

**SOLVING LARGE DOUBLE DIGESTION PROBLEMS FOR  
DNA RESTRICTION MAPPING BY USING  
BRANCH-AND-BOUND INTEGER LINEAR PROGRAMMING**

By

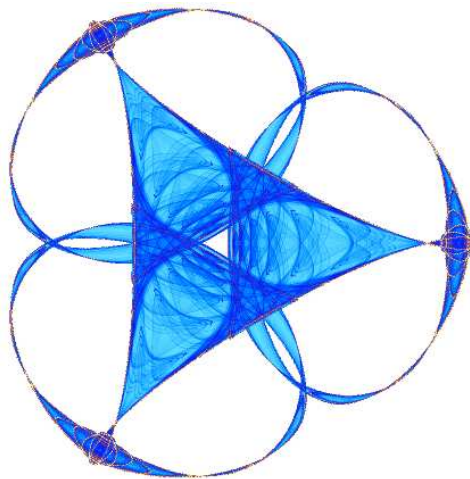
**Zhijun Wu**

and

**Yin Zhang**

**IMA Preprint Series # 2176**

( October 2007 )



**INSTITUTE FOR MATHEMATICS AND ITS APPLICATIONS**

UNIVERSITY OF MINNESOTA  
400 Lind Hall  
207 Church Street S.E.  
Minneapolis, Minnesota 55455-0436

Phone: 612-624-6066 Fax: 612-626-7370

URL: <http://www.ima.umn.edu>

# Solving Large Double Digestion Problems for DNA Restriction Mapping by Using Branch-and-Bound Integer Linear Programming

Zhijun Wu<sup>1</sup> and Yin Zhang<sup>2</sup>

<sup>1</sup>Department of Mathematics  
Iowa State University  
Ames, Iowa, U.S.A.

<sup>2</sup>Department of Computational and Applied Mathematics  
Rice University  
Houston, Texas, U.S.A.

**Abstract.** The double digestion problem for DNA restriction mapping has been proved to be *NP*-complete, and is intractable if the numbers of the DNA fragments generated by the two restriction enzymes are large. Several approaches to the problem have been used, including exhaustive search and simulated annealing, and have proved to be effective only for small problems, as the problem size or in other words, the size of the search space for the problem grows as a factorial function  $(n!)(m!)$  of the numbers  $n$  and  $m$  of the DNA fragments generated by the two restriction enzymes, respectively. In this paper, we formulate the double digestion problem as a mixed-integer linear program by following Waterman 1995 in a slightly different form. We show that with this formulation and by using some state-of-the-art integer programming techniques, we can actually solve double digestion problems for fairly large sizes. In particular, we can solve a set of randomly generated problems of sizes up to  $(242!)(250!)$ , which are many magnitude increases from previously reported  $(16!)(16!)$  testing sizes.

## 1 Introduction

A DNA sequence can be broken down (digested) into a set of small fragments by using some special proteins called restriction enzymes. The sequences of

the fragments usually become easier to determine and the lengths of the fragments can also be measured by using, for example, the electrophoretic gel. Once determined, the fragments can be put back in the original order to obtain the whole sequence. The latter process is called the DNA restriction mapping and is an important step in genomic sequencing [14]. A critical part of DNA restriction mapping is to find the right order of the fragments given only the lengths of the fragments. This problem is called the digestion problem and in particular, the double digestion problem if two different restriction enzymes are used in the experiments [17].

Let the two restriction enzymes be denoted by  $A$  and  $B$ . Let  $a = \{a_1, a_2, \dots, a_n\} \in Z^n$  be the lengths of the fragments obtained after applying enzyme  $A$  to a given sequence, and  $b = \{b_1, b_2, \dots, b_m\} \in Z^m$  the lengths of the fragments obtained after applying enzyme  $B$ , where the lengths are measured in numbers of DNA base pairs. Let  $c = \{c_1, c_2, \dots, c_l\} \in Z^l$  be the lengths of the fragments obtained after first applying enzyme  $A$  and then enzyme  $B$  to the same sequence. The fragments in  $a$ ,  $b$ , and  $c$  are given in arbitrary orders. However, as shown in Figure 1, if we can put the fragments in  $a$  and  $b$  in their original orders, we can immediately obtain all the fragments in  $c$  by simply combining all the restriction sites in  $a$  and  $b$  as if they were digested by a single enzyme. On the other hand, if we put  $a$  and  $b$  in other orders, we will in general obtain a set of fragments different from  $c$ . The double digestion problem is to find the orders for  $a$  and  $b$ , given  $c$  as the result of the sequence being digested by  $A$  and  $B$  combined.

The double digestion problem is solved routinely in molecular biology labs for constructing the restriction maps for newly cloned DNA sequences [17]. The solution of a generalized version of the problem may also be applied to protein determination in proteomics using mass spectrometry [11]. The problem can be difficult to solve, however, if the numbers of the generated fragments are large. In general, it is intractable as shown in Goldstein and Waterman 1987 [4]. Several approaches to the problem have been proposed, such as exhaustive search [12, 15], simulated annealing [4, 5, 18], and fragment matching [3, 8, 16], but they all have proved to be effective only for small problems, as the problem size or in other words, the size of the search space for the problem, grows as a factorial function  $(|a|!)(|b|!)$  of  $|a|$  and  $|b|$ , the numbers of the fragments in  $a$  and  $b$  [4].

In this paper, we formulate the double digestion problem as a mixed-integer linear program by following Waterman 1995 [17] in a slightly different form. We show that with this formulation and by using some state-of-the-art integer programming techniques, we can actually solve double digestion

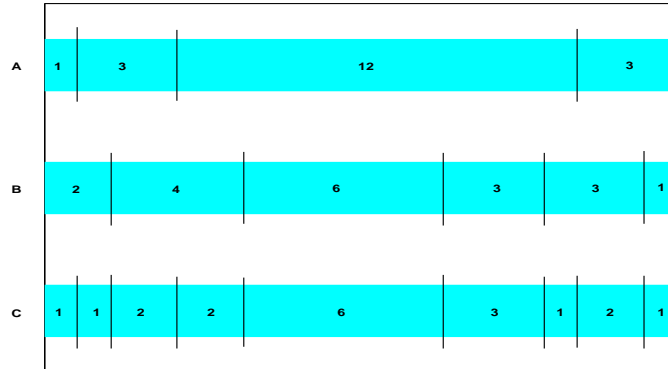


Figure 1: **Double digestion** A sequence of 19 base pairs is digested by a restriction enzyme **A** at sites 1, 4, 16, and is cut into 4 fragments of lengths 1, 3, 12, 3, and by **B** at sites 2, 6, 12, 15, 18, and is cut into 6 fragments of lengths 2, 4, 6, 3, 3, 1. If both **A** and **B** are applied, the sequence will be digested at both **A** and **B** sites, and is cut into 9 fragments of lengths 1, 1, 2, 2, 6, 3, 1, 2, 1 in **C**.

problems for fairly large sizes. In particular, we can solve a set of randomly generated problems of sizes up to  $(242!)(250!)$ , which are many magnitude increases from previously reported  $(16!)(16!)$  testing sizes [4].

More specifically, the double digestion problem can be solved in two steps. First, based only on their lengths, find the fragments in  $c$  that can fit in each of the fragments in  $a$  and  $b$  (like solving some one-dimensional puzzles). This is called the fragment matching. Second, based on the obtained matching relationships of  $a$  and  $b$  with  $c$ , find the orders of the fragments in  $a$  and  $b$ . The first step is the computationally most difficult part of the problem, because the second step can be accomplished easily in order of  $|a| + |b|$  time (Chapter 2 in [17]). We therefore focus on fragment matching and formulate the problem as a mixed-integer linear program that minimizes the matching errors in either  $l_1$  or  $l_\infty$  norm. We then use the CPLEX mixed-integer linear programming software [7] to solve the problem. In this way, we have been able to solve a set of randomly generated test problems with the numbers of the generated fragments in  $a$  and  $b$  ranging from 18 to 242 and 20 to 250, respectively. Since the total number of possible orderings for  $c$  is proportional to  $(|a|!)(|b|!)$  [4], assuming most elements in  $a$  (and  $b$ ) are different, the problem sizes or the search spaces for our test problems are proportional to  $(18!)(20!)$  to  $(242!)(250!)$ .

We describe the mixed-integer linear programming formulation in greater detail in Section 2 and introduce the CPLEX mixed-integer linear programming software in Section 3. Computational results are presented and discussed in Section 4. Comments and remarks are given in Section 5.

## 2 MIP Formulation

Let  $A$  and  $B$  be two restriction enzymes. Let  $a = \{a_1, a_2, \dots, a_n\} \in Z^n$  be the lengths of the fragments obtained after applying enzyme  $A$  to a given DNA sequence and  $b = \{b_1, b_2, \dots, b_m\} \in Z^m$  the lengths of the fragments obtained after applying enzyme  $B$ . Let  $c = \{c_1, c_2, \dots, c_l\} \in Z^l$  be the lengths of the fragments obtained after applying enzyme  $A$  followed by enzyme  $B$ , or in the other way around, applying enzyme  $B$  followed by enzyme  $A$ . In order to find correct orders for  $a$  and  $b$ , we need to find for every element in  $c$  the element in  $a$  it comes from, and the same thing in  $b$ . We call this problem the fragment matching problem. It is an important part of the double digestion problem. Once this problem is solved, a solution to the double digestion problem can be constructed easily using for example a fragment ordering algorithm as described in Chapter 2 of Waterman 1995 [17].

Let  $X$  be a  $n \times l$  matrix and  $x_{i,j}$  the  $i, j$ -element of  $X$ , and

$$x_{i,j} = \begin{cases} 1, & \text{sequence of } c_j \subset \text{sequence of } a_i, \\ 0, & \text{otherwise.} \end{cases}$$

Let  $Y$  be a  $m \times l$  matrix and  $y_{i,j}$  the  $i, j$ -element of  $Y$ , and

$$y_{i,j} = \begin{cases} 1, & \text{sequence of } c_j \subset \text{sequence of } b_i, \\ 0, & \text{otherwise.} \end{cases}$$

Since the elements  $c_j$  in  $c$  that come from the  $i$ th fragment of  $a$  add to  $a_i$  and that from the  $i$ th fragment of  $b$  add to  $b_i$ ,

$$a_i = \sum_{j=1}^l x_{i,j}c_j, \quad i = 1, \dots, n,$$

and similarly,

$$b_i = \sum_{j=1}^l y_{i,j}c_j, \quad i = 1, \dots, m.$$

In matrix form the equations are equivalent to

$$a = Xc, \quad b = Yc.$$

Note that each element in  $c$  may come from only one element in  $a$  or  $b$ . Therefore,

$$\begin{aligned} \sum_{i=1}^n x_{i,j} &= 1, & j &= 1, \dots, l, \\ \sum_{i=1}^m y_{i,j} &= 1, & j &= 1, \dots, l. \end{aligned}$$

Clearly the fragment matching problem for a triple sets  $a$ ,  $b$ , and  $c$  is basically to find  $X$  and  $Y$  so that the above equations hold. As stated in Waterman 1995 [17], the problem can be formulated as a linear integer program,

$$\begin{aligned} \min_{\alpha, \beta, X, Y} \quad & \alpha + \beta \\ \text{subject to} \quad & \alpha \geq \sum_{j=1}^l x_{i,j}c_j - a_i \geq -\alpha, & i = 1, \dots, n \\ & \beta \geq \sum_{j=1}^l y_{i,j}c_j - b_i \geq -\beta, & i = 1, \dots, m \\ & \sum_{i=1}^n x_{i,j} = 1, & j = 1, \dots, l \\ & \sum_{i=1}^m y_{i,j} = 1, & j = 1, \dots, l \\ & x_{i,j}, y_{i,j} \in \{0, 1\}. \end{aligned}$$

Here, we recognize that the program is separable and equivalent to two separate sub-programs,

$$\begin{aligned} \min_{\alpha, X} \quad & \alpha \\ \text{subject to} \quad & \alpha \geq \sum_{j=1}^l x_{i,j}c_j - a_i \geq -\alpha, & i = 1, \dots, n \\ & \sum_{i=1}^n x_{i,j} = 1, & j = 1, \dots, l \\ & x_{i,j} \in \{0, 1\}. \end{aligned}$$

and

$$\begin{aligned}
& \min_{\beta, Y} \quad \beta \\
& \text{subject to} \quad \beta \geq \sum_{j=1}^l y_{i,j} c_j - b_i \geq -\beta, \quad i = 1, \dots, m \\
& \quad \quad \quad \sum_{i=1}^m y_{i,j} = 1, \quad j = 1, \dots, l \\
& \quad \quad \quad y_{i,j} \in \{0, 1\}.
\end{aligned}$$

The two sub-programs have the same mathematical structure. It suffices to consider only one of them. We put it in the following general form,

$$\begin{aligned}
(P_\infty) \quad & \min_{t, X} \quad t \\
& \text{subject to} \quad t \geq \sum_j^l x_{i,j} c_j - a_i \geq -t, \quad i = 1, \dots, n \\
& \quad \quad \quad \sum_i^n x_{i,j} = 1, \quad j = 1, \dots, l \\
& \quad \quad \quad x_{i,j} \in \{0, 1\}.
\end{aligned}$$

This program is to find an appropriate assignment matrix  $X$  such that the errors for the equations

$$a_i = \sum_j^l x_{i,j} c_j, \quad i = 1, \dots, n$$

are minimized in  $l_\infty$  norm. We therefore name the problem  $P_\infty$ . Of course, the errors can also be minimized in  $l_1$  norm, and the program then becomes

$$\begin{aligned}
(P_1) \quad & \min_{t, X} \quad \sum_i^n t_i \\
& \text{subject to} \quad t_i \geq \sum_j^l x_{i,j} c_j - a_i \geq -t_i, \quad i = 1, \dots, n \\
& \quad \quad \quad \sum_i^n x_{i,j} = 1, \quad j = 1, \dots, l \\
& \quad \quad \quad x_{i,j} \in \{0, 1\}.
\end{aligned}$$

Here we name the problem  $P_1$ . Minimizing  $l_1$  or  $l_\infty$  norm may have different effects on the solution to the problem. If the errors are not equal to zero at the solution, minimizing  $l_1$  norm reduces total errors even if there are a few big errors, while minimizing  $l_\infty$  norm keeps the biggest error as small as possible. If an exact solution is found where the errors are all equal to zero, there are no differences in the solution qualities with  $l_1$  or  $l_\infty$  norms. However, the solutions may still be different if there are multiple solutions, and the algorithm may perform quite differently as well. In the following sections, we show how the problems  $P_1$  and  $P_\infty$  can be solved by using the CPLEX mixed-integer linear programming software.

### 3 CPLEX MIP Solver

CPLEX is a widely-used software package for solving integer, linear and quadratic programming problems. It was originally developed by Bixby [2] and later commercialized.

A linear programming problem usually is referred to as an optimization problem with a linear objective function and linear constraints. The following is a typical linear program:

$$\begin{array}{ll}
 \min_{x_1, x_2, \dots, x_n} & c_1x_1 + c_2x_2 \cdots + c_nx_n \\
 \text{subject to} & a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \geq b_1 \\
 & a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n \geq b_2 \\
 & \dots \\
 & a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n \geq b_m \\
 & l_1 \leq x_1 \leq u_1 \\
 & l_2 \leq x_2 \leq u_2 \\
 & \dots \\
 & l_n \leq x_n \leq u_n
 \end{array}$$

where  $x_1, x_2, \dots, x_n$  are continuous variables in general. The problem becomes a mixed-integer linear program if some variables are required to be integers. CPLEX has solvers for linear programs and in particular, a mixed-integer solver for mixed-integer linear programs (MIP). Linear programming problems with continuous variables (LP) can be solved in polynomial time, while mixed-integer linear programs can be hard, including many NP-complete problems. The CPLEX MIP solver is built upon its LP solvers and

uses a branch-and-bound strategy, with the aid of cutting plane techniques, to find integer solutions. Two types of cuts, based on maximal cliques and minimal covers, are generated when necessary to cut off non-integral solutions.

In the branch-and-bound algorithm for solving a MIP problem, a series of LP sub-problems is solved. A tree of subproblems is generated with each subproblem being a node of the tree. The root node is the LP relaxation of the original MIP. If the solution to the relaxation has integer variables taking fractional values, then one of such variables is chosen for branching, and two new subproblems are generated, each with more restrictive bounds for the branching variable. For example, for binary variables, one node will have the variable fixed at zero, and the other at one. The subproblem can result in an all-integer solution, an infeasible problem, or another fractional solution. If the solution is fractional, the process is repeated. A node is “cut off” when the objective function value of the subproblem associated with the node becomes higher than a known upper bound on the optimal value of the original MIP. Note that every all-integer solution gives such an upper bound. Also note that the node is “cut off” because further branching from the node will always generate subproblems with the same or even higher objective function values, and therefore, the whole branch can be removed. For more detailed descriptions on the branch-and-bound algorithm, the readers are referred to [9].

The branch-and-bound algorithm searches for an integer solution for the MIP problem in the subproblem tree. In the worst case, it may need to examine the entire tree, which can be computationally prohibitive since there are exponentially many nodes in the tree with respect to the integer variables. In practice, the performance of the algorithm depends on specific search strategies such as how to select nodes for branching, how to obtain tighter bounds, etc. The CPLEX MIP solver provides several ways for users to choose among best possible search strategies. Readers are referred to the CPLEX Users’ Manual [7] for more detailed descriptions.

For mixed-integer programming, in addition to branch-and-bound, cutting plane is another important solution technique. It is often used in combining with branch-and-bound. A cutting plane, or a cut, is an inequality added to a problem that restricts or cuts off non-integral solutions. Adding cuts usually reduces the number of branches that are needed to solve a mixed-integer program. The CPLEX MIP solver generates two types of cuts derived from maximal cliques and minimal covers:

    Clique inequalities are generated by looking at the relationships between

binary variables before optimization starts. The relationships are given in the all-binary inequalities of the original problem. A graph representing the relationships is constructed, and maximal cliques are found. At most one variable from each clique can be positive in a feasible solution. Inequalities that describe these restrictions are generated. If the solution to a subproblem violates one of these clique inequalities, the inequality is added to the problem.

Cover inequalities are derived by looking at each all-binary inequality with non-unit coefficients. For a constraint with all-negative coefficients, a minimal cover is a subset of the variables in the inequality such that if all variables were set to 1, the inequality would be violated, but if any variable were excluded, the inequality would be satisfied. If a cover inequality is violated by the solution to the current subproblem, it is added to the problem. For more detailed descriptions on clique and cover cuts for linear integer programming, the readers are referred to [6].

## 4 Computational Results

We used the CPLEX MIP solver (CPLEX 6.5.2) to obtain the solutions for problems  $P_1$  and  $P_\infty$ . The program was run on a multi-processor SGI Origin 2000 computer with sixteen 300MHZ R12000 processors and 10GB main memory. However, only one processor was used at a time. In order to see how double digestion problems can be solved by our approach, we generated a set of random test problems by following a procedure suggested by Goldstein and Waterman in [4]. We assumed a set of DNA sequences of 100, 200, 300, 400, and 500 units. The unit here may correspond to one or more than one base pairs. It does not matter what the actual unit is since the difficulty of the problem is correlated to the numbers and lengths of the fragments, which is independent of the size of the unit. For each sequence, we simulate the first enzyme restriction by cutting the sequence at every unit with a probability  $p_1$ , and the second enzyme restriction with a probability  $p_2$ . After combining the two, we obtain a triple set of fragments  $a$ ,  $b$ , and  $c$ . In order to cover certain amount of possibilities, we generated different sets of problems with  $p_1 = p_2 = 0.2$ ,  $p_1 = p_2 = 0.3$ ,  $p_1 = p_2 = 0.4$ , and  $p_1 = p_2 = 0.5$ . Obviously, more fragments are generated with increasing probabilities  $p_1$  and  $p_2$ , and the difficulties of the problems may vary. We list the sizes of the fragment sets for all generated problems in Table 1. Note that since the orders of the fragments in  $a$  and  $b$  determine the order of the

fragments in  $c$  and there are total  $(|a|!)(|b|!)$  possible orderings for  $a$  and  $b$ , the problem size, or the search space, for the double digestion problem is proportional to  $(|a|!)(|b|!)$ . Since the sizes of the generated fragment sets range from  $|a| = 18$  and  $|b| = 20$  to  $|a| = 242$  and  $|b| = 250$ , the search spaces for the generated problems can be as large as  $(242!)(250!)$ , which are many magnitude increases from previously reported  $(16!)(16!)$  testing sizes using simulated annealing [4]. Note also that if the fragments in  $c$  is to be matched to those in  $a$ , the corresponding integer program has  $n \times l$  binary variables, where  $n$  is the number of fragments in  $a$  and  $l$  the number of fragments in  $c$ . In our test problems, the number of fragments in  $a$  or  $b$  ranges from 18 to 250, while that in  $c$  from 35 to 376. This implies that the corresponding integer programming problems can have up to 94,000 binary variables. An integer program with this many binary variables is often considered too large to be solvable even on powerful computers. However, as we can see in Table 2 to 5, the integer programs for our test problems have been solved very well by using the CPLEX MIP solver.

Table 1: Randomly Generated Problems

# units	$p_1 = p_2 = 0.2$	$p_1 = p_2 = 0.3$	$p_1 = p_2 = 0.4$	$p_1 = p_2 = 0.5$
100:	$ a  = 18$ $ b  = 22$ $ c  = 35$	$ a  = 27$ $ b  = 33$ $ c  = 49$	$ a  = 36$ $ b  = 42$ $ c  = 62$	$ a  = 48$ $ b  = 55$ $ c  = 75$
200:	$ a  = 33$ $ b  = 39$ $ c  = 61$	$ a  = 45$ $ b  = 59$ $ c  = 86$	$ a  = 66$ $ b  = 77$ $ c  = 115$	$ a  = 96$ $ b  = 102$ $ c  = 146$
300:	$ a  = 55$ $ b  = 58$ $ c  = 104$	$ a  = 79$ $ b  = 89$ $ c  = 148$	$ a  = 109$ $ b  = 114$ $ c  = 181$	$ a  = 145$ $ b  = 149$ $ c  = 227$
400:	$ a  = 71$ $ b  = 84$ $ c  = 137$	$ a  = 103$ $ b  = 124$ $ c  = 191$	$ a  = 142$ $ b  = 157$ $ c  = 241$	$ a  = 192$ $ b  = 197$ $ c  = 298$
500:	$ a  = 90$ $ b  = 115$ $ c  = 189$	$ a  = 135$ $ b  = 165$ $ c  = 256$	$ a  = 184$ $ b  = 211$ $ c  = 319$	$ a  = 242$ $ b  = 250$ $ c  = 376$

Tables 2 to 5 contain the results obtained from using the CPLEX MIP solver on the test problems we have generated in Table 1. Note that for every triple sets  $a$ ,  $b$ , and  $c$ , there are two integer programs to solve, one for matching  $c$  to  $a$  and another for  $c$  to  $b$ . We denote the two programs by  $c \rightarrow a$  and  $c \rightarrow b$ . For each of these problems, we list in the tables

the number of branch-and-bound nodes used in the CPLEX MIP solver and the time to obtain a solution to the problem. Table 2 contains the results for the problems solved as  $P_1$  type programs. Table 3 contains the results for the problems solved as  $P_\infty$  type programs. These two sets of results were obtained with the given fragment sets in the original orders. However, we have also tested the problems with the given fragment sets sorted in an ascending order. The results are listed in Tables 4 and 5.

From Tables 2 to 5 we can see that the CPLEX MIP solver solved the problems very well although many of their corresponding integer programs have large numbers of integer variables. The problems were solved within reasonable amount of time except for one that took about 10 hours. Comparing the results for  $P_1$  and  $P_\infty$  types problems, the timings are comparable. The choice between solving  $P_1$  or  $P_\infty$  should depend on practical needs. On the other hand, the results for problems with or without sorting the fragment data seem quite different. Most of the problems were solved faster when the fragment sets were sorted except for a couple of cases.

In general, the more nodes are there in the branch-and-bound procedure, the longer time the program will run. However, as we can see in Tables 2 to 5, the number of branch-and-bound nodes does not always correlate with the running time. Some problems required fewer nodes but longer time. The reason is that the cost of solving linear programs corresponding to different nodes can vary greatly.

Table 2: Test Results for  $P_1$  Type Problems

# units	program	$p_1 = p_2 = 0.2$	$p_1 = p_2 = 0.3$	$p_1 = p_2 = 0.4$	$p_1 = p_2 = 0.5$
100	$c \rightarrow a$	45/1s	23/1s	21/2s	21/2s
	$c \rightarrow b$	88/1s	33/1s	38/2s	20/3s
200	$c \rightarrow a$	430/8s	206/11s	73/13s	145/28s
	$c \rightarrow b$	718/12s	762/27s	191/31s	76/24s
300	$c \rightarrow a$	529/32s	400/1m10s	463/2m5s	193/2m35s
	$c \rightarrow b$	4156/2m4s	602/1m56s	2369/7m0s	65/1m9s
400	$c \rightarrow a$	2390/3m13s	1097/6m23s	541/6m49s	160/5m7s
	$c \rightarrow b$	2669/3m36s	1228/6m24s	845/9m49s	324/9m19s
500	$c \rightarrow a$	1235/3m33s	7648/46m36s	662/15m47s	144/7m43s
	$c \rightarrow b$	6616/17m14s	18862/1h30m42s	14661/1h43m13s	1238/38m5s

Table 3: Test Results for  $P_\infty$  Type Problems

# units	program	$p_1 = p_2 = 0.2$	$p_1 = p_2 = 0.3$	$p_1 = p_2 = 0.4$	$p_1 = p_2 = 0.5$
100	$c \rightarrow a$	11/1s	43/1s	28/2s	21/2s
	$c \rightarrow b$	97/2s	50/1s	37/2s	21/3s
200	$c \rightarrow a$	226/4s	225/9s	182/37s	87/35s
	$c \rightarrow b$	4070/38s	1086/39s	148/44s	81/33s
300	$c \rightarrow a$	545/25s	177/1m20s	138/1m16s	97/1m50s
	$c \rightarrow b$	198/18s	237/2m4s	391/4m17s	78/1m16s
400	$c \rightarrow a$	2293/2m27s	269/2m42s	268/5m1s	280/18m35s
	$c \rightarrow b$	3751/4m13s	1727/24m6s	1071/19m36s	324/11m34s
500	$c \rightarrow a$	208/53s	334/10m13s	401/40m31s	144/10m46s
	$c \rightarrow b$	2193/12m10s	1520/1h24m42s	414/17m11s	821/2h44m25s

Table 4: Test Results for  $P_1$  Type Problems with Sorted Fragments

# units	program	$p_1 = p_2 = 0.2$	$p_1 = p_2 = 0.3$	$p_1 = p_2 = 0.4$	$p_1 = p_2 = 0.5$
100	$c \rightarrow a$	232/2s	52/1s	26/2s	33/3s
	$c \rightarrow b$	1198/5s	33/1s	41/3s	19/2s
200	$c \rightarrow a$	446/7s	380/15s	196/30s	42/20s
	$c \rightarrow b$	198/5s	461/21s	162/25s	21/18s
300	$c \rightarrow a$	668/37s	169/50s	83/58s	57/1m12s
	$c \rightarrow b$	1251/1m5s	1115/2m58s	104/1m5s	91/1m40
400	$c \rightarrow a$	1462316/9h59m31s	1345/6m28s	415/5m13s	94/4m14s
	$c \rightarrow b$	4000/4m57s	2341/15m48s	437/5m50s	199/7m0s
500	$c \rightarrow a$	1120/4m8s	1694/1m21s	159/5m24s	119/7m19s
	$c \rightarrow b$	992/4m10s	846/13m32s	342/10m18s	289/15m9s

## 5 Concluding Remarks

In this paper, we have considered a mixed-integer linear programming formulation for the double digestion problem and showed that with this formulation, the problem can be solved efficiently by using state-of-the-art mixed-integer programming techniques. In particular, we used the CPLEX mixed-integer linear programming software and obtained the exact solutions for a set of randomly-generated test problems.

We have focused on the problem of matching the double-digested fragments to the single-digested ones. The problem is formulated as a mixed-integer linear program minimizing the matching errors in either  $l_1$  or  $l_\infty$  norm. The program is separated into two sub-programs. We have used

Table 5: Test Results for  $P_\infty$  Type Problems with Sorted Fragments

# units	program	$p_1 = p_2 = 0.2$	$p_1 = p_2 = 0.3$	$p_1 = p_2 = 0.4$	$p_1 = p_2 = 0.5$
100	$c \rightarrow a$	25/1s	51/1s	21/2s	21/2s
	$c \rightarrow b$	102/1s	77/2s	41/2s	27/3s
200	$c \rightarrow a$	76/2s	191/8s	87/18s	54/19s
	$c \rightarrow b$	152/4s	584/29s	52/13s	50/21s
300	$c \rightarrow a$	212/12s	154/33s	83/51s	70/1m14s
	$c \rightarrow b$	453/27s	254/1m27s	107/51s	106/1m27s
400	$c \rightarrow a$	308/40s	572/5m15s	193/5m39s	69/4m8s
	$c \rightarrow b$	741/1m36s	198/2m47s	175/6m13s	133/6m42s
500	$c \rightarrow a$	222/1m16s	321/11m56s	197/13m33s	124/13m5s
	$c \rightarrow b$	1342/6m32s	261/15m42s	182/9m4s	131/17m8s

the CPLEX MIP solver to solve both  $l_1$  and  $l_\infty$  types of problems. The test problems were generated randomly to cover possible sizes of problems, lengths of fragments, and distributions of random fragments. We have tested problems with hundreds of fragments and hence tens of thousands of binary variables in the corresponding integer programs. Many of the problems were solved in only a few minutes on a single SGI R12000 processor. In terms of the sizes or the search spaces of the problems, they range from  $(18!)(20!)$  to  $(242!)(250!)$  possible orderings, which are many magnitudes larger than previously reported  $(16!)(16!)$  testing sizes.

The reason for such an advancement in solving an NP-hard problem like the double digestion problem is not just because of the improvement on the computer speed over the past ten or twenty years, but is rather because of the development of many new integer programming techniques such as the cutting plane methods and the improved linear programming solvers as implemented in modern software like CPLEX that can exploit the problem structures more effectively. Similar results were obtained for other types of NP-hard problems as well. For instance, in 1987, Padberg and Rinaldi solved a traveling salesman problem with 532 American cities [10], while in 1998, by introducing various integer programming techniques, Applegate et al. were able to solve a problem of same type with 13509 American cities [1], which was considered as a major advance in modern computer science other than just a simple celebration on the computer hardware improvement.

Several issues yet to be resolved in future work. One of them is the handling of the errors in the experimental data. If the measured lengths for the fragments are still integers, the mixed-integer programs will essentially

be at the same level of difficulty as those we have tested. If the lengths are real numbers, however, the solutions to the mixed-integer programs may take longer time to find. We plan to conduct more computational experiments to study this issue.

Another issue is the multiple solutions to the mixed-integer programs and hence the double digestion problem. As discussed in the literature, this problem cannot be solved unless additional knowledge of the sequence is given. In any case, the solutions to the two sub-programs for fragment matching show only how the single and double-digested fragments match but not how they should be ordered. The fragments need to be examined together and put into sensible orders. Fortunately, this can be done in an efficient and systematic fashion [17].

When using the CPLEX MIP solver for the test problems, we have so far not taken advantages that CPLEX provides for improving the performance of the solver, such as using user-supplied initial solutions, heuristic upper bounds, and better node selection orders, etc. We have only used default settings for our test. It is highly probable that by taking these advantages, one should be able to exploit the problem structure more effectively and achieve even better computational results than what we have just observed.

## References

- [1] D. A. Applegate, R. E. Bixby, V. Chvatal, and W. Cook, *On the Solution of Traveling Salesman Problems* Documenta Mathematica Journal der Deutschen Mathematiker-Vereinigung, International Congress of Mathematicians, 1998, 645-656.
- [2] R. E. Bixby, *Implementing the Simplex Method: The Initial Basis*, Technical Report TR90-32, Department of Computational and Applied Mathematics, Rice University, Houston, TX, 1990.
- [3] W. M. Fitch, T. F. Smith, and W. W. Ralph, *Mapping the Order of DNA Restriction Fragments*, Gene, 22, 1983, pp. 19-29.
- [4] L. Goldstein and M. S. Waterman, *Mapping DNA by Stochastic Relaxation*, Adv. Appl. Math., 8, 1987, pp. 194-207.
- [5] A. V. Grigorjev and A. A. Mironov, *Mapping DNA by Stochastic Relaxation: A New Approach to Fragment Sizes*, Applic. Biosci., , pp. 107-111.

- [6] M. Grötschel, L. Lovász, and A. Schrijver, *Geometric Algorithm and Combinatorial Optimization*, Springer, 1987.
- [7] ILOG Inc., *Using the CPLEX Callable Library*, ILOG Inc., 1997.
- [8] M. Krawczak, *Algorithms for the Restriction-Site Mapping of DNA Molecules*, Proc. Natl. Acad. Sci. USA, 85, 1988, pp. 7298-7301.
- [9] G. L. Nemhauser and L. A. Wolsey, *Integer and Combinatorial Optimization*, John Wiley & Sons, Inc., 1988.
- [10] M. Padberg and G. Rinaldi, *Optimization of a 532-City Symmetric Traveling Salesman Problem by Branch and Cut*, Operations Research Letters, 6, 1987, 1-7.
- [11] C. Pandurangan and H. Ramesh, *The Restriction Mapping Problem Revisited*, Department of Computer Science, Brown University, Providence, RI, 2002.
- [12] W. R. Pearson, *Automatic Construction of Restriction Site Maps*, Nucleic Acids Res., 10, 1984, 217-227.
- [13] W. Schmitt and M. S. Waterman, *Multiple Solutions of DNA Restriction Mapping Problems*, Adv. Appl. Math., 12, 1991, pp. 412-427.
- [14] D. P. Snustad and M. J. Simmons, *Principles of Genetics*, John Willey & Sons, Inc., 2003.
- [15] M. Stefik, *Inferring DNA Structures from Segmentation Data*, Artif. Intell., 11, 1978, pp. 85-114.
- [16] P. Tuffery, P. Dessen, C. Mugnier, and S. Hazout, *Restriction Map Construction Using Complete Sentences Compatibility Algorithm*, Comput. Applic. Biosci., 4, pp. 103-110.
- [17] M. S. Waterman, *Introduction to Computational Biology: Sequences, Maps, and Genomes*, Chapman Hall, 1995.
- [18] L. W. Wright, J. B. Lichter, J. Reinitz, M. A. Shifman, K. K. Kidd, and P. L. Miller, *Computer-Assisted Restriction Mapping: An Integrated Approach to Handling Experimental Uncertainty*, CABIOS, Vol. 10, No. 4, 1994, pp. 443-450.