

On Computing Breakpoint Distances for Genomes with Duplicate Genes

Mingfu Shao^(✉) and Bernard M.E. Moret^(✉)

Laboratory for Computational Biology and Bioinformatics,
EPFL, Lausanne, Switzerland
{mingfu.shao,bernard.moret}@epfl.ch

Abstract. A fundamental problem in comparative genomics is to compute the distance between two genomes in terms of its higher-level organization (given by genes or syntenic blocks). For two genomes without duplicate genes, we can easily define (and almost always efficiently compute) a variety of distance measures, but the problem is NP-hard under most models when genomes contain duplicate genes. To tackle duplicate genes, three formulations (exemplar, maximum matching, and any matching) have been proposed, all of which aim to build a matching between homologous genes so as to minimize some distance measure. Of the many distance measures, the breakpoint distance (the number of non-conserved adjacencies) was the first one to be studied and remains of significant interest because of its simplicity and model-free property.

The three breakpoint distance problems corresponding to the three formulations have been widely studied. Although we provided last year a solution for the exemplar problem that runs very fast on full genomes, computing optimal solutions for the other two problems has remained challenging. In this paper, we describe very fast, exact algorithms for these two problems. Our algorithms rely on a compact integer-linear program that we further simplify by developing an algorithm to remove variables, based on new results on the structure of adjacencies and matchings. Through extensive experiments using both simulations and biological datasets, we show that our algorithms run very fast (in seconds) on mammalian genomes and scale well beyond. We also apply these algorithms (as well as the classic orthology tool MSOAR) to create orthology assignment, then compare their quality in terms of both accuracy and coverage. We find that our algorithm for the “any matching” formulation significantly outperforms other methods in terms of accuracy while achieving nearly maximum coverage.

Keywords: Breakpoint distance · Exemplar · Intermediate · Maximum matching · Gene family · ILP · Orthology assignment

1 Introduction

The combinatorics and algorithmics of genomic rearrangements have been the subject of much research in comparative genomics since the problem was formulated in the 1990s (see, e.g., [1]). Perhaps the most fundamental problem is

the computation of some distance measure between two genomes. When the two genomes being compared have no duplicate genes (or synteny blocks, since the basic unit of description need not be restricted to genes), we have linear-time algorithms for most of these distance problems, such as the breakpoint distance, the inversion distance [2], and the DCJ distance [3, 4].

However, gene duplications are widespread events and have long been recognized as a major driving force of evolution [5, 6]. To define the distance in the presence of duplicate genes, three formulations have been proposed, all based on building a matching between duplicate genes and discarding copies not in the matching. This matching leads to treating each matched pair as a new gene family and thus removes duplicates, reducing the distance problem to its simplest version. In all formulations, the goal is to return that matching (from among all those obeying stated constraints) which minimizes the chosen distance. The first formulation is due to Sankoff [7], who proposed the *exemplar* model (E-model): select exactly one matched pair in each gene family. Several years later, Blin *et al.* [8] proposed the *maximum matching* model (M-model): use as many matched pairs as possible—until all genes in the genome with the smaller number of copies in each gene family are matched. Finally, Angibaud *et al.* [9] proposed what we shall term the *intermediate* model (I-model): the matching must contain at least one matched pair per gene family—so that the focus is clearly on minimizing the distance, not on meeting constraints on the matching. Figure 1 illustrates these three formulations. Unfortunately, for almost all distance measures, the corresponding distance problems under the three formulations are NP-hard [10].

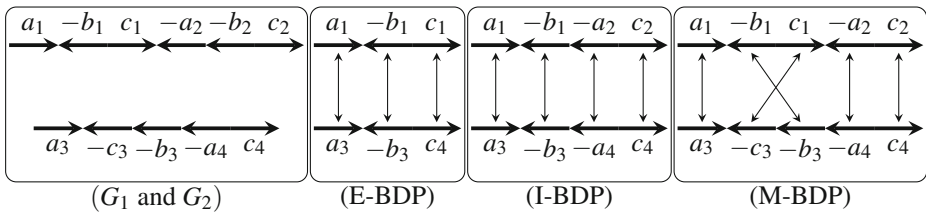


Fig. 1. Examples for E-BDP, I-BDP and M-BDP. Optimal matchings are shown within each model.

In this paper, we focus on the distance problems in the presence of duplicate genes under the breakpoint distance measure. We refer to the three corresponding problems under the three formulations as E-BDP (short for exemplar breakpoint distance problem), M-BDP, and I-BDP, respectively. Many algorithms have been proposed for these problems. For E-BDP, Sankoff gave a first branch-and-bound algorithm in his original paper [7]. Nguyen *et al.* gave a much faster divide-and-conquer approach [11], while Angibaud *et al.* [9] gave an exact algorithm by formulating it as an integer linear program (ILP). In [12], we gave an exact solution through combining ILP and novel preprocessing algorithms and

constraint-generating algorithms, which runs very fast (in a few seconds) and scales up to the largest genomes analyzed on both simulations and biological datasets. For M-BDP, Blin *et al.* [8] had defined (but not tested) a branch-and-bound algorithm, while Swenson *et al.* [13] gave an approximation algorithm. In [14], Angibaud *et al.* gave exact ILP formulations as well as several heuristics for both M-BDP and I-BDP—the heuristics being needed as the ILP formulations did not scale well. The ILP formulations in [14] were tested on 12 bacterial genomes, with numbers of genes varying from a few hundreds to three thousand; among the 66 pairwise comparisons, their ILP for I-BDP took 3 min on 63 of them, but could not terminate on the other 3 within a few hours. It thus remained a challenge to design fast and exact algorithms for I-BDP (and the simpler M-BDP) that could scale easily to mammalian genome sizes and beyond.

From the perspective of biologist, neither the E-model nor the M-model is satisfactory. In real gene families, we expect to find multiple orthologs, not just a single pair, so the E-model is an oversimplification and provides only a small amount of information, which in turn could lead to distortions in the pairwise distances. The M-model suffers from an even more fundamental problem: forcing all genes in the genome with the smaller number of copies to be matched effectively amounts to stating that every one of these genes has an ortholog in the other genome, something that clearly need not always be true. In terms of evolutionary events, this requirement also gives a much larger weight to duplications and losses of genes than to genomic rearrangements, thus implying that duplication and loss events are much less likely than rearrangement events. In other words, the M-model loses the model-free characteristic of the breakpoint distance. In contrast, the I-model retains the model-free property of the breakpoint distance while making full use of the information present in the gene families.

In this paper, we describe fast and exact algorithms for M-BDP and I-BDP. These algorithms employ ILP formulations and use a novel algorithm to reduce the number of key variables in the ILP, based on new results we prove about matchings and adjacencies. We evaluate these algorithms on simulated genomes and on several mammalian genomes. Our results show that our algorithms easily scale beyond the size of mammalian genomes: in all of our testing, almost all instances took a few seconds and none more than 70 s. They also demonstrate that our new algorithm to reduce the number of variables is a crucial contributor to this speed, especially for the more challenging I-BDP. Thus, with our previous (and equally fast) algorithm for E-BDP [12], we now have fast and exact algorithms to compute the breakpoint distance under all three formulations.

We also apply our algorithms to infer orthologs among five mammalian genomes and compare the quality of these assignments both among the three formulations and with the classic orthology tool MSOAR [15–17]. (We use MSOAR as it is one of the few orthology tools based on genome rearrangements and matching and because it has been extensively tested by its authors against other orthology tools.) The results demonstrate that the I-model substantially outperforms all other methods in terms of accuracy while giving very high coverage—close to the M-model (which defines the maximum possible number).

2 Problem Statement

We model each genome as a set of chromosomes and model each chromosome as a linear or circular list of genes. Each gene is represented by a signed (+ or -) symbol, where the sign indicates the transcriptional direction of this gene. Given a chromosome, we can reverse the list of symbols and switch all the signs, which will result in the same chromosome; for instance, $g_1g_2 \cdots g_{n-1}g_n$ and $-g_n - g_{n-1} \cdots - g_2 - g_1$ represent the same linear chromosome. Homologous genes among the given genomes are grouped into *gene families*. In this paper, we assume that the given genomes have the same set of gene families, denoted by \mathcal{F} . For a gene family $f \in \mathcal{F}$ and a genome G , we use $F(G, f)$ to denote the set of genes in G that come from f . We say a gene family $f \in \mathcal{F}$ is a *singleton* in G if $|F(G, f)| = 1$; otherwise we say f is a *multi-gene family*.

Two consecutive genes g and h on the same chromosome, with g ahead of h along the chromosome, form an *adjacency*, written as gh . Given two genomes G_1 and G_2 , we say two adjacencies $g_1h_1 \in G_1$ and $g_2h_2 \in G_2$ form a *pair of shared adjacencies* or a *PSA*, written as $\langle g_1h_1, g_2h_2 \rangle$, if g_1 and g_2 (and also h_1 and h_2) have the same sign and come from the same gene family, or g_1 and h_2 (and also h_1 and g_2) have opposite signs and come from the same gene family. If two given genomes G_1 and G_2 contain only singletons (have at most one gene each per gene family), then, for each adjacency in G_1 (resp. G_2), there exists at most one adjacency in G_2 (resp. G_1) shared with it.

Given two genomes G_1 and G_2 that may contain multi-gene families, we define a *matching* between them as a one-to-one correspondence between a subset of genes in G_1 and a subset of genes in G_2 , such that each element of the matching is a pair of homologous genes. We denote by \mathcal{M} the set of all possible matchings between G_1 and G_2 . For a matching $M \in \mathcal{M}$ and a gene family $f \in \mathcal{F}$, we use $M(f)$ to denote the set of gene pairs in M that come from f . We say a gene is *covered* by M if it appears in some pair in M . Given $M \in \mathcal{M}$, we can modify G_1 and G_2 as follows: we first remove all genes that are not covered by M , then set up a new distinct gene family for each pair of genes in M . Clearly, these two new genomes contain only singletons; we denote by $\mathcal{S}(M)$ the set of PSAs between such two new genomes induced by M .

Given two genomes, the *breakpoint distance problem* is to compute a matching (satisfying some requirements depending on certain models) such that the number of shared adjacencies between the new genomes induced by this matching is maximized. Define the following sets of matchings:

$$\begin{aligned} \mathcal{M}_e &= \{M \in \mathcal{M} : |M(f)| = 1, \forall f \in \mathcal{F}\}; \\ \mathcal{M}_i &= \{M \in \mathcal{M} : |M(f)| \geq 1, \forall f \in \mathcal{F}\}; \\ \mathcal{M}_m &= \{M \in \mathcal{M} : |M(f)| = \min\{|F(G_1, f)|, |F(G_2, f)|\}, \forall f \in \mathcal{F}\}. \end{aligned}$$

Then the three problems corresponding to the three formulations can be written as:

$$\max_{M \in \mathcal{M}_e} |\mathcal{S}(M)| \quad (\text{E-BDP})$$

$$\max_{M \in \mathcal{M}_i} |\mathcal{S}(M)| \quad (\text{I-BDP})$$

$$\max_{M \in \mathcal{M}_m} |\mathcal{S}(M)| \quad (\text{M-BDP})$$

Notice that our formulations are essentially based on *conserved adjacencies*, not breakpoints. For the E-model and the M-model, the two are equivalent; for the I-model, they are different (although clearly similar). This is due to the fact that the number of breakpoints equals the total number of adjacencies minus the number of conserved adjacencies, and that for the E-model and the M-model, the number of adjacencies is a fixed value, while for the I-model it is variant. Using conserved adjacencies rather than breakpoints has a significant advantage when the size of the matching is not fixed, as in the I-model the measure rewards both increased coverage and increased structural similarity.

3 Algorithms

We described a fast and exact algorithm for E-BDP in previous work [12]. In this section, we describe fast and exact algorithms for I-BDP and M-BDP. Both algorithms use the same framework, consisting of an integer linear program (ILP), described in Sect. 3.1, and an algorithm to reduce the number of variables in the ILP, described in Sect. 3.2. In each of them, we first describe the formulation or algorithm for I-BDP, then we state them for M-BDP, mainly focusing on clarifying their differences.

3.1 ILP Formulations

We first generalize the definition of adjacency. For two genes g and h on the same chromosome (g is ahead of h), we use $[g, h]$ to represent the genes from g to h along the chromosome (including g and h), and use (g, h) to represent the genes between g and h (excluding g and h). We say $[g, h]$ forms a *potential adjacency* for I-BDP, if we can remove all genes in (g, h) such that at least one gene remains in each gene family (as required by I-BDP). We say two potential adjacencies $[g_1, h_1] \in G_1$ and $[g_2, h_2] \in G_2$ form a *pair of shared potential adjacencies* (PSPA), written as $\langle [g_1, h_1], [g_2, h_2] \rangle$, if g_1 and g_2 (and also h_1 and h_2) have the same sign and come from the same gene family (case 1), or g_1 and h_2 (and also h_1 and g_2) have opposite signs and come from the same gene family (case 2). Without loss of generality, in the following we always assume that a PSPA belongs to case 1—all corresponding results for PSPAs belonging to case 2 can be derived directly in a symmetric way. We denote by \mathcal{P}_i the set of all

PSPAs between G_1 and G_2 for I-BDP. Given a matching $M \in \mathcal{M}_i$, we say a PSPA $p = \langle [g_1, h_1], [g_2, h_2] \rangle \in \mathcal{P}_i$ survives w.r.t. M if we have $\langle g_1 h_1, g_2 h_2 \rangle \in \mathcal{S}(M)$.

Let $M_i^* \in \mathcal{M}_i$ be an optimal matching for I-BDP. Our ILP formulation to compute M_i^* has three types of variables. First, for every gene g in the two given genomes, we use one binary variable x_g to indicate whether g is covered by M_i^* . Second, for each pair of homologous genes $g_1 \in G_1$ and $g_2 \in G_2$, we use one binary variable y_{g_1, g_2} to indicate whether pair $\langle g_1, g_2 \rangle$ is in M_i^* . Third, for every PSPA $p \in \mathcal{P}_i$, we use one binary variable z_p to indicate whether p survives w.r.t. M_i^* .

Our ILP to compute M_i^* has three types of constraints. First, we require that for each gene family in each genome, at least one gene is covered by M_i^* :

$$\begin{aligned} \sum_{g \in F(G_1, f)} x_g &\geq 1, & \forall f \in \mathcal{F}; \\ \sum_{g \in F(G_2, f)} x_g &\geq 1, & \forall f \in \mathcal{F}. \end{aligned} \tag{1}$$

Second, we use the following constraints to guarantee that gene g is covered by M_i^* if and only if there exists a pair in M_i^* that includes g :

$$\begin{aligned} \sum_{g_2 \in F(G_2, f)} y_{g_1, g_2} &= x_{g_1}, & \forall g_1 \in F(G_1, f), \forall f \in \mathcal{F}; \\ \sum_{g_1 \in F(G_1, f)} y_{g_1, g_2} &= x_{g_2}, & \forall g_2 \in F(G_2, f), \forall f \in \mathcal{F}. \end{aligned}$$

Third, for each PSPA $p = \langle [g_1, h_1], [g_2, h_2] \rangle \in \mathcal{P}_i$, we require that, if p survives w.r.t. M_i^* , then we must have $\langle g_1, g_2 \rangle \in M_i^*$, $\langle h_1, h_2 \rangle \in M_i^*$, and that all genes in $(g_1, h_1) \cup (g_2, h_2)$ cannot be covered by M_i^* :

$$\begin{aligned} y_{g_1, g_2}, y_{h_1, h_2} &\geq z_p, & \forall p = \langle [g_1, h_1], [g_2, h_2] \rangle \in \mathcal{P}_i; \\ 1 - z_p &\geq x_g, & \forall g \in (g_1, h_1) \cup (g_2, h_2). \end{aligned} \tag{2}$$

The objective of the ILP is to maximize the sum of the variables for PSPAs:

$$\max \sum_{p \in \mathcal{P}_i} z_p.$$

A few modifications to this ILP for I-BDP provide an ILP for M-BDP. We say that $[g, h]$ is a potential adjacency for M-BDP if at least $\min\{|F(G_1, f)|, |F(G_2, f)|\}$ genes remain in gene family $f \in \mathcal{F}$ after genes in (g, h) are removed. Let \mathcal{P}_m be the set of PSPAs for M-BDP and let $M_m^* \in \mathcal{P}_m$ be one optimal matching. We can construct the ILP to compute M_m^* from the above one to compute M_i^* by replacing M_i^* with M_m^* , \mathcal{P}_i with \mathcal{P}_m , and constraints (1) with the following ones:

$$\begin{aligned} \sum_{g \in F(G_1, f)} x_g &= \min\{|F(G_1, f)|, |F(G_2, f)|\}, & \forall f \in \mathcal{F}; \\ \sum_{g \in F(G_2, f)} x_g &= \min\{|F(G_1, f)|, |F(G_2, f)|\}, & \forall f \in \mathcal{F}. \end{aligned}$$

Our ILP formulations use the same ideas we used for the exact algorithm to solve E-BDP in [12]. They are similar to those proposed in [14], but have fewer variables. In fact, the variables in our formulations are a subset of those in the formulations of [14]: the authors use two additional binary variables, c_{g_1, h_1} and c_{g_2, h_2} , for each PSPA $\langle [g_1, h_1], [g_2, h_2] \rangle$ to indicate whether $g_1 h_1$ and $g_2 h_2$ form adjacencies in the resulting genomes. In our formulations, the functions of these variables are carried out by constraints (2). ILP solvers typically do better with fewer variables and do not suffer (and often benefit) from the addition of extra constraints.

3.2 Building a Sufficient Subset

The number of binary variables corresponding to the PSPAs is the most important factor affecting the efficiency of our ILP formulations. We now describe an algorithm to reduce the number of these variables while preserving the optimality of the ILP. Formally, we say $\mathcal{P} \subset \mathcal{P}_i$ is *sufficient* if there exists an optimal matching $M_i^* \in \mathcal{M}_i$ such that for any PSA $\langle g_1 h_1, g_2 h_2 \rangle \in \mathcal{S}(M_i^*)$, its corresponding PSPA $\langle [g_1, h_1], [g_2, h_2] \rangle$ is in \mathcal{P} . Clearly, if we replace \mathcal{P}_i by a sufficient subset \mathcal{P} in the ILP formulation described in Sect. 3.1, an optimal matching can be still obtained. In the following we first prove two conditions to reduce a single PSPA, then use these results to devise an iterative algorithm to compute a sufficient subset of reduced cardinality.

Intuitively, the first lemma states that we can remove a PSPA without breaking sufficiency if it can be split into two PSPAs.

Lemma 1. *Let $\mathcal{P} \subset \mathcal{P}_i$ be a sufficient subset. Let $p = \langle [g_1, h_1], [g_2, h_2] \rangle \in \mathcal{P}$. If there exists gene $x_1 \in (g_1, h_1)$ and gene $x_2 \in (g_2, h_2)$ that have the same sign and come from the same gene family, then we have that $\mathcal{P} \setminus \{p\}$ is sufficient.*

Proof. We prove this lemma by contradiction. Suppose that $\mathcal{P} \setminus \{p\}$ is not sufficient. Let M_i^* be any optimal matching. Since \mathcal{P} is sufficient but $\mathcal{P} \setminus \{p\}$ is not, p survives *w.r.t.* M_i^* , i.e., we have $\langle g_1 h_1, g_2 h_2 \rangle \in \mathcal{S}(M_i^*)$. Let $M'_i = M_i^* \cup \{\langle x_1, x_2 \rangle\}$. Clearly, we have that $\mathcal{S}(M'_i) = \mathcal{S}(M_i^*) \setminus \{\langle g_1 h_1, g_2 h_2 \rangle\} \cup \{\langle g_1 x_1, g_2 x_2 \rangle, \langle x_1 h_1, x_2 h_2 \rangle\}$. Thus we have that $|\mathcal{S}(M'_i)| = |\mathcal{S}(M_i^*)| + 1$, contradicting with the assumption that M_i^* is optimal. \square

We now give another condition to reduce the size of a sufficient subset. Intuitively, the following lemma states that we can remove a PSPA without breaking sufficiency if there exists a gene inside that can play the same role as one of its four boundary genes.

Lemma 2. *Let $\mathcal{P} \subset \mathcal{P}_i$ be a sufficient subset. Set $p = \langle [g_1, h_1], [g_2, h_2] \rangle \in \mathcal{P}$ and $X = \{\langle x_1, x_2 \rangle : \langle [h_1, x_1], [h_2, x_2] \rangle \in \mathcal{P}\}$. If there exists $h \in (g_1, h_1)$ satisfying $\langle [g_1, h], [g_2, h_2] \rangle \in \mathcal{P}$, such that for all $\langle x_1, x_2 \rangle \in X$, we have $\langle [h, x_1], [h_2, x_2] \rangle \in \mathcal{P}$, then $\mathcal{P} \setminus \{p\}$ is sufficient.*

Proof. Let M_i^* be an optimal matching. If p does not survive *w.r.t.* M_i^* , then we can conclude that $\mathcal{P} \setminus \{p\}$ is sufficient. Now consider the case where p survives

w.r.t. M_i^* , i.e., $\langle g_1h_1, g_2h_2 \rangle \in \mathcal{S}(M_i^*)$. Set $M'_i = M_i^* \setminus \{\langle h_1, h_2 \rangle\} \cup \{\langle h, h_2 \rangle\}$. Consider the difference between $\mathcal{S}(M_i^*)$ and $\mathcal{S}(M'_i)$. Clearly, at most two PSAs in $\mathcal{S}(M_i^*)$ can be affected by the removal of $\langle h_1, h_2 \rangle$. One of them is $\langle g_1h_1, g_2h_2 \rangle$, while we have $\langle g_1h, g_2h_2 \rangle \in \mathcal{S}(M'_i)$ in return. If we further have $\langle h_1x_1, h_2x_2 \rangle \in \mathcal{S}(M_i^*)$ for some $\langle x_1, x_2 \rangle$, we must also have $\langle hx_1, h_2x_2 \rangle \in \mathcal{S}(M'_i)$ in return. Thus, we have $|\mathcal{S}(M_i^*)| = |\mathcal{S}(M'_i)|$, which implies that M'_i is also optimal. On the other hand, we can show that all PSAs in $\mathcal{S}(M'_i)$ have their corresponding PSPAs in $\mathcal{P} \setminus \{p\}$. In fact, p does not survive *w.r.t.* M'_i , and we have $\langle [g_1, h], [g_2, h_2] \rangle \in \mathcal{P}$ and $\langle [h, x_1], [h_2, x_2] \rangle \in \mathcal{P}$; finally, the validity of all other PSAs in $\mathcal{S}(M'_i)$ is guaranteed by the sufficiency of \mathcal{P} . Thus, in addition to the optimality of M'_i , we have that $\mathcal{P} \setminus \{p\}$ is sufficient. \square

Given a PSPA $\langle [g_1, h_1], [g_2, h_2] \rangle$, Lemma 2 only proves the condition related to h_1 . We can derive the other three conditions for g_1 , g_2 and h_2 analogously. If a PSPA passes any of these four conditions, it can be removed while keeping sufficiency.

Based on the above two lemmas, our algorithm to build a sufficient subset (BSS) of reduced cardinality proceeds as follows. The algorithm initializes the current set of PSPAs \mathcal{P} as \mathcal{P}_i , and then it iteratively removes redundant PSPAs in \mathcal{P} . In each iteration, the algorithm examines each PSPA $p \in \mathcal{P}$ using the above two lemmas: if p passes either of the tests, then we update \mathcal{P} as $\mathcal{P} \setminus \{p\}$ and starts a new iteration. The algorithm terminates if no PSPA passes any tests in an iteration.

We give two examples to illustrate the effect of the BSS algorithm. First, suppose that two genomes contain a pair of nontrivial ($n > 1$) shared segments $\langle g^1g^2 \cdots g^n, h^1h^2 \cdots h^n \rangle$, i.e., g^k and h^k have the same sign and come from the same gene family, for $1 \leq k \leq n$. Then there are $n \cdot (n - 1)/2$ PSPAs generated from this shared segment, namely, $\langle [g^i, g^j], [h^i, h^j] \rangle$ for all $1 \leq i < j \leq n$. After running BSS on these PSPAs, we can get a sufficient subset with only $(n - 1)$ PSPAs, namely, $\langle [g^k, g^{k+1}], [h^k, h^{k+1}] \rangle$, $1 \leq k \leq n - 1$. The removal of PSPAs is due here to Lemma 1—Lemma 1 is very effective when genomes contain duplicate segments. Secondly, suppose that G_1 contains a segment $a_1b_1c_1$ and G_2 contains a segment $a_2b_2b_3 \cdots b_nc_2$. Here $2 \cdot (n - 1)$ PSPAs can be generated from these two segments, namely, $\langle [a_1, b_1], [a_2, b_k] \rangle$ and $\langle [b_1, c_1], [b_k, c_2] \rangle$, $2 \leq k \leq n$. After applying BSS on these PSPAs, we can see that a sufficient subset with at most 4 PSPAs is returned, namely, $\{\langle [a_1, b_1], [a_2, b_2] \rangle, \langle [b_1, c_1], [b_n, c_2] \rangle\} \cup \{\langle [a_1, b_1], [a_2, b_k] \rangle, \langle [b_1, c_1], [b_k, c_2] \rangle\}$ for some k , $2 \leq k \leq n$, where the value of k depends on the order in which these PSPAs are tested in BSS. The removal of PSPAs here is due to Lemma 2—Lemma 2 is very effective when genomes contain locally duplicate genes.

Since M-BDP requires keeping the maximum number of pairs, we have that a PSPA $\langle [g_1, h_1], [g_2, h_2] \rangle \in \mathcal{P}_m$ cannot have a pair of genes $x_1 \in (g_1, h_1)$ and $x_2 \in (g_2, h_2)$ within the same gene family. In other words, the PSPAs that are tested for removal in Lemma 1 do not appear in \mathcal{P}_m by definition. In fact, it is clear that \mathcal{P}_m is a (usually very small) subset of \mathcal{P}_i , making M-BDP significantly easier to solve than I-BDP (as is also apparent from the results of Sects. 4 and 5).

Lemma 2 holds unchanged for M-BDP. Thus the BSS algorithm for M-BDP uses only Lemma 2.

4 Simulation Results

We refer to the full algorithms (ILP plus BSS) for I-BDP and M-BDP as I-A1 and M-A1, respectively, and refer to the baseline algorithms (just the ILP) for I-BDP and M-BDP as I-A0 and M-A0, respectively. We use simulated data to evaluate these four algorithms in order to illustrate both the performance of the full algorithms and the effectiveness of the BSS reduction algorithm. We do not explicitly compare with the algorithms proposed in [14] for I-BDP and M-BDP, because their performance is similar to and bounded by that of I-A0 and M-A0 respectively, as discussed in Sect. 3.1.

We simulate a pair of genomes as follows. We start from a genome with only one linear chromosome consisting of N singletons. We then perform S_1 segmental duplications to make some multi-gene families, which then forms the ancestor genome. A segmental duplication randomly chooses a segment of length L and inserts its copy to another random position. The two extant genomes then speciate independently from this ancestor genome. After that on each branch it happens randomly mixed I inversions and S_2 segmental duplications. An inversion randomly chooses two positions in the genome and then reverses the segment in between. We make sure that the expected number of genes per gene family in each extant genome is 2 (the average copy number of each gene family in human, mouse and rat genomes, are 1.46, 1.55 and 1.28, respectively; thus in terms of duplicated genes, our simulated genomes are more complicated than typical mammalian genomes). Therefore, we have that $(S_1 + S_2) \cdot L = N$. Let $r_1 = S_1/(S_1 + S_2)$ be the percentage of segmental duplications before speciation for each genome. Let $r_2 = I/(I + S_2)$ be the percentage of inversions after speciation. We can calculate that $S_1 = r_1 \cdot N/L$, $S_2 = (1 - r_1) \cdot N/L$, and $I = r_2 \cdot (1 - r_1) \cdot N/(L \cdot (1 - r_2))$. Thus, a simulation configuration is determined by parameters (N, L, r_1, r_2) . For each parameter combination, we randomly simulate 10 independent instances and run the four algorithms (I-A1, I-A0, M-A1 and M-A0) to calculate the average running time over these 10 instances. Since ILP might take very long time, we give a time limit of 30 min for each instance—an algorithm will be terminated when it exceeds such time limit.

In the following experiments, all our ILP instances are solved with GUROBI [18]. We first test parameters $(N = 10000, L, r_1 = 0, r_2)$, where $r_2 \in \{0.0, 0.1, \dots, 0.9\}$ and $L \in \{1, 5, 10\}$. In this setting, all segmental duplications appear after speciation ($r_1 = 0$), and the expected number of genes in each extant genome is 20000. The results on these parameters are shown in Table 1. First, we can observe that as L increases, all algorithms take more time, indicating that the simulated instances become harder with larger L . This is because longer shared segments create drastically more PSPAs (in order of $O(n^2)$, where n is the length of this shared segments). Second, all the four algorithms take less time as r_2 increases. This is because larger r_2 means more inversions after

Table 1. Comparison of the four algorithms on parameters ($N = 10000, L, r_1 = 0, r_2$). For each combination, if all 10 instances finish in 30 min, the average running time (in seconds) is recorded; otherwise, the number of finished instances is recorded in parentheses. All four programs were run on an 8-core (2.1GHz) PC with 16GB memory.

r_2	$L = 1$				$L = 5$				$L = 10$			
	I-A1	I-A0	M-A1	M-A0	I-A1	I-A0	M-A1	M-A0	I-A1	I-A0	M-A1	M-A0
0.0	1	2	0	0	25	(0)	8	8	55	(0)	10	12
0.1	1	1	0	0	21	(0)	7	9	43	(0)	10	13
0.2	0	0	0	0	14	(0)	7	7	33	(0)	9	11
0.3	0	0	0	0	11	(0)	5	4	26	(0)	8	8
0.4	0	0	0	0	6	554	3	3	22	(0)	7	8
0.5	0	0	0	0	4	123	2	2	15	(0)	6	5
0.6	0	0	0	0	2	38	1	1	10	(0)	4	4
0.7	0	0	0	0	1	3	1	0	7	430	3	2
0.8	0	0	0	0	0	0	0	0	4	70	1	1
0.9	0	0	0	0	0	0	0	0	0	0	0	0

speciation, which can destroy existing PSPAs. Third, we can see I-A1 can finish in less than 1 min for all parameters, while I-A0 will exceed time limit for large L and small r_2 —this proves that the BSS algorithm is very crucial for I-BDP. Fourth, we can see both M-A1 and M-A0 can finish in a very short time for all parameters, indicating that the ILP formulation for M-BDP is already very efficient. This is because the definition of M-BDP determines that the number of PSPAs is reasonably small, as we analyzed theoretically in Sect. 3.2. For the same reason, we can also observe that M-BDP is easier to solve than I-BDP.

We then test parameters ($N = 20000, L, r_1 = 0.5, r_2$), where $r_2 \in \{0.0, 0.1, \dots, 0.9\}$ and $L \in \{1, 5, 10\}$. In this setting, the expected number of genes in the ancestor genome is 30000, and the expected number of genes in each extant genome is 40000. The results are shown in Table 2. Again, we can observe that the instances become harder to solve as L increases and r_2 decreases. Also observe that M-A1 and M-A0 takes similar (and very small) amount of time, since the ILP formulation for M-BDP is already very efficient. We emphasize that I-A0 can only finish within the time limit for small L and large r_2 , while I-A1 can get the optimal solution for all parameters very fast, showing that the BSS algorithm plays a key role in producing a fast algorithm for I-BDP. Notice that in these simulations, the size of extant genomes exceeds that of typical mammalian genomes and our full algorithms (I-A1 and M-A1) return optimal solutions in a very short time.

Table 2. Comparison of the four algorithms on parameters ($N = 20000, L, r_1 = 0.5, r_2$). The setup is the same as in Table 1.

r_2	$L = 1$				$L = 5$				$L = 10$			
	I-A1	I-A0	M-A1	M-A0	I-A1	I-A0	M-A1	M-A0	I-A1	I-A0	M-A1	M-A0
0.0	6	(0)	2	2	48	(0)	22	25	104	(0)	44	36
0.1	8	(0)	2	2	45	(0)	26	22	93	(0)	39	39
0.2	5	521	1	1	42	(0)	22	24	89	(0)	38	37
0.3	6	200	1	1	37	(0)	20	19	83	(0)	36	32
0.4	8	28	1	1	31	(0)	18	17	65	(0)	31	29
0.5	4	3	2	1	29	(0)	16	17	57	(0)	35	31
0.6	2	1	1	1	25	(0)	13	14	52	(0)	29	26
0.7	1	1	1	0	18	(0)	11	15	42	(0)	27	22
0.8	1	1	1	0	8	336	5	7	27	(0)	18	15
0.9	1	1	2	0	2	3	1	1	11	437	8	5

5 Biological Results

We study five well-annotated genomes, human (*H.s.*), gorilla (*G.g.*), orangutan (*P.a.*), mouse (*M.m.*), and rat (*R.n.*). For each species, we collect all the protein-coding genes and download their positions on the chromosomes and their Ensembl gene family names from Ensembl (<http://www.ensembl.org>). Genes are grouped into the same gene family if they have the same Ensembl gene family name. For each species, we merge each group of tandemly arrayed genes into a single gene by keeping only the first gene and discarding the following ones in the group.

We perform pairwise comparisons among these five species. For each pair of species, we compare the running time for the four algorithms (I-A1, I-A0, M-A1 and M-A0). We also run another two algorithms for comparison, namely, the exact algorithm to solve E-BDP described in [12] (referred to as E-A1), and MSOAR [16], which uses heuristics to compute a matching such that the inversion distance induced by this matching is minimized. The results are shown in Table 3. Note that I-A0 cannot finish within the time limit for any pair, while I-A1 can finish in a very short time (less than 30s except one taking 72s) for all pairs, once again indicating that the BSS algorithm is indispensable for solving I-BDP. The BSS algorithm improves the solution for M-BDP as well, allowing M-A1 to finish within 2s for all pairs. (MSOAR takes quite long time for these pairs, ranging from half an hour to a few hours.) Thus we now have exact and very fast algorithms (E-A1, I-A1 and M-A1) for all three formulations.

All of these algorithms give a matching between the homologous genes for each pair of species, a matching that defines a subset of orthologs under a parsimonious evolutionary assumption. We thus apply our full algorithms (E-A1, I-A1 and M-A1) together with MSOAR to infer orthologs among these 5 species.

Table 3. The running time (in seconds) of the algorithms on comparing five genomes. I-A0 cannot finish in 30 min for any species pair.

Species pairs	I-A1	I-A0	M-A1	M-A0	E-A1	MSOAR
<i>G.g.&H.s.</i>	8	N/A	1	1	2	1770
<i>G.g.&M.m.</i>	20	N/A	1	59	3	5298
<i>G.g.&P.a.</i>	16	N/A	1	1	2	3191
<i>G.g.&R.n.</i>	22	N/A	1	4	3	12555
<i>H.s.&M.m.</i>	23	N/A	1	2	2	3660
<i>H.s.&P.a.</i>	9	N/A	1	2	2	1585
<i>H.s.&R.n.</i>	22	N/A	1	6	2	7328
<i>M.m.&P.a.</i>	18	N/A	1	1	2	4287
<i>M.m.&R.n.</i>	72	N/A	2	66	2	5627
<i>R.n. &P.a.</i>	18	N/A	1	2	3	6009

The quality of the inferred orthologs is evaluated using two measures, the coverage (the number of orthologous pairs identified) and the accuracy. To compute the accuracy of a matching, we use the gene symbols (HGNC symbols for primate genes, MGI symbols for mouse genes, and RGD symbols for rat genes, downloaded from Ensembl): those gene pairs that have the same gene symbol form the set of *true orthology pairs*. We say a pair in a matching is *trivial* if it consists of two singletons. Notice that by definition all trivial pairs must appear in the matchings returned by these four algorithms. We thus exclude trivial pairs for comparison. Among the non-trivial pairs in a matching, we say a pair is *assessable* if at least one gene in this pair is covered by some true orthology pair. Then the *accuracy* of a matching is defined as the ratio between the number of non-trivial true orthology pairs in this matching and the number of (non-trivial) assessable pairs in this matching.

The quality of the orthologs inferred by these four algorithms is shown in Table 4. First, we can observe that the coverage of I-A1 (by definition, it is between E-A1 and M-A1) is much closer to M-A1 (much higher than E-A1). Second, notice that the accuracy of I-A1 significantly outperforms the other three algorithms—it is 1.08% higher than E-A1, 0.89% higher than M-A1, and 1% higher than MSOAR, on average over the 10 pairs. Third, the quality of MSOAR is very close to M-A1, in terms of both the coverage and accuracy (the accuracy of M-A1 is 0.11% higher than that of MSOAR on average). Thus, we believe that I-A1 is an excellent choice for inferring orthologs, which outperforms in terms of both coverage and accuracy.

We do not evaluate the accuracy of these algorithms on simulation data. The reason is that orthology assignment requires a biologically credible model for generating simulation data—a model that need combine duplications, losses, rearrangements, and sequence mutations and indels. Such a model that takes into account all these evolutionary events in a biologically reasonable way is

Table 4. The performance of the algorithms on inferring orthologs among the five genomes.

Species pairs	Trivial	Non-trivial				Accuracy (%)			
		E-A1	I-A1	M-A1	MSOAR	E-A1	I-A1	M-A1	MSOAR
<i>G.g. & H.s.</i>	8331	3448	7830	8131	8051	97.51	98.18	97.56	97.60
<i>G.g. & M.m.</i>	7304	3478	7572	8025	7858	97.33	98.29	97.61	97.48
<i>G.g. & P.a.</i>	7737	3399	7466	7893	7720	97.26	98.10	97.37	97.17
<i>G.g. & R.n.</i>	6915	3787	7826	8610	8317	94.65	96.12	94.56	94.43
<i>H.s. & M.m.</i>	7932	3250	7546	7834	7722	97.89	98.71	98.16	98.20
<i>H.s. & P.a.</i>	8091	3223	7355	7585	7501	98.27	98.73	98.03	98.05
<i>H.s. & R.n.</i>	7436	3638	7808	8196	8072	94.64	96.10	94.94	94.93
<i>M.m. & P.a.</i>	7311	3276	7208	7537	7408	97.91	98.61	97.95	97.81
<i>M.m. & R.n.</i>	7953	3706	8376	8837	8671	94.82	96.52	95.63	95.20
<i>R.n. & P.a.</i>	6911	3610	7466	7987	7789	94.90	96.63	95.27	95.08

currently not available. Besides, we have high-quality fully assembled genomes with curated annotations to assess these methods, which leaves little necessity to perform comparison on simulation data.

6 Conclusion and Discussion

We have described exact and very fast algorithms for the three breakpoint distance problems, by formulating them as integer linear programs and designing additional algorithms to improve their efficiency through proving key properties of these problems. Using extensive experiments on both simulations and biological datasets, we have demonstrated that these algorithms scale beyond the size of mammalian genomes, and also achieve very high accuracy when applied to infer orthologs. We conclude that among those orthology assignment tools, I-A1 is the most suitable choice, since it gives highest accuracy and nearly highest coverage.

We also help understanding the structures of these problems through proving several properties. As we have already illustrated, these properties are crucial in designing efficient algorithms for the corresponding problems. Notice that some properties are common among the three problems, while some others only hold for one or two of them. In Sect. 3.2, we give some theoretical analysis between I-BDP and M-BDP through showing that Lemma 1 only applies for I-MDP, while Lemma 2 holds for both of them. Now we further state the relationship between E-BDP and I-BDP/M-BDP. We can easily prove that both lemmas in this paper also apply to E-BDP. On the other side, in [12], we proved a lemma for E-BDP stating that if one pair in a PSA is in an optimal matching, then the other pair can be also fixed optimally. However, this property does not hold for I-BDP and M-BDP. To see that, consider the example of $G_1 = a_1b_1c_1b_2d_1$

and $G_2 = a_2b_3d_2b_4c_2$. Notice that $\langle a_1, a_2 \rangle$ is a singleton pair and thus must be in any optimal matching. If we apply this property, then $\langle b_1, b_3 \rangle$ is also in some optimal matching. However, we can easily verify that for both I-BDP and M-BDP, $\langle b_1, b_3 \rangle$ is not in any optimal matching.

Acknowledgements. We thank Daniel Dörr for helpful discussions.

References

1. Fertin, G., Labarre, A., Rusu, I., Tannier, E., Vialette, S.: *Combinatorics of Genome Rearrangements*. MIT Press, Cambridge (2009)
2. Bader, D.A., Moret, B.M.E., Yan, M.: A fast linear-time algorithm for inversion distance with an experimental comparison. *J. Comput. Biol.* **8**(5), 483–491 (2001)
3. Yancopoulos, S., Attie, O., Friedberg, R.: Efficient sorting of genomic permutations by translocation, inversion and block interchange. *Bioinformatics* **21**(16), 3340–3346 (2005)
4. Bergeron, A., Mixtacki, J., Stoye, J.: A unifying view of genome rearrangements. In: Bücher, P., Moret, B.M.E. (eds.) *WABI 2006*. LNCS (LNBI), vol. 4175, pp. 163–173. Springer, Heidelberg (2006)
5. Bailey, J.A., Eichler, E.E.: Primate segmental duplications: crucibles of evolution, diversity and disease. *Nat. Rev. Genet.* **7**(7), 552–564 (2006)
6. Lynch, M.: *The Origins of Genome Architecture*. Sinauer, Sunderland (2007)
7. Sankoff, D.: Genome rearrangement with gene families. *Bioinformatics* **15**(11), 909–917 (1999)
8. Blin, G., Chauve, C., Fertin, G.: The breakpoint distance for signed sequences. In: *Proceedings of the 1st Conference on Algorithms and Computational Methods for Biochemical and Evolutionary Networks (CompBioNets 2004)*, vol. 3, pp. 3–16 (2004)
9. Angibaud, S., Fertin, G., Rusu, I., Thévenin, A., Vialette, S.: A pseudo-boolean programming approach for computing the breakpoint distance between two genomes with duplicate genes. In: Tesler, G., Durand, D. (eds.) *RECMOB-CG 2007*. LNCS (LNBI), vol. 4751, pp. 16–29. Springer, Heidelberg (2007)
10. Blin, G., Chauve, C., Fertin, G., Rizzi, R., Vialette, S.: Comparing genomes with duplications: a computational complexity point of view. *ACM/IEEE Trans. Comput. Bio. Bioinf.* **14**, 523–534 (2007)
11. Nguyen, C.T., Tay, Y.C., Zhang, L.: Divide-and-conquer approach for the exemplar breakpoint distance. *Bioinformatics* **21**(10), 2171–2176 (2005)
12. Shao, M., Moret, B.M.E.: A fast and exact algorithm for the exemplar breakpoint distance. In: Przytycka, T.M. (ed.) *RECOMB 2015*. LNCS, vol. 9029, pp. 309–322. Springer, Heidelberg (2015)
13. Swenson, K.M., Marron, M., Earnest-DeYoung, J.V., Moret, B.M.E.: Approximating the true evolutionary distance between genomes. In: *Proceedings of the 7th SIAM Workshop on Algorithm Engineering and Experiments (ALENEX 2005)*, pp. 121–129. SIAM Press (2005)
14. Angibaud, S., Fertin, G., Rusu, I., Thévenin, A., Vialette, S.: Efficient tools for computing the number of breakpoints and the number of adjacencies between two genomes with duplicate genes. *J. Comput. Biol.* **15**(8), 1093–1115 (2008)

15. Chen, X., Zheng, J., Fu, Z., Nan, P., Zhong, Y., Lonardi, S., Jiang, T.: Assignment of orthologous genes via genome rearrangement. *ACM/IEEE Trans. Comput. Bio. Bioinf.* **2**(4), 302–315 (2005)
16. Fu, Z., Chen, X., Vacic, V., Nan, P., Zhong, Y., Jiang, T.: MSOAR: a high-throughput ortholog assignment system based on genome rearrangement. *J. Comput. Biol.* **14**(9), 1160–1175 (2007)
17. Shi, G., Zhang, L., Jiang, T.: MSOAR 2.0: incorporating tandem duplications into ortholog assignment based on genome rearrangement. *BMC Bioinform.* **11**(1), 10 (2010)
18. Gurobi Optimization Inc.: Gurobi optimizer reference manual (2013)