

Dartmouth College

## Dartmouth Digital Commons

---

Computer Science Technical Reports

Computer Science

---

11-20-1995

### A 2-2/3 Approximation for the Shortest Superstring Problem

Chris Armen

*Dartmouth College*

Clifford Stein

*Dartmouth College*

Follow this and additional works at: [https://digitalcommons.dartmouth.edu/cs\\_tr](https://digitalcommons.dartmouth.edu/cs_tr)



Part of the [Computer Sciences Commons](#)

---

#### Dartmouth Digital Commons Citation

Armen, Chris and Stein, Clifford, "A 2-2/3 Approximation for the Shortest Superstring Problem" (1995).  
Computer Science Technical Report PCS-TR95-262. [https://digitalcommons.dartmouth.edu/cs\\_tr/117](https://digitalcommons.dartmouth.edu/cs_tr/117)

This Technical Report is brought to you for free and open access by the Computer Science at Dartmouth Digital Commons. It has been accepted for inclusion in Computer Science Technical Reports by an authorized administrator of Dartmouth Digital Commons. For more information, please contact [dartmouthdigitalcommons@groups.dartmouth.edu](mailto:dartmouthdigitalcommons@groups.dartmouth.edu).

# A $2\frac{2}{3}$ -Approximation Algorithm for the Shortest Superstring Problem

Chris Armen and Clifford Stein  
Dartmouth College Technical Report PCS-TR95-262

November 20, 1995

## Abstract

Given a collection of strings  $S = \{s_1, \dots, s_n\}$  over an alphabet  $\Sigma$ , a *superstring*  $\alpha$  of  $S$  is a string containing each  $s_i$  as a substring; that is, for each  $i$ ,  $1 \leq i \leq n$ ,  $\alpha$  contains a block of  $|s_i|$  consecutive characters that match  $s_i$  exactly. The *shortest superstring problem* is the problem of finding a superstring  $\alpha$  of minimum length.

The shortest superstring problem has applications in both data compression and computational biology. In data compression, the problem is a part of a general model of string compression proposed by Gallant, Maier and Storer (JCSS '80). Much of the recent interest in the problem is due to its application to DNA sequence assembly.

The problem has been shown to be NP-hard; in fact, it was shown by Blum et al. (JACM '94) to be MAX SNP-hard. The first  $O(1)$ -approximation was also due to Blum et al., who gave an algorithm that always returns a superstring no more than 3 times the length of an optimal solution. Several researchers have published results that improve on the approximation ratio; of these, the best previous result is our algorithm SHORTSTRING, which achieves a  $2\frac{3}{4}$ -approximation (WADS '95).

We present our new algorithm, G-SHORTSTRING, which achieves a ratio of  $2\frac{2}{3}$ . It generalizes the SHORTSTRING algorithm, but the analysis differs substantially from that of SHORTSTRING. Our previous work identified classes of strings that have a nested periodic structure, and which must be present in the worst case for our algorithms. We introduced machinery to describe these strings and proved strong structural properties about them. In this paper we extend this study to strings that exhibit a more relaxed form of the same structure, and we use this understanding to obtain our improved result.

## 1 Introduction

The shortest superstring problem has applications in both computational biology [7, 16, 18] and data compression [10, 20]. We begin by describing the former. DNA sequencing is the task of determining the sequence of nucleotides in a molecule of DNA. These nucleotides are one of adenine, cytosine, guanine, and thymine, and are typically represented by the alphabet  $\{a, c, g, t\}$ . A molecule of human DNA is about  $10^8$  nucleotides long. Current laboratory procedures can directly determine the nucleotides of a fragment of DNA up to about 600 nucleotides long. In *shotgun sequencing*, several copies of a DNA molecule are fragmented using various restriction enzymes.

Once the nucleotides of all of the fragments have been determined, the *sequence assembly problem* is the computational task of reconstructing the original molecule from the overlapping fragments. The shortest superstring problem is an abstraction of this problem, in which the shortest reconstruction is assumed to be the most likely on the grounds that it is the most parsimonious. We state the problem as follows.

Given a collection of strings  $S = \{s_1, \dots, s_n\}$  over an alphabet  $\Sigma$ , a *superstring*  $\alpha$  of  $S$  is a string containing each  $s_i$  as a substring, that is, for each  $i$ ,  $1 \leq i \leq n$ ,  $\alpha$  contains  $|s_i|$  consecutive characters that match  $s_i$  exactly. The *shortest superstring problem* is the problem of finding a superstring  $\alpha$  of minimum length.

The shortest superstring problem is MAX SNP-hard [4]; several heuristics and approximation algorithms have been proposed. One often used algorithm is a greedy algorithm that repeatedly merges the pair of strings with the maximum amount of overlap. Turner [23] and Tarhio and Ukkonen [21] independently proved that the greedy algorithm constructs a superstring that achieves at least half as much overlap as an optimal superstring. However, this does not guarantee a constant approximation with respect to the length of the resulting superstring.

The first bound on the length approximation of the greedy algorithm was provided by Blum et al.[4], who showed that the greedy algorithm returns a string that is no longer than four times optimal; they also give a modified greedy algorithm that returns a string that is within three times optimal. Teng and Yao [22] gave a nongreedy algorithm that finds a string that is within  $2\frac{8}{9}$  of optimal. Subsequently, three results appeared that achieved better approximation ratios using very different techniques. Czumaj et al.[6] refined the algorithm of [22] to achieve a  $2\frac{5}{6}$  approximation. Kosaraju et al. obtained an improved result for the maximum traveling salesman problem; this more general result can be used by the algorithm of [4] to obtain an approximation slightly better than 2.8 [15]. Our result of  $2\frac{3}{4}$  [2, 1] was the best known until recently, and in fact can be combined with the algorithm of [15] to obtain an approximation ratio of about 2.725.

In this report we describe our  $2\frac{2}{3}$ -approximation algorithm for the shortest superstring problem, which also appears in [3]. Algorithmically, the approach is a generalization of the one taken in [2], but the analysis is very different.

We now give a brief overview of our approach. All of the above mentioned algorithms begin by finding a minimum-weight cycle cover on a graph which has a node for every string and an edge between string  $u$  and  $v$  of length  $|u| - ov(u, v)$ , where  $ov(u, v)$  is the amount of overlap that can be obtained by merging  $u$  and  $v$ . This cycle cover partitions the strings into cycles; the remaining work is in patching the cycles together to form one cycle covering the whole graph. The key to our new algorithm is to exploit the periodic structure of the cycles of strings that arise in this problem. In particular, the 3-approximation of [4] uses a theorem about infinite periodic functions [8], and the correspondence between periodic functions and strings in cycles. However, the particular instances of cycle patching that appear to be difficult actually involve short periodic strings, that is, strings that are periodic, but whose period may repeat only slightly more than once. We prove several new properties about such strings, allowing us to answer questions of the following form: given a string with some periodic structure, characterize *all* the possible periodic strings that can have a large amount of overlap with the first string. Given this understanding, we will be able to predict the ways in which overlap between certain strings can occur, and thus plan for it algorithmically.

## 2 Preliminaries

For consistency, we use some notation and definitions of [4] and [22]. We assume, without loss of generality, that the set  $S$  of strings is *substring free*, i.e. no  $s_j$  is a substring of  $s_i$ ,  $i \neq j$ . We use  $|s_i|$  to denote the length of string  $s_i$ ,  $|S|$  to denote the sum of the lengths of all the strings, and  $\text{opt}(S)$  to denote the length of the shortest superstring of  $S$ .

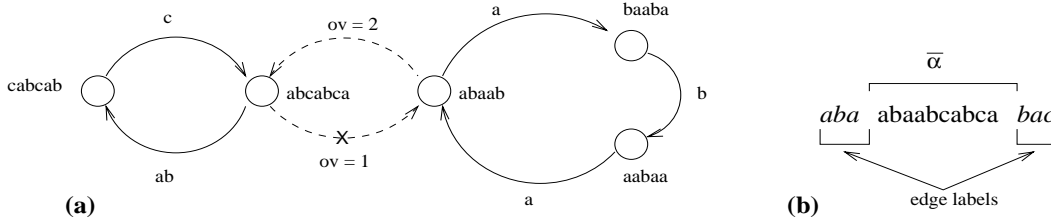


Figure 1: Execution of GENERIC SUPERSTRING ALGORITHM. Nodes are labeled with strings, edges with prefixes. (a) The graph after Step (4). Solid edges are in  $C$ , dashed edges in  $CC$ . The edge with an X is the one discarded in Step (4). (b) The final string consisting of  $\bar{\alpha}$  (the merge of  $abaab$  and  $abcabca$ ) along with the labels from the edges of the cycles.

Given two strings  $s$  and  $t$ , we define  $ov(s, t)$ , the *overlap* between  $s$  and  $t$ , to be the length of the longest string  $x$ , such that there exist non-empty  $u$  and  $v$  with  $s = ux$  and  $t = xv$ . We call  $u$  the *prefix* of  $s$  with respect to  $t$ ,  $\text{pref}(s, t)$ , and refer to  $|u|$  as the distance from  $s$  to  $t$ ,  $d(s, t)$ . Observe that for any  $s$  and  $t$ ,  $ov(s, t) + d(s, t) = |s|$ . String  $uxv$ , the shortest superstring of  $s$  and  $t$  in which  $s$  appears before  $t$  is denoted by  $\langle s, t \rangle$ , and  $|\langle s, t \rangle| = |s| + |t| - ov(s, t)$ .

We can map the superstring problem to a graph problem by defining the *distance graph*. We create a graph  $G = (V, E)$  with a vertex  $v_i \in V$  for each string  $s_i \in S$ . For every ordered pair of vertices  $v_i, v_j$ , we place a directed edge of length  $d(s_i, s_j)$  and label the edge with  $\text{pref}(s_i, s_j)$ . We can now observe that a minimum length hamiltonian cycle (traveling salesman tour)  $v_{\pi_1}, \dots, v_{\pi_n}, v_{\pi_1}$ , in  $G$ , with edge  $i, j$  labeled by  $\text{pref}(s_{\pi_i}, s_{\pi_j})$ , almost corresponds to a superstring in  $S$ , the only difference being that we must replace  $\text{pref}(s_{\pi_n}, s_{\pi_1})$  with  $s_{\pi_n}$ . Since  $\text{pref}(s_i, s_j) \leq |s|$ , we can conclude that  $\text{opt}(TSP) \leq \text{opt}(S)$ , where  $\text{opt}(TSP)$  is the optimal solution to TSP defined above. This TSP is directed (sometimes called *asymmetric*); thus the best known approximation [9] is only within a factor of  $O(\log n)$ . Therefore, we must exploit more of the structure of the problem in order to achieve better bounds.

Given a directed graph  $G$ , with weights on the edges, a *cycle cover*  $C$  is a set of cycles such that each vertex is in exactly one cycle. A minimum-cost cycle cover is a cycle cover such that the sum of the weights of the edges in all the cycles is minimized. A minimum-cost cycle cover can be computed in  $O(n^3)$  time by a well-known reduction to the assignment problem [17]. Since a tour is a cycle cover,  $\text{opt}(C) \leq \text{opt}(TSP)$ . When we say that a string  $s_i$  is in some cycle  $c$  of cycle cover  $C$ , we mean that the vertex  $v_i$  with which  $s_i$  is associated is in cycle  $c$ . Throughout the paper, when we refer to a cycle, we will be referring to a cycle that is in a minimum-cost cycle cover in the distance graph.

Because  $ov(s_i, s_j) + d(s_i, s_j) = |s_i|$ , one could also weight the edges by their overlap, find a maximum-cost cycle cover and obtain the same solution. A superstring which has minimum length, or distance, also has maximum overlap. However, this correspondence breaks down for approximations; approximating the largest overlap appears to be an easier problem (cf. [23, 22, 15]) than approximating the shortest superstring.

We now describe a generic superstring algorithm used, in some form, by [4],[22] and [6]. An execution of the algorithm appears as Fig. 1.

#### GENERIC SUPERSTRING ALGORITHM

- 1) Find a minimum cost cycle cover  $C$  in the distance graph  $G$ .

- 2) For each cycle  $c \in C$ , choose one string to be a representative  $r_c$ .  
Let  $G'$  be the subgraph induced by the representative set  $R$ .
- 3) Compute a cycle cover  $CC$  on  $G'$ .
- 4) Break each cycle  $\gamma \in CC$  by deleting one edge.
- 5) Concatenate the remaining strings arbitrarily.
- 6) Extend each representative  $r_c$  by the concatenation of the prefixes around  $c$ .

The first cycle cover identifies sets of strings that have large amounts of overlap. This allows us to form the second cycle cover, in which approximating overlap and the string length are roughly comparable, so stronger bounds apply. Step (6) correctly extends the superstring for  $R$  into a superstring for  $S$ , as proved in [22].

We now analyze the generic algorithm in a way that anticipates our improvements. A more detailed analysis appears in [4]. Let  $d(C')$  be the sum of the distances and  $ov(C')$  be the sum of the overlaps of the edges in a cycle cover  $C'$ . Consider the second cycle cover  $CC$ . Let  $opt(R)$  be the optimal superstring on the strings in  $r_c \in R$  and observe that  $opt(R) \leq opt(S)$ . Let  $\bar{\alpha}$  be the string produced in Step 5, a superstring of  $R$ , and let  $opt(ov(R)) = |R| - opt(R)$  be the amount of overlap in the optimal superstring for  $R$ . Since the shortest superstring for  $R$  is a cycle cover for  $G'$ ,  $ov(CC) \geq opt(ov(R))$ . However, the superstring  $\bar{\alpha}$  does not have as much overlap as  $CC$ , since we delete one edge from each cycle.

For a cycle  $\gamma$ , let  $ov_\gamma^n$  denote the overlap in the edge deleted in Step 4, and let  $ov_\gamma^t$  denote the remaining overlap in  $\gamma$ . Let  $ov^t = \sum_{\gamma \in CC} ov_\gamma^t$  and  $ov^n = \sum_{\gamma \in CC} ov_\gamma^n$ . Thus,  $|\bar{\alpha}| \leq |R| - ov^t$ . By definition,  $|R| \leq opt(R) + opt(ov(R)) \leq opt(R) + ov(CC)$ . Combining these two inequalities with  $ov(CC) = ov^n + ov^t$ , gives that  $\bar{\alpha} \leq opt(R) + ov^n$ . We then must extend each cycle, in Step 6. Let  $Ext(\gamma)$  be the cost of extending all cycles  $c \in C$  s.t.  $r_c \in \gamma$ . Then we can express the length of  $\alpha$ , the string obtained, as

$$|\alpha| \leq opt(R) + \sum_{\gamma \in CC} (ov_\gamma^n + Ext(\gamma)) . \quad (1)$$

Let  $d(c)$  be the sum of the weights of the edges of a cycle  $c \in C$ ; so  $d(C) = \sum_{c \in C} d(c)$ . To obtain a 3-approximation, observe that the set of edges which contribute to  $ov^n$  form a matching  $M$  on  $G'$ . Now we employ a key lemma from [4]:

**Lemma 2.1 ([4])** *Let  $c, c'$  be cycles in a minimum cycle cover  $C$  with strings  $s \in c$  and  $s' \in c'$ . Then the overlap between  $s, s'$  is less than  $d(c) + d(c')$ .*

Since  $M$  is a matching, each cycle  $c$  is at an endpoint of a string at most once, and hence  $ov^n \leq d(C)$ . Now, we extend  $\bar{\alpha}$  by the edge labels on each cycle, adding a total of  $d(C)$  to the length of the string. Let  $\alpha$  be the resulting string. We conclude that

$$|\alpha| \leq opt(R) + \sum_{\gamma \in CC} ov_\gamma^n + Ext(\gamma) \leq opt(R) + d(C) + d(C) \leq 3opt(S) , \quad (2)$$

since both  $d(C)$  and  $opt(R)$  are lower bounds on  $opt(S)$ .

The analysis above makes it clear that the cycle cover  $CC$  actually partitions the cycles in the cycle cover  $C$ , and hence each cycle in  $CC$  can be analyzed separately. As was observed by [22] in their  $2\frac{8}{9}$  algorithm, if  $\gamma$  has three or more vertices, then  $ov_\gamma^n \leq \frac{2}{3} \sum_{c \in \gamma} d(c)$ .

Thus we can restrict our attention to 2-cycles in  $CC$ . We will analyze each 2-cycle in  $CC$  separately, and obtain a  $2\frac{2}{3}$  bound by proving structural properties of these cycles.

Given a representative  $v = r_c$  for some cycle  $c$ , we use  $c_v$  to denote the cycle  $c$  of which  $v$  is a representative. We summarize this discussion with the following lemma:

**Lemma 2.2** *An algorithm following the framework of the generic algorithm above, that, for each 2-cycle  $\gamma$  in  $CC$  consisting of vertices  $v$  and  $t$ , attains a bound of  $ov_\gamma^n + Ext(\gamma) \leq \beta(d(c_v) + d(c_t))$ , for some  $\beta \geq \frac{5}{3}$ , is a  $(1 + \beta)$ -approximation algorithm for the shortest superstring problem.*

We define a few terms describing the structure of cycles. The reader is referred to [4] for a more complete discussion. We call a string  $s$  *irreducible* if all cyclic shifts of  $s$  yield unique strings, and *reducible* otherwise. We say that  $s$  has *periodicity*  $x$  if there exists a string  $t$  with  $|t| = x$  such that  $s$  is substring of  $t^\infty$ . Let  $per(c)$  be the string formed by concatenating all the labels on the edges of a cycle  $c$ . Then for each string  $s \in c$ ,  $s$  is a substring of  $per(c)^\infty$ . Note that  $per(c)$  must be irreducible; otherwise a cycle with less total distance could generate the same strings, contradicting the minimality of the cycle cover. The irreducibility of the periods of cycles in a minimum cycle cover will figure prominently in many of our proofs.

We can now state a corollary to Lemma 2.1 that we will also use frequently in our proofs.

**Corollary 2.3** ([4]) *Let  $w$  be a substring of both  $(\sigma_j)^\infty$  and  $(\sigma_k)^\infty$ . Then if  $|w| \geq |\sigma_j| + |\sigma_k|$ , either  $\sigma_j$  or  $\sigma_k$  is reducible.*

### 3 Repeaters and their Characteristics

In the previous section, we saw that in order to obtain a better approximation for the shortest superstring problem it is sufficient to consider 2-cycles in the second cycle cover of the generic superstring algorithm. In this section we describe the machinery for describing 2-cycles developed in [2].

Suppose we choose  $v$  and  $t$  as representatives of two cycles of the first cycle cover  $C$ , and they form a 2-cycle in  $CC$  in which one of  $ov(v, t)$  or  $ov(t, v)$  is large but the other is small. In Step 4 we will break the 2-cycle to form a string, and since we are trying to maximize overlap, the obvious choice is to keep the high-overlap edge and discard the other. But if both edges have high overlap, we must discard one of them. In a 2-cycle this will cost us up to half of the overlap, which is the “worst case” of the generic algorithm. We observe that both edges in such a 2-cycle cannot participate in an optimal solution; in this sense the second cycle cover has achieved “false overlap”. We formalize the idea of a “high-overlap 2-cycle” as follows:

**Definition 3.1** Let  $\gamma$  be a 2-cycle in the second cycle cover  $CC$  of the GENERIC algorithm, consisting of vertices  $r_j$  and  $r_k$ , the representatives of cycles  $c_j$  and  $c_k$  in  $C$ . Without loss of generality assume that  $d(c_j) \geq d(c_k)$ . Then  $\gamma$  is a  $(g, h)$ -HO2-cycle if  $\min\{ov(r_j, r_k), ov(r_k, r_j)\} \geq gd(c_j) + hd(c_k)$ .

Our strategy is to anticipate, when we choose representatives, the potential of each string to participate in a  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycle. In particular we evaluate the potential of each string to play the role of the larger-period string in the 2-cycle. Such a string must have a very specific structure; if we find a string without such a structure, we use it as representative. Otherwise we know a great deal about the structure of the entire

$ \begin{array}{l} z = \quad \underline{\text{abababrstababab}} \\ \phantom{z = } \phantom{\underline{\text{abababrstababab}}} \phantom{)} \phantom{(} \\ y = y_\ell = \text{ababab} \\ y_r = \phantom{\text{ababab}} \quad \text{ababab} \\ \sigma = \quad \text{ab} \\ \phantom{\sigma = } \phantom{\text{ab}} \phantom{(a)} \\ \phantom{\sigma = } \phantom{\text{ab}} \phantom{(a)} \phantom{(b)} \end{array} $	$ \begin{array}{l} z = \quad \underline{\text{ababadababadabababadababadabab}} \\ \phantom{z = } \phantom{\underline{\text{ababadababadabababadababadabab}}} \phantom{(} \phantom{(} \phantom{(} \\ y = y_\ell = \text{ababadababadababa} \\ y_r = \phantom{\text{ababadababadababa}} \quad \text{ababadababadabab} \\ \sigma = \quad \text{ababad} \\ \phantom{\sigma = } \phantom{\text{ababad}} \phantom{(b)} \\ \phantom{\sigma = } \phantom{\text{ababad}} \phantom{(b)} \phantom{(a)} \end{array} $
---	---

Figure 2: Positive and Negative Characteristics.  $\text{per}(c)$  is underlined. (a) shows a negative characteristic. (b) shows a positive characteristic and  $\sigma$  are also shown.

cycle and can trade off the amount of two-way overlap against the cost of extending the representative to include the rest of the cycle.

In order to have the potential to be the larger-period string in a high-overlap 2-cycle, a string  $z$  must have a significant prefix that has some smaller period. This smaller period might correspond to the period of another cycle in the cover, and hence some other representative  $w$  such that  $\text{ov}(w, z)$  would be large. The suffix of  $z$  must similarly have the same smaller period, so that  $\text{ov}(z, w)$  would be large. We require some notation to describe this potential.

**Definition 3.2** Let  $z$  be a string in cycle  $c$  and let  $\sigma$  be an irreducible string with  $|\sigma| < d(c)$ . Then  $\sigma$  is a  $(g, h)$ -repeater of  $z$  if there exist witnesses  $y_\ell$  and  $y_r$ , such that

1.  $y_\ell$  is a prefix of  $z$  and  $y_r$  is a suffix of  $z$ .
2.  $y_\ell$  and  $y_r$  are substrings of  $(\sigma)^\infty$ .
3.  $|y_\ell|, |y_r| > gd(c) + h|\sigma|$ .

We will always choose  $y_\ell$  and  $y_r$  to be the maximum length prefix and suffix that satisfy conditions 1–3 above.

Consider the string  $z$  in Fig. 2b and let  $g = h = \frac{2}{3}$ . Here  $\text{per}(c) = \text{ababadababadab}$ ,  $\sigma = \text{ababad}$ ,  $y_\ell = \text{ababadababadababa}$  and  $y_r = \text{ababadababadabab}$ . So  $|y_\ell|, |y_r| > \frac{2}{3}d(c) + \frac{2}{3}|\sigma|$ , and we say that  $\sigma$  is a  $(\frac{2}{3}, \frac{2}{3})$ -repeater of  $z$ .

Note that in our example  $y_\ell$  and  $y_r$  are almost the same; this is not a complete coincidence. All the repeaters we will be considering in this paper will have  $g \geq \frac{1}{2}$  and hence  $y_\ell$  and  $y_r$  *must* overlap, often significantly (as in this example). For convenience we will define one witness  $y_\sigma$  which contains both  $y_\ell$  and  $y_r$ ; that is, we define  $y_\sigma$  to be the maximum-length substring of  $(\sigma)^\infty$  that is also a substring of  $\text{per}(c)^\infty$ . In other words, if you took  $\sigma$  and tried to repeat it as many times as possible, in both directions, while being consistent with  $c$ , you get  $y_\sigma$ . In the example above  $y_\sigma = y_\ell$ . When the context is clear, we will drop the  $\sigma$  and just refer to witness  $y$ .

Henceforth when discussing and proving properties of cycles, we will refer to the maximal witness  $y_\sigma$  rather than to the underlying pair of witnesses  $y_\ell$  and  $y_r$ . This simplification is conservative.

The idea behind  $(g, h)$ -repeaters is to identify periodic substrings of the period of a cycle in  $C$ . We will also be interested in identifying that portion of a cycle that is *not* consistent with some  $(g, h)$ -repeater  $\sigma$ . Note that a copy of  $y_\sigma$  begins every  $d(c)$  in  $\text{per}(c)^\infty$ , and that  $|y| < 2d(c)$ , since by Corollary 2.3,  $|y| < d(c) + |\sigma| \leq 2d(c)$ .

**Definition 3.3** Let  $c$  be a cycle with  $(g, h)$ -repeater  $\sigma$  and maximal witness  $y$ . Fix a copy of  $y$  in  $\text{per}(c)^\infty$ . The point just to the left of the first character of  $y$  is the *head of  $y$* . Index this point as 0 and continue the indices between each character leftward and rightward to cover the interval  $[-d(c)..d(c)]$ . Now mark the point  $|y| - d(c)$  and call it the *tail of  $\sigma$* . The *characteristic* of  $\sigma$ ,  $X_\sigma$ , is the interval from the head to the tail. If  $|y| - d(c) > 0$  we call  $[0..|y| - d(c)]$  a *positive characteristic*  $X_\sigma$ . If  $|y| - d(c) \leq 0$  we call  $[|y| - d(c)..0]$  a *negative characteristic*  $X_\sigma$ .

We can picture the characteristics of the repeaters of a cycle  $c$  in terms of parentheses. Fig. 2b illustrates this idea for positive characteristics. The left and right ends of  $y_\sigma$  are marked with left and right parentheses; these correspond to the head and tail of adjacent copies of  $X_\sigma$ .

A negative characteristic appears in Fig. 2a and can be pictured as a single solid entity (perhaps of size zero) which spans the gap between copies of  $y$ . In this example  $rst$  is the negative characteristic. Each characteristic appears once every  $d(c)$ . Intuitively, the characteristic of a repeater borders the portion of  $\text{per}(c)$  which must be included as a prefix and suffix of some string  $z$  if  $z$  is to participate in a high-overlap 2-cycle. Recall that we defined  $(g, h)$ -repeaters (Def. 3.2) in terms of some string  $z$  in a cycle  $c$  which contained witnesses  $y_\ell$  and  $y_r$  as a prefix and suffix. In general there might be several such strings in  $c$  which could satisfy the definition. We say that  $\sigma$  is *active* in each of these strings. We say that two characteristics  $X_{\sigma_i}, X_{\sigma_j}$  are *nested* if  $X_{\sigma_i}$  is a positive characteristic and  $X_{\sigma_j}$  falls within  $X_{\sigma_i}$ . We say that two characteristics  $X_{\sigma_i}, X_{\sigma_j}$  are *disjoint* if their intervals are disjoint. Otherwise we say that  $X_\sigma$  and  $X_{\sigma'}$  are *linked*.

We will frequently be interested in the relationship between two substrings of  $\text{per}(c)^\infty$ , for instance between two witness strings  $y$  and  $y'$ . As noted above, a copy of any substring of  $\text{per}(c)^\infty$  occurs every  $d(c)$  in  $\text{per}(c)^\infty$ . We overload our notation for  $d(\cdot)$  and  $\text{ov}(\cdot)$  in the obvious way to refer to prefix distance  $d(y, y')$  and overlap  $\text{ov}(y, y')$ . We also define the *suffix distance*  $\tilde{d}(y, y')$  to be the distance from the last character of a copy of  $y$  to the last character of the first copy of  $y'$  that ends after  $y$ .

We will require the following bound on the length of a witness string:

**Lemma 3.4 ([2])** *Let  $y_\sigma$  be a maximal witness for some  $(g, h)$ -repeater  $\sigma$  in a cycle  $c$ . Then  $|y_\sigma| < d(c) + |\sigma| < 2d(c)$ .*

## 4 The Algorithm

We present our algorithm G-SHORTSTRING, which is a  $2\frac{2}{3}$ -approximation algorithm for the shortest superstring problem. We describe the algorithm in Section 4.1. In order to prove our bound on its approximation ratio in Section 4.3, we present some technical lemmas on the structure of cycles with  $(\frac{2}{3}, \frac{2}{3})$ -repeaters in Section 4.2.

### 4.1 Algorithm G-SHORTSTRING

In order to achieve a bound of  $2\frac{2}{3}$  within the framework of GENERIC, Lemma 2.2 states that we need to concentrate on  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycles. As in [2], our strategy is to anticipate, when we select a representative  $r_j$ , the possible involvement of  $r_j$  as the larger-period string in a  $(g, h)$ -HO2-cycle. In [2] we used criteria for doing so which were based on our detailed knowledge of the structure of  $(\frac{3}{4}, \frac{3}{4})$ -repeaters. G-SHORTSTRING does not depend on such knowledge.



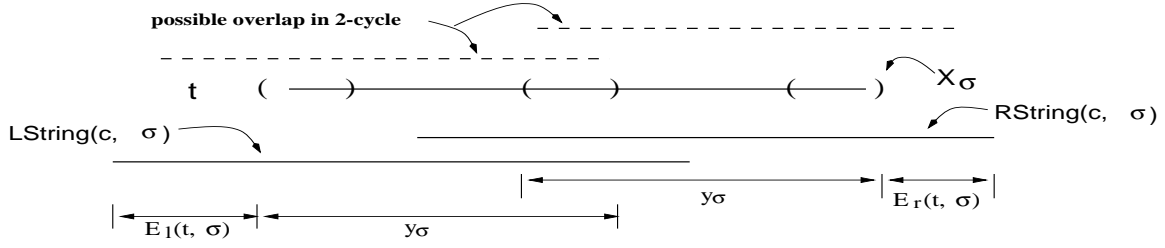


Figure 3: Definitions 4.2 and 4.3.

Our new procedure for selecting representatives is to evaluate a cost function for each string in a cycle, and to select the string with the best *worst-case* cost. We identify a cost function which resembles the desired bounds, and we explicitly attempt to minimize this function in the algorithm. We achieve our improved bound by more careful extension of each representative  $r_j$  of a cycle  $c_j$  that is also the larger-period string in a  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycle. We therefore need some new ideas about extension and some notation for expressing it.

**Definition 4.1** Let  $\sigma$  be a  $(g, h)$ -repeater with maximal witness  $y_\sigma$  in an  $m$ -cycle  $c$ . Index the strings  $s_i$  such that  $d(y_\sigma, s_i) < d(y_\sigma, s_{i+1})$ ,  $1 \leq i < m$ . Then we define the *right string of  $\sigma$  in  $c$* ,  $\text{RString}(c, \sigma) = s_m$ . The *left string of  $\sigma$  in  $c$* ,  $\text{LString}(c, \sigma)$  is defined symmetrically; reindex the strings  $s_i$  such that  $\tilde{d}(s_i, y_\sigma) > \tilde{d}(s_{i+1}, y_\sigma)$ ,  $1 \leq i < m$ . Then we define  $\text{LString}(c, \sigma) = s_1$ .

In other words, if we align a copy of each of the strings in  $c$  in such a way that the first one begins as soon after a copy of  $y_\sigma$  as possible, then the rightmost string is  $\text{RString}(c, \sigma)$ . The idea is that if we choose as representative a string  $t$  in which  $\sigma$  is active, and  $t$  becomes the larger-period string in a  $(g, h)$ -HO2-cycle, then  $\text{RString}(c, \sigma)$  is the rightmost string which we will have to include if we extend to the right. Figure 3 illustrates Definitions 4.1 and 4.2.

**Definition 4.2** Let  $\sigma$  be a  $(g, h)$ -repeater which is active in a string  $t$  in cycle  $c$ . Then the *right  $\sigma$ -extension with respect to  $t$* ,  $E_r(t, \sigma) = \tilde{d}(y_\sigma, \text{RString}(c, \sigma))$ . The *left  $\sigma$ -extension with respect to  $t$* ,  $E_\ell(t, \sigma) = d(\text{LString}(c, \sigma), y_\sigma)$ .

Given a  $(\frac{2}{3}, \frac{2}{3})$ -repeater  $\sigma$  which is active in a string  $t$  in cycle  $c_t$ , we wish to calculate the cost of choosing  $t$  and having  $t$  involved in a  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycle  $\gamma$  with some string  $v$  such that  $\text{per}(c_v) = \sigma$ . In particular, by Lemma 2.2 we are interested in anticipating  $\text{ov}_\gamma^n + \text{Ext}(c_t) + \text{Ext}(c_v)$ . Consider without loss of generality the right end of  $t$ , and let  $y_r$  be the suffix of  $t$  which is the witness string for  $\sigma$ . Then we know that  $\text{ov}_\gamma^n \leq |y_r| \leq |y_\sigma|$ . If there is slack in either of these inequalities, then we use the slack as part of our upper bound on extension cost. We have to extend to the right well beyond the end of  $y_\sigma$  in any case, so it does not matter whether we charge  $|y_\sigma| - |y_r|$  to  $\text{ov}_\gamma^n$  or to  $\text{Ext}(c_t)$ . From the end of  $y_\sigma$ , we need to extend to the right only as far as necessary to include  $\text{RString}(c, \sigma)$ . We also have to extend  $v$  to include the remaining strings in  $c_v$ ; we assume the cost of full extension. This motivates the following definition.

**Definition 4.3** Let  $\sigma$  be a  $(g, h)$ -repeater that is active in string  $t$  in cycle  $c$ . Then the *anticipated cost* of choosing  $t$  as representative and forming a 2-cycle with a string with

period  $\sigma$  is

$$\text{Cost}(t, \sigma) = |y_\sigma| + \min\{E_\ell(t, \sigma), E_r(t, \sigma)\} + |\sigma|.$$

What we seek, then, is to *minimize*, in our choice of representative  $t$ , the *maximum* over all  $(\frac{2}{3}, \frac{2}{3})$ -repeaters active in  $t$ , the anticipated cost  $\text{Cost}(t, \sigma)$ . Allowing  $\sigma \in t$  to mean “ $\sigma$  active in  $t$ ”, we seek .

$$\text{BestRep}(c) = \operatorname{argmin}_{t \in c} \left\{ \max_{\sigma \in t} \{\text{Cost}(t, \sigma)\} \right\}$$

Procedure G-FINDREPS( $c$ ), shown below, calculates the anticipated cost for each pair  $(t, \sigma)$  such that  $t$  is a string in  $c$  and  $\sigma$  is active in  $t$ .

**Procedure G-FINDREPS( $c_j$ )**

- 1) Find all  $(\frac{2}{3}, \frac{2}{3})$ -repeaters and associated characteristics in  $c_j$ .
- 2) **If** any string  $t$  has no  $(\frac{2}{3}, \frac{2}{3})$ -repeaters  
     **Then**  $r_j = t$ ;
- 3) **Else**  
      $r_j = \text{BestRep}(c_j)$ ;
- 4) **Return**  $r_j$ .

The main body of G-SHORTSTRING is exactly SHORTSTRING, except that representatives are selected in Step 2) by a call to procedure G-FINDREPS( $c$ ), and the parameters of the  $(g, h)$ -HO2-cycle in Step 4) are different.

**Algorithm G-SHORTSTRING**

- (1) Form minimum cycle cover  $C$  on distance graph  $G$ .
- (2) For each cycle  $c \in C$   
     Call G-FINDREPS( $c$ ) to choose representative  $r_c$ .  
     Add  $r_c$  to  $R$ .  
     Let  $G'$  be the subgraph induced by  $R$ .
- (3) Form minimum cycle cover  $CC$  on  $G'$ .
- (4) For each cycle  $\gamma$  in  $CC$ :  
     **if**  $\gamma$  is a  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycle  $(v, t)$   
     (a)     **then if**  $\text{ov}(t, v) + E_r(t, \text{per}(c_v)) \leq \text{ov}(v, t) + E_\ell(t, \text{per}(c_v))$   
             **then** Extend  $\langle v, t \rangle$ ;  
             **else** Extend  $\langle t, v \rangle$ ;  
     (b) **else** discard edge of cycle  $\gamma$  with least overlap; Extend each vertex  $w$  by  $d(c_w)$
- (5) Concatenate strings from (4) to form superstring  $\alpha$

In Step (4a) above, the instruction “Extend  $\langle v, t \rangle$ ” is shorthand for the following idea. We extend  $v$  to the left to include all of the strings in  $c_v$ ; we assume in the analysis in Section 4.3 that this length is  $d(c_v)$ . We extend  $t$  to the right, as far as is necessary to include  $\text{RString}(c, \text{per}(c_v))$ . Extending  $\langle t, v \rangle$  is done symmetrically.

The algorithm G-SHORTSTRING correctly computes a superstring of the set of strings  $S$ . This follows from the correctness of GENERIC. Our method of choosing representatives for each cycle is a special case of the method of GENERIC, which chooses an arbitrary string as representative. In step (4b), we do exactly what GENERIC does. In step (4a), we use a different criterion for breaking a cycle  $\gamma \in CC$ , and we only extend each

representative far enough to “cover” all of the strings in its cycle. Each string is therefore included in the solution  $\alpha$ .

G-SHORTSTRING runs in polynomial time. The distance graph  $G$  can be built in  $O(|S| + n^2)$  time [11], and the cycle cover computations take  $O(n^3)$  time [17]. These two results determine the running time of GENERIC. In addition, our algorithm must find all of the  $(\frac{2}{3}, \frac{2}{3})$ -repeaters in each cycle  $c \in C$  in G-FINDREPS( $c$ ). This can be done naively in polynomial time by examining a prefix and suffix of each string, and determining whether the prefix and suffix have periodicity  $2 \leq j < d(c)$ .

In order to analyze the approximation ratio achieved by G-SHORTSTRING, we require a few technical lemmas pertaining to  $(\frac{2}{3}, \frac{2}{3})$ -repeaters.

## 4.2 Properties of Strings with $(\frac{2}{3}, \frac{2}{3})$ -Repeaters

A *small*  $(g, h)$ -repeater in a cycle  $c$  is one whose minimum witness length is less than  $d(c)$ . A  $(\frac{2}{3}, \frac{2}{3})$ -repeater  $\sigma$  is small if  $|\sigma| < \frac{1}{2}d(c)$ . Generally we are interested in avoiding small  $(g, h)$ -repeaters. To see why this is so, suppose that we choose a representative  $r_j$  for cycle  $c_j$ , and  $r_j$  is involved in a  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycle with another representative  $r_k$  of cycle  $c_k$ , and  $\text{per}(c_k) = \sigma$ . Then we will want to bound the extension cost we incur,  $\text{Ext}(c_j)$ , in terms of  $d(c_j) + d(c_k) = d(c_j) + |\sigma|$ . So if  $\sigma$  is larger, then our extension cost, as a fraction of  $d(c_j) + d(c_k)$ , is smaller.

There may be several small  $(\frac{2}{3}, \frac{2}{3})$ -repeaters in a cycle, but we are able to bound the number of small  $(\frac{2}{3}, \frac{2}{3})$ -repeaters in a string.

**Lemma 4.4** *Let  $s$  be a string in a cycle  $c$ . Then at most one small  $(\frac{2}{3}, \frac{2}{3})$ -repeater can be active in  $s$ .*

**Proof:** Suppose for purpose of contradiction that there exist two such  $(\frac{2}{3}, \frac{2}{3})$ -repeaters  $\sigma$  and  $\sigma'$ . Let  $y_\ell(\sigma)$  and  $y_\ell(\sigma')$  be the prefixes of  $s$  which are the left witness strings of  $\sigma$  and  $\sigma'$  respectively. Let  $y_\ell = \text{argmin}\{|y_\ell(\sigma)|, |y_\ell(\sigma')|\}$  be the prefix of  $s$  which is periodic in both  $\sigma$  and  $\sigma'$ . Applying Corollary 2.3, Definition 3.2, and the fact that  $|\sigma'| < \frac{1}{2}d(c)$ , we get

$$\begin{aligned} |\sigma| &> |y_\ell| - |\sigma'| \\ &> \frac{2}{3}d(c) - \frac{1}{3}|\sigma'| \\ &> \frac{1}{2}d(c), \end{aligned}$$

a contradiction since  $\sigma$  is a small  $(\frac{2}{3}, \frac{2}{3})$ -repeater. ■

The following lemma gives us a lower bound on the size of a  $(\frac{2}{3}, \frac{2}{3})$ -repeater whose characteristic has the characteristic of another  $(\frac{2}{3}, \frac{2}{3})$ -repeater nested within it.

**Lemma 4.5** *Let  $X_\sigma$  be a positive characteristic in cycle  $c$  and  $X_{\sigma'}$  a characteristic nested within  $X_\sigma$  with  $|\sigma| > |\sigma'|$ . Then  $|\sigma| > \frac{1}{2}d(c)$ .*

**Proof:** In this case the witness  $y'$  is completely contained within the witness  $y$ . We apply Corollary 2.3 and the definition of  $(\frac{2}{3}, \frac{2}{3})$ -repeater to get

$$\begin{aligned} |\sigma| + |\sigma'| &> \text{ov}(y, y') \\ &> |y'| \\ &> \frac{2}{3}(d(c) + |\sigma'|) \\ \Rightarrow |\sigma| &> \frac{2}{3}d(c) - \frac{1}{3}|\sigma'|, \end{aligned}$$

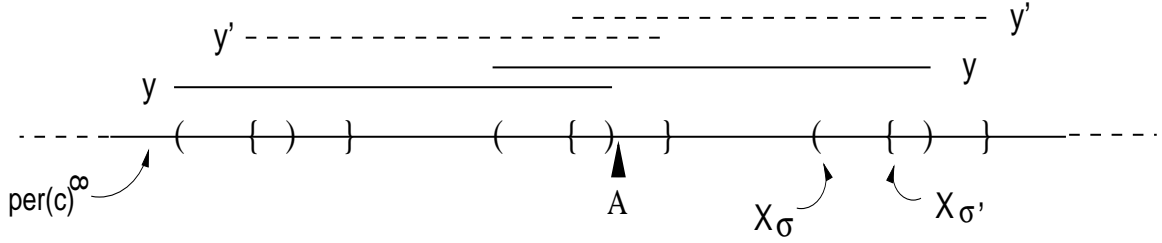


Figure 4: The characteristics  $X_\sigma$  and  $X_{\sigma'}$  are linked, as in Lemma 4.7.

which implies  $|\sigma| > \frac{1}{2}d(c)$  because  $|\sigma| > |\sigma'|$ . ■

Because  $(\frac{2}{3}, \frac{2}{3})$ -repeaters may not be well parenthesized, we will often be faced in our analysis with situations in which two positive characteristics are linked, as pictured in Figure 4. (Recall that two positive characteristics are linked if they overlap, but neither contains the other.) The following lemma and its corollary gives us strong bounds on the size of the two  $(\frac{2}{3}, \frac{2}{3})$ -repeaters and on their difference. In order to prove the lemma, we require a proof technique introduced in [2], the *shift argument*. We describe this technique below.

We apply the shift argument to cycles that include two or more repeaters. We are generally interested in proving that some property holds; we assume that it does not, and use the shift argument to derive a contradiction. We begin with the following observation, which can easily be verified by the definition of maximal witness.

**Observation 4.6** *Let  $y$  be the maximal witness for a  $(g, h)$ -repeater  $\sigma$  in a cycle  $c$ , and fix a copy  $y^*$  of  $y$  in  $\text{per}(c)^\infty$ . Index the character positions of  $\text{per}(c)^\infty$  with the character to the left of  $y^*$  as 0, the first character of  $y^*$  as 1, and continuing to the right beyond the end of  $y^*$ . Let  $\text{Char}(i)$  be the character in position  $i$ . Then*

$$a) \text{Char}(0) \neq \text{Char}(|\sigma|)$$

and

$$b) \text{Char}(|y^*| - |\sigma| + 1) \neq \text{Char}(|y^*| + 1).$$

In each shift argument our goal will be to show that either inequality a) or b) in Observation 4.6 is violated and the terms are indeed equal. We will do so by making a series of *shifts* between characters, which we know to be identical, by the periodic structure of the strings. In particular, within any  $y_\sigma$ , any two characters that are  $\sigma$  apart are identical, and in  $\text{per}(c)^\infty$ , any two characters that are  $d(c)$  apart are identical. We call such shifts *valid*. We will begin at either the character immediately preceding or following a copy of  $y$  or  $y'$ , and perform a series of shifts which will bring us to the position whose character is supposed to be unequal. If these shifts are valid, then the two characters must be equal, contradicting our initial assumption that the characteristics  $X_\sigma$  and  $X_{\sigma'}$  could overlap.

We introduce notation to describe the sequence of shifts. We give a starting position and a position at which we wish to arrive, relative to the starting position. We also give the series of moves and a set of requirements, that is, conditions on the various parameters that must be met in order for the moves to all be valid. Below the box, we show that the conditions for validity are indeed satisfied, which gives us a contradiction for the region in which the shifts are valid.

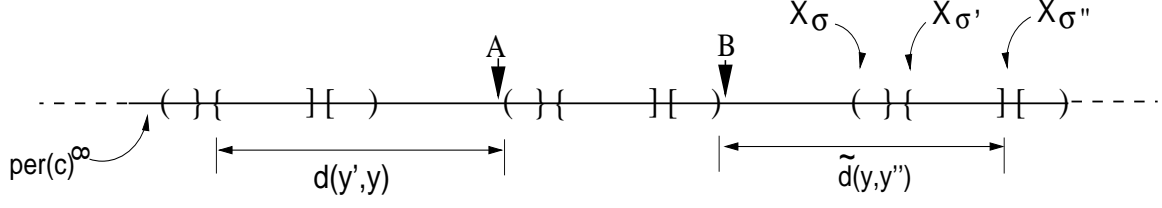


Figure 5: Proof of Lemma 4.9.

**Lemma 4.7** *Let  $\sigma$  and  $\sigma'$  be two  $(\frac{2}{3}, \frac{2}{3})$ -repeaters with positive characteristics in a cycle  $c$ , with  $|\sigma| > |\sigma'|$ , and  $X_\sigma$  and  $X_{\sigma'}$  linked. Let  $k = \lfloor \frac{|\sigma|}{|\sigma'|} \rfloor$ . Then  $|\sigma| - k|\sigma'| > |y_\sigma| - d(c)$ .*

**Proof:** We apply the following shift argument, using start position (A) in Figure 4:

Start: (A)			Goal: $- \sigma $
No.	Move	Requirement	Comments
1.	$+k \sigma' $	$k \sigma'  < d(c)$	$k \sigma'  <  \sigma  < d(c)$
2.	$- \sigma $	$ \sigma  - k \sigma'  \leq  y_\sigma  - d(c)$	See below.
3	$-k \sigma' $	$ \sigma  < d(c)$	Def. Repeater.

Because only move #2 is the only one whose validity is conditional, we conclude the negation of that condition, i.e.  $|\sigma| - k|\sigma'| > |y_\sigma| - d(c)$ . ■

**Corollary 4.8** *Let  $\sigma$  and  $\sigma'$  be two  $(\frac{2}{3}, \frac{2}{3})$ -repeaters with positive characteristics in a cycle  $c$ , with  $|\sigma| > |\sigma'|$ , and  $X_\sigma$  and  $X_{\sigma'}$  linked. Let  $k = \lfloor \frac{|\sigma|}{|\sigma'|} \rfloor$ . Then  $|\sigma'| > |y_\sigma| - d(c)$ .*

**Proof:** By the choice of  $k$  and Lemma 4.7,

$$|\sigma'| > |\sigma| - k|\sigma'| > |y_\sigma| - d(c).$$

■

In our analysis, we will be interested in *lower bounds* on the size of potentially small  $(\frac{2}{3}, \frac{2}{3})$ -repeaters in terms of some measure of distance which will correspond to extension cost. The following two lemmas provides such bounds for two important cases in which three characteristics are involved. Our choice of dimensions for identifying the relative positions of the three characteristics will seem unnatural now, but will simplify our task in Section 4.3.

**Lemma 4.9** *Let  $\sigma, \sigma'$  and  $\sigma''$  be  $(\frac{2}{3}, \frac{2}{3})$ -repeaters in cycle  $c$ , with  $X_{\sigma'}$  and  $X_{\sigma''}$  disjoint, and with  $X_{\sigma'}$  nested to the left of  $X_{\sigma''}$  within  $X_\sigma$ , Then  $|\sigma'| > d(y', y) + |y| - 2d(c)$  and  $|\sigma''| > \tilde{d}(y, y'') + |y| - 2d(c)$ .*

**Proof:** Figure 5 illustrates the start positions of our shift arguments.

Start: (A)			Goal: $+ \sigma $
No.	Move	Requirement	Comments
1.	$+ \sigma'' $	$ \sigma''  \leq \tilde{d}(y, y'') +  y  - 2d(c)$	See below.
2.	$+ \sigma $	$ \sigma  +  \sigma''  <  y  - d(c) + \tilde{d}(y, y'')$	See below.
3	$- \sigma'' $	$ \sigma  +  \sigma''  >  y  - d(c) +  \sigma'' $	See below.

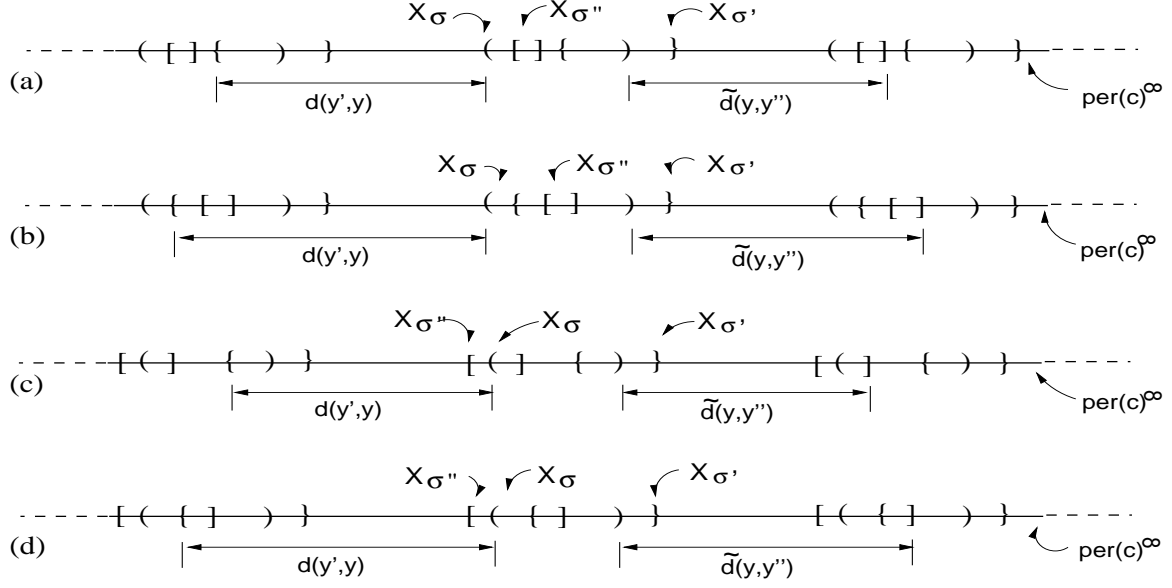


Figure 6: Lemma 4.10:  $X_{\sigma'}$  is linked with  $X_{\sigma}$ .  $X_{\sigma''}$  may be nested within  $X_{\sigma'}$  (b), or only nested within  $X_{\sigma}$  (a). If  $X_{\sigma''}$  is linked with  $X_{\sigma}$ , it may or may not also be linked with  $X_{\sigma'}$  ((d) and (c) respectively).

Start: (B)			Goal: $- \sigma $
No.	Move	Requirement	Comments
1.	$- \sigma' $	$ \sigma'  \leq d(y', y) +  y  - 2d(c)$	See below.
2.	$- \sigma $	$ \sigma  +  \sigma'  <  y  - d(c) + d(y', y)$	See below.
3.	$+ \sigma' $	$ \sigma  +  \sigma'  >  y  - d(c) +  \sigma' $	See below.

Requirement #3 for both of the above is always true because  $|\sigma| > |y| - d(c)$  by Lemma 3.4. Because  $|\sigma| < d(c)$ , #1 implies #2; therefore #1 must be false:

$$\begin{aligned}
|\sigma''| &> \tilde{d}(y, y'') + |y| - 2d(c) \\
|\sigma'| &> d(y', y) + |y| - 2d(c).
\end{aligned}$$

■

**Lemma 4.10** *Let  $\sigma$ ,  $\sigma'$  and  $\sigma''$  be  $(\frac{2}{3}, \frac{2}{3})$ -repeaters in cycle  $c$ , with maximal witnesses  $y$ ,  $y'$  and  $y''$ . Let  $|\sigma| > |\sigma'| > |\sigma''|$ , and  $X_{\sigma}$  and  $X_{\sigma'}$  positive. If  $X_{\sigma'}$  is linked with  $X_{\sigma}$ , then  $\min\{d(y', y), \tilde{d}(y, y'')\} < \frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y_{\sigma}|$ .*

**Proof:** By the condition of the lemma we know that  $X_{\sigma}$  and  $X_{\sigma'}$  are linked, but we do not know the relationship between between  $X_{\sigma''}$  and the other two characteristics. The characteristic  $X_{\sigma''}$  may be nested within one or both of  $X_{\sigma}$  and  $X_{\sigma'}$  as in Figure 6(a) or (b), or it may be linked with one or both of  $X_{\sigma}$  and  $X_{\sigma'}$  as in Figure 6 (c) or (d). In any of these cases we can apply Corollary 2.3 to the overlap between  $y'$  and  $y''$ :

$$|\sigma'| + |\sigma''| > d(y', y) + \tilde{d}(y, y'') + |y| - 2d(c),$$

which implies

$$|\sigma'| > \frac{1}{2}(d(y', y) + \tilde{d}(y, y'')) + \frac{1}{2}|y| - d(c). \quad (3)$$

We now use Lemma 4.7 and Equation 3 to obtain

$$\begin{aligned} |\sigma| &> |\sigma'| + |y| - d(c) \\ &> \frac{1}{2}(d(y', y) + \tilde{d}(y, y'')) + \frac{1}{2}|y| - d(c) + |y| - d(c) \\ &= \frac{1}{2}(d(y', y) + \tilde{d}(y, y'')) + \frac{3}{2}|y| - 2d(c). \end{aligned}$$

Solving for  $\frac{1}{2}(d(y', y) + \tilde{d}(y, y''))$  and using Definition 3.2 gives us our result:

$$\begin{aligned} \frac{1}{2}(d(y', y) + \tilde{d}(y, y'')) &< |\sigma| - \frac{3}{2}|y| + 2d(c) \\ &< |\sigma| - |y| + 2d(c) - \frac{1}{2}\left(\frac{2}{3}d(c) + \frac{2}{3}|\sigma|\right) \\ &= \frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y|. \end{aligned}$$

■

### 4.3 Analysis of the Algorithm

We now analyze our algorithm G-SHORTSTRING. The structure of our approach is similar to that of [2], though the analysis we use in each case is completely different than that used for SHORTSTRING. We relate the performance of our algorithm to that of GENERIC; the case of interest is when a cycle in  $CC$  is a  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycle.

**Lemma 4.11** *For each cycle  $\gamma \in CC$  which is not a  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycle, Algorithm G-SHORTSTRING produces a superstring no longer than GENERIC would produce on the same cycle  $\gamma$ .*

**Proof:** We observe that step 4b) of G-SHORTSTRING handles any cycle  $\gamma \in CC$  which is not a  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycle. It selects an edge  $e$  and extends the cycle  $\gamma$  in exactly the same way as GENERIC. It then fully extends each representative  $r_\ell \in \gamma$  to cover the remaining strings in each cycle  $c_\ell$ . The only difference between the two algorithms in their handling of these cycles is that we perform full extension before concatenating with the strings from other cycles in  $CC$ . This does not affect the length of the resulting string. ■

We now must show, according to Lemma 2.2, that for each  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycle, we attain the bound specified by Lemma 2.2.

**Lemma 4.12** *Let  $\gamma$  be a  $(\frac{2}{3}, \frac{2}{3})$ -HO2-cycle in  $CC$  with  $r_j$  the representative of cycle  $c_j$  and  $r_k$  the representative of  $c_k$ . Then  $ov_\gamma^n + Ext(\gamma) \leq \frac{5}{3}(d(c_j) + d(c_k))$ .*

**Proof:** Assume without loss of generality  $d(c_j) \geq d(c_k)$ . Because  $r_j$  has high overlap at both ends with  $r_k$ , there must be at least one  $(\frac{2}{3}, \frac{2}{3})$ -repeater  $\sigma'$  in  $c_i$ , with  $\sigma' = \text{per}(c_j)$ . All strings in  $c_j$  must have at least one  $(\frac{2}{3}, \frac{2}{3})$ -repeater, otherwise we would not have chosen  $r_j$  as representative.

We consider two cases:

1. All strings in  $c_j$  have a small  $(\frac{2}{3}, \frac{2}{3})$ -repeater.
2. At least one string in  $c_j$  has no small  $(\frac{2}{3}, \frac{2}{3})$ -repeaters.

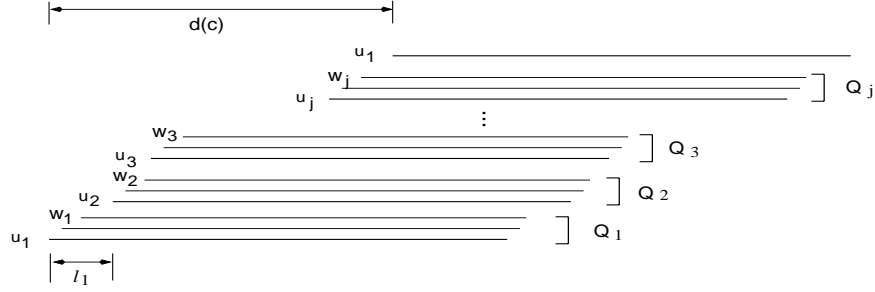


Figure 7: Case 1 of Lemma 4.12.

In each case, we must show that *some* string in  $c_j$  must have been able to achieve the bound. Because the representative is selected by comparing the worst case costs of each string, the existence of such a string is sufficient.

**Case 1:** All strings in  $c_j$  have a small  $(\frac{2}{3}, \frac{2}{3})$ -repeater.

The proof of Lemma 4.4 suggests our strategy: if two strings with different  $(\frac{2}{3}, \frac{2}{3})$ -repeaters begin near each other, then the sum of their periods must be close to  $d(c)$ . If they do not begin near each other, then we can save on extension by the amount of this gap.

Because we're in Case 1, each string has at least one small  $(\frac{2}{3}, \frac{2}{3})$ -repeater. No string has more than one small  $(\frac{2}{3}, \frac{2}{3})$ -repeater by Lemma 4.4, and so each string has exactly one small  $(\frac{2}{3}, \frac{2}{3})$ -repeater. More than one string may have the same small  $(\frac{2}{3}, \frac{2}{3})$ -repeater active.

**Claim 4.13** *Let  $\sigma$  and  $\sigma'$  be small  $(\frac{2}{3}, \frac{2}{3})$ -repeaters in cycle  $c$ . Let  $Q$  be the set of strings in which  $\sigma$  is active, and let  $Q'$  be the set of strings in which  $\sigma'$  is active. Then there is a rotation of the cyclic ordering of the strings in  $c$  such that all of the strings in  $Q$  appear before all of the strings in  $Q'$ .*

**Proof:** For purpose of contradiction let  $t$  and  $v$  be two strings in  $Q$  and let  $t'$  and  $v'$  be two strings in  $Q'$  such that they appear in the cyclic order  $t, t', v, v'$ . Without loss of generality let  $d(t, v) \leq \frac{1}{2}d(c)$ ; otherwise  $d(v, t) \leq \frac{1}{2}d(c)$  and the same argument follows. Consider the prefixes of  $t$  and  $v$  which are the left witness for  $\sigma$ ; both prefixes must be substrings of the same copy of  $y_\sigma$ . Since  $t'$  is between  $t$  and  $v$ , then it also must have a prefix  $y'_t$  which has period  $\sigma$ . The same argument holds for the suffixes of  $t, v$  and  $t'$ , so  $\sigma$  must be active in  $t'$ . But then  $t'$  has both  $\sigma$  and  $\sigma'$  active, contradicting Lemma 4.4. ■

We resume our analysis of Case 1. Let  $\sigma_1$  be the largest of the small  $(\frac{2}{3}, \frac{2}{3})$ -repeaters in  $c$ , and let  $Q_1$  be the set of strings in which  $\sigma_1$  is active. Number the remaining small  $(\frac{2}{3}, \frac{2}{3})$ -repeaters  $s_2, \dots, s_m$ , and let  $Q_i$ ,  $1 \leq i \leq m$ , be the set of strings in which  $\sigma_i$  is active. The  $Q_i$  partition the strings of the cycle, and by Claim 4.13 the  $Q_i$  form a cyclic ordering. Let  $u_i$ ,  $1 \leq i \leq m$  be the leftmost string in each group  $Q_i$ , and let  $w_i$ ,  $1 \leq i \leq m$  be the rightmost string in each group  $Q_i$ . Let  $\ell_i = d(v_i, v_{i+1})$ ,  $1 \leq i < j$ . (See Figure 7.)

First we apply Corollary 2.3 to derive a lower bound on the distance  $\ell_1$  between  $u_1$  and  $u_2$ .

$$\begin{aligned}
|\sigma_1| + |\sigma_2| &> \text{ov}(|y_{\sigma_1}|, |y_{\sigma_2}|) \\
&\geq |y_{\sigma_1}| - \ell_1 \\
\Rightarrow 2|\sigma_1| &> |y_{\sigma_1}| - \ell_1
\end{aligned}$$



$$\Rightarrow \ell_1 > |y_{\sigma_1}| - 2|\sigma_1|. \quad (4)$$

Now we bound the anticipated extension cost  $\text{Cost}(u_1, \sigma_1)$ ,

$$\begin{aligned} \text{Cost}(u_1, \sigma_1) &= |y_{\sigma_1}| + \min\{E_\ell(u_1, \sigma_1), E_r(u_1, \sigma_1)\} + |\sigma_1| \\ &\leq |y_{\sigma_1}| + E_\ell(u_1, \sigma_1) + |\sigma_1|. \end{aligned}$$

If we extend  $u_1$  to the left, the last string we will have to cover will be  $u_2$ , so  $E_\ell(u_1, \sigma_1) = d(c) - \ell_1$ , and then we use Equation 4:

$$\begin{aligned} &= |y_{\sigma_1}| + d(c) - \ell_1 + |\sigma_1| \\ &\leq d(c) + 3|\sigma_1| \\ &\leq \frac{5}{3}(d(c) + |\sigma_1|). \end{aligned}$$

The last inequality follows from the fact that  $\sigma_1$  is a small  $(\frac{2}{3}, \frac{2}{3})$ -repeater, so  $|\sigma_1| < \frac{1}{2}d(c) < \frac{1}{3}(d(c) + |\sigma_1|)$ .

**Case 2:** At least one string in  $c_j$  has no small  $(\frac{2}{3}, \frac{2}{3})$ -repeaters.

Throughout the proof of this case, we fix  $s$  to be a particular string; in some cases, but not all,  $s$  will prove to be a good choice of representative. When it does not, we will show that there is another string whose anticipated cost is small enough.

Let  $A$  be the set of  $m'$  strings which do not have a small repeater; there is at least one such string because we are in Case 2. For the purpose of identifying  $s$ , rename the strings in  $A$ ,  $a_1, \dots, a_{m'}$ . Let  $\sigma_i$  be the smallest repeater which is active in each of the strings  $a_i$ . Then let  $s = a_k$ , with  $k$  chosen such that  $|\sigma_k| \geq |\sigma_i|$ ,  $1 \leq i \leq m'$ . In other words,  $s$  is the string whose smallest  $(\frac{2}{3}, \frac{2}{3})$ -repeater is the largest, over all the strings in  $c$ .

By our choice of  $s$ , we know that for any other string  $t$  in  $c_j$ ,  $t$  has at least one  $(\frac{2}{3}, \frac{2}{3})$ -repeater  $\sigma'$  such that  $|\sigma'| \leq |\sigma|$ .

Our strategy will be to show that either  $s$  can be extended to include any other strings in  $c$  within our bound, or that there is some particular string  $t$  whose position and length causes the extension of  $s$  to be too costly. In the latter case we show that  $t$  can be extended within our bounds.

We will consider four cases, which depend on the the composition of the cycle  $c$ .

**Case 2A:**  $\min\{E_\ell(s, \sigma), E_r(s, \sigma)\} \leq \frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y_\sigma|$ .

In this case we can extend  $s$  either to the left if  $E_\ell(s, \sigma) \leq E_r(s, \sigma)$ , or to the right otherwise to cover the remaining strings in  $c$ . We bound  $\text{Cost}(s, \sigma)$ :

$$\begin{aligned} \text{Cost}(s, \sigma) &\leq |y_\sigma| + \min\{E_\ell(s, \sigma), E_r(s, \sigma)\} + |\sigma| \\ &\leq |y_\sigma| + \frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y_\sigma| + |\sigma| \\ &= \frac{5}{3}(d(c) + |\sigma|). \end{aligned}$$

This concludes the analysis of Case 2A. If Case 2A does not apply, then as in Figure 8 there must be a string  $t = \text{LString}(c, \sigma)$  and a string  $u = \text{RString}(c, \sigma)$ , not necessarily distinct, which extend to the left and right, respectively, too far for  $s$  to be extended within the bounds of Case 2A. In particular, let  $X_\sigma^\ell$  and  $X_\sigma^r$  be the copies of  $X_\sigma$  in which  $s$  begins and ends. Then  $t$  must extend into  $X_\sigma^\ell$ , because otherwise  $E_\ell(s, \sigma) \leq$

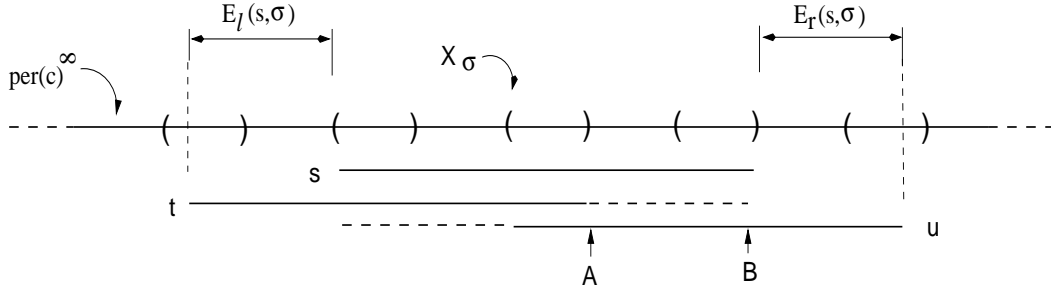


Figure 8: Case 2 of Lemma 4.12. Determining the range of possible  $t$  and  $u$ .

$2d(c) - |y_\sigma| \leq \frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y_\sigma|$ , since  $|\sigma| > \frac{1}{2}d(c)$ . We also note that  $t$  cannot extend to the left beyond  $X_\sigma^\ell$ , or we could simply shift it over  $d(c)$  to the right. Therefore the left end of  $t$  is in  $X_\sigma^\ell$ . The right end of  $t$  must also be within  $d(c)$  of the right end of  $s$ , or between points A and B marked in Figure 8. Similarly, the right end of  $u$  is in  $X_\sigma^r$ , and the left end may be anywhere within  $d(c)$  to the right of the left end of  $s$ .

Because each string in  $c$  must have at least one  $(\frac{2}{3}, \frac{2}{3})$ -repeater active, let  $\sigma'$  be the smallest  $(\frac{2}{3}, \frac{2}{3})$ -repeater active in  $t$ , and  $\sigma''$  the smallest  $(\frac{2}{3}, \frac{2}{3})$ -repeater active in  $u$ . The position of the right end of  $t$  (left end of  $u$ ) will determine whether  $X_{\sigma'}$  ( $X_{\sigma''}$ ) is nested within  $X_\sigma$  or linked with it. The remaining cases which we consider all have  $\min\{E_\ell(s, \sigma), E_r(s, \sigma)\} > \frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y_\sigma|$  and are determined by whether  $t = u$  and whether  $X_{\sigma'}$  and  $X_{\sigma''}$  are linked with or nested within  $X_\sigma$ .

In order to simplify our analysis, we will often assume that a string with an active repeater  $\sigma$  extends from the left end of one copy of  $y_\sigma$  to the right end of another copy of  $y_\sigma$ . This assumption is pessimistic in two ways; first, we may be over-charging for extension, if a string does not go as far as the right end of  $y_\sigma$  and we assume it does. Second, witnesses longer than the minimum for  $(\frac{2}{3}, \frac{2}{3})$ -repeaters give us stronger results when we apply Corollary 2.3.

**Case 2B:**  $\min\{E_\ell(s, \sigma), E_r(s, \sigma)\} > \frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y_\sigma|$ ,  $t = u$ .

We will show that  $t$  can be extended within the desired bounds. Recall that  $\sigma'$  is the smallest  $(\frac{2}{3}, \frac{2}{3})$ -repeater active in  $t$ . Observe that  $E_\ell(s, \sigma)$  and  $E_r(s, \sigma)$  span the length of a single copy of  $y_{\sigma'}$  with some overlap between two copies of  $X_\sigma$ . This observation gives rise to the following identity:

$$E_\ell(s, \sigma) + E_r(s, \sigma) = |y_{\sigma'}| + 2d(c) - |y_\sigma|. \quad (5)$$

Now consider extending  $t$  to the right. Any string  $t'$  which begins within  $d(t, s)$  of the beginning of  $t$  must end before  $s$  due to the no-substring assumption, and we will only need to extend  $t$  by  $\bar{d}(y_{\sigma'}, y_\sigma)$ , to the end of  $X_\sigma$ . (See Figure 9(a).) We will also have to consider the case where a string  $v$  begins to the right of  $s$  and extends beyond the right end of  $s$ . We call  $v$  an *interloper*. We first consider the case where there are no interlopers, then when there is an interloper on one side, and finally when there is an interloper on each side.

If there are no interlopers, then by the definition of interloper, we only have to extend  $t$  left or right to the end of string  $s$ . Therefore  $E_\ell(t, \sigma') \leq d(c) - E_\ell(s, \sigma)$  and  $E_r(t, \sigma') \leq$

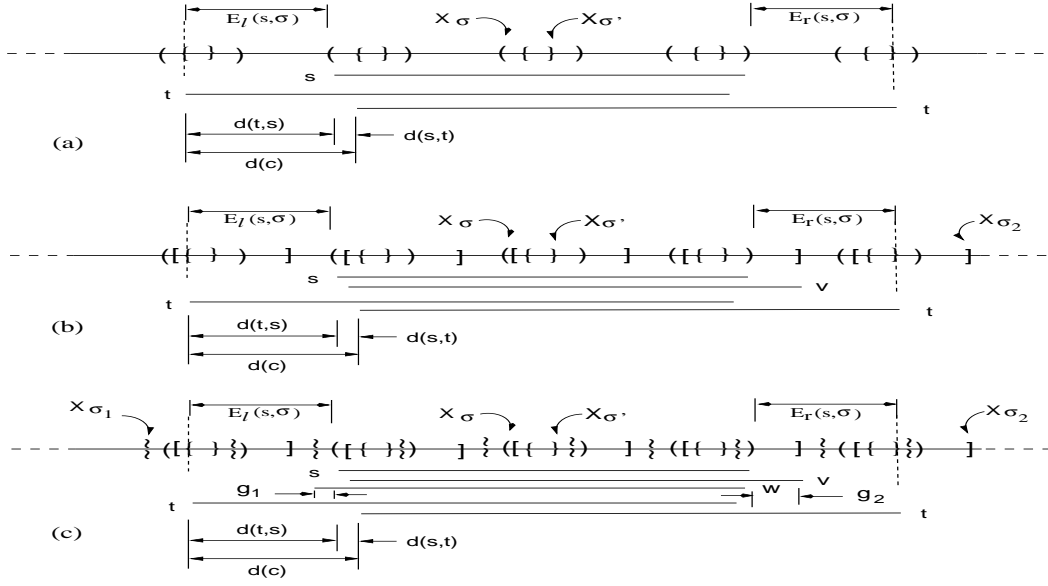


Figure 9: Case 2B of Lemma 4.12. (a) Without an interloper. (b) With one interloper  $v$ . (c) With two interlopers  $v$  and  $w$ .

$d(c) - E_r(s, \sigma)$ :

$$\begin{aligned}
\text{Cost}(t, \sigma') &= |y_{\sigma'}| + \min\{d(c) - E_\ell(s, \sigma), d(c) - E_r(s, \sigma)\} + |\sigma'| \\
&\leq |y_{\sigma'}| + d(c) - \frac{1}{2}(E_\ell(s, \sigma) + E_r(s, \sigma)) + |\sigma'| \\
&= \frac{1}{2}|y_{\sigma'}| + \frac{1}{2}|y_\sigma| + |\sigma'| && \text{(Eq. 5.)} \\
&< d(c) + \frac{1}{2}|\sigma| + \frac{3}{2}|\sigma'| && \text{(Lemma 3.4.)} \\
&< \frac{3}{2}(d(c) + |\sigma'|).
\end{aligned}$$

Suppose there is an interloper on one side. Let  $v$  be the interloper which extends the furthest to the right as in Figure 9(b). Because all strings must have an active  $(\frac{2}{3}, \frac{2}{3})$ -repeater, let  $\sigma_2$  be the smallest  $(\frac{2}{3}, \frac{2}{3})$ -repeater active in  $v$ . By our conditions on where  $v$  starts and ends,  $X_{\sigma_2}$  must be linked with  $X_\sigma$  and contain  $X_{\sigma'}$  as shown. We know by our choice of  $s$  and  $\sigma$  that  $|\sigma| > |\sigma_2|$ . By Lemma 4.5,  $|\sigma_2| > \frac{1}{2}d(c)$ , so we apply Lemma 4.7 to conclude that

$$|y_\sigma| < \frac{3}{2}d(c). \quad (6)$$

If  $v$  goes beyond  $X_\sigma$  to the right as in the Figure, we will extend  $t$  to the left. As above when there were no interlopers, we use  $E_\ell(t, \sigma') = d(c) - E_\ell(s, \sigma)$ :

$$\begin{aligned}
\text{Cost}(t, \sigma') &= |y_{\sigma'}| + d(c) - E_\ell(s, \sigma) + |\sigma'| \\
&= |y_{\sigma'}| + d(c) - (|y_{\sigma'}| + 2d(c) - |y_\sigma| - E_r(s, \sigma)) + |\sigma'| && \text{(Eq. 5.)} \\
&= |y_\sigma| + E_r(s, \sigma) - d(c) + |\sigma'| \\
&< |y_\sigma| + |\sigma'| && (E_\ell(s, \sigma) < d(c).) \\
&< \frac{3}{2}(d(c) + |\sigma'|). && \text{(Eq. 6.)}
\end{aligned}$$

Finally, suppose that there is an interloper in each direction, say  $w$  and  $v$  with  $(\frac{2}{3}, \frac{2}{3})$ -repeaters  $\sigma_1$  and  $\sigma_2$  respectively, as in Figure 9(c). Although this seems to present some difficulties, the situation also gives us stronger bounds because multiple characteristics are linked and we can employ Lemma 4.7.

Note that  $X_{\sigma_1}$  and  $X_{\sigma_2}$  are linked, as are  $X_{\sigma_1}$  and  $X_\sigma$ . Let  $g_1 = d(y_{\sigma_1}, y_\sigma)$  be the amount that  $w$  extends to the left beyond  $X_\sigma$ , and let  $g_2 = \tilde{d}(y_\sigma, y_{\sigma_2})$  be the amount that  $v$  extends to the right beyond  $X_\sigma$ . We derive a lower bound on  $|\sigma|$ :

$$\begin{aligned}
|\sigma| &> |\sigma_1| + |y_\sigma| - d(c) && \text{(Lemma 4.7.)} \\
&> |\sigma_2| + |y_{\sigma_1}| - d(c) + |y_\sigma| - d(c) && \text{(Lemma 4.7.)} \\
&> |y_{\sigma_1}| + |y_\sigma| - \frac{3}{2}d(c) && \text{(Lemma 4.5.)} \\
\Rightarrow \frac{1}{3}|\sigma| &> \frac{2}{3}|\sigma_1| - \frac{1}{6}d(c) && \text{(Def. 3.2.)} \\
\Rightarrow |\sigma_1| &< \frac{1}{4}d(c) + \frac{1}{2}|\sigma|.
\end{aligned}$$

Without loss of generality let  $|\sigma_1| > |\sigma_2|$ . We will choose to extend in the direction of the larger of  $\sigma_1$  and  $\sigma_2$ , so in this case we will extend  $t$  to the left. Since  $g_1 = d(y_{\sigma_1}, y_\sigma)$  and  $X_{\sigma_1}$  and  $X_\sigma$  are linked, we conclude that

$$g_1 < |y_{\sigma_1}| - d(c). \quad (7)$$

We use Equation 7, Lemma 4.7, and Equation 4.3 to bound  $g_1$ :

$$g_1 < |y_{\sigma_1}| - d(c) < |\sigma_1| - |\sigma_2| < \frac{1}{4}d(c) + \frac{1}{2}|\sigma| - |\sigma_2| < \frac{1}{2}|\sigma| - \frac{1}{4}d(c). \quad (8)$$

We now calculate the anticipated cost of extending  $t$  to the left (in the direction of  $\sigma_1$ , the larger of  $\sigma_1$  and  $\sigma_2$ ):

$$\begin{aligned}
\text{Cost}(t, \sigma') &\leq |y_{\sigma'}| + E_\ell(t, \sigma') + |\sigma'| \\
&\leq |y_{\sigma'}| + d(c) - E_\ell(s, \sigma) + g_1 + |\sigma'| \\
&< |y_{\sigma'}| + d(c) - (\frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y_\sigma|) + g_1 + |\sigma'| \quad \text{(Case bound.)} \\
&= |y_{\sigma'}| - \frac{11}{12}d(c) - \frac{1}{6}|\sigma| + |y_\sigma| + |\sigma'| \quad \text{(Eq. 8.)} \\
&< |y_{\sigma'}| - \frac{11}{12}d(c) - \frac{1}{6}|\sigma| + (|\sigma| - |\sigma_1| + d(c)) + |\sigma'| \quad \text{(Lemma 4.7.)} \\
&= |y_{\sigma'}| + \frac{1}{12}d(c) + \frac{5}{6}|\sigma| - |\sigma_1| + |\sigma'|.
\end{aligned}$$

In the last inequality above we were able to apply Lemma 4.7 because  $X_\sigma$  and  $X_{\sigma_1}$  are linked; now we can apply it again, because  $X_{\sigma_1}$  and  $X_{\sigma_2}$  are also linked.

$$\begin{aligned}
&< |y_{\sigma'}| + \frac{1}{12}d(c) + \frac{5}{6}|\sigma| - (|\sigma_2| + |y_{\sigma_1}| - d(c)) + |\sigma'| \quad \text{(Lemma 4.7.)} \\
&= |y_{\sigma'}| + \frac{13}{12}d(c) + \frac{5}{6}|\sigma| - |\sigma_2| - |y_{\sigma_1}| + |\sigma'| \\
&< |y_{\sigma'}| + \frac{13}{12}d(c) + \frac{5}{6}|\sigma| - |\sigma_2| - (\frac{2}{3}d(c) + \frac{2}{3}|\sigma_1|) + |\sigma'| \quad \text{(Def. 3.2.)} \\
&< |y_{\sigma'}| + \frac{5}{12}d(c) + \frac{5}{6}|\sigma| - \frac{5}{3}|\sigma_2| + |\sigma'| \quad (|\sigma_1| > |\sigma_2|. ) \\
&< \frac{9}{4}d(c) + 2|\sigma'| - \frac{5}{3}|\sigma_2| \quad \text{(Lemma 3.4.)} \\
&< \frac{17}{12}d(c) + 2|\sigma'| \\
&< \frac{29}{18}(d(c) + |\sigma'|). \quad (|\sigma'| < \frac{1}{2}d(c).)
\end{aligned}$$

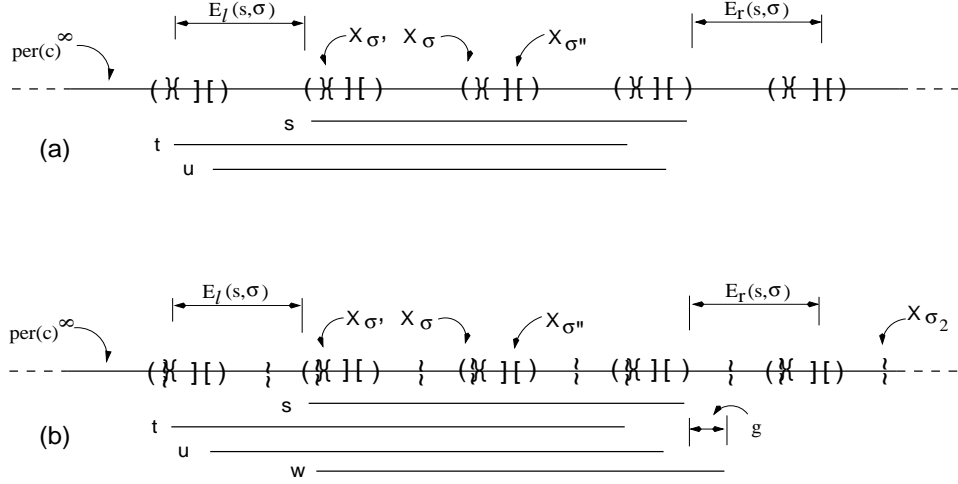


Figure 10: Case 2C of Lemma 4.12. (a) Without an interloper. (b) With an interloper  $w$ .

This concludes the analysis of Case 2B. In the remaining two cases,  $t \neq u$ ; that is,  $\text{LString}(c, \sigma) \neq \text{RString}(c, \sigma)$ . Let  $\sigma'$  be the smallest  $(\frac{2}{3}, \frac{2}{3})$ -repeater active in  $t$  and  $\sigma''$  be the smallest  $(\frac{2}{3}, \frac{2}{3})$ -repeater active in  $u$ , and without loss of generality let  $|\sigma'| > |\sigma''|$ . By our choice of  $s$  we know that  $|\sigma| > |\sigma'| > |\sigma''|$ .

If  $X_{\sigma'}$  is linked with  $X_\sigma$ , we observe that  $E_\ell(s, \sigma) = d(y', y)$  and  $E_r(s, \sigma) = \tilde{d}(y'', y)$ , so we can apply Lemma 4.10 and conclude that  $\min\{E_\ell(s, \sigma), E_r(s, \sigma)\} < \frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y_\sigma|$ . This satisfies the bound for Case 2A. We therefore only need to consider two remaining cases: when neither  $X_{\sigma'}$  nor  $X_{\sigma''}$  is linked with  $X_\sigma$  (Case 2C), and when only  $X_{\sigma''}$  is linked with  $X_\sigma$  (Case 2D).

**Case 2C:**  $\min\{E_\ell(s, \sigma), E_r(s, \sigma)\} > \frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y_\sigma|$ ,  $X_{\sigma'}$  and  $X_{\sigma''}$  both nested.

We show that  $t$  can be extended to the right within our bounds. (See Figure 10a.) Here again interlopers are possible, so we will first consider the case without an interloper, and then the case with an interloper on at least one side.

If there is no interloper, then we only have to extend  $t$  to the right as far as the end of  $X_\sigma$ . We use Lemma 4.9, the Case bound on  $E_\ell(s, \sigma)$  and  $E_r(s, \sigma)$ , and the fact that  $E_\ell(s, \sigma) + E_r(s, \sigma) = |y_{\sigma'}| + 2d(c) - |y_\sigma|$ :

$$\begin{aligned}
\text{Cost}(t, \sigma') &\leq |y_{\sigma'}| + E_\ell(t, \sigma') + |\sigma'| \\
&\leq |y_{\sigma'}| + \tilde{d}(y', y) + |\sigma'| \\
&= |y_{\sigma'}| + (|y_\sigma| - |y_{\sigma'}| - d(y', y) + |\sigma'|) \\
&= |y_\sigma| + E_\ell(s, \sigma) - d(c) + |\sigma'| \\
&= \frac{5}{3}|\sigma'| + |y_\sigma| + E_\ell(s, \sigma) - d(c) - \frac{2}{3}|\sigma'|.
\end{aligned}$$

We apply Lemma 4.9 and the fact that  $E_\ell(s, \sigma) = d(y', y)$  to bound the last term above,

$$\begin{aligned}
&\leq \frac{5}{3}|\sigma'| + |y_\sigma| + E_\ell(s, \sigma) - d(c) - \frac{2}{3}(E_\ell(s, \sigma) + |y_\sigma| - 2d(c)) \\
&= \frac{5}{3}|\sigma'| + \frac{1}{3}|y_\sigma| + \frac{1}{3}E_\ell(s, \sigma) + \frac{1}{3}d(c) \\
&< \frac{5}{3}|\sigma'| + \frac{4}{3}d(c) && \text{(Lemma 3.4.)} \\
&< \frac{5}{3}(d(c) + |\sigma'|).
\end{aligned}$$

Because  $|\sigma'| > |\sigma''|$  and  $\sigma'$  is active in  $t$ ,  $t$  was our choice of representative and we elected to extend to the right. Therefore the only interloper which concerns us is one like  $w$  in Figure 10b. Let  $\sigma_2$  be the smallest  $(\frac{2}{3}, \frac{2}{3})$ -repeater active in  $w$ . Due to our choice of  $s$ ,  $|\sigma| > |\sigma_2|$ , we can apply Lemma 4.7 to obtain:

$$|\sigma_2| < |\sigma| - |y_\sigma| + d(c). \quad (9)$$

Let  $g = \tilde{d}(y, y_{\sigma_2})$  be the distance that the interloper  $w$  extends beyond  $X_\sigma$ . We observe that

$$g < |y_{\sigma_2}| - d(c) - (E_\ell(s, \sigma) - 2d(c) - |y_\sigma|) = |y_{\sigma_2}| + d(c) - E_\ell(s, \sigma) - |y_\sigma|. \quad (10)$$

We now calculate the cost of extending  $t$  to the right:

$$\begin{aligned}
\text{Cost}(t, \sigma') &\leq |y_{\sigma'}| + E_r(t, \sigma') + |\sigma'| \\
&\leq |y_{\sigma'}| + \tilde{d}(y_{\sigma'}, y_{\sigma_2}) + g + |\sigma'| \\
&= |y_{\sigma'}| + (|y_\sigma| - (d(c) - E_\ell(s, \sigma)) - |y_{\sigma'}|) + g + |\sigma'| \\
&= |y_\sigma| - d(c) + E_\ell(s, \sigma) + g + |\sigma'| \\
&< |y_\sigma| - d(c) + E_\ell(s, \sigma) \\
&\quad + (|y_{\sigma_2}| + d(c) - E_\ell(s, \sigma) - |y_\sigma|) + |\sigma'| && \text{(Eq. 10.)} \\
&= |y_{\sigma_2}| + |\sigma'| \\
&< d(c) + |\sigma_2| + |\sigma'| && \text{(Lemma 3.4.)} \\
&< d(c) + |\sigma'| + (|\sigma| - |y_\sigma| + d(c)) && \text{(Eq. 9.)} \\
&< d(c) + |\sigma'| + \frac{1}{3}|\sigma| + \frac{1}{3}d(c) && \text{(Def. 3.2.)} \\
&< \frac{5}{3}(d(c) + |\sigma'|).
\end{aligned}$$

**Case 2D:**  $\min\{E_\ell(s, \sigma), E_r(s, \sigma)\} > \frac{5}{3}d(c) + \frac{2}{3}|\sigma| - |y_\sigma|$ ,  $X_{\sigma''}$  (but not  $X_{\sigma'}$ ) linked with  $X_\sigma$ .

In this case  $X_{\sigma'}$  might be nested within  $X_{\sigma''}$  (Figure 11 a), or not (Figure 11b). It is an unlikely case to give us trouble, because here the smaller  $(\frac{2}{3}, \frac{2}{3})$ -repeater has the larger characteristic, and it turns out that we achieve a stronger bound than in other cases.

Subcase (i). Because  $X_{\sigma''}$  contains  $X_{\sigma'}$ , Lemma 4.5 applies, so  $|\sigma''| > \frac{1}{2}d(c)$ . Since Lemma 4.7 also applies we have

$$|y_\sigma| < |\sigma| - |\sigma''| + d(c) < \frac{3}{2}d(c). \quad (11)$$

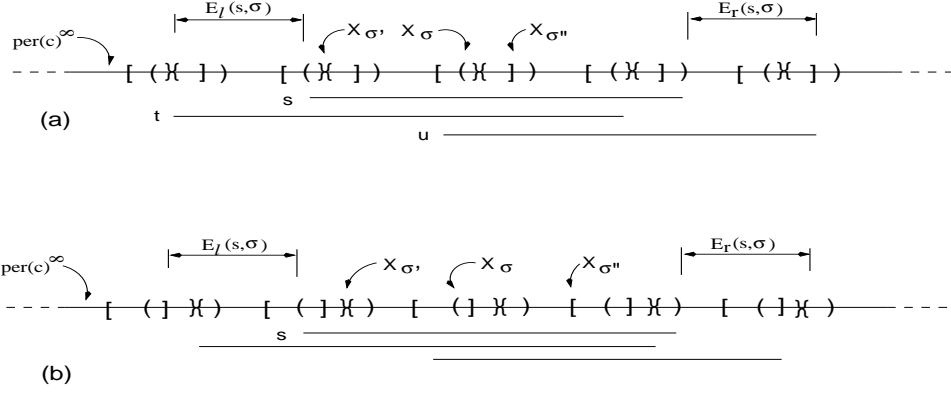


Figure 11: Case 2D of Lemma 4.12. (a)  $X_{\sigma'}$  nested within  $X_{\sigma''}$ . (b)  $X_{\sigma'}$  not nested within  $X_{\sigma''}$ .

If there are no interlopers, we can now bound the anticipated cost of extending  $t$  to the right as follows:

$$\begin{aligned}
\text{Cost}(t, \sigma') &\leq |y_{\sigma'}| + E_\ell(t, \sigma') + |\sigma'| \\
&\leq |y_{\sigma'}| + \tilde{d}(y_{\sigma'}, y_\sigma) + |\sigma'| \\
&= |y_{\sigma'}| + (|y_\sigma| - |y_{\sigma'}| - (d(c) - E_\ell(s, \sigma))) + |\sigma'| \\
&< |y_\sigma| + |\sigma'| && (E_\ell(s, \sigma) < d(c)) \\
&< \frac{3}{2}(d(c) + |\sigma'|). && (\text{Eq. 11})
\end{aligned}$$

Now suppose there was an interloper  $v$  with smallest  $(\frac{2}{3}, \frac{2}{3})$ -repeater  $\sigma_2$ . Then  $X_{\sigma_2}$  would be linked with  $X_{\sigma''}$  and  $X_\sigma$ , and Lemma 4.10 would apply as in Figure 6(d), and we would once again be in Case 2A.

Subcase (ii). Now  $X_{\sigma'}$  is not nested within  $X_{\sigma''}$ , as in Figure 11b. If there are no interlopers, then we only have to extend  $t$  to the right to the end of  $X_\sigma$ :

$$\begin{aligned}
\text{Cost}(t, \sigma') &\leq |y_{\sigma'}| + E_\ell(t, \sigma') + |\sigma'| \\
&\leq |y_{\sigma'}| + E_\ell(s, \sigma) + d(c) - |y_{\sigma'}| - (2d(c) - |y_\sigma|) + |\sigma'| \\
&= |y_\sigma| + E_\ell(s, \sigma) - d(c) + |\sigma'|.
\end{aligned}$$

We apply Lemma 4.9 to complete the analysis:

$$\begin{aligned}
\text{Cost}(t, \sigma') &< |y_\sigma| + (|\sigma'| - |y_\sigma| + 2d(c)) - d(c) + |\sigma'| \\
&= 2|\sigma'| + d(c) \\
&< \frac{3}{2}(d(c) + |\sigma'|).
\end{aligned}$$

As in Case (i), if there is an interloper then Lemma 4.10 will apply (Figures 6c or d), and we have Case 2A.

This completes the proof of Case IId, which completes the proof of the lemma.  $\blacksquare$

We now combine Lemmas 2.2, 4.11, and 4.12 to obtain:

**Theorem 4.14** *Algorithm G-SHORTSTRING( $S$ ) is a  $2\frac{2}{3}$ -approximation for the shortest superstring problem.*

## References

- [1] C. Armen and C. Stein. Improved length bounds for the shortest superstring problem. In *Proceedings of Workshop on Algorithms and Data Structures*, pages 494–505, 1995.
- [2] C. Armen and C. Stein. Short superstrings and the structure of overlapping strings. To appear in *J. of Computational Biology*, 1995.
- [3] Chris Armen. *Approximation Algorithms for the Shortest Superstring Problem*. PhD thesis, Dartmouth College, July 1995.
- [4] A. Blum, T. Jiang, M. Li, J. Tromp, and M. Yannakakis. Linear approximation of shortest superstrings. In *Proceedings of the 23rd Annual ACM Symposium on Theory of Computing*, pages 328–336, 1991.
- [5] Thomas H. Cormen, Charles E. Leiserson, and Ronald L. Rivest. *Introduction to Algorithms*. MIT Press/McGraw-Hill, 1990.
- [6] A. Czumaj, L. Gasieniec, M. Piotrow, and W. Rytter. Parallel and sequential approximations of shortest superstrings. In *Proceedings of Fourth Scandinavian Workshop on Algorithm Theory*, pages 95–106, 1994.
- [7] A. Lesk (edited). *Computational Molecular Biology, Sources and Methods for Sequence Analysis*. Oxford University Press, 1988.
- [8] N. Fine and H. Wilf. Uniqueness theorems for periodic functions. *Proceedings of the American Mathematical Society*, 16:109–114, 1965.
- [9] A.M. Frieze, G. Galbiati, and F. Maffoli. On the worst case performance of some algorithms for the asymmetric travelling salesman problem. *Networks*, 12:23–39, 1982.
- [10] J. Gallant, D. Maier, and J. Storer. On finding minimal length superstrings. *Journal of Computer and System Sciences*, 20:50–58, 1980.
- [11] D. Gusfield, G. Landau, and B. Schieber. An efficient algorithm for the all pairs suffix-prefix problem. *Information Processing Letters*, (41):181–185, March 1992.
- [12] Tao Jiang and Ming Li. Approximating shortest superstrings with constraints. *Theoretical Computer Science*, (134):473–491, 1994.
- [13] J.D. Kececioğlu and E.W. Myers. Combinatorial algorithms for dna sequence assembly. *Algorithmica*, 13(1/2):7–51, 1995.
- [14] John D. Kececioğlu. *Exact and approximation algorithms for DNA sequence reconstruction*. PhD thesis, University of Arizona, 1991.
- [15] R. Kosaraju, J. Park, and C. Stein. Long tours and short superstrings. In *Proceedings of the 35th Annual Symposium on Foundations of Computer Science*, November 1994.
- [16] M. Li. Towards a DNA sequencing theory (learning a string). In *Proceedings of the 31st Annual Symposium on Foundations of Computer Science*, pages 125–134, 1990.
- [17] Christos H. Papadimitriou and Kenneth Steiglitz. *Combinatorial Optimization, Algorithms and Complexity*. Prentice-Hall, Englewood Cliffs, NJ, 1982.
- [18] H. Peltola, H. Soderlund, J. Tarjio, and E. Ukkonen. Algorithms for some string matching problems arising in molecular genetics. In *Proceedings of the IFIP Congress*, pages 53–64, 1983.
- [19] Graham A. Stephen. *String searching algorithms*. World Scientific, 1994.



- [20] J. Storer. *Data compression: methods and theory*. Computer Science Press, 1988.
- [21] J. Tarhio and E. Ukkonen. A greedy approximation algorithm for constructing shortest common superstrings. *Theoretical Computer Science*, 57:131–145, 1988.
- [22] Shang-Hua Teng and Frances Yao. Approximating shortest superstrings. In *Proceedings of the 34th Annual Symposium on Foundations of Computer Science*, pages 158–165, November 1993.
- [23] J. Turner. Approximation algorithms for the shortest common superstring problem. *Information and Computation*, 83:1–20, 1989.