

On the complexity of the DNA Simplified Partial Digest Problem

Jacek Blazewicz

Marta Kasprzak

Institute of Computing Science, Poznan University of Technology, Poznan,
and Institute of Bioorganic Chemistry, Polish Academy of Sciences, Poznan, Poland.
Emails: jblazewicz@cs.put.poznan.pl, marta@cs.put.poznan.pl

Abstract

The problem to be addressed is one of the genome mapping of DNA molecules. The new approach — the Simplified Partial Digest Problem (SPDP), is analyzed. This approach is easy in laboratory implementation and robust with respect to measurement errors. In the paper, it is formulated in terms of a combinatorial search problem and proved to be strongly NP-hard for the general error-free case. For a subproblem of the SPDP, a simple $O(n \log n)$ -time algorithm is given, where n is a number of restriction sites.

Keywords: combinatorial optimization, DNA restriction mapping, partial digest, computational complexity.

1 Introduction

A construction of a DNA physical map is an important step in the genome sequencing process (cf. Waterman (1995), Setubal & Meidanis (1997), Pevzner (2000), and Fogel & Corne (2003) as excellent sources of algorithmic ideas used in computational biology). This map of a DNA molecule contains the information about locations of short, specific subsequences called markers and in turn places longer DNA subchains on the chromosome. One approach to create the map is based upon splitting the molecule into many shorter ones and *hybridizing* them with the *markers* (known very short DNA sequences). This approach results in the *interval graph model* used successfully in the past (Benzer 1959).

A more recent approach to the genome map construction relies on a digestion of a DNA molecule with restriction enzymes (Waterman 1995). These enzymes cut DNA molecules within specific, short patterns of nucleotides called restriction sites. After the digestion, the lengths of obtained fragments are measured and the original ordering of these fragments must be reconstructed, and this is the place where combinatorial optimization methods come to the effect. In practice, several variants of this approach are used. Two of the best known are *the double digest* and *the partial digest* (with its many variants).

The research has been supported by KBN grant 3T11F00227. The research of the first author was sponsored by the ESPRC Fellowship at the University of Nottingham.

Corresponding author: Marta Kasprzak, Institute of Computing Science, Poznan University of Technology, Piotrowo 3A, 60-965 Poznan, Poland.

Copyright ©2006, Australian Computer Society, Inc. This paper appeared at Computing: The Australasian Theory Symposium (CATS2006), Hobart, Australia. Conferences in Research and Practice in Information Technology, Vol. 51. Barry Jay and Joachim Gudmundsson, Ed. Reproduction for academic, not-for profit purposes permitted provided this text is included.

In the double digest approach two restriction enzymes are used. A target DNA is amplified, e.g. using a PCR reaction, and the copies are divided into three sets. Molecules from the first set are digested by one enzyme, molecules from the second set are digested by the other enzyme and molecules from the third set are cut by both enzymes. All digestions are complete for the time span of each reaction is long enough to allow the enzyme to cut the target strand at each occurrence of the restriction site. As the result one obtains three collections of short DNA fragments that correspond to three digestion processes. The lengths of these fragments are measured during a gel electrophoresis process and recorded as three multisets. On the basis of this data locations of restriction sites in the target DNA are reconstructed. Unfortunately, from the combinatorial point of view the *Double Digest Problem* (DDP) is NP-hard even in an ideal case involving no errors (Goldstein & Waterman 1987), thus unlikely to admit a polynomial time algorithm. Another difficulty which must be taken into account when dealing with this approach is an exponential number of possible solutions (as measured with respect to a number of restriction sites — n) (Schmitt & Waterman 1991). As an alternative *the partial digest approach* has been proposed (Skiena, Smith & Lemke 1990, Skiena & Sundaram 1994). Here, one enzyme only is applied to cut the DNA molecule into fragments of different lengths by using the enzyme for different time periods. In the error-free case one gets here all $\binom{n}{2}$ fragment lengths between all pairs of cuts (including two ends of the DNA molecule). The resulting *Partial Digest Problem* (PDP) consists in reconstructing the original positions of the cuts. This problem is known also in discrete geometry (Skiena, Smith & Lemke 1990, Skiena 1997), where having all interpoint distances, one reconstructs the positions of a set of points on a line. For PDP two backtracking algorithms with an exponential worst case complexity have been proposed by Skiena and co-workers (Skiena, Smith & Lemke 1990, Skiena & Sundaram 1994). In general, the complexity of the PDP remains an open question, although measurement errors and noisy data result in the strong NP-hardness of the problem ((Cieliebak & Eidenbenz 2004) and (Cieliebak, Eidenbenz & Penna 2003), respectively). It is worth stressing that the number of possible solutions in the PDP is bounded from above by a polynomial in the number of restriction sites (Skiena, Smith & Lemke 1990). Other known approaches in that area include *optical mapping* (Schwartz et al. 1993, Karp & Shamir 1998), *probed partial digest mapping* (Newberg & Naor 1993), and *labeled partial digest* (Pandurangan & Ramesh 2002). Let us especially comment on the labeled partial digest approach, which although more complicated from the implementation point of view (labeling the ends of a DNA molecule by using radioactive labeling), results in a

polynomial algorithm for the related variant of the PDP and produces a unique solution (Pandurangan & Ramesh 2002).

In this paper, we study yet another variant of the PDP, called the *Simplified Partial Digest Problem* (SPDP), which is very simple for the laboratory implementation (Blazewicz et al. 2001). In gel approach only two digestions are performed. We will call them, respectively, a short digestion and a complete (long) digestion. After amplification, the copies of a target strand are splitted into two sets. The goal of a short digestion is to have all molecules from one of the sets cut in at most one occurrence of the restriction site. This is assured by properly chosen time span of the reaction. Molecules from the other set are cut in all occurrences of the restriction site due to the long reaction time span (a complete digestion). Then, as in other methods, the lengths of restriction fragments obtained, are measured during a gel electrophoresis process. Although, the algorithms proposed for the SPDP (Blazewicz et al. 2001, Blazewicz & Jaroszewski 2003) are very efficient on the average and robust with respect to measurement errors, the complexity of the problem in the error-free case was open for several years.

In this paper, we present the proof that the general error-free variant of the SPDP is strongly NP-hard (of course in its *search version*). On the other hand, if the complete digestion yields a multiset composed of 1s and 2s only, the problem is solvable in $O(n \log n)$ time, where n is a number of digestion sites. An organization of the paper is as follows. Section 2 contains the formal definition of the SPDP and the proof of its strong NP-hardness. Section 3 presents a polynomial time algorithm for the restricted case of the problem. Conclusions in Section 4 summarize the work and indicate further research problems.

2 Computational complexity of the SPDP

In this section we prove that the combinatorial search version of the Simplified Partial Digest Problem without errors is strongly NP-hard. It is worth stressing that we do not do it in a standard way by a transformation from the decision version of a known intractable problem to the decision version of our problem, since the decision counterpart of the SPDP is polynomially solvable. This is because — on the assumption that the instance of our problem has no errors — the restriction fragments from both digestion reactions exactly match those existing in the real restriction map, therefore the solution always exists. However, the process of finding this solution (i.e. the map) is not easy from the computational complexity point of view.

A similar situation — where the decision version of a problem is easy but its search version is hard — is encountered in the DNA Sequencing problem (Blazewicz & Kasprzak 2003). In this problem, as well as in SPDP, the additional information about instances — the lack of errors or some type of errors — results in the answer “yes” for every instance of the decision version. Therefore, the process of proving strong NP-hardness of the search problem must go through an artificial decision problem, allowing also instances with the answer “no”.

We based our reasoning on a deduction formed in Johnson (1985) for the Hamiltonian Circuit Problem (being NP-complete in its decision version). As it has been stated there, even if we knew a graph contains a Hamiltonian circuit (in this case the decision version is, of course, trivially easy), we could not find it in polynomial time unless $P=NP$. For “if we had such an algorithm A , we could use it to tell in polynomial

time whether an arbitrary graph G has a Hamiltonian circuit. Let p be the polynomial that bounds A ’s running time on graphs with Hamiltonian circuits. Apply A to G . If G has a Hamiltonian circuit, A will find one in time $p(|G|)$. If G does not have such a circuit, then after $p(|G|)$ steps A could not have found one, and we will know that none exists”.

2.1 Formulations of related problems

The innovation of the well-known Partial Digest Problem, studied here, was based on a simplified biochemical experiment and proposed for the first time by Blazewicz et al. (2001). Its search version without any error in instances can be formulated as below.

Problem 1 *Simplified Partial Digest Problem* (Π_{SPDP}) — *search version*.

Instance: Multiset $A = \{a_1, a_2, \dots, a_{2n}\}$ of lengths of restriction fragments coming from short digestion and multiset $B = \{b_1, b_2, \dots, b_{n+1}\}$ of lengths of restriction fragments coming from long digestion, A and B containing no errors.

Answer: A map of n restriction sites of a DNA chain consistent with multisets A and B , i.e. such an order of elements of B which allows to cover the indicated cuts by some order of elements of A .

(See Appendix for original formulation from Blazewicz et al. (2001).) We see that $D_{\Pi_{\text{SPDP}}} = Y_{\Pi_{\text{SPDP}}}$, where $D_{\Pi_{\text{SPDP}}}$ is the set of all instances of Π_{SPDP} and $Y_{\Pi_{\text{SPDP}}}$ is the set of all instances of Π_{SPDP} for which such a solution exists. (This follows from the fact, that both multisets come from a real error-free digestion experiment.) Example 1 shows an instance of this problem together with a possible solution.

Example 1 Let the output of the digestion reactions performed without errors be: $A = \{3, 6, 8, 9, 10, 11, 13, 16\}$ from the short reaction and $B = \{2, 3, 3, 5, 6\}$ from the long one. The feasible solution for the problem is presented in Figure 1 and can be written as the following ordered list of B : $[3, 5, 2, 3, 6]$. \square

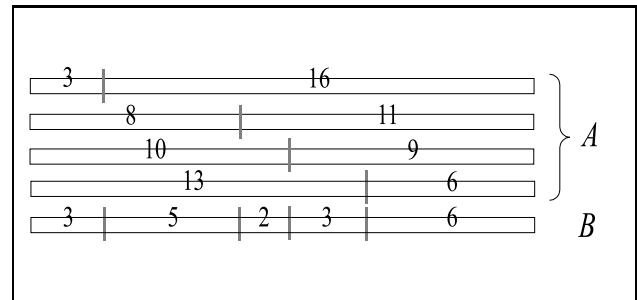


Figure 1: An example for the Simplified Partial Digest Problem without errors.

To study the computational complexity of problem Π_{SPDP} , we must introduce an additional decision *quasi-mapping* problem Π_{QMD} . In the quasi-mapping problem the multisets contain arbitrary positive integers instead of the ones coming from an errorless experiment, thus the answer to the new problem is not always “yes”.

Problem 2 *Quasi-Mapping Problem* (Π_{QMD}) — *decision version*.

Instance: Multisets $A = \{a_1, a_2, \dots, a_{2n}\}$ and $B = \{b_1, b_2, \dots, b_{n+1}\}$ of positive integers.

Question: Does there exist a map of n “restriction

sites” consistent with multisets A and B , i.e. such an ordering of elements of B which allows to cover the indicated “cuts” by some order of elements of A ?

The two problems Π_{SMs} and Π_{QMd} have the same sets of instances with the positive answer: $Y_{\Pi_{\text{SMs}}} = Y_{\Pi_{\text{QMd}}}$, because any instance of Π_{QMd} resulting in the “yes” answer belongs also to $D_{\Pi_{\text{SMs}}}$. However, problem Π_{QMd} contains also instances with the negative answer ($D_{\Pi_{\text{SMs}}} \subset D_{\Pi_{\text{QMd}}}$).

In the next subsection we construct a pseudo-polynomial transformation from the problem Numerical Matching With Target Sums, described below, to the considered SPDP.

Problem 3 Numerical Matching With Target Sums (Π_{NMd}) — decision version (Garey & Johnson 1979).

Instance: Disjoint sets X and Y , each containing m elements, sizes $s(x_i)$ and $s(y_i)$ for every $x_i \in X$ and $y_i \in Y$, and a target vector $[z_1, z_2, \dots, z_m]$; $s(x_i), s(y_i)$, and z_i being positive integers, $i = 1..m$.

Question: Can $X \cup Y$ be partitioned into m disjoint sets W_1, W_2, \dots, W_m , each containing exactly one element from X and one element from Y , such that $\sum_{w \in W_i} s(w) = z_i, i = 1..m$?

This problem is strongly NP-complete (Garey & Johnson 1979).

2.2 The transformation

The first stage of proving strong NP-hardness of the Simplified Partial Digest Problem without errors consists in a simple modification of problem Π_{NMd} . We are interested in a variant with the ranges of variables shifted by some added values. The modified ranges will have several properties, in particular, they will be disjoint. Initially, the ranges of values of variables $s(x_i), s(y_i)$, and $z_i, i = 1..m$, can be written as follows:

$$\begin{aligned} s(x_i) &\in \langle x_L, x_R \rangle, \\ s(y_i) &\in \langle y_L, y_R \rangle, \\ z_i &\in \langle z_L, z_R \rangle, \end{aligned}$$

where the variables with index L mean the smallest values and the ones with index R mean the largest values in the respective collections. These ranges are arbitrary, however, we assume here that they satisfy few obvious conditions:

$$\begin{aligned} z_R &\leq x_R + y_R, \\ z_L &\geq x_L + y_L, \\ z_R &> x_R, \\ z_R &> y_R. \end{aligned}$$

The above assumptions do not change the computational complexity of problem Π_{NMd} — if one of them is not satisfied, the problem becomes easy because the answer is obviously “no”.

We modify the initial ranges of the variables in the following way (where “:=” assigns the right-hand side value to the left-hand side variable):

$$\begin{aligned} s(x_i) &:= s(x_i) + z_R, & i &= 1..m \\ s(y_i) &:= s(y_i) + 2x_R + 2z_R, & i &= 1..m \\ z_i &:= z_i + 2x_R + 3z_R, & i &= 1..m. \end{aligned}$$

The new ranges have the following form ($i = 1..m$):

$$\begin{aligned} s(x_i) &\in \langle x_L + z_R, x_R + z_R \rangle, \\ s(y_i) &\in \langle y_L + 2x_R + 2z_R, y_R + 2x_R + 2z_R \rangle, \\ z_i &\in \langle z_L + 2x_R + 3z_R, 4z_R + 2x_R \rangle. \end{aligned}$$

They are visualized in Figure 2.

Lemma 1 Problem Π_{NMd} and its version with the ranges of variable values shifted as above, are equivalent.

PROOF

Both problems have the same sets of instances (at the beginning). The only difference is in the equation from the question, where, instead of $s(x_i) + s(y_j) = z_k$ in problem Π_{NMd} , we have $s(x_i) + z_R + s(y_j) + 2x_R + 2z_R = z_k + 2x_R + 3z_R$ in the new problem with the shifted ranges. We see that the two problems are the same. \square

The variables from problem Π_{NMd} increased by the proposed values have several useful properties.

Lemma 2 None of the modified $z_i, i = 1..m$, can be equal to some $s(x_j)$ or to a sum of any $s(x_j)$ and $s(x_k)$.

PROOF

The first statement is obvious, see Figure 2. As to the second one, the sum of two largest possible values of $s(x_j)$ and $s(x_k)$, i.e. $2x_R + 2z_R$, is still smaller than $z_L + 2x_R + 3z_R$ being the smallest z_i . \square

Lemma 3 None of the modified $z_i, i = 1..m$, can be equal to some $s(y_j)$ or to a sum of any $s(y_j)$ and $s(y_k)$.

PROOF

Also here the first statement is obvious, see Figure 2. On the other hand, the sum of two smallest possible values of $s(y_j)$ and $s(y_k)$, i.e. $2y_L + 4x_R + 4z_R$, is larger than $2x_R + 4z_R$ being the largest possible z_i . \square

Lemma 4 None of the modified $z_i, i = 1..m$, can be equal to a sum of any $s(y_j), s(x_k)$, and $s(x_l)$.

PROOF

The sum of one smallest $s(y_j)$ and two smallest $s(x_k)$ and $s(x_l)$, i.e. $y_L + 2x_L + 2x_R + 4z_R$, is larger than the largest $z_i = 2x_R + 4z_R$. \square

Now, we can define a transformation from problem Π_{NMd} to problem Π_{QMd} , which is given below.

The transformation

Given an instance of problem Π_{NMd} , the corresponding instance of Π_{QMd} is constructed as follows.

- (1) Shift the ranges of numbers in the problem Π_{NMd} as specified above, i.e.

$$\begin{aligned} s(x_i) &:= s(x_i) + z_R, & i &= 1..m \\ s(y_i) &:= s(y_i) + 2x_R + 2z_R, & i &= 1..m \\ z_i &:= z_i + 2x_R + 3z_R, & i &= 1..m. \end{aligned}$$

From now on all the variables have these modified values, if not stated otherwise.

- (2) Add values $s(x_i)$ and $s(y_i), i = 1..m$, to (initially empty) multiset B . Also add to B ($\sum_{i=1..m} z_i - 2m$) times value 1.
- (3) Set L to $2 \sum_{i=1..m} z_i - 2m$ and n to $\sum_{i=1..m} z_i - 1$.
- (4) For all $i = 1.. \frac{L}{2}$ add values i and $L - i$ to (initially empty) multiset A . Also add to A values $\sum_{j=1..i} z_j - i$ and $L - \sum_{j=1..i} z_j - i$ for all $i = 1..m - 1$. \square

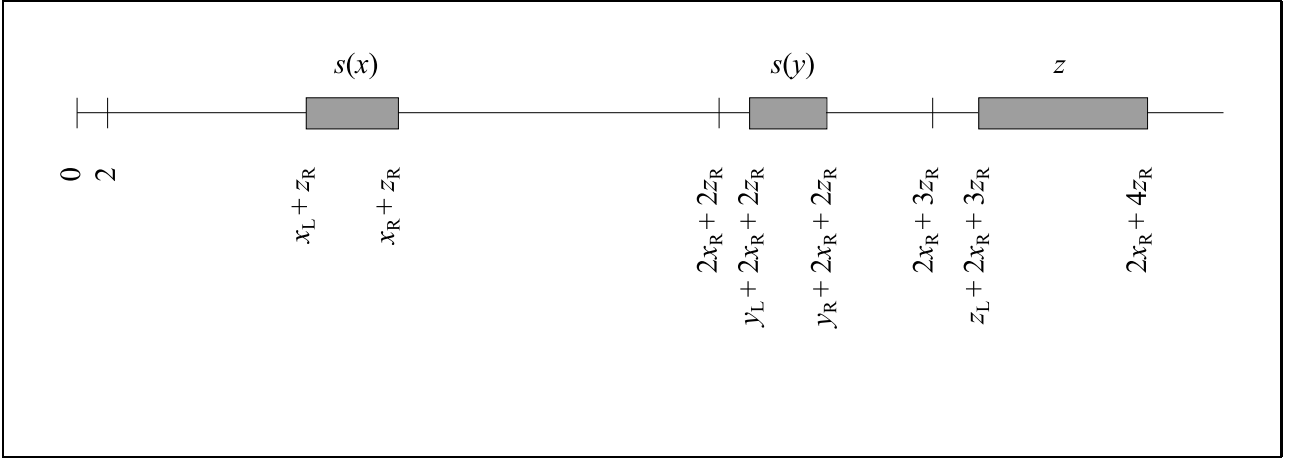


Figure 2: The ranges of values of variables $s(x_i)$, $s(y_i)$, and z_i , $i = 1..m$, after the modification.

Lemma 5 *The proposed transformation can be computed in time bounded by a polynomial in the length of the instance of Π_{NMd} (Len_{NMd}) and the maximal number appearing in this instance (Max_{NMd}).*

PROOF

Len_{NMd} is $O(m \lceil \log \text{Max}_{\text{NMd}} \rceil)$. In the first step of the transformation we make $O(\text{Len}_{\text{NMd}})$ operations. New values of the variables do not change Len_{NMd} and Max_{NMd} substantially: m is not changed, new Max_{NMd} is up to 6 times larger than previously.

Filling multiset B requires $O(\text{Len}_{\text{NMd}} \text{Max}_{\text{NMd}})$ operations. Step (3) is $O(\text{Len}_{\text{NMd}})$ and filling A takes $O(\text{Len}_{\text{NMd}} \lceil \log \text{Len}_{\text{NMd}} \rceil \text{Max}_{\text{NMd}})$ operations. Taking the above functions together, we have $O(\text{Len}_{\text{NMd}} \lceil \log \text{Len}_{\text{NMd}} \rceil \text{Max}_{\text{NMd}})$ as the complexity of the proposed transformation. Thus, it is pseudo-polynomial in time. \square

We can now prove the following main theorem.

Theorem 1 *Quasi-mapping problem Π_{QMd} is strongly NP-complete.*

PROOF

The proof uses the proposed transformation, which is pseudo-polynomial one (see Lemma 5). Below we prove that the transformation is correct. For all $I \in D_{\Pi_{\text{NMd}}}$, $I \in Y_{\Pi_{\text{NMd}}}$ if and only if $t(I) \in Y_{\Pi_{\text{QMd}}}$, where t means the transformation.

The first step of the transformation slightly modifies problem Π_{NMd} , but both versions are equivalent (see Lemma 1). Thus, in the following the shifted ranges of variable values are used.

Let us assume that there exists a solution for problem Π_{NMd} . It means that there is a partition of $X \cup Y$ such that every disjoint subset W_i , $i = 1..m$, contains one x_j and one y_k and $s(x_j) + s(y_k) = z_i$, for some $j, k \in \langle 1, m \rangle$. Then, the solution of problem Π_{QMd} can be constructed by ordering elements of B in the following way.

```

for  $i := 1$  to  $m$  with step 1
begin
  take  $s(x_j) : x_j \in W_i$ ;
  for  $j := 1$  to  $z_i - s(x_j) - 1$  with step 1
    take 1;
end
for  $i := m$  to 1 with step -1
begin
  take  $s(y_k) : y_k \in W_i$ ;
  for  $j := 1$  to  $z_i - s(y_k) - 1$  with step 1
    take 1;
end

```

We take to the solution m times some $s(x_j)$, m times some $s(y_k)$, and $\sum_{i=1..m} (z_i - s(x_j) - 1 + z_i - s(y_k) - 1)$ times element 1, $x_j, y_k \in W_i$. All x_j and y_k , $j, k = 1..m$, are in the partitioning and the sums of sizes of these pairs cover the whole vector $[z_1, z_2, \dots, z_m]$. Thus, all the elements of B are used after finishing the above procedure.

Moreover, we can cover n indicated “cuts” by an order of pairs of all elements of A . The pairs can be easily determined by summing the elements up to L (which is the length of the solution of Π_{QMd}). The procedure of ordering them (placing an element on the left or on the right side of the solution, i.e. calculating the distance from the left or from the right end of the restriction map, respectively), is shown below.

```

for  $i := 1$  to  $m$  with step 1
  for  $j := 1$  to  $z_i - 1$  with step 1
    begin
      if  $j < s(x_k) : x_k \in W_i$  then
        begin
          place  $\sum_{l=1..i-1} (z_l - 1) + j$  on the
            right side;
          place  $L - \sum_{l=1..i-1} (z_l - 1) + j$  on
            the left side;
        end
      else
        begin
          place  $\sum_{l=1..i-1} (z_l - 1) + j$  on the
            left side;
          place  $L - \sum_{l=1..i-1} (z_l - 1) + j$  on
            the right side;
        end
      end
    end
  for  $i := 1$  to  $m - 1$  with step 1
    begin
      place  $\sum_{j=1..i} (z_j - 1)$  on the right side;
      place  $L - \sum_{j=1..i} (z_j - 1)$  on the left
        side;
    end
end

```

We use in the procedure all elements of A : values increasing from 1 to $\sum_{i=1..m} (z_i - 1)$ with step 1, as well as the elements $\sum_{j=1..i} (z_j - 1)$, $i = 1..m-1$, together with their complements. Moreover, their ordering agrees with the ordering within B in the sense of “cuts”, what follows the procedures. Therefore, a feasible solution for problem Π_{QMd} exists.

Now, let us assume, that there exists a solution for problem Π_{QMd} . Thus, there is an order of $n + 1$ elements of B which allows to cover the indicated “cuts”

by some order of n pairs of elements of A . The solution for the corresponding instance of problem Π_{NMD} can be then constructed as follows.

Having the solution for problem Π_{QMD} , we know that several “cuts” appear for sure. Those are the ones corresponding to duplicated pairs of complements: if we have in A value d twice (together with $L-d$ twice, of course), then in any solution to problem Π_{QMD} one d will be placed on the left and the second d will be placed on the right. (Two “cuts” cannot appear on the same side, because we assume the lack of any errors in the data of the problem. Duplication of a “cut” would be such an error.) The guaranteed “cuts” are determined by pairs $\sum_{j=1..i} z_j - i$ and $L - \sum_{j=1..i} z_j - i$ for all $i = 1..m - 1$, together with their mirrors. Also the “cut” in the middle of L is guaranteed by the pair of two $\frac{L}{2}$. These “cuts” become ends of fragments, which correspond to disjoint subsets from a solution of problem Π_{NMD} . Every pair of mirroring fragments F_{iL} and F_{iR} of length $z_i - 1$ corresponds to set W_i of size z_i .

How are elements of B arranged in the solution? Except for many 1s, B contains also all $s(x_i)$ and $s(y_i)$, $i = 1..m$. These sizes are significant, especially for Y . Summing up Lemmae 2, 3, and 4, which are still true if values z_i will be decreased by 1, we deduce that every pair F_{iL} , F_{iR} can contain at most one $s(y_j)$ of any value. Because B includes m sizes of elements of Y and there are m pairs of fragments, all $s(y_i)$, $i = 1..m$, must be placed in separate pairs. In addition, none two (or more) sizes of elements of X can be added to a pair of fragments (already containing some $s(y_j)$), so all $s(x_i)$, $i = 1..m$, must be also placed in separate pairs.

At first glance, it is not obvious why two mirroring fragments F_{iL} and F_{iR} cannot include two sizes of elements of Y . It is displayed in Figure 3. If some size is assigned to one fragment, the corresponding place in the second fragment must be filled by “cuts”. Because sizes of elements of Y are greater than halves of lengths of fragments plus 1, none two such sizes can be inserted into a mirroring pair.

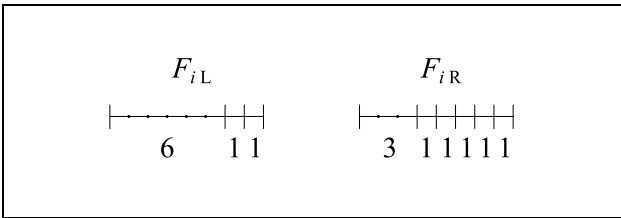


Figure 3: Two mirroring fragments F_{iL} and F_{iR} with some components.

If we have exactly one $s(x_j)$ and exactly one $s(y_k)$ assigned to a pair F_{iL} , F_{iR} , for some $i, j, k \in \langle 1, m \rangle$, plus in addition only 1s, there is only one possible placement. One of the longer sizes is shifted to the left (or right) end of F_{iL} and the second one is shifted to the left (or right, respectively) end of F_{iR} . The sum of both sizes is equal to z_i . The remaining places are filled by 1s (see Figure 3). Another placement would result in the occurrence of an additional element of a value greater than 1, what is in contradiction to previous proofs.

Having the solution of problem Π_{QMD} we can immediately read the solution of problem Π_{NMD} : those elements x_j and y_k , which sizes appear in pair F_{iL} , F_{iR} , will compose set W_i . All the requirements for a solution of problem Π_{NMD} are satisfied. This ends the proof. \square

The proposed transformation of an instance of problem Π_{NMD} to an instance of problem Π_{QMD} is illustrated by the following example.

Example 2 Let the example instance of problem Π_{NMD} be:

$$\begin{aligned} m &= 3, \\ X &= \{x_1, x_2, x_3\}, \\ Y &= \{y_1, y_2, y_3\}, \\ s(x_1) &= 2, s(x_2) = 3, s(x_3) = 5, \\ s(y_1) &= 4, s(y_2) = 5, s(y_3) = 6, \\ z_1 &= 7, z_2 = 9, z_3 = 9. \end{aligned}$$

After shifting the ranges of the variables we get the values:

$$\begin{aligned} s(x_1) &= 11, s(x_2) = 12, s(x_3) = 14, \\ s(y_1) &= 32, s(y_2) = 33, s(y_3) = 34, \\ z_1 &= 44, z_2 = 46, z_3 = 46. \end{aligned}$$

The construction of the instance of problem Π_{QMD} ends with the following result:

$$\begin{aligned} n &= 135, L = 266, \\ A &= \bigcup_{i=1..133} \{i, 266 - i\} \cup \{43, 223, 88, 178\}, \\ B &= \bigcup_{i=1..130} \{1\} \cup \{11, 12, 14, 32, 33, 34\}. \end{aligned}$$

The feasible solution for the instance after the transformation (i.e. of problem Π_{QMD}) is shown in Figure 4.

This solution can be easily translated to a feasible solution of problem Π_{NMD} :

$$\begin{aligned} W_1 &= \{x_1 + y_2\} \quad (11 + 33 = 44 \rightarrow 2 + 5 = 7), \\ W_2 &= \{x_2 + y_3\} \quad (12 + 34 = 46 \rightarrow 3 + 6 = 9), \\ W_3 &= \{x_3 + y_1\} \quad (14 + 32 = 46 \rightarrow 5 + 4 = 9). \quad \square \end{aligned}$$

Finally, we prove the computational hardness of the Simplified Partial Digest Problem without errors in the search version Π_{SMS} .

Theorem 2 The Simplified Partial Digest Problem without errors Π_{SMS} (search version) is strongly NP-hard.

PROOF

Proving strong NP-completeness of quasi-mapping problem Π_{QMD} (Theorem 1) directly leads to proving strong NP-hardness of the corresponding problem Π_{SMS} . For, if we had an algorithm solving Π_{SMS} in polynomial time, we could use it to solve problem Π_{QMD} in polynomial time in the following way. We could apply the algorithm to the instance of Π_{QMD} , and after a number of steps bounded by the polynomial function known for the algorithm we could have the answer for Π_{QMD} . Either the algorithm would find the solution and the answer would be “yes”, or the algorithm would not find one and the answer would be “no”. As pointed out earlier, a similar reasoning has been used in Johnson (1985), where the problem of looking for a Hamiltonian circuit in a graph has been considered. \square

3 A polynomially solvable subproblem of SPDP

On the other hand, a special version of the Simplified Partial Digest Problem without errors (a combinatorial search problem), in which multiset B is composed of only the elements of values 1 and 2, is polynomially solvable. In this section we present the algorithm solving the subproblem, defined below.

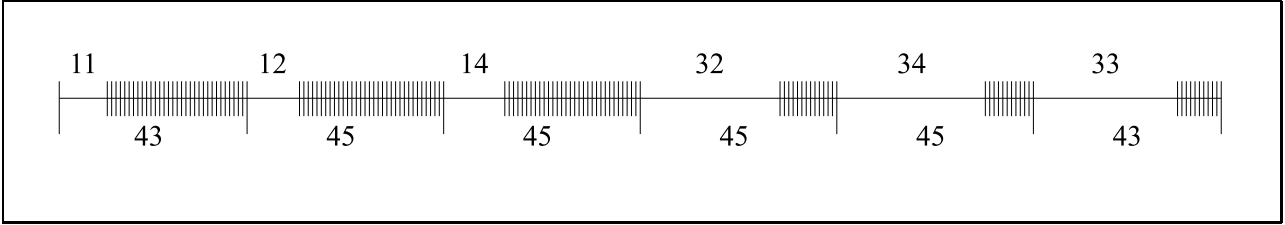


Figure 4: An example solution of problem Π_{QMd} .

Problem 4 *Simplified Partial Digest Problem with long-digestion fragments of lengths only 1 or 2 (Π_{SM12s}) — search version.*

Instance: Multiset $A = \{a_1, a_2, \dots, a_{2n}\}$ of lengths of restriction fragments coming from short digestion and multiset $B = \{b_1, b_2, \dots, b_{n+1}\}$ of lengths of restriction fragments coming from long digestion, A and B containing no errors and $b_i \in \{1, 2\}$, $i = 1..n + 1$.

Answer: A map of n restriction sites of a DNA chain consistent with multisets A and B , i.e. such an order of elements of B which allows to cover the indicated cuts by some order of the elements of A .

The algorithm solving Π_{SM12s}

- (1) Calculate $L = \sum_{i=1..n+1} b_i$. Initially the solution is one fragment of length L without cuts.
- (2) Match complementary pairs of lengths in A , i.e. the ones summing up to L . Let the pairs be denoted by A_i , $i = 1..n$, and satisfying $\bigcup_{i=1..n} A_i = A$ and $\bigcap_{i=1..n} A_i = \emptyset$.
- (3) Select all identical pairs A_i and A_j , $i, j = 1..n$, and mark the cuts corresponding to them in the solution symmetrically (from both sides). Remove all selected A_i and A_j from A .
- (4) If A is empty, go to step (7). Otherwise, select such A_i in A , that contains the shortest length l in current A . Mark the cut placed l units from the left end of the solution and remove A_i from A .
- (5) If A is empty, go to step (7). Otherwise, increase l by 1. If some cut already exists in the solution l units from the right end, go to step (4). Otherwise, mark the cut placed l units from the right end of the solution and remove A_i which contains l from A .
- (6) If A is empty, go to step (7). Otherwise, increase l by 1. If some cut already exists in the solution l units from the left end, go to step (4). Otherwise, mark the cut placed l units from the left end of the solution, remove A_i which contains l from A , and go to step (5).
- (7) The current solution contains all cuts, i.e. it corresponds to an order of elements of multiset B .

Let us note that the algorithm uses only information coming from multiset A and does not check it with multiset B . This is because the applied moves are the only possible ones, and in the case where we have no errors they have to be correct. This is demonstrated in the following proof.

Theorem 3 *Problem Π_{SM12s} is solvable in time $O(n \log n)$.*

PROOF

Let us analyze the consecutive steps of the algorithm.

The first two steps are obvious. The third one consists in placing the symmetric cuts indicated by

pairs from A . Of course, if we have in A two identical subsets $\{x, L-x\}$, then in the problem without errors both cuts must appear. Thus, they must be placed in opposite parts of the solution, i.e. one cut x units from the left and the second one x units from the right end.

All remaining cuts are asymmetric. Because the current solution is symmetric, no matter what side we place the next cut. So, the left side can be chosen, as well as the strategy of taking the cut nearest the left end first. Therefore, the fourth step is correct. It also concerns further calling of step (4), when we fill a part between symmetric cuts and go to the next slot (i.e. a fragment between two pairs of symmetric cuts from step (3)) — then there is no difference from what side we start to place next cuts.

After placing the previous cut in step (4) we must look at the opposite side of the solution. In the symmetric place we cannot have a cut, since all cuts added after step (3) are asymmetric. And because in multiset B there are only lengths equal to 1 or 2, we must place a cut to get a fragment of length at most 2 after one unit on the right side without a cut. So, step (5) is correct as well as step (6) with a similar reasoning. The only exception, when we reach an existing cut, concerns the situation when we fill the whole current slot and jump to the next one (the slot closer to the middle of the solution).

Summing up the above reasoning, we may conclude that the algorithm is correct. As to its complexity, if the elements of A are sorted in step (2), the steps can be done in time, respectively, $O(n \log n)$, $O(n \log n)$, $O(n \log n)$, $O(\log n)$, $O(\log n)$, and $O(\log n)$. Steps (4)–(6) are performed up to n times, so the whole algorithm is of the complexity $O(n \log n)$. \square

Example 3 presents in detail the proposed algorithm.

Example 3 *Let our instance of problem Π_{SM12s} be: $A = \{1, 2, 3, 4, 5, 5, 6, 7, 7, 8, 8, 9, 10, 10, 11, 12, 13, 14\}$ and $B = \{1, 1, 1, 1, 1, 2, 2, 2, 2, 2\}$. Thus, we have 9 restriction places within the solution of length 15. The complementary pairs of lengths in A are: $\{1, 14\}$, $\{2, 13\}$, $\{3, 12\}$, $\{4, 11\}$, $\{5, 10\}$, $\{5, 10\}$, $\{6, 9\}$, $\{7, 8\}$, and $\{7, 8\}$. To the initially empty solution we add all symmetric cuts determined by the instance and we get the sequence of fragments as in Figure 5.*

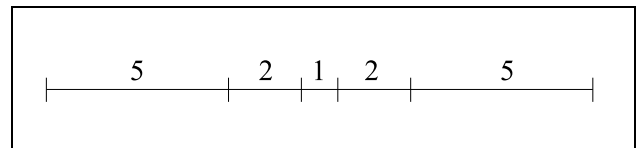


Figure 5: A symmetric part of an example solution of problem Π_{SM12s} .

After that $A = \{\{1, 14\}, \{2, 13\}, \{3, 12\}, \{4, 11\}, \{6, 9\}\}$. We add the remaining cuts starting from the

one nearest the left end of the solution, the result is shown in Figure 6.

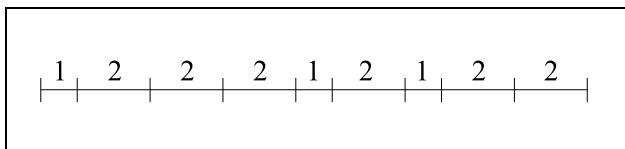


Figure 6: An example solution of problem Π_{SM12s} with the first slot filled by cuts.

Finally, we jump to the next slot and place the only remaining cut $\{6, 9\}$. The whole solution is in Figure 7. \square

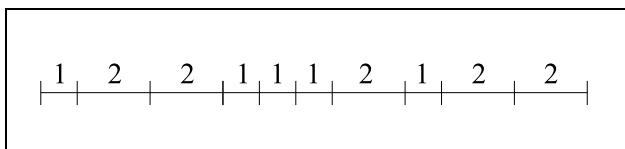


Figure 7: A complete solution for the example problem Π_{SM12s} .

4 Conclusions

In the paper, the DNA Simplified Partial Digest Problem was formulated and analyzed as a combinatorial search problem. The general error-free case was proved to be strongly NP-hard. On the other hand, a special form of the multiset resulting from the long digestion and composed of only 1s and 2s results in a simple, polynomial-time algorithm. The question remains open, whether or not, a similar approach can be used for the multisets composed of other restricted values (e.g. 1 and 3). An interesting problem is a question of approximability, i.e. an existence of polynomial-time algorithm constructing solutions being not far from the optimum one with respect to a suitably chosen optimality criterion.

Appendix

In the following, the original formulation of the Simplified Partial Digest Problem without errors in instances, from Blazewicz *et al.* (2001), is presented. We need some additional definitions.

$C = \langle c_1, c_2, \dots, c_{2n} \rangle$ — multiset A sorted in non-decreasing order,

$L = c_i + c_{2n-i+1}$, for any i — the length of the analyzed strand,

$\mathcal{P}_i = \langle c_i, c_{2n-i+1} \rangle$, $i = 1..2n$ — an ordered pair of complementary fragments (i.e. summing up to L), where c_i is called the predecessor,

$\mathcal{Q} = \{q_1, q_2, \dots, q_n\}$ — a set of ordered pairs, where $q_i = \mathcal{P}_i$ or $q_i = \mathcal{P}_{2n-i+1}$, $i = 1..n$,

$\mathcal{O} = \langle o_1, o_2, \dots, o_n \rangle$ — the list of predecessors from every q_i , $i = 1..n$, sorted in non-decreasing order,

$\mathcal{R} = \{r_1, r_2, \dots, r_{n+1}\}$ — a multiset derived from \mathcal{O} according to the following rule: $r_1 = o_1$, $r_i = o_i - o_{i-1}$, $i = 2..n$, $r_{n+1} = L - o_n$.

Now, the problem is:

Problem 1' *Simplified Partial Digest Problem — search version.*

Instance: Multiset $A = \{a_1, a_2, \dots, a_{2n}\}$ of lengths of restriction fragments coming from short digestion and multiset $B = \{b_1, b_2, \dots, b_{n+1}\}$ of lengths of restriction fragments coming from long digestion, A and B containing no errors.

Answer: The set \mathcal{Q} such that the corresponding multiset \mathcal{R} is equal to B .

References

- Benzer, S. (1959), 'On the topology of the genetic fine structure', *Proceedings of the National Academy of Sciences of the USA* **45**, 1607–1620.
- Blazewicz, J., Formanowicz P., Kasprzak M., Jaroszewski M. & Markiewicz W.T. (2001), 'Construction of DNA restriction maps based on a simplified experiment', *Bioinformatics* **17**, 398–404.
- Blazewicz, J. & Jaroszewski M. (2003), 'New algorithm for the Simplified Partial Digest Problem', *Lecture Notes in Bioinformatics* **2812**, 95–110.
- Blazewicz, J. & Kasprzak M. (2003), 'Complexity of DNA sequencing by hybridization', *Theoretical Computer Science* **290**, 1459–1473.
- Cieliebak, M. & Eidenbenz S. (2004), 'Measurement errors make the partial digest problem NP-hard', *Lecture Notes in Computer Science* **2976**, 379–390.
- Cieliebak, M., Eidenbenz S. & Penna P. (2003), 'Noisy data make the partial digest problem NP-hard', *Lecture Notes in Bioinformatics* **2812**, 111–123.
- Fogel, G.B. & Corne D.W., eds. (2003), *Evolutionary Computations in Bioinformatics*, Morgan Kaufman: Boston.
- Garey, M.R. & Johnson D.S. (1979), *Computers and Intractability. A Guide to the Theory of NP-Completeness*, W.H. Freeman and Company: San Francisco.
- Goldstein, L. & Waterman M.S. (1987), 'Mapping DNA by stochastic relaxation', *Advanced Applied Mathematics* **8**, 194–207.
- Johnson, D.S. (1985), 'The NP-completeness column: an ongoing guide', *Journal of Algorithms* **6**, 291–305.
- Karp, R.M. & Shamir R. (1998), Algorithms for optical mapping, in 'Proceedings of the Second International Conference on Computational Biology (RECOMB)', New York, 117–124.
- Newberg, L. & Naor D. (1993), 'A lower bound on the number of solutions of the probed partial digest problem', *Advanced Applied Mathematics* **14**, 172–183.
- Pandurangan, G. & Ramesh H. (2002), 'The restriction mapping problem revisited', *Journal of Computer and System Sciences* **65**, 526–544.
- Pevzner, P.A. (2000), *Computational Molecular Biology: An Algorithmic Approach*, MIT Press: Cambridge.
- Schmitt, W. & Waterman M.S. (1991), 'Multiple solutions of DNA restriction mapping problem', *Advanced Applied Mathematics* **12**, 412–427.

- Schwartz, D.C., Li X., Hernandez L.I., Ramnarain S.P., Huff E.J. & Wang Y.K. (1993), 'Ordered restriction maps of *Saccharomyces cerevisiae* chromosomes constructed by optical mapping', *Science* **262**, 110–114.
- Setubal, J. & Meidanis J. (1997), *Introduction to Computational Biology*, PWS Publishing Company: Boston.
- Skiena, S.S. (1997), Geometric reconstruction problems, in J.E. Goodman & J. O'Rourke, eds, 'Handbook of Discrete and Computational Geometry', CRC Press: Boca Raton, 481–490.
- Skiena, S.S., Smith W.D. & Lemke P. (1990), Reconstructing sets from interpoint distances, in 'Proceedings of Sixth ACM Symposium on Computational Geometry', 332–339.
- Skiena, S.S. & Sundaram G. (1994), 'A partial digest approach to restriction site mapping', *Bulletin of Mathematical Biology* **56**, 275–294.
- Waterman, M.S. (1995), *Introduction to Computational Biology. Maps, Sequences and Genomes*, Chapman and Hall: London.