

On the Greedy Algorithm for the Shortest Common Superstring Problem with Reversals

Gabriele Fici

Dipartimento di Matematica e Informatica, Università di Palermo, Italy

Tomasz Kociumaka

Faculty of Mathematics, Informatics and Mechanics, University of Warsaw, Poland

Jakub Radoszewski

*Faculty of Mathematics, Informatics and Mechanics, University of Warsaw, Poland
and Newton International Fellow at Department of Informatics, King's College London, UK*

Wojciech Rytter

Faculty of Mathematics, Informatics and Mechanics, University of Warsaw, Poland

Tomasz Waleń

Faculty of Mathematics, Informatics and Mechanics, University of Warsaw, Poland

Abstract

We study a variation of the classical Shortest Common Superstring (SCS) problem in which a shortest superstring of a finite set of strings S is sought containing as a factor every string of S or its reversal. We call this problem Shortest Common Superstring with Reversals (SCS-R). This problem has been introduced by Jiang et al. [9], who designed a greedy-like algorithm with length approximation ratio 4. In this paper, we show that a natural adaptation of the classical greedy algorithm for SCS has (optimal) *compression ratio* $\frac{1}{2}$, i.e., the sum of the overlaps in the output string is at least half the sum of the overlaps in an optimal solution. We also provide a linear-time implementation of our algorithm.

Keywords: Shortest Common Superstring, reversal, greedy algorithm.

1. Introduction

The Shortest Common Superstring (SCS) problem is a classical combinatorial problem on strings with applications in many domains, e.g. DNA fragment assembly, data compression, etc. (see [6] for a recent survey). It consists, given a finite set of strings S , in finding a shortest string containing

Email addresses: `gabriele.fici@unipa.it` (Gabriele Fici), `kociumaka@mimuw.edu.pl` (Tomasz Kociumaka), `jrad@mimuw.edu.pl` (Jakub Radoszewski), `rytter@mimuw.edu.pl` (Wojciech Rytter), `walen@mimuw.edu.pl` (Tomasz Waleń)

as factors (substrings) all the strings in S . The decision version of the problem is known to be NP-complete [13, 5, 4], even under several restrictions on the structure of S (see again [6]). However, a particularly simple greedy algorithm introduced by Gallant in his Ph.D. thesis [5] is widely used in applications since it has very good performance in practice (see for instance [12] and references therein). It consists in repetitively replacing a pair of strings with maximum overlap with the string obtained by overlapping the two strings, until one string remains. The greedy algorithm can be implemented using Aho-Corasick automaton in $\mathcal{O}(n)$ randomized time (with hashing on the symbols of the alphabet) or $\mathcal{O}(n \min(\log m, \log |\Sigma|))$ deterministic time (see [17]), where n is the sum of the lengths of the strings in S and m its cardinality.

The approximation of the greedy algorithm is usually measured in two different ways: one consists in taking into account the *approximation ratio* (also known as the *length ratio*) k_g/k_{min} , where k_g is the length of the output string of greedy and k_{min} the length of a shortest superstring, the other consists in taking into account the *compression ratio* $(n - k_g)/(n - k_{min})$.

For the approximation ratio, Turner [16] proved that there is no constant $c < 2$ such that $k_g/k_{min} \leq c$. The *greedy conjecture* states that this approximation ratio is in fact 2 [1]. The best bound currently known is 3.5 due to Kaplan and Shafir [10]. Algorithms with better approximation ratio are known; the best one is due to Mucha, with an approximation ratio of $2\frac{11}{23}$ [14].

For the compression ratio, Tarhio and Ukkonen [15] proved that $(n - k_g)/(n - k_{min}) \geq 1/2$ and this bound is tight, since it is achieved for the set $S = \{ab^h, b^h a, b^{h+1}\}$ when greedy makes the first choice merging the first two strings together.

Let us formally state the SCS problem:

SHORTEST COMMON SUPERSTRING (SCS)

Input: strings $S = \{s_1, \dots, s_m\}$ of total length n .

Output: a shortest string u that contains s_i for each $i = 1, \dots, m$ as a factor.

Several variations of SCS have been considered in literature. For example, shortest common superstring problem with reverse complements was considered in [11]. In this setting the alphabet is $\Sigma = \{\mathbf{a}, \mathbf{t}, \mathbf{g}, \mathbf{c}\}$ and the complement of a string s is \bar{s}^R , where $\bar{\cdot}$ is defined by $\bar{\mathbf{a}} = \mathbf{t}$, $\bar{\mathbf{t}} = \mathbf{a}$, $\bar{\mathbf{g}} = \mathbf{c}$, $\bar{\mathbf{c}} = \mathbf{g}$, and t^R denotes the reversal of t , that is the string obtained reading t backwards. In particular, this problem was shown to be NP-complete.

Other variations of the SCS problem can be found in [8, 3, 7, 2].

In this paper, we address the problem of searching for a string u of minimal length such that for every $s_i \in S$, u contains as a factor s_i or its reversal s_i^R .

SHORTEST COMMON SUPERSTRING WITH REVERSALS (SCS-R)

Input: strings $S = \{s_1, \dots, s_m\}$ of total length n .

Output: a shortest string u that contains for each $i = 1, \dots, m$ at least one of the strings s_i or s_i^R as a factor.

For example, if $S = \{aabb, aaac, abbb\}$, then a solution of SCS-R for S is $caaabbb$. Notice that a shortest superstring with reversals can be much shorter than a classical shortest superstring. An extremal example is given by an input set of the form $S = \{ab^h, cb^h\}$.

The SCS-R problem was already considered by Jiang et al. [9], who observed that the problem is still NP-hard. We provide a proof of this fact for the sake of completeness (see Proposition 1).

In [9], the authors proposed a greedy 4-approximation algorithm. Here, we show that an adaptation of the classical greedy algorithm can be used for solving the SCS-R problem with an (optimal) compression ratio $\frac{1}{2}$.

2. Basics and Notation

Let Σ be a finite alphabet. We assume that Σ is linearly sortable, e.g., $\Sigma = \{0, \dots, n^{\mathcal{O}(1)}\}$. The *length* of a string s over Σ is denoted by $|s|$. The *empty string*, denoted by ε , is the unique string of length zero. A string t *occurs* in a string s if $s = vtz$ for some strings v, z . In this case we say that t is a *factor* of s . In particular, we say that t is a *prefix* of s when $v = \varepsilon$ and a *suffix* of s when $z = \varepsilon$. We say that a factor t is *proper* if $s \neq t$.

The string s^R obtained by reading s from right to left is called the *reversal* (or *mirror image*) of s . Given a set of strings $S = \{s_1, \dots, s_m\}$, we define the set $S^R = \{s_1^R, \dots, s_m^R\}$ and the set $\tilde{S} = S \cup S^R$.

Given two strings u, v , we define the (maximum) overlap between u and v , denoted by $ov(u, v)$, as the length of the longest suffix of u that is also a prefix of v . Sometimes we abuse the notation and also say that the suffix of u of length $ov(u, v)$ is the overlap of u and v . In general $ov(u, v)$ is not equal to $ov(v, u)$, but it is readily verified that $ov(u, v) = ov(v^R, u^R)$. Additionally, we define $pr(u, v)$ as the prefix of u obtained by removing the suffix of length $ov(u, v)$ and denote $u \otimes v = pr(u, v)v$. Note that the \otimes operation is in general neither symmetric nor associative.

A set of strings S is called *factor-free* if no string in S is a factor of another string in S . We say that S is *reverse-factor-free* if there are no distinct strings $u, v \in S$ such that u is a factor of v or v^R .

Given a factor-free set of strings $S = \{s_1, \dots, s_m\}$, the SCS problem for S is known to be equivalent to that of finding a maximum-weight Hamiltonian path π in the *overlap graph* G_S , which is a directed weighted graph (S, E, w) with arcs $E = \{(s_i, s_j) \mid i \neq j\}$ of weights $w(s_i, s_j) = ov(s_i, s_j)$ (cf. Theorem 2.3 in [15]). In this setting, a path $\pi = v_1, \dots, v_k$ corresponds to a string $str(\pi) := pr(v_1, v_2) \cdots pr(v_{k-1}, v_k)v_k$. By $ov(\pi)$ we denote the total weight of arcs in the path π .

To accommodate reversals we extend the notion of an overlap graph to $\tilde{G}_S = (V, E, w)$. Here $V = \{v_s : s \in S\} \cup \{v_s^R : s \in S\}$ so every $s \in S$ corresponds to exactly two vertices, v_s and v_s^R . We define $str(v_s) = s$ and $str(v_s^R) = s^R$. For a vertex $\alpha \in \tilde{G}_S$ we define α^R as v_s^R if $\alpha = v_s$ for some s or as v_s if $\alpha = v_s^R$ for some s . Note that $str(\alpha^R) = str(\alpha)^R$. For every $\alpha, \beta \in V$, $\alpha \neq \beta$, we introduce an arc from α to β with weight $ov(str(\alpha), str(\beta))$. For an arc $e = (\alpha, \beta)$ we define $e^R = (\beta^R, \alpha^R)$. Note that the weight of e^R is the same as the weight of e .

For paths π in \tilde{G}_S we also use the notions of $str(\pi)$ and $ov(\pi)$. We say that a path π in \tilde{G}_S is *semi-Hamiltonian* if π contains, for every vertex $\alpha \in \tilde{G}_S$, exactly one of the two vertices α, α^R . Observe that a solution to SCS-R problem for a reverse-factor-free set S corresponds to a maximum-weight semi-Hamiltonian path π in the overlap graph \tilde{G}_S .

3. NP-completeness of the SCS-R Problem

Let $\Gamma = \Sigma \cup \{\$, \#\}$, where $\$, \# \notin \Sigma$. Let h be the morphism from Σ^* to Γ^* defined by $h(c) = \$ \# c$, for every $c \in \Sigma$. Given $k \geq 1$, we also define the morphism $g_k(c) = c^k$, for every $c \in \Sigma$.

Observation 1. For every nonempty strings u, v , the strings $h(u)$ and $h(v)^R$ have an overlap of length at most one.

Observation 2. For every nonempty strings u, v , one has

$$|h(g_k(u))| = 3k|u| \quad \text{and} \quad ov(h(g_k(u)), h(g_k(v))) = 3k \cdot ov(u, v).$$

In the decision version of the SCS-R problem we need to check if a given set of strings admits a common superstring with reversals of length at most ℓ , where ℓ is an additional input parameter.

Proposition 1. *The decision version of the SCS-R problem is NP-complete.*

Proof. Let s_1, \dots, s_m be an instance of the SCS problem. Set $y_i = h(g_m(s_i))$, for $i = 1, \dots, m$. We have the following claim.

Claim 1. There exists a common superstring of s_1, \dots, s_m of length at most ℓ if and only if there exists a common superstring with reversal of y_1, \dots, y_m of length at most $3m\ell$.

Proof. (\Rightarrow) If u is a common superstring of s_1, \dots, s_m , then $h(g_m(u))$ is a common superstring of y_1, \dots, y_m . Hence, $h(g_m(u))$ is also a common superstring with reversals of y_1, \dots, y_m . If $|u| = \ell$, then $|h(g_m(u))| = 3m\ell$.

(\Leftarrow) Let u be a common superstring with reversals of y_1, \dots, y_m . We will show that, thanks to the special form of the strings, there exists a common superstring *without* reversals of y_1, \dots, y_m of length not much greater than $|u|$.

Let $\pi = z_{i_1}, \dots, z_{i_m}$ be the sequence of nodes on the path in the overlap graph G that corresponds to u ; here $\{i_1, \dots, i_m\} = \{1, \dots, m\}$ and each z_i is either y_i or y_i^R . Let us construct a new sequence of nodes π' , that first contains all nodes from π of the form y_i in the same order as in π , and then all nodes from π of the form y_i^R , but given in the reverse order and taken without reversal. Let u' be the common superstring corresponding to π' . By Observation 1, $|u'| \leq |u| + m - 1$. Note that u' is an ordinary common superstring of y_1, \dots, y_m .

By Observation 2, $|u'|$ is a multiple of $3m$ and u' corresponds to a common superstring v of s_1, \dots, s_m of length $|u'|/(3m)$. If $|u| \leq 3m\ell$, then

$$|v| = \frac{|u'|}{3m} \leq \left\lfloor \frac{|u| + m - 1}{3m} \right\rfloor \leq \ell.$$

□

The claim provides a reduction of the decision version of the SCS problem to the decision version of the SCS-R problem. This shows that the latter is NP-hard, hence NP-complete, as it is obviously in NP. □

4. Greedy Algorithm and its Linear-Time Implementation

We define an auxiliary procedure MAKE-REVERSE-FACTOR-FREE(S) that removes from S all strings u which are contained as a factor in v or v^R for some $v \in S$, $v \neq u$. Note that the resulting set S' is reverse-factor-free and, moreover, a string is a common superstring with reversals for S' if and only if it is a common superstring with reversals for S .

Example 1. Let $S = \{ab, aaa, aab, baa\}$. Then MAKE-REVERSE-FACTOR-FREE(S) produces $S' = \{aaa, aab\}$ or $S' = \{aaa, baa\}$.

Algorithm Greedy-R(S)**Input:** a non-empty set of strings S **Output:** a superstring of S that approximates a solution of SCS-R problem for S **begin** $S := \text{MAKE-REVERSE-FACTOR-FREE}(S)$ **while** $|S| > 1$ **do** $P := \{(u, v) : u, v \in \tilde{S}, u \notin \{v, v^R\}\}$ $\{ S \text{ is reverse-factor-free and } |S| > 1, \text{ so } |P| \geq 1 \}$ take $(u, v) \in P$ with the maximal value of $ov(u, v)$ $S := S \cup \{u \otimes v\}$ $S := S \setminus \{u, v, u^R, v^R\}$ **return** the only element of S **end**

The Greedy-R algorithm works as follows: while $|S| > 1$, choose $u, v \in \tilde{S}$ (excluding $u = v$ and $u = v^R$) with largest overlap, insert into S the string $u \otimes v$, and remove from S all strings among u, v, u^R, v^R that belong to S ; see the pseudocode.

Let us state two properties of this algorithm useful for its efficient implementation.

Lemma 2. *The set S stays reverse-factor-free after each iteration of the **while** loop.*

Proof. Suppose that at some point S ceases to be reverse-factor-free. This might only be due to the fact that, when $w = u \otimes v$ is introduced to S , \tilde{S} contains a string $w' \notin \{u, u^R, v, v^R\}$ such that w' is a factor of w but not of u or v . The latter, however, implies $ov(u, w') > ov(u, v)$ and $ov(w', v) > ov(u, v)$. That contradicts the choice of $(u, v) \in P$ maximizing $ov(u, v)$. \square

Lemma 3. *Before $u \otimes v$ is inserted to S , we have $ov(w, u \otimes v) = ov(w, u)$ and $ov(u \otimes v, w) = ov(v, w)$ for every $w \in \tilde{S}$.*

Proof. Clearly, $ov(w, u \otimes v) \geq ov(w, u)$ and $ov(u \otimes v, w) \geq ov(v, w)$. Moreover, one of these inequalities might be strict only if $ov(w, u \otimes v) > |u|$ or $ov(u \otimes v, w) > |v|$, in particular, only if w contains respectively u or v as a proper factor. This, however, contradicts Lemma 2. \square

4.1. Interpretation on the Overlap Graph

Tarhio and Ukkonen [15] work with the following interpretation of the greedy algorithm on the overlap graph G_S : They maintain a set of arcs $F \subseteq E(G_S)$ forming a collection of disjoint paths and at each step they add to F an arc $e \in E(G_S) \setminus F$ of maximum weight so that the resulting set still forms a collection of disjoint paths. Here, paths correspond to strings in the collection S maintained by the original implementation. Insertion of an arc (s_i, s_j) to F results in merging two paths into one and this corresponds to replacing strings $u, v \in S$ with $u \otimes v$. It turns out that $ov(u, v) = ov(s_i, s_j)$ and thus both implementations of the greedy algorithm are equivalent.

Here, we provide an analogous interpretation of Greedy-R on the overlap graph \tilde{G}_S and, for completeness, explicitly prove its equivalence to the original implementation. We also maintain a set of arcs $F \subseteq E(\tilde{G}_S)$ forming a collection of disjoint paths. In each step we extend F with a pair of arcs $\{e, e^R\}$ of maximum (common) weight so that the resulting set still forms a collection of disjoint paths; see the pseudocode of algorithm Greedy-R2. Observe that this way paths formed by F come in pairs π, π^R such that $\pi = \alpha_1, \dots, \alpha_p$ and $\pi^R = \alpha_p^R, \dots, \alpha_1^R$ (in particular, $str(\pi^R) = str(\pi)^R$).

Algorithm Greedy-R2(S)**Input:** a non-empty set of strings S **Output:** a superstring of S that approximates a solution of SCS-R problem for S **begin** $S := \text{MAKE-REVERSE-FACTOR-FREE}(S)$ Construct \tilde{G}_S $F := \emptyset$ **for** $i := 1$ **to** $|S| - 1$ **do** $P := \{e \in E(\tilde{G}_S) : F \cup \{e, e^R\} \text{ forms a collection of disjoint paths in } \tilde{G}_S\}$ take $e \in P$ with the maximal weight $w(e)$ $F := F \cup \{e, e^R\}$ **return** $\text{str}(\pi)$ for one of the maximal paths π formed by F in \tilde{G}_S **end**

We claim that these pairs of paths correspond to strings $s \in S$ of the original implementation (with $s = \text{str}(\pi)$ or $s = \text{str}(\pi^R)$). In particular, each string $s \in \tilde{S}$ corresponds to a single path formed by F unless $s = s^R$ when it corresponds to two reverse paths. The claim is certainly true at the beginning of the algorithm, so let us argue that single iterations of the main loops in both algorithms perform analogous operations.

First, consider an arc $e = (\alpha, \beta)$ such that there exists a path π that ends at α and a path π' that starts at β . (Otherwise, $F \cup \{e, e^R\}$ does not form a collection of disjoint paths). Observe that strings $u = \text{str}(\pi)$ and $v \in \text{str}(\pi')$ belong to \tilde{S} in Greedy-R and that $u \in \{v, v^R\}$ if and only if $F \cup \{e, e^R\}$ yields a cycle. Thus, both algorithms consider essentially the same set of possibilities (the only caveat is that if u or v is a palindrome, then several arcs e correspond to the same pair of strings from \tilde{S}). By Lemma 3 the weight of an arc (α, β) is $ov(\text{str}(\alpha), \text{str}(\beta)) = ov(u, v)$. Thus, the maximum-weight arcs correspond to pairs $(u, v) \in P$ with largest overlap.

Clearly, setting $F := F \cup \{e, e^R\}$ results in merging π' with π and π^R with π'^R . These paths represent $u \otimes v$ and $v^R \otimes u^R$. Since the set S in Greedy-R stays reverse-factor-free due to Lemma 2 and, in particular, S does not contain any pair of strings v, v^R such that v is not a palindrome, the “ $S := S \setminus \{u, u^R, v, v^R\}$ ” instruction always removes exactly two elements of S . This means that the bijection between S and pairs of paths formed by F is preserved.

Finally, observe that after exactly $|S| - 1$ steps F forms two semi-Hamiltonian paths. Either of them can be returned as a solution.

4.2. Linear-Time Implementation

Ukkonen [17] showed that the Greedy algorithm for the original SCS problem can be implemented in linear time based on the overlap-graph interpretation. Our linear-time implementation of the Greedy-R2 algorithm is quite similar.

First, let us show how to efficiently implement the MAKE-REVERSE-FACTOR-FREE operation. It is actually slightly easier to compute it for \tilde{S} instead of S . However, this is not an issue: if we substitute S with \tilde{S} in the very beginning of the greedy algorithm, then the overlap graph will stay the same.

Lemma 4. *The result of MAKE-REVERSE-FACTOR-FREE(\tilde{S}) can be computed in time linear in the total length of strings in S .*

Proof. Let us introduce an auxiliary procedure $\text{MAKE-FACTOR-FREE}(X)$ that removes from X all strings u for which there exists a string v in X such that u is a proper factor of v .

Ukkonen [17] applied the following result for the preprocessing phase of the greedy algorithm for the ordinary SCS problem.

Claim 2 ([17]). $\text{MAKE-FACTOR-FREE}(X)$ can be implemented in time linear in the total length of the strings in X .

Observe that in order to compute $\text{MAKE-REVERSE-FACTOR-FREE}(\tilde{S})$ it suffices to determine $S' = \text{MAKE-FACTOR-FREE}(\tilde{S})$ and then for every pair of strings $u, u^R \in S'$ leave exactly one of these strings, e.g., the lexicographically smaller one. Note that $u \in S'$ if and only if $u^R \in S'$, so for the latter it suffices to iterate through S' and report a string u if and only if $u \leq u^R$.

The whole procedure works in linear time in the total length of strings in \tilde{S} , which is at most twice the total length of strings in S . \square

Now, we show the main result.

Theorem 5. *Greedy-R algorithm can be implemented in time linear in the total length of strings in S .*

Proof. By Lemma 4, we can make the input set reverse-factor-free in linear time. Let $S = \{s_1, \dots, s_m\}$ be the set of remaining strings and n the total length of strings in S .

We actually implement the equivalent algorithm Greedy-R2. We cannot store the whole graph \tilde{G}_S , since this would take too much space. Therefore we only store its vertex set, whereas we will be considering the edges of the graph in an indirect way.

Denote by $\text{Pref}(X)$ the set of all different prefixes of strings in X . Each element of this set can be represented as a state of the Aho-Corasick automaton constructed for X , thus using $\mathcal{O}(1)$ space per element. Further, given a string w , denote by $\text{PrefSet}(w, X)$, $\text{SufSet}(w, X)$ the sets (represented as lists of identifiers) of strings in X having w as a prefix, suffix, respectively. Ukkonen [17] applied the Aho-Corasick automaton for X to show the following fact:

Claim 3 ([17]). $\text{PrefSet}(w, X)$, $\text{SufSet}(w, X)$ for all $w \in \text{Pref}(X)$ can be computed in time linear in the total length of the strings in X .

In the implementation we use sets of the form $\text{Pref}(\tilde{S})$, $\text{PrefSet}(w, \tilde{S})$ and $\text{SufSet}(w, \tilde{S})$. Our implementation actually requires PrefSet and SufSet sets to consist of vertices $\alpha \in V(\tilde{G}_S)$ instead of strings $\text{str}(\alpha)$. Thus, we replace every string identifier with one or two vertex identifiers if the string is not a palindrome or is a palindrome, respectively.

Instead of simulating the main loop directly, we will consider all overlaps w of $\text{str}(\alpha)$ and $\text{str}(\beta)$ to find the maximum-weight arc $e = (\alpha, \beta) \in E(\tilde{G}_S)$. Observe that at subsequent iterations of the loop, the weight $w(e)$ may only decrease. Hence, we iterate over all $w \in \text{Pref}(\tilde{S})$ in decreasing length order and for each w check if there exists an appropriate arc $e = (\alpha, \beta)$ such that $F \cup \{e, e^R\}$ forms a collection of disjoint paths. More formally, we seek vertices $\alpha, \beta \in V(\tilde{G}_S)$ such that:

- (a) a path formed by F ends in α and a *different* path formed by F starts at β ,
- (b) $\beta \neq \alpha^R$, and
- (c) w is the overlap of $\text{str}(\alpha)$ and $\text{str}(\beta)$.

Once we find such an arc e , we set $F := F \cup \{e, e^R\}$.

Let us explain how this approach can be implemented efficiently. To iterate through all $w \in \text{Pref}(\tilde{S})$ in decreasing length order we simply traverse all the states of the Aho-Corasick automaton in reverse-BFS order, breaking ties arbitrarily. To check the conditions (a)-(c), we could iterate through pairs of elements $\beta \in \text{PrefSet}(w, \tilde{S})$ and $\alpha \in \text{SufSet}(w, \tilde{S})$ and verify if they satisfy the conditions. However, to avoid repetitively scanning *redundant* elements, we remove $\beta \in \text{PrefSet}(w, \tilde{S})$ and $\alpha \in \text{SufSet}(w, \tilde{S})$ if no path starts in β and no path ends in α , respectively. For each vertex we will remember if a path starts or ends there, and, if so, what is the other end-point of the path. Observe that for each $\alpha \in \text{SufSet}(w, \tilde{S})$ there are at most two non-redundant elements $\beta \in \text{PrefSet}(w, \tilde{S})$ for which the arc $e = (\alpha, \beta)$ is not valid. Indeed, these might only be α^R and the starting vertex of the path ending at α .

Hence, with amortized constant-time overhead, either an arc $e = (\alpha, \beta)$ satisfying (a)-(c) is found, or it can be verified that no such arc exists for this particular string w and then we can continue iterating through strings in $\text{Pref}(\tilde{S})$. If an arc is found, we start the next search with the same string w , since there could still be arcs satisfying conditions (a)-(c) for the same string w .

To conclude: in every step of the simulation, in amortized constant time we either find a pair of arcs to be introduced to F or discard the given candidate $w \in \text{Pref}(\tilde{S})$. The former situation takes place at most $m - 1$ times and the latter happens at most n times. The whole algorithm thus works in $\mathcal{O}(n)$ time. \square

5. Compression Ratio

In this section we prove that the compression ratio of the Greedy-R algorithm is always at least $\frac{1}{2}$, and that this value is effectively achieved, so the bound is tight.

Let $S = \{s_1, \dots, s_m\}$ be the input set of strings. We assume that S is already reverse-factor-free. Let $\text{opt}(S)$ be the length of a longest semi-Hamiltonian path in $G = \tilde{G}_S$. Let $\text{pgreedy}(S)$ denote the length of the semi-Hamiltonian path produced by the Greedy-R algorithm for S . We will show that $\text{pgreedy}(S) \geq \frac{1}{2} \text{opt}(S)$.

In the proof we use as a tool the following fact from [15] (see Lemma 3.1 in [15]):

Lemma 6. *If strings x_1, x_2, x_3, x_4 satisfy*

$$\max(\text{ov}(x_1, x_4), \text{ov}(x_2, x_3)) \leq \text{ov}(x_1, x_3),$$

then

$$\text{ov}(x_1, x_4) + \text{ov}(x_2, x_3) \leq \text{ov}(x_1, x_3) + \text{ov}(x_2, x_4).$$

We proceed with the following crucial lemma.

Lemma 7. *Let S be a set of strings and let $u, v \in \tilde{S}$ be two elements for which $\text{ov}(u, v)$ is maximal ($u \notin \{v, v^R\}$). Set $OV = \text{ov}(u, v)$. Let $\text{opt}(S)$ be the length of a longest semi-Hamiltonian path in G and let $\text{opt}'(S)$ be the length of a longest semi-Hamiltonian path in G that contains the arc (u, v) . Then:*

$$\text{opt}'(S) \geq \text{opt}(S) - OV.$$

Proof. We consider the path π corresponding to $\text{opt}(S)$ and show how it can be modified without losing its semi-Hamiltonicity so that the arc (u, v) occurs in the path and the length of the path decreases by at most OV .

Obviously, if π already contains the arc (u, v) , nothing is to be done. If both u and v occur in π , we perform transformations as in the proof of a similar fact from [15] related to the ordinary SCS problem. If u occurs in π before v then we select π' as in Fig. 1 and:

$$ov(\pi') \geq ov(\pi) - ov(u, b) - ov(c, v) + ov(u, v) \geq ov(\pi) - ov(u, b) \geq ov(\pi) - OV.$$

Note that both nodes b, c exist (they could be the same node, though). If any of the remaining nodes does not exist, it is simply skipped on the path.

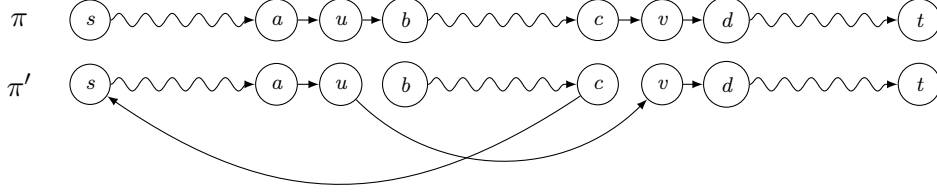


Figure 1: Proof of Lemma 7, first case: u occurs in π before v .

If v occurs before u , then by applying the inequality of Lemma 6 (with $x_1 = u$, $x_2 = a$, $x_3 = v$, $x_4 = d$) we have

$$ov(u, d) + ov(a, v) \leq ov(u, v) + ov(a, d).$$

By this inequality, for the path π' defined in Fig. 2 we have:

$$\begin{aligned} ov(\pi') &\geq ov(\pi) - ov(a, v) - ov(u, d) + ov(u, v) + ov(a, d) - ov(v, b) \\ &\geq ov(\pi) - ov(v, b) \geq ov(\pi) - OV. \end{aligned}$$

As before, if any of the depicted nodes does not exist, we simply skip the corresponding part of the path. In particular, if any of the nodes u, v is an endpoint of the path π , we do not need to use the aforementioned inequality to show that $ov(\pi') \geq ov(\pi) - OV$.

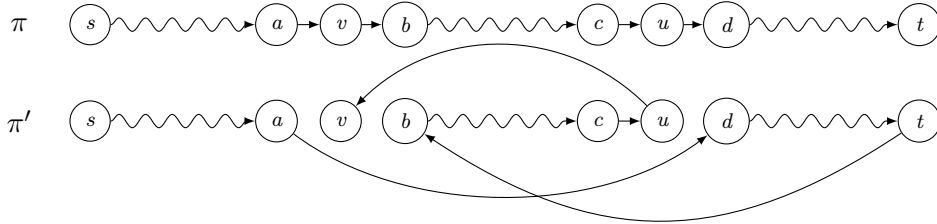


Figure 2: Proof of Lemma 7, second case: u occurs in π after v .

Differently from the original SCS problem considered in [15], it might not be the case that u and v are in π . If none of them is, then π contains both u^R and v^R and by reversing π (that is, taking the path π^R) we obtain a semi-Hamiltonian path that contains both u and v , which was the case considered before. Thus, we can assume that u and v^R occur in π (the case of u^R and v is symmetric). Again, we have two cases, depending on which of the two nodes comes first in π . If u occurs before v^R in π , then we have (note that $ov(c, v^R) = ov(v, c^R)$):

$$ov(\pi') \geq ov(\pi) - ov(u, b) - ov(v^R, d) + ov(u, v) \geq ov(\pi) - ov(u, b) \geq ov(\pi) - OV,$$

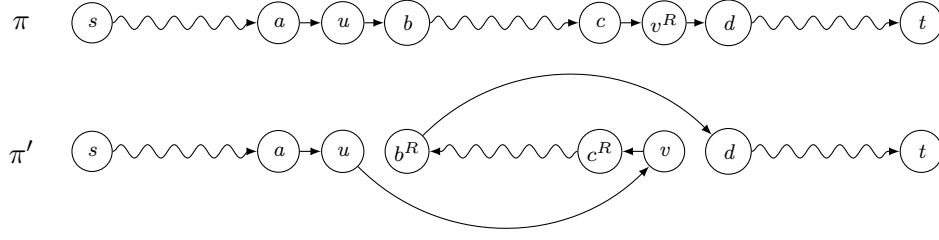


Figure 3: Proof of Lemma 7, third case: u occurs in π before v^R .

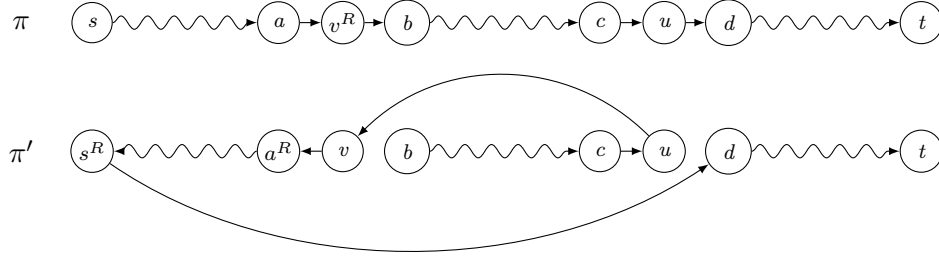


Figure 4: Proof of Lemma 7, third case: u occurs in π after v^R .

see also Fig. 3.

Finally, if v^R occurs before u , then (see Fig. 4):

$$ov(\pi') \geq ov(\pi) - ov(v^R, b) - ov(u, d) + ov(u, v) \geq ov(\pi) - ov(v^R, b) \geq ov(\pi) - OV.$$

This completes the proof of the lemma. \square

To conclude the proof of the compression ratio of the Greedy-R algorithm we use an inductive argument. Assume that $pgreedy(S') \geq \frac{1}{2} \text{opt}(S')$ for all S' such that $|S'| < |S|$. By Lemma 7, for $S' = S \setminus \{u, v, u^R, v^R\} \cup \{u \otimes v\}$:

$$\text{opt}(S) - OV \leq \text{opt}'(S) \leq \text{opt}(S') + OV,$$

so that $\text{opt}(S') + 2OV \geq \text{opt}(S)$. Moreover, $pgreedy(S) = pgreedy(S') + OV$ and, by the inductive hypothesis, $pgreedy(S') \geq \frac{1}{2} \text{opt}(S')$. Consequently:

$$pgreedy(S) = pgreedy(S') + OV \geq \frac{1}{2}(\text{opt}(S') + 2OV) \geq \frac{1}{2} \text{opt}(S).$$

We arrive at the following theorem.

Theorem 8. *Greedy-R has compression ratio $\frac{1}{2}$.*

It turns out that the bound on the compression ratio of the Greedy-R algorithm is tight. It suffices to consider the set of strings $\{ab^h, b^h c, b^{h+1}\}$ (this is actually the same example as from the analysis of the Greedy algorithm [1]). The output of Greedy-R can be the string $ab^h c b^{h+1}$ with total overlap h , whereas an optimal solution to SCS-R is $ab^{h+1}c$ of total overlap $2h$.

6. Acknowledgments

The authors thank anonymous referees for a number of helpful comments and remarks.

Gabriele Fici acknowledges the support of the PRIN 2010/2011 project “Automi e Linguaggi Formali: Aspetti Matematici e Applicativi” of the Italian Ministry of Education (MIUR). Tomasz Kociumaka is supported by Polish budget funds for science in 2013-2017 as a research project under the ‘Diamond Grant’ program. Jakub Radoszewski, Wojciech Rytter and Tomasz Waleń are supported by the Polish National Science Center, grant no 2014/13/B/ST6/00770.

- [1] Avrim Blum, Tao Jiang, Ming Li, John Tromp, and Mihalis Yannakakis. Linear approximation of shortest superstrings. *Journal of the ACM*, 41(4):630–647, 1994.
- [2] Raphael Clifford, Zvi Gotthilf, Moshe Lewenstein, and Alexandru Popa. Restricted common superstring and restricted common supersequence. In *Combinatorial Pattern Matching*, volume 6661 of *Lecture Notes in Computer Science*, pages 467–478. Springer Berlin Heidelberg, 2011.
- [3] Maxime Crochemore, Marek Cygan, Costas Iliopoulos, Marcin Kubica, Jakub Radoszewski, Wojciech Rytter, and Tomasz Waleń. Algorithms for three versions of the shortest common superstring problem. In *Combinatorial Pattern Matching*, volume 6129 of *Lecture Notes in Computer Science*, pages 299–309. Springer Berlin Heidelberg, 2010.
- [4] John Gallant, David Maier, and James A. Storer. On finding minimal length superstrings. *Journal of Computer and System Sciences*, 20(1):50–58, 1980.
- [5] John K. Gallant. *String compression algorithms*. PhD thesis, Princeton University, 1982.
- [6] Theodoros P. Gevezes and Leonidas S. Pitsoulis. The shortest superstring problem. In *Optimization in Science and Engineering*, pages 189–227. Springer New York, 2014.
- [7] Zvi Gotthilf, Moshe Lewenstein, and Alexandru Popa. On shortest common superstring and swap permutations. In *String Processing and Information Retrieval*, volume 6393 of *Lecture Notes in Computer Science*, pages 270–278. Springer Berlin Heidelberg, 2010.
- [8] Tao Jiang and Ming Li. Approximating shortest superstrings with constraints. *Theoretical Computer Science*, 134(2):473 – 491, 1994.
- [9] Tao Jiang, Ming Li, and Ding-Zhu Du. A note on shortest superstrings with flipping. *Information Processing Letters*, 44(4):195–199, 1992.
- [10] Haim Kaplan and Nira Shafrir. The greedy algorithm for shortest superstrings. *Information Processing Letters*, 93(1):13–17, 2005.
- [11] John D. Kececiloglu and Eugene W. Myers. Combinatorial algorithms for DNA sequence assembly. *Algorithmica*, 13(1/2):7–51, 1995.
- [12] Bin Ma. Why greed works for shortest common superstring problem. *Theoretical Computer Science*, 410(51):5374–5381, 2009.
- [13] David Maier and James A. Storer. A note on the complexity of the superstring problem. Technical Report 233, Princeton University, 1977.
- [14] Marcin Mucha. Lyndon words and short superstrings. In *Proceedings of SODA 2013*, pages 958–972, 2013.
- [15] Jorma Tarhio and Esko Ukkonen. A greedy approximation algorithm for constructing shortest common superstrings. *Theoretical Computer Science*, 57:131–145, 1988.
- [16] Jonathan S. Turner. Approximation algorithms for the shortest common superstring problem. *Information and Computation*, 83(1):1–20, 1989.
- [17] Esko Ukkonen. A linear-time algorithm for finding approximate shortest common superstrings. *Algorithmica*, 5(3):313–323, 1990.