

Genome Rearrangements with Partially Ordered Chromosomes

Chunfang Zheng¹ and David Sankoff²

¹ Department of Biology
University of Ottawa, Canada K1N 6N5
czhen033@uottawa.ca

² Department of Mathematics and Statistics
University of Ottawa, Canada K1N 6N5
sankoff@uottawa.ca

Abstract. Genomic maps often do not specify the order within some groups of two or more markers. The synthesis of a master map from several sources introduces additional order ambiguity due to markers missing from some sources. We represent each chromosome as a partial order, summarized by a directed acyclic graph (DAG), to account for poor resolution and of missing data. The genome rearrangement problem is then to infer a minimum number of translocations and reversals for transforming a set of linearizations, one for each chromosomal DAG in the genome of one species, to linearizations of the DAGs of another species. We augment each DAG to a directed graph (DG) in which all possible linearizations are embedded. The chromosomal DGs representing two genomes are combined to produce a single bicoloured graph. From this we extract a maximal decomposition into alternating coloured cycles, determining an optimal sequence of rearrangements. We test this approach on simulated partially ordered genomes.

1 Introduction

1.1 Formalizing Genomes and Their Evolutionary Mechanisms

The genetic structure of a genome can be modeled as a set $\chi = \{\chi_1, \dots, \chi_k\}$ of $k \geq 1$ chromosomes, where each chromosome χ_i consists of $n_i > 0$ genes or other genetic markers, signed and totally ordered: e.g., $g_1 < \dots < g_{n_i}$, also written simply as $g_1 \cdots g_{n_i}$, where each g is a positive or negative integer, and where each g appears only once in all of χ . Without loss of generality, we may assume each integer between 1 and $n = n_1 + \dots + n_k$ appears exactly once in χ , either with a plus or minus sign. n is the size of the genome. An example of a genome of size 8 is thus $\{\chi_1, \chi_2, \chi_3\}$, where $\chi_1 = 5 - 1 - 4$, $\chi_2 = 7 8 3$, and $\chi_3 = -2 6$. Henceforward we will use the term gene to refer to either genes or genetic markers of any kind. The order relation abstracts the position of the gene on the linear chromosome, and the sign carries information about which of the two DNA strands the gene is located.

Genome evolution can be modeled by the transformation of a genome χ of size n to a genome ψ of the same size, by means of reversals within a chromosome and translocations between chromosomes. A reversal transforms any contiguous part of an order to its reverse order, changing the polarity of all genes in its scope, e.g., $g h i j k \rightarrow g - j - i - h k$ is a reversal of the “segment” $h i j$. A (reciprocal) translocation exchanges any prefixes of two chromosomes (or equivalently, any suffixes of two chromosomes) or the reversed prefix of one with the reversed suffix of the other, for example $g h i, x y z \rightarrow g h y z, x i$; $g h i, x y z \rightarrow g - x, -i - h y z$. There are two special cases: $g h i, x y z \rightarrow g h i x y z$; $g h i x y z \rightarrow g h i, x y z$ where a null prefix of one chromosome $x y z$ is exchanged with the largest prefix of the other $g h i$ (chromosome fusion) and where a null chromosome translocates with $g h i x y z$ (chromosome fission), respectively. To make biological sense, since a chromosome does not change its nature when it is moved around in space, a reversal of an entire chromosome is considered to leave it unchanged: $g_1 \cdots g_{n_i} = -g_{n_i} \cdots -g_1$.

The Hannenhalli-Pevzner algorithms [2, 3] for comparing genomes, and their improvements (e.g., [1], [4]), infer $d(\chi, \psi)$, the smallest number of reversals and translocations necessary to transform genome χ into ψ , as well as a sequence of such operations that actually achieves this minimum.

1.2 Partially Ordered Genomes

The representation of a genome as a set of totally ordered chromosomes must often be weakened in the case of real data, where mapping information only suffices to partially order the set of genes on a chromosome. The concepts and methods of genome rearrangement, however, pertain only to totally ordered sets of genes or markers, and are meaningless in the context of partial orders.

Our approach is to extend genome rearrangement theory to the more general context where all the chromosomes are general DAGs rather than total orders [5]. The use of DAGs reflects uncertainty of the gene order on chromosomes in the genomes of most advanced organisms. This may be due to lack of resolution, where several genes are mapped to the same chromosomal position, to missing data from some of the datasets used to compile a gene order, and/or to conflicts between these datasets.

We construct the chromosomal DAGs for each species from two or more incomplete data sets, or from a single low-resolution data set. The frequent lack of order information in each data set, due to missing genes or missing order information, is converted into parallel subpaths within each chromosomal DAG in a straightforward manner.

Outright conflicts of order create cycles that must be broken to preserve a DAG structure. We suggest a number of reasonable alternative conventions for breaking cycles. This is not the focus of our analysis, however; whatever convention is adopted does not affect our subsequent analysis.

The rearrangement problem is then TO INFER A TRANSFORMATION SEQUENCE (TRANSLOCATIONS AND/OR REVERSALS) FOR TRANSFORMING A SET OF LINEARIZATIONS (TOPOLOGICAL SORTS), ONE FOR EACH CHROMOSOMAL

DAG IN THE GENOME OF ONE SPECIES, TO A SET OF LINEARIZATIONS OF THE CHROMOSOMAL DAGS IN THE GENOME OF ANOTHER SPECIES, MINIMIZING THE NUMBER OF TRANSLOCATIONS AND REVERSALS REQUIRED. To do this, we embed the set of all possible linearizations in each DAG by appropriately augmenting the edge set, so that it becomes a general directed graph (DG). We combine the two sets of chromosomal DGs representing two genomes to produce a single large bicoloured graph from which we extract a maximal decomposition into alternating coloured cycles, so that a Hannenhalli-Pevzner type of procedure can then generate an optimal sequence of rearrangements. We focus here on obtaining the cycle decomposition; this is equivalent to optimally linearizing the partial orders, so that finding the rearrangements themselves can be done using the previously available algorithms.

2 Gene Order Data

2.1 The Methodological Origins of Incomplete Maps

Maps of genes or other markers produced by recombination analysis, physical imaging and other methods, no matter how highly resolved, inevitably are missing some (and usually most) genes or markers and fail to order some pairs of neighbouring genes with respect to each other. Even at the ultimate level of resolution, that of genome sequences, the application of different gene-finding protocols usually gives maps with different gene content.

Moreover, experimental methodologies and statistical mapping procedures inevitably give rise to some small proportion of errors, two neighbouring genes incorrectly ordered, a gene mapped to the wrong chromosome, a gene incorrectly named or annotated. However it is not these errors we focus on in this paper, but the more widespread issues of lack of resolution and genes missing from a map. These should not be considered errors; they are normal and inherent in all ways of constructing of a map except for highly polished genome sequencing with accurate gene identification (something that has not yet been achieved in the higher eukaryotes, even for humans).

2.2 Simulating Incomplete Maps of Pairs of Two Related Genomes

How incomplete maps arise may perhaps be best understood through a description of how we simulate them.

Simulating the Genomes. For a given n , we pick a small integer k , as well as positive n_1, \dots, n_k with the constraint $n = n_1 + \dots + n_k$. Then we define

$$\chi = \{m_1 \cdots n_1, m_2 \cdots n_2, \dots, m_k \cdots n_k\}, \quad (1)$$

where $m_1 = 1, n_k = n$ and the remaining $m_i = n_{i-1} + 1$. It is well known that the genes in two genomes being compared through translocation and reversal

3 Constructing the Chromosomal DAGs

A linear map of a chromosome that has several genes or markers at the same position π , because their order has not been resolved, can be reformulated as a partial order, where all the genes before π are ordered before all the genes at π and all the genes at π are ordered before all the genes following π , but the genes at π are not ordered amongst themselves. We call this procedure **make_po**.

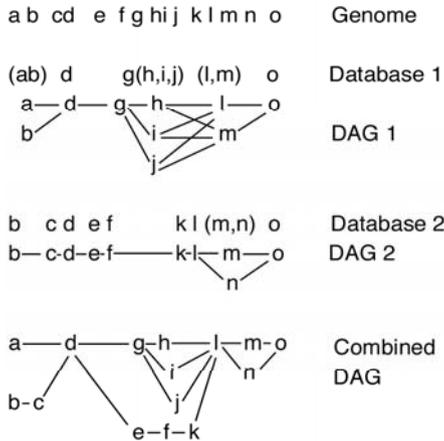


Fig. 1. Construction of DAGs from individual databases each containing partial information on genome, due to missing genes and missing order information, followed by construction of combined DAG representing all known information on the genome. All edges directed from left to right.

For genomes with two or more gene maps constructed from different kinds of data or using different methodologies, there is only one meaningful way of combining the order information on two (partially ordered) maps of the same chromosome containing different subsets of genes. Assuming there are no conflicting order relations ($a < b, b < a$) nor conflicting assignments of genes to chromosomes among the data sets (as in the data sets on our simulated genomes), for each chromosome we simply take the union of the partial orders, and extend this set through transitivity. This procedure is **combine_po**.

All the partial order data on a chromosome can be represented in a minimal DAG whose vertex set is the union of all gene sets on that chromosome in the contributing data sets, and whose edges correspond to just those order relations that cannot be derived from other order relations by transitivity. The outcome of this construction, **dagger**, is illustrated in Figure 1.

In real applications, different maps of the same genome do occasionally conflict, either because $b < a$ in one data set while $a < b$ in the other or because a gene is assigned to different chromosomes in the two data sets. There are a variety of possible ways of resolving order conflicts or, equivalently, of avoiding

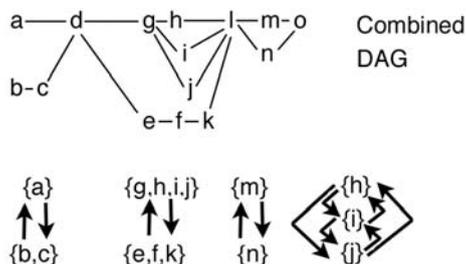


Fig. 2. Edges added to DAG to obtain DG containing all linearization as paths (though not all paths in the DG are linearizations of the DAG!). Each arrow represents a set of directed edges, one from each element in one set to each element of the other set.

any cycles in the construction of the DAG. One way is to delete all order relations that conflict with at least one other order relation. Another is to delete a minimal set of order relations so that all conflicts can be resolved. Still another is to ignore a minimum set of genes that will accomplish the same end. The latter method also resolves conflicts due to gene assignment to different chromosomes. Any of these approaches, or others, which we denote by the generic routine name **resolve**, will produce results appropriate for our subsequent analysis.

4 The DG Embedding of Topological Sorts

A DAG can generally be linearized in many different ways, all derivable from a topological sorting routine. All the possible adjacencies in these linear sorts can be represented by the edges of a directed graph (DG) containing all the edges of the DAG plus two edges of opposite directions between all pairs of vertices, which are not ordered by the DAG. This is illustrated in Figure 2. The routine for constructing this graph is **dgger**.

5 The Algorithm

5.1 Background

Before discussing our algorithm for comparing DGs derived from our DAG representations, we review the existing technology for the special case when the DAG and the associated DG represent a total order, which is the traditional subject of computational comparative genomics.

Hannenhalli and Pevzner [2] showed how to find a shortest sequence of reversals and translocations that transform one genome χ with n genes on k chromosomes into another genome ψ of the same size but with h chromosomes, in polynomial time. As described in [4], this construction begins by combining the signed order representations of all the chromosomes in the two genomes. The following procedure, **make-bicoloured**, produces a bicoloured graph on $2n + 2k$

vertices that decomposes uniquely into a set of alternating-coloured cycles and set of $h + k$ alternating-colour paths. First, each gene or marker x in χ determines two vertices, x_t and x_h , to which two additional dummy vertices e_{i_1} and e_{i_2} are added to the ends of each chromosome χ_i . One colour edge, say red, is determined by the adjacencies in χ . If x is the left-hand neighbour of y in χ , and both have positive polarity, then x_h is connected by a red edge to y_t . If they both have negative polarity, it is x_t that is joined to y_h . If x is positive and y negative, or x is negative and y positive, x_h is joined to y_h , or x_t is joined to y_t , respectively. If x is the first gene in χ_i , then e_{i_1} is joined to x_t or x_h depending on whether x has positive or negative polarity, respectively. If x is the last gene, then e_{i_2} is joined to x_t or x_h depending on whether x is negative or positive.

Black edges are added according to the same rules, based on the adjacencies in genome ψ , though no dummy vertices are added in this genome.

It can be seen that each vertex is incident to exactly one red edge and one black edge, except for the dummy vertices in χ , which are each incident to only a red edge, plus the two (non-dummy) vertices at the ends of each chromosome in ψ , which are also each incident only to a red edge. The bicoloured graph decomposes uniquely into a number of alternating cycles plus $h + k$ alternating paths terminating in either the dummy vertices of χ or the end vertices of ψ , or one of each. Suppose the number of these paths that terminate in at least one dummy vertex (good paths) is $j \leq h + k$. If the number of cycles is c , then the minimum number of reversals r and translocations t necessary to convert χ into ψ is given by the Hannenhalli- Pevzner equation:

$$r + t = n - j - c + \theta \quad (4)$$

where θ is a correction term that is usually zero for simulated or empirical data. For simplicity of exposition, we ignore this correction here, though an eventual full-scale program will incorporate it with little or no computational cost.

5.2 Generalization to Partial Orders

The routine **make_bicoloured** can also be applied to the set of edges in the DGs for two partially ordered genomes. In the resulting graph, each of the DAG edges and both of the edges connecting each of the unordered pairs in the DG for each chromosome represent potential adjacencies in our eventual linearization of a genome. The n genes or markers and $2k$ dummies determine $2n + 2k$ vertices and the potential adjacencies determine the red and black edges, based on the polarity of the genes or markers. Where the construction for the totally ordered genomes contains exactly $n + k$ red edges and $n - h$ black edges, in our construction in the presence of uncertainty there are more potential edges of each colour, but only $2n + k - h$ can be chosen in our construction of the cycle graph, which equivalent to the simultaneous linearization by topological sorting of each chromosome in each genome. IT IS THIS PROBLEM OF SELECTING THE RIGHT SUBSET OF EDGES THAT MAKES THE PROBLEM DIFFICULT (AND, WE CONJECTURE, NP-HARD.

The choice of certain edges generally excludes the choice of certain other ones. This is not just a question of avoiding multiple edges of the same colour incident to a single vertex. There are more subtle conflicts particularly involving the non-DAG edges, as illustrated in Figure 3. Our approach to this problem is a depth-first branch and bound search, `find_cycle_decomp`, in the environment of $h + k$ continually updated partial orders, one for each chromosome in each genome. The strategy is to build cycles and paths one at a time.



Fig. 3. If ab and cd are DAG edges, then the two non-DAG edges da and bc are mutually exclusive, since using them both leads to the wrong order for a and b .

Initially all edges in the DG for each chromosome are “eligible” and all vertices are “unused”. We choose any initiating vertex u in the bicoloured graph and an edge ϵ connecting it to another vertex v . All remaining edges of the same colour incident to either u or v then become ineligible. For certain choices of ϵ , the partial order associated with that chromosome must be updated through the addition of the order relation represented by ϵ , plus all others involving one vertex ordered before u and one ordered after v .

At each successive stage of the search we add an eligible edge ϵ that does not conflict with the current partial order, incident to the most recently included vertex u to extend the current cycle or path to some as yet unused vertex v or, preferably, to close a cycle or complete a cycle or path. A complication is that when the construction reaches a potential end vertex in a chromosome of ψ , it is not always clearly the termination of a path since the DAG may contain several competing end vertices. (This is not a problem with the chromosomes in χ because the dummies e_{i_1} and e_{i_2} are always at the ends of the chromosomes.) In this case, the choice of path termination or not becomes one of the branches to explore in the branch and bound.

When an edge ϵ is added, the partial order for the chromosome containing ϵ is updated, if necessary, including whenever ϵ is not a DAG edge. If ϵ is in the DAG, no update is necessary (since the initial partial orders for the branch and bound are determined by the DAGs) unless u or v is incident to more than one eligible path of the same colour as ϵ , in which case additional order is imposed by the choice of ϵ . All remaining edges of the same colour incident to u or v are made ineligible and u is now “used”.

When a cycle is completed, the initiating vertex also becomes “used”. When a path is complete, both the initiating and terminating vertices become “used”. Any unused vertex can then be chosen as to initiate a new cycle or path. (In our implementation, we increase the efficiency of the search by choosing a dummy or an end vertex whenever possible, including the very first choice of u .)

The search is bounded by using the fact that a cycle has at least two edges, and that a complete solution, representing some linearization, optimal or not,

always has $2n + k - h$ edges. Suppose the current best solution has c^* cycles and (necessarily) $h + k$ paths of which j^* are “good” as in (4). Suppose further the construction now underway is at a point where there are c' cycles and l paths, with j' good ones, and this has used m edges. This means there are only $2n + k - h - m$ edges left to chose of which at least $h + k - l$ must be in paths. Then the final number of cycles when the current construction is terminated will be no more than $c' + (2n + k - h - m - h - k + l)/2 = c' + n - h - (m - l)/2$. The final number of “good” paths will be no more than $j' + h + k - l$. So if

$$\begin{aligned} c' + n - h - (m - l)/2 + j' + h + k - l &= c' + j' + n + k - (m + l)/2 \\ &< c^* + j^*, \end{aligned} \tag{5}$$

this branch of the search is abandoned and backtracking begins.

Backtracking is also invoked if no cycles or paths can be made up of the unused vertices. During backtracking, when an edge is removed, so are the extra partial order relations it induced. The “eligible” and “unused” status it annulled are restored. An initial value of c^* can be found using any linearizations of the chromosomal DAGs of the two genomes or simply by running the depth-first algorithm until a first complete decomposition of the bicoloured graph is found.

6 Summary of the Analysis

The steps in our analysis, starting from several sets of incomplete chromosomal orders for each of two genomes, and outputting two genomes with totally ordered chromosomes, as well as a minimum number of reversals and translocations necessary to convert one to the other, are as follows.

Input: A number of incomplete maps for each genome

Remove: Genes or markers that do not appear in at least one map for each genome

For each chromosome in each map,

make_po

For each genome,

resolve

For each chromosome,

combine_po

dagger

dgger

make_bicoloured

find_cycle_decomp

Output: Optimal cycles, paths and linearizations

The major time and space costs of our method are of course due to the branch and bound procedure in **find_cycle_decomp**. The number of potential edges to be considered for inclusion in the decomposition can grow as $O(n^2 S^2)$, where S is the maximum number of parallel paths through the DAGS, but the depth of our search tree remains $O(n)$. The costs at each step are dominated by the necessity

of checking and updating a partial order matrix of size $O(n^2/h^2)$, assuming $h = k$, and all chromosomes are about the same size.

7 Analyzing the Simulated Incomplete Data

We submitted the data in (2) to our analysis. In (6) we compare the results to the original genomes in (2).

$$\begin{array}{l}
 \text{True genomes:} \\
 \{1\ 2\ 3\ 4\ 5\ 6\ 7, \\
 = 8\ 9\ 10\ 11\ 12\ 13\ 14\ 15, \\
 16\ 17\ 18\ 19\ 20, \\
 21\ 22\ 23\ 24\ 25\} \\
 \\
 \text{Reconstructed genomes:} \\
 \{1\ 3\ 2\ 5\ 4\ 6\ 7, \\
 = 8\ 9\ 10\ 11\ 12\ 13\ 14\ 15, \\
 16\ 17\ 18\ 19\ 20, \\
 22\ 21\ 23\ 24\ 25\} \\
 \\
 \{-17\ -3\ -2\ -1\ 18\ -5\ -4\ -16, \\
 = 8\ -13\ -12\ -11\ -22\ -21, \\
 -7\ -20\ 10\ -24\ -23\ 25, \\
 14\ 15, \\
 9\ -19\ 6\} \\
 \\
 \{-17\ -2\ -3\ -1\ 18\ -4\ -5\ -16, \\
 = 8\ -13\ -12\ -11\ -21\ -22, \\
 -7\ -20\ 10\ -24\ -23\ 25, \\
 14\ 15, \\
 9\ -19\ 6\}
 \end{array} \tag{6}$$

The only reconstruction errors appear to be the reverse order of genes 2 and 3, 4 and 5, and 21 and 22 in both reconstructed genomes. An inspection of the simulated incomplete data in (3), however, shows that no information is present in any of the data sets for either genome that bears on the ordering within any of these pairs.

In addition, the reconstructed linearized chromosomes of χ in (6) require 4 translocations, 1 fission and four inversions to be transformed into the reconstructed ψ , exactly how ψ was originally obtained from χ in (2).

8 Performance and Future Work

The experimental version of our program can handle moderate size maps. Tests on simulated 6-chromosome maps with 100 genes, of which 10 % were missing, and 20% unresolved from each of three datasets for both of the two genomes being compared ($d(\chi, \psi) = 20$), executed in less than a second on a Macintosh G4. With 20% missing and 40% unresolved, an analysis usually required about 3 minutes. Increasing uncertainty beyond this quickly led to run times of several hours.

The current implementation is fairly straightforward, and there are a number of promising possibilities for increasing efficiency. Since many comparative maps have only a few hundred genes our method seems quite practical.

The one-sided version of our problem, where the chromosomes of one of the genomes being compared are totally ordered is of great interest. If one genome is known in great detail, we can then resolve many of the uncertainties in less densely mapped species, despite somewhat rearranged genomes, using our technique.

Acknowledgements

Research supported in part by grants from the Natural Sciences and Engineering Research Council of Canada (NSERC). DS holds the Canada Research Chair in Mathematical Genomics and is a Fellow of the Evolutionary Biology Program of the Canadian Institute for Advanced Research.

References

1. Bader, D. A., Bernard M.E. Moret, B. M. E. and Yan, M. 2001. A linear-time algorithm for computing inversion distance between signed permutations with an experimental study. *Journal of Computational Biology* 8, 483–91.
2. Hannenhalli, S. and Pevzner, P.A. 1995. Transforming men into mice (polynomial algorithm for genomic distance problem. *Proceedings of the IEEE 36th Annual Symposium on Foundations of Computer Science*. 581–92.
3. Hannenhalli, S. and Pevzner, P.A. 1999. Transforming cabbage into turnip (polynomial algorithm for sorting signed permutations by reversals). *Journal of the ACM* 48, 1–27.
4. Tesler, G. 2002. Efficient algorithms for multichromosomal genome rearrangements. *Journal of Computer and System Sciences* 65, 587–609.
5. Zheng, C., Lenert, A. and Sankoff, D. 2005. Reversal distance for partially ordered genomes. *Bioinformatics* 21, supplementary issue, *Proceedings of ISMB 2005*.