \geq \min\{|\text{ov}(f)|, |\text{ov}(f')|\} - w(c).$$

322 It follows that

$$323 \quad |\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')|$$

$$324 \quad > |\text{ov}(e)| + \min\{|\text{ov}(f)|, |\text{ov}(f')|\} - w(c) - |\text{ov}(f)| - |\text{ov}(f')|$$

$$325 \quad = |\text{ov}(e)| - \max\{|\text{ov}(f)|, |\text{ov}(f')|\} - w(c).$$

327

328 The following lemma is, also, due to [7]. Here, we state it in a slightly different way, but
 329 the proof is essentially the same and included in Appendix C for completeness.

330 ► **Lemma 11.** Consider the edges $e = (u, v)$, $f = (v', v)$, $f' = (u, u')$, and $e' = (v', u')$
 331 between (not necessarily different) nodes u, u', v, v' in G_{ov} . Suppose u' and v' are strings in
 332 the same cycle c' of CC and that whichever of f or f' has larger overlap connects a string
 333 from cycle c and a string from cycle $c' \neq c$ (if $|\text{ov}(f)| = |\text{ov}(f')|$ then it is sufficient if one of
 334 them satisfies this). If $|\text{ov}(e)| \geq w(c) + w(c')$, then

$$335 \quad |\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')| > |\text{ov}(e')| - w(c').$$

336 The following lemma draws conclusions from the previous ones in a way that will be
 337 useful later for our analysis.

338 ► **Lemma 12.** Consider the edges $e = (u, v)$, $f = (v', v)$, $f' = (u, u')$, and $e' = (v', u')$
 339 between (not necessarily different) nodes u, u', v, v' in G_{ov} . Suppose e is an edge in a cycle c
 340 of CC. Suppose further that $|\text{ov}(e)| \geq \max\{|\text{ov}(f)|, |\text{ov}(f')|\}$ and the edge of f and f' that has
 341 larger overlap connects a string of cycle c and a string of cycle $c' \neq c$ (if $|\text{ov}(f)| = |\text{ov}(f')|$,
 342 then either one of f and f' may satisfy this condition). All of the following statements hold:

- 343 (a) $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')| \geq 0$.
- 344 (b) If $w(c) \geq (\beta/2 - 1/6) \cdot w(c')$, then $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')| \geq |\text{ov}(e)| - \gamma w(c)$.
- 345 (c) If $w(c) \geq (\beta - 1) \cdot w(c')$, then $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')| \geq |\text{ov}(e)| - \gamma w(c) -$
 346 $w(c')/2 + w(c)/2$.
- 347 (d) Furthermore, if v' and u' are strings in the same cycle in CC, then also $|\text{ov}(e)| + |\text{ov}(e')| -$
 348 $|\text{ov}(f)| - |\text{ov}(f')| \geq \max\{|\text{ov}(e')| - \gamma w(c'), |\text{ov}(e)| - \gamma w(c) + |\text{ov}(e')| - \gamma w(c')\}$.

349 **Proof.** We show the relevant lower bounds on $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')|$ separately.

- 350 (a) Due to Lemma 8, we have $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')| \geq 0$.

351 (b) If $w(c) \geq (\beta/2 - 1/6) \cdot w(c')$, due to Lemma 9, we have

$$\begin{aligned}
 352 \quad & |\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')| \geq |\text{ov}(e)| - \max\{|\text{ov}(f)|, |\text{ov}(f')|\} - w(c) \\
 353 \quad & \geq |\text{ov}(e)| - 2w(c) - w(c') \geq |\text{ov}(e)| - \gamma w(c),
 \end{aligned}$$

355 where the second step uses Lemma 7 and the last inequality follows from $2 + 6/(3\beta - 1) = \gamma$.

356 (c) If $w(c) \geq (\beta - 1) \cdot w(c')$, we have due to Lemma 9 that

$$\begin{aligned}
 357 \quad & |\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')| \geq |\text{ov}(e)| - \max\{|\text{ov}(f)|, |\text{ov}(f')|\} - w(c) \\
 358 \quad & \geq |\text{ov}(e)| - 2w(c) - w(c') \\
 359 \quad & = |\text{ov}(e)| - \frac{5}{2}w(c) - w(c')/2 - w(c')/2 + w(c)/2 \\
 360 \quad & \geq |\text{ov}(e)| - \gamma w(c) - w(c')/2 + w(c)/2,
 \end{aligned}$$

362 where the second step uses Lemma 7 and the last inequality follows from $5/2 + 1/(2(\beta - 1)) \leq \gamma$.

364 (d) \blacksquare Suppose v' and u' are strings in the same cycle in CC. If $|\text{ov}(e)| \geq w(c) + w(c')$, we apply
 365 Lemma 11 to get $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')| \geq |\text{ov}(e')| - w(c') \geq |\text{ov}(e')| - \gamma w(c')$.
 366 Otherwise, we have $|\text{ov}(e)| < w(c) + w(c')$ and hence,

$$\begin{aligned}
 367 \quad & |\text{ov}(f)| \leq w(c) + w(c') = w(c') + \gamma w(c) - (\gamma - 1)w(c) \\
 368 \quad & \leq w(c') + (\gamma - 1)o(c) - (\gamma - 1)w(c) \\
 369 \quad & \leq w(c') + (\gamma - 1)|\text{ov}(e)| - (\gamma - 1)w(c) < \gamma w(c'),
 \end{aligned}$$

371 since it holds $\beta \geq \frac{\gamma}{\gamma - 1}$ and $o(c) > \beta w(c)$ for any large or small cycle c (recall that we
 372 assume that CC contains no extra large cycle). We get $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| -$
 373 $|\text{ov}(f')| \geq |\text{ov}(e')| - |\text{ov}(f)| \geq |\text{ov}(e')| - \gamma w(c')$.

374 \blacksquare Suppose v' and u' are strings in the same cycle in CC. Due to Lemma 7,

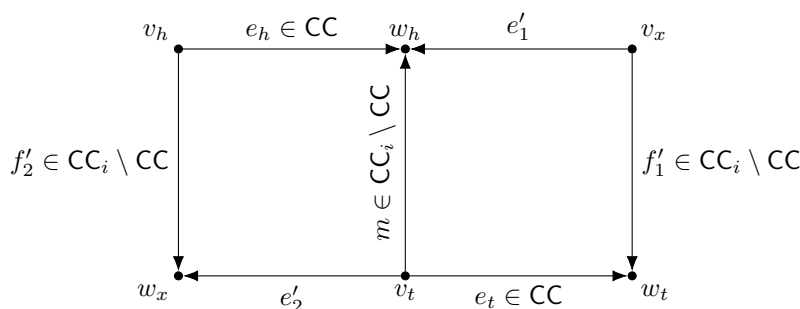
$$\begin{aligned}
 375 \quad & |\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(f)| - |\text{ov}(f')| \geq |\text{ov}(e)| + |\text{ov}(e')| - 2 \max\{|\text{ov}(f)|, |\text{ov}(f')|\} \\
 376 \quad & \geq |\text{ov}(e)| - 2w(c) + |\text{ov}(e')| - 2w(c') \\
 377 \quad & \geq |\text{ov}(e)| - \gamma w(c) + |\text{ov}(e')| - \gamma w(c').
 \end{aligned}$$

379 ◀

380 3.3 The Induction Step

381 We specify how an edge swap is made at a fixed step i in which we obtain cycle cover CC_{i+1}
 382 from CC_i . We start by identifying the largest-overlap edge $m = (v_t, w_h)$ in $\text{CC}_i \setminus \text{CC}$, breaking
 383 ties arbitrarily. Six further edges will be important. First, let $e_h = (v_h, w_h)$ and $e_t = (v_t, w_t)$
 384 be the edges in CC that share heads and tails with m , respectively. Further, let $f'_1 = (v_x, w_t)$
 385 and $f'_2 = (v_h, w_x)$ be the two edges in $\text{CC}_i \setminus \text{CC}$ that share heads with e_t and tails with e_h ,
 386 respectively. Lastly, define $e'_1 = (v_x, w_h)$ and $e'_2 = (v_t, w_x)$. See Figure 2 for a summary of
 387 this notation. It is important to note that the six strings v_h, w_h, v_x, w_x, v_t , and w_t are not
 388 necessarily different.

389 With this, we can define two potential edge swaps. In the first one, we add e_t and e'_1 to
 390 the cycle cover and instead remove m and f'_1 . In the second one, we add e_h and e'_2 to the cycle
 391 cover and instead remove m and f'_2 . Which one of these two swaps we will perform depends
 392 on a few properties of the edges involved. First of all, we assume that $|\text{ov}(e_h)| \geq |\text{ov}(e_t)|$.



■ **Figure 2** Illustration of the notation. Note that we also allow nodes to be equal to one another here, e.g., it could be that $w_t = w_x$, in which case $e_t = e'_2$, $v_h = v_x$, $e_h = e'_1$, and $f'_1 = f'_2$.

393 Otherwise, all the remaining arguments follow symmetrically by considering e_t instead of e_h
 394 and vice versa. Furthermore, we have that

$$395 \quad |\text{ov}(e_h)| \geq |\text{ov}(m)|, \quad (11)$$

396 since otherwise $|\text{ov}(m)| > |\text{ov}(e_h)| \geq |\text{ov}(e_t)|$ and m would be added to CC by the greedy
 397 algorithm for the optimal cycle cover before e_h and e_t , contradicting the choice of m as an
 398 edge of largest overlap in $\text{CC}_i \setminus \text{CC}$.

399 We observe that there are two reasons why $\phi(i+1)$ may be larger than $\phi(i)$.

400 ■ The first potential reason is a difference between the sets $\mathcal{M}(\text{CC}_{i+1})$ and $\mathcal{M}(\text{CC}_i)$. We
 401 know that $\mathcal{M}(\text{CC}_{i+1}) \supseteq \mathcal{M}(\text{CC}_i)$, because if a cycle c is in $\mathcal{M}(\text{CC}_i)$, then there is no edge
 402 in CC_i connecting a string of c to a string of another cycle. That means that the edges f
 403 and f' that we remove from CC_i in the process of constructing CC_{i+1} either have both
 404 their endpoints in c or both their endpoints not in c . If both endpoints of both edges f
 405 and f' are part of c , then also the two edges that are swapped in to obtain CC_{i+1} from
 406 CC_i have their endpoints entirely in c . Therefore, c would still be in $\mathcal{M}(\text{CC}_{i+1})$ after the
 407 swap. If both endpoints of both edges f and f' are outside of c , then also the two edges
 408 that are swapped in to obtain CC_{i+1} from CC_i have their endpoints entirely outside of
 409 c . Again, c would still be in $\mathcal{M}(\text{CC}_{i+1})$ after the swap in this case. Finally, if one of f
 410 and f' has both endpoints in c and the other one has both endpoints outside of c , then
 411 the two edges that are swapped in both have one endpoint in c and the other endpoint
 412 outside of c . However, this is not possible because one of the edges we swap in is e_h or e_t
 413 and must therefore be part of the optimal cycle cover CC .

414 We can further observe that $\mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i)$ must either be equal to \emptyset , $\{c\}$, $\{c'\}$, or
 415 $\{c, c'\}$, where c and c' are the cycles that e_h and e_t belong to in CC , respectively. (It is
 416 possible that $c = c'$.) To see this, observe that one edge being swapped out to obtain
 417 CC_{i+1} from CC_i is m and that m has one endpoint (w_h) in c and the other endpoint (v_t)
 418 in c' . However, for each cycle of CC , it is clear from a parity argument that the number
 419 of edges of CC_i connecting the cycle to other cycles must be even. Hence, for a cycle c'' to
 420 be in $\mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i)$, each of the edges being swapped out must have a string from
 421 cycle c'' as an endpoint. This can only be true for c or c' and not for any other cycle.

422 Overall, if this reason for the difference between $\phi(i+1)$ and $\phi(i)$ applies, we have that

$$\begin{aligned}
 423 \quad \phi(i+1) - \phi(i) &= \sum_{c \in \mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i)} \left(\min\{|\text{ov}(\hat{e})| \mid \hat{e} \in \text{CC}_i \text{ connects two strings of } c\} - \gamma \cdot w(c) \right. \\
 424 \quad &\quad \left. - \sum_{c': (c, c') \in R} \left(w(c') - \frac{o(c')}{2} \right) \right). \\
 425
 \end{aligned}$$

- The second potential reason why $\phi(i+1)$ may be larger than $\phi(i)$ is that for a cycle $c \in \mathcal{M}(\text{CC}_i)$ the term $\min\{|\text{ov}(\hat{e})| \mid \hat{e} \in \text{CC}_i \text{ connects two strings of } c\}$ could change. However, this can only happen if $\mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i) = \emptyset$ and, furthermore, it can only happen for a cycle c when both edges f and f' that are swapped out have both their endpoints in cycle c . In this case, all four strings involved in the swap (either $v_h, w_x, w_h,$ and v_t or $v_x, w_t, w_h,$ and v_t), must be part of the same cycle in CC . If the value $\min\{|\text{ov}(\hat{e})| \mid \hat{e} \in \text{CC}_{i+1} \text{ connects two strings of } c\}$ is larger than the value $\min\{|\text{ov}(\hat{e})| \mid \hat{e} \in \text{CC}_i \text{ connects two strings of } c\}$, then an edge in $\arg \min\{|\text{ov}(\hat{e})| \mid \hat{e} \in \text{CC}_i \text{ connects two strings of } c\}$ must have been swapped out. This means, that if f and f' are the edges being swapped out to obtain CC_{i+1} from CC_i , then $\min\{|\text{ov}(\hat{e})| \mid \hat{e} \in \text{CC}_i \text{ connects two strings of } c\} = \min\{|\text{ov}(f)|, |\text{ov}(f')|\}$. If e and e' are the two edges being swapped in, the new value of $\min\{|\text{ov}(\hat{e})| \mid \hat{e} \in \text{CC}_{i+1} \text{ connects two strings of } c\}$ can be at most $\min\{|\text{ov}(e)|, |\text{ov}(e')|\}$ because e and e' are in CC_{i+1} and satisfy the condition that they connect two strings of c . So overall, in this situation,

$$\phi(i+1) - \phi(i) \leq \min\{|\text{ov}(e)|, |\text{ov}(e')|\} - \min\{|\text{ov}(f)|, |\text{ov}(f')|\}.$$

426 In summary, we note that only one of the two reasons can apply for any fixed step i . If there
 427 is an increase of $\phi(i+1)$ over $\phi(i)$ due to the first reason (a change in the set $\mathcal{M}(\text{CC}_{i+1})$
 428 compared to $\mathcal{M}(\text{CC}_i)$), then there is no increase due to the second reason and vice versa.

429 We are now ready to complete the proof by showing how to select one of the two identified
 430 swap operations such that the total overlap increases by at least $\phi(i+1) - \phi(i)$.

- 431 ■ If m connects two strings of the same cycle in CC , then observe that $\mathcal{M}(\text{CC}_{i+1}) = \mathcal{M}(\text{CC}_i)$.
 432 We swap in e_h and e'_2 and swap out f'_2 and m . Since $|\text{ov}(e_h)| \geq |\text{ov}(m)|$ by (11), we can
 433 apply Lemma 8 and establish that the total overlap does not decrease when this swap is
 434 performed.

Furthermore, if $v_h, w_h, v_t,$ and w_x all belong to the same cycle of CC , then the total
 overlap increases by $|\text{ov}(e_h)| + |\text{ov}(e'_2)| - |\text{ov}(f'_2)| - |\text{ov}(m)| \geq |\text{ov}(e'_2)| - |\text{ov}(f'_2)| \geq$
 $\min\{|\text{ov}(e_h)|, |\text{ov}(e'_2)|\} - \min\{|\text{ov}(f'_2)|, |\text{ov}(m)|\}$, where the second inequality uses $|\text{ov}(f'_2)| \leq$
 $|\text{ov}(m)|$ by the definition of m . This is the only case in which

$$\min\{|\text{ov}(e)| \mid e \text{ is edge of } \text{CC}_i \text{ connecting two strings of cycle } c\}$$

435 can change for a cycle in $c \in \mathcal{M}(\text{CC}_i)$ and the increase is at least $\min\{|\text{ov}(e_h)|, |\text{ov}(e'_2)|\} -$
 436 $\min\{|\text{ov}(f'_2)|, |\text{ov}(m)|\} \geq \phi(i+1) - \phi(i)$, as required.

- 437 ■ If m connects strings of two different cycles in CC and $|\text{ov}(e_t)| \geq |\text{ov}(m)|$. Let c be the
 438 cycle of e_h and c' be the cycle of e_t . If $w(c) \geq w(c')$, we swap in $e = e_h$ and $e' = e'_2$ and
 439 swap out $f' = f'_2$ and m . Otherwise, we swap in $e = e_t$ and $e' = e'_1$ and swap out $f' = f'_1$
 440 and m .

441 We distinguish between these two cases:

- 442 ■ Suppose $w(c) \geq w(c')$.

Then, if $c' \in \mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i)$, Lemma 12(d) applies and we know that the increase in total overlap due to the swap is $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(m)| - |\text{ov}(f')| \geq \max\{|\text{ov}(e')| - \gamma w(c'), |\text{ov}(e)| - \gamma w(c) + |\text{ov}(e')| - \gamma w(c')\} \geq \phi(i+1) - \phi(i)$, as required since $\phi(i+1) - \phi(i)$ is either equal to $|\text{ov}(e')| - \gamma w(c')$ or equal to $|\text{ov}(e)| - \gamma w(c) + |\text{ov}(e')| - \gamma w(c')$ depending on whether $\mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i) = \{c'\}$ or $\mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i) = \{c', c\}$.

Otherwise, if $c' \notin \mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i)$, Lemma 12(a) and (b) both apply and we know that the increase in total overlap due to the swap is $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(m)| - |\text{ov}(f')| \geq \max\{0, |\text{ov}(e)| - \gamma w(c)\} \geq \phi(i+1) - \phi(i)$, as required since $\phi(i+1) - \phi(i)$ is either equal to 0 or equal to $|\text{ov}(e)| - \gamma w(c)$ depending on whether $\mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i) = \emptyset$ or $\mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i) = \{c\}$.

- Suppose $w(c) < w(c')$.

Then, the same argument as above holds with the only difference being that the roles of e and e' and of c and c' are reversed. Specifically, if $c \in \mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i)$, Lemma 12(d) applies with the roles of e and e' and the roles of c and c' reversed. It follows that the increase in total overlap due to the swap is $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(m)| - |\text{ov}(f')| \geq \max\{|\text{ov}(e)| - \gamma w(c), |\text{ov}(e')| - \gamma w(c') + |\text{ov}(e)| - \gamma w(c)\} \geq \phi(i+1) - \phi(i)$, as required. Otherwise, if $c \notin \mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i)$, Lemma 12(a) and (b) both apply (again with the roles of e and e' and c and c' reversed) and we know that the increase in total overlap due to the swap is $|\text{ov}(e)| + |\text{ov}(e')| - |\text{ov}(m)| - |\text{ov}(f')| \geq \max\{0, |\text{ov}(e')| - \gamma w(c')\} \geq \phi(i+1) - \phi(i)$, as required.

- If m connects strings of two different cycles in CC and $|\text{ov}(e_t)| < |\text{ov}(m)|$, then we swap in e_h and e'_2 and swap out f'_2 and m . Let c be the cycle of e_h and c' be the cycle of e_t .

- If $\mathcal{M}(\text{CC}_{i+1}) = \mathcal{M}(\text{CC}_i)$, then Lemma 12(a) shows that the total overlap does not decrease, while the potential $\phi(i)$ does not increase.
- If $c' \in \mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i)$, then w_x and v_t must both be strings in cycle c' as otherwise, v' is a string of cycle c' and w_x is a string of a different cycle and thus e'_2 , which is an edge in CC_{i+1} , would connect a string of cycle c' to a string of another cycle. Thus, by Lemma 12(d), $|\text{ov}(e_h)| + |\text{ov}(e'_2)| - |\text{ov}(m)| - |\text{ov}(f'_2)| \geq \max\{|\text{ov}(e'_2)| - \gamma w(c'), |\text{ov}(e_h)| - \gamma w(c) + |\text{ov}(e'_2)| - \gamma w(c')\} \geq \phi(i+1) - \phi(i)$, as required.
- If $\mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i) = \{c\}$ and $(c, c') \in R$, we first observe

$$w(c) \geq |\text{ov}(m)| - w(c') > |\text{ov}(e_t)| - w(c') \geq o(c') - w(c') \geq (\beta - 1) \cdot w(c'),$$

where the third inequality follows from the fact that e_t is an edge of the cycle c' and the last step follows because c' is not extra large. Therefore, we can apply Lemma 12(c) which is sufficient because $w(c)/2 - w(c')/2 = w(c)/2 + w(c')/2 - w(c') \geq |\text{ov}(m)|/2 - w(c') \geq o(c')/2 - w(c')$ and therefore, $|\text{ov}(e_h)| + |\text{ov}(e'_2)| - |\text{ov}(m)| - |\text{ov}(f'_2)| \geq |\text{ov}(e)| - \gamma w(c) - w(c')/2 + w(c)/2 \geq |\text{ov}(e)| - \gamma w(c) - w(c') + o(c')/2 \geq \phi(i+1) - \phi(i)$, as required.

- If $\mathcal{M}(\text{CC}_{i+1}) \setminus \mathcal{M}(\text{CC}_i) = \{c\}$ and $(c, c') \notin R$, there are two possibilities.
 1. If c' is a small cycle, then $w(c') \leq o(c') - w(c') \leq |\text{ov}(e_t)| - w(c') < |\text{ov}(m)| - w(c') \leq w(c)$, where the first step uses the definition of a small cycle and the last step uses Lemma 7.
 2. If c' is a large cycle and $(c, c') \notin R$, then, because $|\text{ov}(m)| > |\text{ov}(e_t)| \geq \beta w(c')$ by the definition of related cycles, $w(c) > (\beta/2 - 1/6) \cdot w(c')$.

Either way $w(c) > (\beta/2 - 1/6) \cdot w(c')$, which means that Lemma 12(b) implies $|\text{ov}(e_h)| + |\text{ov}(e'_2)| - |\text{ov}(m)| - |\text{ov}(f'_2)| \geq |\text{ov}(e_h)| - \gamma w(c) \geq \phi(i+1) - \phi(i)$, as required.

489 — References

- 490 1 Chris Armen and Clifford Stein. Improved length bounds for the shortest superstring problem.
491 In *Proceedings of the 4th International Workshop on Algorithms and Data Structures (WADS)*,
492 pages 494–505, 1995. doi:10.1007/3-540-60220-8_88.
- 493 2 Chris Armen and Clifford Stein. A $2\frac{2}{3}$ superstring approximation algorithm. *Discret. Appl.*
494 *Math.*, 88(1-3):29–57, 1998. doi:10.1016/S0166-218X(98)00065-1.
- 495 3 Avrim Blum, Tao Jiang, Ming Li, John Tromp, and Mihalis Yannakakis. Linear approximation
496 of shortest superstrings. *Journal of the ACM*, 41(4):630–647, 1994. doi:10.1145/179812.
497 179818.
- 498 4 Dany Breslau, Tao Jiang, and Zhigen Jiang. Rotations of periodic strings and short
499 superstrings. *J. Algorithms*, 24(2):340–353, 1997. doi:10.1006/jagm.1997.0861.
- 500 5 M. Crochemore and W. Rytter. *Text Algorithms*. Oxford University Press, 1994.
- 501 6 Artur Czumaj, Leszek Gasieniec, Marek Piotrów, and Wojciech Rytter. Sequential and
502 parallel approximation of shortest superstrings. *J. Algorithms*, 23(1):74–100, 1997. doi:
503 10.1006/jagm.1996.0823.
- 504 7 Matthias Englert, Nicolaos Matsakis, and Pavel Veselý. Improved approximation guarantees
505 for shortest superstrings using cycle classification by overlap to length ratios. In *Proceedings*
506 *of the 54th ACM Symposium on Theory of Computing (STOC)*, pages 317–330. ACM, 2022.
507 doi:10.1145/3519935.3520001.
- 508 8 Alan M. Frieze and Wojciech Szpankowski. Greedy algorithms for the shortest common
509 superstring that are asymptotically optimal. *Algorithmica*, 21(1):21–36, 1998.
- 510 9 Dan Gusfield. *Algorithms on Strings, Trees, and Sequences: Computer Science and Computa-*
511 *tional Biology*. Cambridge University Press, 1997. doi:10.1017/CB09780511574931.
- 512 10 Haim Kaplan, Moshe Lewenstein, Nira Shafir, and Maxim Sviridenko. Approximation
513 algorithms for asymmetric TSP by decomposing directed regular multigraphs. *Journal of the*
514 *ACM*, 52(4):602–626, 2005. doi:10.1145/1082036.1082041.
- 515 11 Haim Kaplan and Nira Shafir. The greedy algorithm for shortest superstrings. *Inf. Process.*
516 *Lett.*, 93(1):13–17, 2005. doi:10.1016/j.ipl.2004.09.012.
- 517 12 Marek Karpinski and Richard Schmed. Improved inapproximability results for the shortest
518 superstring and related problems. In *Proceedings of the 19th Computing: The Australasian*
519 *Theory Symposium (CATS)*, pages 27–36, 2013.
- 520 13 S. Rao Kosaraju, James K. Park, and Clifford Stein. Long tours and short superstrings. In
521 *Proceedings of the 35th IEEE Symposium on Foundations of Computer Science (FOCS)*, pages
522 166–177, 1994. doi:10.1109/SFCS.1994.365696.
- 523 14 Bin Ma. Why greed works for shortest common superstring problem. *Theor. Comput. Sci.*,
524 410(51):5374–5381, 2009.
- 525 15 Marcin Mucha. A tutorial on shortest superstring approximation. [https://www.mimuw.edu.](https://www.mimuw.edu.pl/~much/teaching/aa2008/ss.pdf)
526 [pl/~much/teaching/aa2008/ss.pdf](https://www.mimuw.edu.pl/~much/teaching/aa2008/ss.pdf), 2007. [Accessed 15-June-2023].
- 527 16 Marcin Mucha. Lyndon words and short superstrings. In *Proceedings of the 24th ACM-*
528 *SIAM Symposium on Discrete Algorithms (SODA)*, pages 958–972, 2013. doi:10.1137/1.
529 9781611973105.69.
- 530 17 Katarzyna Paluch, Khaled Elbassioni, and Anke van Zuylen. Simpler approximation of the
531 maximum asymmetric traveling salesman problem. In *Proceedings of the 29th Symposium*
532 *on Theoretical Aspects of Computer Science (STACS)*, pages 501–506, 2012. doi:10.4230/
533 LIPICs.STACS.2012.501.
- 534 18 Steven Skiena. *The Algorithm Design Manual, Third Edition*. Texts in Computer Science.
535 Springer, 2020.
- 536 19 Ondřej Sladký, Pavel Veselý, and Karel Břinda. Masked superstrings as a unified framework for
537 textual k-mer set representations. *bioRxiv*, 2023. URL: [https://www.biorxiv.org/content/](https://www.biorxiv.org/content/early/2023/02/03/2023.02.01.526717)
538 [early/2023/02/03/2023.02.01.526717](https://www.biorxiv.org/content/early/2023/02/03/2023.02.01.526717), doi:10.1101/2023.02.01.526717.
- 539 20 Z. Sweedyk. A $2\frac{1}{2}$ -approximation algorithm for shortest superstring. *SIAM J. Comput.*,
540 29(3):954–986, 1999. doi:10.1137/S0097539796324661.

- 541 21 Jorma Tarhio and Esko Ukkonen. A greedy approximation algorithm for constructing shortest
542 common superstrings. *Theor. Comput. Sci.*, 57:131–145, 1988. doi:10.1016/0304-3975(88)
543 90167-3.
- 544 22 Shang-Hua Teng and Frances Yao. Approximating shortest superstrings. *SIAM Journal on*
545 *Computing*, 26(2):410–417, 1997. doi:10.1137/S0097539794286125.
- 546 23 Jonathan S. Turner. Approximation algorithms for the shortest common superstring problem.
547 *Inf. Comput.*, 83(1):1–20, 1989. doi:10.1016/0890-5401(89)90044-8.
- 548 24 Virginia Vassilevska. Explicit inapproximability bounds for the shortest superstring problem.
549 In *30th International Symposium, MFCS, Gdansk, Poland*, volume 3618 of *Lecture Notes in*
550 *Computer Science*, pages 793–800. Springer, 2005.
- 551 25 Vijay Vazirani. *Approximation algorithms*. Springer, 2001.

552 **A** Deriving Approximation Guarantees from Theorem 3

553 The technical contribution of the paper is proving Theorem 3 that shows an improved
 554 inequality for overlaps of cycle-closing edges in terms of the optimal superstring length n
 555 and the length w of the optimal cycle cover CC . In the next two subsections, we explain how
 556 our improved approximation guarantees follow, using essentially the same arguments (and
 557 algorithms) as in previous work.

558 **A.1** The GREEDY Algorithm for SSP

559 The $|S|^2$ edges of the overlap graph G_{ov} are assumed to be ordered by non-increasing overlap
 560 length. The GREEDY algorithm for SSP chooses edges from this order, unless an edge shares
 561 an endpoint with an already chosen edge or closes a cycle. The edges corresponding to the
 562 latter case are called *bad back edges*. As proven in [3], bad back edges do not intersect each
 563 other, forming a laminar family of edges. Each inner-most bad back edge forms a cycle
 564 in the output of GREEDY and each such cycle is called *culprit*. The sum of lengths of all
 565 culprit cycles is denoted by w_c and the sum of overlap lengths of the cycle-closing edges of
 566 all culprits is denoted by o_c .

567 Blum et al. have shown the following two inequalities (Section 5 in [3]):

$$568 \quad |\text{GREEDY}(S)| \leq 2n + o_c - w_c \quad (12)$$

$$569 \quad o \leq n + 2w \quad (13)$$

571 Moreover, the application of the GREEDY algorithm for the optimal cycle cover CC on
 572 the set of strings comprising the culprit cycles only, outputs the exact same set of culprit
 573 cycles (Lemma 15 in [3]). By this and (13) it follows that $o_c \leq n + 2w_c$, which by (12) gives
 574 $|\text{GREEDY}(S)| \leq 4n$, completing their proof.

575 Theorem 3 shows that $o \leq n + \frac{\sqrt{67}-4}{3}w$ which implies that $o_c \leq n + \frac{\sqrt{67}-4}{3}w_c$ using the same
 576 syllogism (Lemma 15 in [3]). By this and (12), we have $|\text{GREEDY}(S)| \leq \frac{\sqrt{67}+2}{3}n \approx 3.396 \cdot n$,
 577 completing our proof.

578 **A.2** SSP Algorithms Based on Max-ATSP Approximations

579 Blum et al. proposed the following 4-approximate SSP algorithm, called MGREEDY:

- 580 1. Apply GREEDY to find an optimal cycle cover CC .
- 581 2. Open all cycle-closing edges in CC to obtain a set of strings called *representatives*.
- 582 3. Concatenate the representatives in an arbitrary order.

583 If instead of concatenating the representatives in the third step, we merge them using a
 584 Max-ATSP approximation algorithm (executed on the overlap graph of the representatives),
 585 then we will obtain an SSP approximation algorithm which, obviously, cannot perform
 586 worse. This is the idea behind the 3-approximate TGREEDY algorithm [3]. The Max-ATSP
 587 algorithm utilized as a black-box within TGREEDY is GREEDY, which had been already
 588 shown [21, 23] to be a $\frac{1}{2}$ -approximate Max-ATSP algorithm for the overlap graphs.

589 We will need the following theorem from [7], which has already appeared in similar forms
 590 in literature (e.g., [3, 4, 15]).

591 **► Theorem 13.** *If MGREEDY is a $(2 + \zeta)$ -approximate SSP algorithm and there exists a*
 592 *δ -approximate algorithm for Max-ATSP then there exists a $(2 + (1 - \delta) \cdot \zeta)$ -approximate SSP*
 593 *algorithm.*

594 Showing that $o \leq n + (\sqrt{67} - 4)w/3 \approx n + 1.396w$ implies that MGREEDY is a 3.396-
 595 approximate SSP algorithm, since $|\text{MGREEDY}(S)| = w + o \leq w + n + (\sqrt{67} - 4)w/3 < 3.396n$.
 596 Moreover, the currently best Max-ATSP approximation algorithms are $\frac{2}{3}$ -approximate, due
 597 to Kaplan et al. [10] or due to Paluch et al. [17]. Setting $\delta = \frac{2}{3}$ and $\zeta = (\sqrt{67} - 4)/3 \approx 1.396$
 598 in Theorem 13, we obtain an SSP algorithm with approximation guarantee $\frac{\sqrt{67}+14}{9} \approx 2.466$.

599 Finally, regarding TGREEDY, setting $\delta = \frac{1}{2}$ and $\zeta = (\sqrt{67} - 4)/3 \approx 1.396$ in Theorem 13,
 600 we improve the approximation guarantee of TGREEDY to $(\sqrt{67} + 8)/6 \approx 2.698$, from
 601 $(25 + \sqrt{57})/12 \approx 2.712$ as shown in [7].

602 **B** Dealing with extra large cycles (as in [7])

603 Let $\bar{S} \subseteq S$ be the subset of strings that belong to all small and large cycles of CC. Observation
 604 5.1 in [7] implies that the optimal cycle cover for \bar{S} (in short $\text{CC}(\bar{S})$) consists of all small and
 605 large cycles of the optimal cycle cover for S (for simplicity denoted by $\text{CC}(S) = \text{CC}$), while
 606 the optimal cycle cover for $S - \bar{S}$ (in short $\text{CC}(S - \bar{S})$) consists of all extra large cycles of
 607 $\text{CC}(S)$.

608 Let \hat{w} denote the sum of lengths of the (extra large) cycles in $\text{CC}(S - \bar{S})$ and let \hat{o} be the
 609 sum of overlap lengths of the cycle-closing edges of the cycles in $\text{CC}(S - \bar{S})$. Similarly, let \bar{o}
 610 be the sum of overlap lengths of the cycle-closing edges in $\text{CC}(\bar{S})$ and let \bar{w} be the sum of
 611 lengths of the cycles in $\text{CC}(\bar{S})$.

612 Proving $o \leq n + \beta \cdot w$ for input \bar{S} implies that $\bar{o} \leq |\text{OPT}(\bar{S})| + \beta \cdot \bar{w}$, and assuming this,
 613 we show $o \leq n + \beta \cdot w$. Indeed, we take the sum of inequality $\bar{o} \leq |\text{OPT}(\bar{S})| + \beta \cdot \bar{w}$ with
 614 inequality $\hat{o} \leq \beta \cdot \hat{w}$ (which holds by the definition of extra large cycles) and obtain:

$$615 \quad o = \bar{o} + \hat{o} \leq |\text{OPT}(\bar{S})| + \beta \cdot \bar{w} + \beta \cdot \hat{w} = |\text{OPT}(\bar{S})| + \beta \cdot w \leq n + \beta \cdot w$$

616 where the penultimate step uses $w = \bar{w} + \hat{w}$ and the last inequality uses $|\text{OPT}(\bar{S})| \leq$
 617 $|\text{OPT}(S)| = n$, which follows from $\bar{S} \subseteq S$. Therefore, for proving $o \leq n + \beta \cdot w$, we assume
 618 w.l.o.g. that $\text{CC}(S) = \text{CC}$ has no extra large cycle.

619 **C** Lemma 11 (slightly modified from [7])

620 For completeness, we include a proof of Lemma 11. The proof is almost identical to the one
 621 in [7] with only very minor changes to make it more general.

622 We start by stating a corollary, a version of which is already stated in [7] and in slight
 623 variations has been known already before (e.g. see Lemma 9 in [3] and Lemma 7 in [15]).

624 **► Corollary 14.** *Let c and c' be any two cycles of CC. Any string h , which is a substring of*
 625 *both $s(c)^\infty$ and $s(c')^\infty$,¹ satisfies $|h| < w(c) + w(c')$.*

626 This enables us to restate the proof of Lemma 11.

627 **Proof of Lemma 11.** We show that $|\text{ov}(e)| > |\text{ov}(f)| + |\text{ov}(f')| - w(c')$, which implies the
 628 lemma. If $\min\{|\text{ov}(f)|, |\text{ov}(f')|\} \leq w(c')$, this inequality holds because by using Lemma 7,

¹ The definitions of $s(c)$ and s^∞ appear below Lemma 9.

26:18 Approximation Guarantees for Shortest Superstrings: Simpler and Better

629 we get

$$\begin{aligned}
 630 \quad |\text{ov}(e)| &\geq w(c) + w(c') \\
 631 \quad &> \max\{|\text{ov}(f)|, |\text{ov}(f')|\} \\
 632 \quad &\geq \max\{|\text{ov}(f)|, |\text{ov}(f')|\} + \min\{|\text{ov}(f)|, |\text{ov}(f')|\} - w(c') \\
 633 \quad &= |\text{ov}(f)| + |\text{ov}(f')| - w(c') . \\
 634
 \end{aligned}$$

635 Hence, for the remainder of the proof, we assume that we have $\min\{|\text{ov}(f)|, |\text{ov}(f')|\} > w(c')$.

Now, assume for contradiction that $|\text{ov}(e)| \leq |\text{ov}(f)| + |\text{ov}(f')| - w(c')$. We claim that in this case $\text{ov}(e)$ has a periodicity of length $w(c')$, i.e., $\text{ov}(e)$ is a prefix of x^∞ for some string x with $|x| = w(c')$. To show this, first recall that $|\text{ov}(e)| \geq w(c) + w(c') > \max\{|\text{ov}(f')|, |\text{ov}(f)|\}$ by Lemma 7. Since $\text{ov}(f)$ is a prefix of v and a suffix of v' and since $\text{ov}(e)$ is a prefix of v , the first $|\text{ov}(f)|$ characters of $\text{ov}(e)$ are also a suffix of v' , i.e.,

$$\text{ov}(e)[1, |\text{ov}(f)|] = \text{ov}(f) = v' [|v'| - |\text{ov}(f)| + 1, |v'|] .$$

Similarly, since $\text{ov}(f')$ is a prefix of u' and a suffix of u and since $\text{ov}(e)$ is a suffix of u , we get that

$$\text{ov}(e)[|\text{ov}(e)| - |\text{ov}(f')| + 1, |\text{ov}(e)|] = \text{ov}(f') = u'[1, |\text{ov}(f')|] .$$

636 Observe that for all $1 \leq i \leq |\text{ov}(e)| - w(c')$, a character at position i of $\text{ov}(e)$ must be the same
 637 as the character at position $i + w(c')$ of $\text{ov}(e)$. Indeed, if $i + w(c') \leq |\text{ov}(f)|$, this is true as v'
 638 has a periodicity of length $w(c')$. If $i > |\text{ov}(e)| - |\text{ov}(f')|$, it is true because u' has a periodicity
 639 of length $w(c')$. One of these two cases must apply because otherwise, $i + w(c') > |\text{ov}(f)|$
 640 and $i \leq |\text{ov}(e)| - |\text{ov}(f')|$, which implies $|\text{ov}(f)| - w(c') < i \leq |\text{ov}(e)| - |\text{ov}(f')|$, contradicting
 641 our assumption that $|\text{ov}(f')| + |\text{ov}(f)| \geq |\text{ov}(e)| + w(c')$. Hence, $\text{ov}(e)$ has a periodicity of
 642 length $w(c')$ (in particular, $\text{period}(\text{ov}(e)) \leq w(c')$).

643 Next, we show that $\text{ov}(e)$ is a substring of the semi-infinite string $s(c')^\infty$. Because
 644 $\text{ov}(e)$ has a periodicity of length $w(c')$ and $s(c')^\infty$ has period $w(c')$, it is sufficient to argue
 645 that the first $w(c')$ characters of $\text{ov}(e)$ are a substring of $s(c')^\infty$. This is indeed the case
 646 since $\text{ov}(e)[1, |\text{ov}(f)|]$ is a substring of v' which is a substring of $s(c')^\infty$ and we assume that
 647 $|\text{ov}(f)| > w(c')$.

648 Since $\text{ov}(e)$ is a substring of $s(c')^\infty$ as well as of $s(c)^\infty$ (because $\text{ov}(e)$ is a substring of a
 649 string that is part of c), Corollary 14 implies $|\text{ov}(e)| < w(c) + w(c')$ which contradicts the
 650 assumption of the lemma. ◀