

# The Complexity of the Shortest Common Matching String Problem

Jonathan S. Turner \*

April 30, 1986  
WUCS-86-9

## Abstract

This paper describes the *shortest common matching string problem*, which arises from a data analysis problem in molecular genetics, and shows that it is NP-complete.

---

\*The research described here was supported in part by the National Science Foundation (grant #DCR-8409435) and the National Institutes of Health (grant #RR-01380)

# The Complexity of the Shortest Common Matching String Problem

Jonathan S. Turner

Let  $s_1 = a_1 \dots a_r$  and  $s_2 = b_1 \dots b_s$  be strings over some finite alphabet  $\Sigma$ . We say that  $s_1$  is a substring of  $s_2$  if there is an integer  $i \in [0, s - r]$  such that  $a_j = b_{i+j}$  for  $1 \leq j \leq r$ . We also say in this case that  $s_2$  is a *superstring* of  $s_1$ .

A *bag*  $b = \langle a_1, \dots, a_r \rangle$  is an unordered collection of symbols from some alphabet  $\Sigma$  in which the same symbol may appear more than once. (A bag is often referred to as a *multi-set*.) If  $s = b_1 \dots b_s$  is a string we define  $\langle s \rangle$  to be the bag  $\langle b_1, \dots, b_s \rangle$ . We say that a bag  $b$  *matches* a string  $s$  if  $s$  contains some substring  $s'$  such that  $\langle s' \rangle = b$ . We also say that  $s$  matches  $b$  or that  $s$  is a *matching string* of  $b$ . For example, the string `debcabf` is a matching string of the bag  $\langle a, b, b, c \rangle$ .

An instance of the *shortest common matching string problem* (SCMS) is a set of bags  $B = \{b_1, \dots, b_n\}$  over a finite alphabet  $\Sigma$  and an integer  $m$ . The object of the problem is to determine if there is a string of length  $\leq m$  that matches every bag in  $B$ . Alternatively, we can view the object as being to find a minimum length string that matches every bag in  $B$ . We let  $\chi^*(B)$  denote the length of a minimum length matching string for  $B$ .

EXAMPLE. If  $B = \{\langle \text{aceghi} \rangle, \langle \text{abfgik} \rangle, \langle \text{adfhki} \rangle, \langle \text{defghi} \rangle, \langle \text{afghik} \rangle\}$ , the string `bfgiakhfdegiach` is a minimum length solution.

This problem has applications to molecular genetics. In particular, it arises in the analysis of experimental data used to map restriction enzyme sites in DNA from complex organisms. This connection is explained fully in [4]. The problem does not appear to have been studied previously, although a related problem, the *shortest common superstring problem* (SCS) has been [1,2,3].

The purpose of this paper is to introduce the shortest common matching string problem and prove that it is NP-complete. The transformation is from the shortest common superstring problem [1]. To

simplify the presentation, we introduce an intermediate problem and use a two step transformation from SCMS to SCS.

Let  $s = a_1 \dots a_r$  be a string. The notation  $s[i]$  denotes the symbol  $a_i$  if  $i > 0$  and  $a_{r+i+1}$  if  $i < 0$ . The notation  $s[i, j]$  denotes the string  $s[i] \dots s[j]$ .

If  $s = a_1 \dots a_r$  is a string, we let  $rev(s) = a_r \dots a_1$ .

The number of symbols in a string  $s$  is denoted  $|s|$  and for any set of strings  $S$ ,  $\|S\| = \sum_{s \in S} |s|$ .

An instance of the *shortest common superstring problem* is a set of strings  $S = \{s_1, \dots, s_n\}$  over a finite alphabet  $\Sigma$  and an integer  $m$ . The object of the problem is to determine if there is a string of length  $\leq m$  that is a superstring of every  $s_i \in S$ .

EXAMPLE. If  $S = \{\text{egiach, bfgiak, hfdegi, iakhfd, fgiakh}\}$ , the string  $\text{bfgiakhfdegiach}$  is a minimum length solution.

The NP-completeness of SCS is proved in reference [1].

An instance of the *reversible shortest common superstring problem* is also a set of strings  $S = \{s_1, \dots, s_n\}$  over a finite alphabet  $\Sigma$  and an integer  $m$ . The object in this case, is to determine if there a string of length  $\leq m$  that for all  $s \in S$  contains either  $s$  or  $rev(s)$ .

EXAMPLE. If  $S = \{\text{hcaige, kaigfb, igedfh, iakhfd, fgiakh}\}$ , the string  $\text{bfgiakhfdegiach}$  is a minimum length solution.

**THEOREM 1.** *RSCS is NP-complete.*

*Proof.* Clearly  $RSCS \in NP$  since a nondeterministic Turing machine can guess a string of length  $\leq m$  and check in polynomial time that it is a superstring of either  $s$  or  $rev(s)$  for all  $s \in S$ .

We now show how to transform an instance  $(S = \{s_1, \dots, s_n\}, \Sigma, m)$  of SCS to an instance  $(S' = \{s'_1, \dots, s'_n\}, \Sigma', m')$  of RSCS. We assume without loss of generality that no string in  $S$  is a substring of another.

First, define  $\Sigma' = \Sigma \cup \{0, 1\}$  where 0 and 1 are not in  $\Sigma$ . For any string  $s = a_1 \dots a_r$ , define  $f(s) = 0a_110a_21 \dots 0a_r10$ . We now define  $s'_i = f(s_i)$  and let  $m' = 3m + 1$ .

For example, if  $\Sigma = \{a, b, c, d\}$ ,  $S = \{\text{dccbda, bacbad, bdaabc, cbadcc}\}$  and  $m = 14$  then  $\Sigma' = \{a, b, c, d, 0, 1\}$ ,  $m' = 43$  and

$$S' = \{ \text{0d10c10c10b10d10a10, 0b10a10c10b10a10d10,} \\ \text{0b10d10a10a10b10c10, 0c10b10a10d10c10c10} \}$$

The original problem has the string `bacbadccbdaabc` as a solution. The corresponding solution to the transformed problem is

`Ob10a10c10b10a10d10c10c10b10d10a10a10b10c10`

We claim that in general  $S$  has a solution of size  $\leq m$  if and only if  $S'$  has a solution of size  $\leq m'$ . First, assume that  $\sigma$  is a superstring of all  $s_i \in S$  and that  $|\sigma| \leq m$ . Renumber the  $s_i$  in order of their first appearance in  $\sigma$  and let  $\pi_i$  be the smallest  $j$  such that  $s_i = \sigma[j, j + |s_i| - 1]$ . We will assume without loss of generality that  $\pi_1 = 1$ ,  $\pi_{i+1} \leq \pi_i + |s_i|$  for  $1 \leq i \leq n-1$  and  $\pi_n = |\sigma| - |s_n| + 1$ . Now, let  $\psi_i = \pi_i + |s_i| - \pi_{i+1}$  for  $1 \leq i \leq n-1$  be the amount of overlap between consecutive strings in  $\sigma$  and note that for  $1 \leq i \leq n-1$ ,  $s'_i[-(3\psi_i + 1), -1] = s'_{i+1}[1, (3\psi_i + 1)]$ . Hence, the string

$$\sigma' = s'_1 s'_2 [3\psi_1 + 1, -1] s'_3 [3\psi_2 + 1, -1] \cdots s'_n [3\psi_{n-1} + 1, -1]$$

is a superstring of all strings in  $S'$  and  $|\sigma'| = 3|\sigma| + 1 \leq 3m + 1 = m'$ . Hence, if  $S$  has a solution of size  $\leq m$ ,  $S'$  has a solution of size  $\leq m'$ .

We now show that if  $S'$  has a solution of size  $\leq m'$ , then  $S$  must have a solution of size  $\leq m$ . Let  $\sigma'$  be any string which for all  $s'_i \in S'$  is a superstring of either  $s'_i$  or  $\text{rev}(s'_i)$  and let  $|\sigma'| \leq m'$ . Renumber the  $s_i$  in order of the first appearance of either  $s_i$  or  $\text{rev}(s_i)$  in  $\sigma'$  and let  $\pi'_i$  be the smallest  $j$  such that either  $s'_i = \sigma'[j, j + |s'_i| - 1]$  or  $\text{rev}(s'_i) = \sigma'[j, j + |s'_i| - 1]$ . We will assume without loss of generality that  $\pi'_1 = 1$ ,  $\pi'_{i+1} \leq \pi'_i + |s'_i|$  for  $1 \leq i \leq n-1$  and  $\pi'_n = |\sigma'| - |s'_n| + 1$ . Now, let  $\psi'_i = \pi'_i + |s'_i| - \pi'_{i+1}$  for  $1 \leq i \leq n-1$  be the amount of overlap between consecutive strings. Note that if  $\psi'_i > 1$  then either  $\sigma'[\pi'_i + 2, \pi'_i + 3] = \sigma'[\pi'_{i+1} + 2, \pi'_{i+1} + 3] = 10$  or  $\sigma'[\pi'_i, \pi'_i + 1] = \sigma'[\pi'_{i+1}, \pi'_{i+1} + 1] = 01$ . That is, either both  $s'_i$  and  $s'_{i+1}$  are reversed in  $\sigma'$  or neither one is. Consequently, there is a string  $\sigma''$  of the same length as  $\sigma'$  which is a superstring of all the strings in  $S'$  (that is, none of the strings is reversed in  $\sigma''$ ). We will assume therefore, that each  $s' \in S$  is a substring of  $\sigma'$ . Hence

$$\sigma = s_1 s_2 [((\psi'_1 - 1)/3) + 1, -1] s_3 [((\psi'_2 - 1)/3) + 1, -1] \cdots s_n [((\psi'_{n-1} - 1)/3) + 1, -1]$$

is a superstring of all strings in  $S$  and  $|\sigma| = (|\sigma'| - 1)/3 \leq (m' - 1)/3 = m$ . To complete the proof, we note that  $(S', \Sigma', m')$  can be computed deterministically in time polynomial in  $\|S\|$ .  $\square$

*Remark.* In [1], it is shown that SCS is NP-complete, even when the alphabet is limited to two symbols. Since the proof of Theorem 1 adds just two symbols to the alphabet, it follows that RSCS is NP-complete when the alphabet is limited to four symbols.

**THEOREM 2.** *SCMS is NP-complete.*

*Proof.* Clearly SCMS is in NP, since a nondeterministic Turing machine can guess a string of length  $\leq m$  and check in polynomial time that it matches each of the bags.

We now show how to transform an instance  $(S = \{s_1, \dots, s_n\}, \Sigma, m)$  of RSCS to an instance  $(B = \{b_1, \dots, b_p\}, \Sigma', m')$  of SCMS, where  $\Sigma' = \Sigma$  together with the new symbols  $\{L, R, x_1, \dots, x_n\}$ ,  $m' = n(r+2) + m + \|S\|$  and  $r = 1 + 4\|S\|$ .  $B = B_1 \cup \dots \cup B_n$  where

$$B_i = \{ \langle x_i^r \rangle, \langle Lx_i^r \rangle, \langle x_i^r R \rangle \} \cup \{ \langle x_i^r R s_i[1, j] \rangle \mid 1 \leq j \leq |s_i| \} \\ \cup \{ \langle s_i[j, -1] Lx_i^r \rangle \mid 1 \leq j \leq |s_i| \}$$

For example, if  $\Sigma = \{a, b, c, d\}$ ,  $S = \{bcd, dc, abc\}$  and  $m = 7$  then  $\Sigma' = \{a, b, c, d, L, R, x_1, x_2, x_3\}$ ,  $m' = 172$  and  $B = B_1 \cup B_2 \cup B_3$  where

$$B_1 = \{ \langle x_1^{40} \rangle, \langle Lx_1^{40} \rangle, \langle x_1^{40} R \rangle, \langle x_1^{40} R b \rangle, \langle x_1^{40} R bc \rangle, \langle x_1^{40} R bcd \rangle, \langle x_1^{40} R bcd b \rangle, \\ \langle b Lx_1^{40} \rangle, \langle db Lx_1^{40} \rangle, \langle cdb Lx_1^{40} \rangle, \langle bcd b Lx_1^{40} \rangle \}$$

The sets  $B_2$  and  $B_3$  are similar. The original problem has the string  $abcdbcd$  as a solution. The corresponding solution to the transformed problem is

$$abc b Lx_3^{40} R a b c b c d R x_2^{40} L c b c d b Lx_1^{40} R b c d b$$

Note that its length is 172.

We claim that in general  $S$  has a solution string of length  $\leq m$  if and only if  $B$  has a solution string of length  $\leq m'$ . First, assume that  $\sigma$  is a superstring of either  $s_i$  or  $rev(s_i)$  for all  $s_i \in S$  and that  $|\sigma| \leq m$ . Renumber the  $s_i$  in order of the first appearance of either  $s_i$  or  $rev(s_i)$  in  $\sigma$  and let  $\pi_i$  be the smallest  $j$  such that either  $s_i = \sigma[j, j + |s_i| - 1]$  or  $rev(s_i) = \sigma[j, j + |s_i| - 1]$ . We will assume without loss of generality that  $\pi_1 = 1$ ,  $\pi_{i+1} \leq \pi_i + |s_i|$  for  $1 \leq i \leq n-1$  and  $\pi_n = |\sigma| - |s_n| + 1$ .

Define

$$s'_i = \begin{cases} s_i Lx_i^r R s_i & \text{if } \sigma[\pi_i, \pi_i + |s_i| - 1] = s_i \\ rev(s_i) R x_i^r L rev(s_i) & \text{if } \sigma[\pi_i, \pi_i + |s_i| - 1] = rev(s_i) \end{cases}$$

for  $1 \leq i \leq n$  and note that  $s'_i$  is a matching string for all the bags in  $B_i$ . Now, let  $\psi_i = \pi_i + |s_i| - \pi_{i+1}$  for  $1 \leq i \leq n-1$  be the amount of overlap between consecutive strings in  $\sigma$  and note that the string

$$\sigma' = s'_1 s'_2 [\psi_1 + 1, -1] s'_3 [\psi_2 + 1, -1] \cdots s'_n [\psi_{n-1} + 1, -1]$$

is a matching string for all the bags in  $B$  and

$$|\sigma'| = n(r+2) + 2\|S\| - \sum_{i=1}^{n-1} \psi_i \leq n(r+2) + 2\|S\| - (\|S\| - m) = n(r+2) + m + \|S\| = m'$$

Hence, if  $S$  has a solution of length  $\leq m$ , then  $B$  has a solution of length  $\leq m'$ .

We now show that if  $B$  has a solution of length  $\leq m'$  then  $S$  has a solution of length  $\leq m$ . Let  $\sigma'$  be a shortest matching string for  $B$  and assume that  $|\sigma'| \leq m'$ . Let  $\pi'_i$  be the smallest  $j$  such that  $b_i = \langle \sigma'[j, j + |b_i| - 1] \rangle$ . Define

$$\psi'_{ij} = |\{\pi'_i, \dots, \pi'_i + |b_i| - 1\} \cap \{\pi'_j, \dots, \pi'_j + |b_j| - 1\}|$$

That is,  $\psi'_{ij}$  is the amount of overlap between bags  $b_i$  and  $b_j$  in  $\sigma'$ .

Now, note that  $nr \leq |\sigma'| \leq m' < (n+1)r$ . Consequently for any  $h \in [1, n]$  if  $b_i = \langle x_h^r \rangle$  and  $b_j = \langle Lx_h^r \rangle$  then  $\psi'_{ij} \geq 1$ . If  $\pi'_j < \pi'_i - 1$  then the string  $\sigma'[\pi'_j, \pi'_i]$  has the form  $x_h^r Lx_h^r$ , where  $0 \leq s \leq \pi'_i - \pi'_j - 1$  and  $s + t = \pi'_i - \pi'_j$ . Hence,  $\sigma'[1, \pi'_j - 1]L\sigma'[\pi'_i, |\sigma'|]$  is also a matching string of  $B$  and is shorter than  $\sigma'$ . Since we assumed that  $\sigma'$  was a shortest matching string for  $B$  it follows that  $\pi'_j \geq \pi'_i - 1$ . Similarly, we can show that  $\pi'_j \leq \pi'_i$  and consequently  $\psi'_{ij} = r$ . This argument can be extended to show that  $\psi'_{ij} = r$  for any string  $b_j \in B_h$  and  $b_i = \langle x_h^r \rangle$ .

The above observations imply that for all  $h \in [1, n]$ ,  $\sigma'$  contains a string  $s'_h$  of the form  $s_h Lx_h^r R s_h$  or  $rev(s_h) R x_h^r L rev(s_h)$ . Renumber the  $s'_h$  in order of their first appearance in  $\sigma'$  and note that for all  $i \in [1, n-1]$  the overlapping portion of  $s'_i$  and  $s'_{i+1}$  is also a valid overlap for  $s_i$  and  $s_{i+1}$ . Now, redefine  $\pi'_i$  to be the smallest  $j$  such that  $\sigma'[j, j + |s'_i| - 1] \in \{s'_i, rev(s'_i)\}$  for  $1 \leq i \leq n$ , and redefine  $\psi'_i = \pi'_i + |s'_i| - \pi'_{i+1}$  for  $1 \leq i \leq n-1$  and note that the string

$$\sigma = s_1 s_2 [\psi'_1 + 1, -1] s_3 [\psi'_2 + 1, -1] \dots s_n [\psi'_{n-1} + 1, -1]$$

is a solution to the original RSCS instance and that

$$|\sigma| = \|S\| - \sum_{i=1}^{n-1} \psi'_i \leq \|S\| - \left[ \sum_{i=1}^n |s'_i| - m' \right] = m' - n(r+2) - \|S\| = m$$

Hence, whenever  $B$  has a solution of length  $\leq m'$ ,  $S$  has a solution of length  $\leq m$ . To complete the proof, we note that  $(B, \Sigma', m')$  can be computed deterministically in time polynomial in  $\|S\|$ .  $\square$

The NP-completeness of SCMS makes it unlikely that there exists an efficient algorithm to solve it exactly. In a separate paper [5] we address the issue of good approximation algorithms.

## References

- [1] Gallant, John K. "On finding Minimal Length Superstrings," *Journal of Computer and System Sciences*, vol. 20, no. 1, 50-58, 2/80.

- 
- [2] Gallant, John K. "String Compression Algorithms," Ph.D. Dissertation, Princeton University, Department of Electrical Engineering and Computer Science, June 1982.
  - [3] Maier, David and James A. Storer. "A Note on the Complexity of the Superstring Problem," Princeton University Technical Report 233, Department of Electrical Engineering and Computer Science, October 1977.
  - [4] Turner, Jonathan S. "The DNA Mapping Problem," Washington University, Computer Science Department technical report, in preparation.
  - [5] Turner, Jonathan S. "Algorithms for the Shortest Common Matching String Problem," Washington University, Computer Science Department technical report, in preparation.