

The Footprint Sorting Problem[†]

Claudia Fried,^{‡,§} Wim Hordijk,^{||} Sonja J. Prohaska,^{‡,§} Claus R. Stadler,[⊥] and Peter F. Stadler^{*,‡,§}

Bioinformatics, Department of Computer Science, University of Leipzig, Germany, Institute for Theoretical Chemistry and Structural Biology, University of Vienna, Austria, Department of Mathematics and Statistics, University of Canterbury, Christchurch, New Zealand, and Department of Computer Science, University of Leipzig, Germany

Received July 23, 2003

Phylogenetic footprints are short pieces of noncoding DNA sequence in the vicinity of a gene that are conserved between evolutionary distant species. A seemingly simple problem is to sort footprints in their order along the genomes. It is complicated by the fact that not all footprints are collinear: they may cross each other. The problem thus becomes the identification of the crossing footprints, the sorting of the remaining collinear cliques, and finally the insertion of the noncollinear ones at “reasonable” positions. We show that solving the footprint sorting problem requires the solution of the “Minimum Weight Vertex Feedback Set Problem”, which is known to be NP-complete and APX-hard. Nevertheless good approximations can be obtained for data sets of interest. The remaining steps of the sorting process are straightforward: computation of the transitive closure of an acyclic graph, linear extension of the resulting partial order, and finally sorting w.r.t. the linear extension. Alternatively, the footprint sorting problem can be rephrased as a combinatorial optimization problem for which approximate solutions can be obtained by means of general purpose heuristics. Footprint sortings obtained with different methods can be compared using a version of multiple sequence alignment that allows the identification of unambiguously ordered sublists. As an application we show that the rat has a slightly increased insertion/deletion rate in comparison to the mouse genome.

1. INTRODUCTION

Phylogenetic footprints are short pieces of noncoding DNA sequence in the vicinity of a gene that are conserved between evolutionarily distant species.¹ They are conserved over time scales of sometimes hundreds of millions of years because their function is crucial for the survival of the organism: Phylogenetic footprints are (predominantly) binding sites for transcription factors that regulate the expression of the associated genes.^{2–5} A common methodology for detecting phylogenetic footprints is the comparison of the DNA sequences from two or more organisms with suitable pairwise distances so that the conserved sequence pieces can be discriminated from the intervening sequences that are randomized by the accumulation of mutations, see e.g. ref 6. Automatic procedures for phylogenetic footprinting such as **footprinter**⁷ or **tracker**⁸ can produce large amounts of data that require automatized analysis tools.

A seemingly simple problem is to sort the detected footprints from a multispecies comparison sort in their order along the genomes. This task is complicated by the fact that not all footprints are collinear: they may cross each other. Often, a properly sorted list is mostly a convenience for presenting the data; below we will present an application where a sorted list is a necessary prerequisite. When one

studies the fate of ancestral footprints, those that violate collinearity are neglected because they are unlikely to be true homologues.⁹ Nevertheless they can well be real binding sites, see e.g., ref 10, that might have arisen from duplications and translocations.¹¹ Hence we cannot simply discard non-collinear footprints. The problem thus becomes to identify the crossing footprints, to sort the remaining collinear cliques, and finally to insert the noncollinear ones at “reasonable” positions.

In this contribution we show that the footprint sorting problem is in fact a hard combinatorial problem. Identification of those footprints that violate collinearity can be formulated as a Minimum Feedback Vertex Set Problem.^{12,13} Alternatively, we can directly search for a suitable sorting by assigning a cost to each collinearity violation of a given permutation (sorting) and then using a heuristic to minimize this cost function.

Mathematically speaking, we are given N intervals \mathcal{X}^i , $i = 1, \dots, N$ representing the DNA sequences. Let us denote by $[i; a, l]$ the subinterval $[a, a + l - 1] \subseteq \mathcal{X}^i$ where i identifies the DNA sequence, a is the initial position of the subinterval, and l is the length of the interval. A *footprint clique* J is a collection of subintervals with the property that $\alpha = [i; a, l] \in J$ and $\alpha' = [i'; a', l'] \in J$ implies either $\alpha = \alpha'$ or $i \neq i'$, i.e., a footprint clique contains at most one subinterval from each sequence \mathcal{X}^i . The output of a footprinting program is a collection \mathcal{J} of M footprint cliques J_k , $k = 1, \dots, M$.

Since not all footprint cliques are of equal importance (or have been determined with equal certainty), it is useful to assign a weight $w: \mathcal{J} \rightarrow [0, 1]$ to each footprint clique. For

* Corresponding author phone: ++49 341 14951 20; fax: ++49 341 14951 19; e-mail: peter.stadler@bioinf.uni-leipzig.de.

[†] Dedicated to George W. A. Milne, a former long-term Editor-in-Chief of *JCICS*.

[‡] Bioinformatics, Department of Computer Science, University of Leipzig.

[§] University of Vienna.

^{||} University of Canterbury.

[⊥] Department of Computer Science, University of Leipzig.

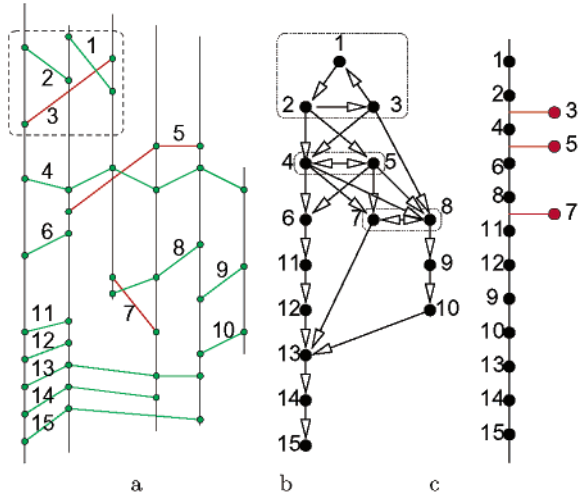


Figure 1. (a) Phylogenetic footprints (small balls) of different sequences (vertical lines) belonging to the common footprint clique are connected by lines. (b) Directed graph G describing their relative locations (not all arcs are shown for clarity). 2-connected components, i.e., obstructions to partial ordering are shown by boxes. (c) A well-ordering on a maximal set of collinear cliques. The diagram also indicates the obstructing cliques at positions with a minimal number of conflicts.

instance, one might use

$$\omega([i; a, l]) = \frac{l}{\max_j l_j} \quad (1)$$

where the maximum is taken over all intervals in all footprint cliques. More sophisticated weight functions, that e.g. take the sequence conservation into account, might also be useful. For each subset $\mathcal{J} \subseteq \mathcal{J}$ we define the weight of the subset

$$\omega(\mathcal{J}) = \sum_{\alpha \in \mathcal{J}} \omega(\alpha) \quad (2)$$

For two subintervals $\alpha = [i; a, l]$ and $\alpha' = [i; a', l']$ on the same sequence i we have the (trivial) order relation

$$\alpha < \alpha' \leftrightarrow \begin{cases} a < a' \\ l < l' \end{cases} \text{ and } a = a' \quad (3)$$

Clearly, the ordering (3) implies an order-like relation $<^*$ on \mathcal{J} :

Definition 1. For $J, J' \in \mathcal{J}$ we set $J <^* J'$ [Due to limitations with the composition system the symbols “ $<$ ” and “ $<^*$ ” are used to mean “less than” and “less than or equal”, respectively, in the context of a partial order instead of the usual mathematical symbols for these relations.] if and only if for all $i \in \{1, \dots, N\}$ for which there is an $\alpha = [i; a, l] \in J$ and an $\alpha' = [i; a', l'] \in J'$ we have $\alpha \leq \alpha'$.

Our task is hence to find a well-order on \mathcal{J} that is consistent with the order of the subintervals on each sequence, i.e., that is an extension of $<^*$ on \mathcal{J} .

Recall that a relation \leq on a set X is a partial order if the following three axioms are satisfied

- (O1) $x \leq x$ for all $x \in X$ (reflexivity).
- (O2) $x \leq y$ and $y \leq x$ implies $x = y$ (antisymmetry).
- (O3) $x \leq y$ and $y \leq z$ implies $x \leq z$ (transitivity).

A partial order is totally ordered if in addition we have (O4) $x \leq y$ or $y \leq x$ for all $x, y \in X$.

A total order consistent with any given partial order (a so-called *linear extension*) can be computed efficiently, see e.g. ref 14. As a necessary condition, the transitive closure $\overline{<^*}$ of $<^*$ must therefore be a partial order. In general, however, this is not the case for realistic data. Consider the following two simple examples:

(1) $J = \{[1; a_1, l_1], [2; a_2, l_2], [3; a_3, l_3]\}$ and $J' = \{[1; a_1, l_1], [2; a_2, l_2], [4; a_4, l_4]\}$. In this case we have $J <^* J'$ and $J' <^* J$ but $J \neq J'$, i.e., antisymmetry (O2) is violated.

(2) $J_1 = \{[1; a_1, l_1], [2; a_2, l_2]\}$, $J_2 = \{[1; a'_1, l'_1], [3; a_3, l_3]\}$, and $J_3 = \{[2; a'_2, l'_2], [3; a'_3, l'_3]\}$, such that $a_1 < a'_1$, $a'_2 < a_2$, $a_3 < a'_3$. This implies by definition $J_1 <^* J_2$, $J_3 <^* J_1$, and $J_2 <^* J_3$. For the transitive closure, hence, we have $J_1 <^* J_3$ and $J_3 <^* J_1$ but $J_1 \neq J_3$, again violating (O2).

The relation $<^*$ therefore is not antisymmetric in general. Our task therefore becomes to identify a maximal subset $\mathcal{J} \subseteq \mathcal{J}$ of footprint cliques that can be well-ordered. Maximality is defined conveniently w.r.t. some weight function such as eq 2. The remaining footprint cliques that have to be removed from \mathcal{J} are those that are called “noncollinear” in ref 8. We remark that in the case of just two sequences the maximum increasing subsequence algorithm¹⁵ [12.5.1] can be used.

2. MINIMUM FEEDBACK VERTEX SETS

The set \mathcal{J} can be regarded as the vertex set of a directed graph \vec{G} with arcs $\alpha \rightarrow \beta$ if and only if $\alpha <^* \beta$ and $\alpha \neq \beta$. The following result is obvious from the definition of a partial order:

Lemma 1. The transitive closure of the relation \leq^* is a partial order if and only if the associated graph \vec{G} is acyclic.

Proof. The relation \leq^* is antisymmetric if and only if for any two vertices x and y with a direct path from x to y there is no directed path from y to x , i.e., no two vertices in \vec{G} are contained in a circuit (q.e.d.).

The problem of detecting noncollinear footprint cliques can therefore be rephrased as follows:

Given the vertex-weighted directed graph $\vec{G} = (V, A)$ with vertex set V and arc-set A , find a maximal (w.r.t. ω) subset of vertices $U \subseteq V$ such that the induced subgraph $\vec{G}[U]$ is acyclic.

In other words, we want to remove a set $W = V \setminus U$ of *noncollinear* footprints with minimal total weight. This problem is known as the *Minimum Feedback Vertex Set Problem*,¹² see ref 13 for a recent review. It is known to be NP-hard¹⁶ and has applications in many diverse areas, including program verification¹⁷ and Bayesian inference.¹⁸

Unfortunately, the exact algorithm described in ref 19 is not fast enough for large examples (with thousands of nodes) with a larger number, say 10%, of noncollinear footprints. Well-tested heuristics based on a greedy randomized adaptive search procedure are available.²⁰ For the purpose of this study we have reimplemented both approaches in C. Recently a branch-and-cut algorithm was proposed²¹ that might be a

useful alternative. A general approximation algorithm for node-deletion problems is described in ref 22.

3. TOPOLOGICAL SORTING

For the next step we have two options: (1) We may sort the acyclic subset U and then reinsert the feedback set W at appropriate positions. (2) Alternatively, we may modify $\tilde{G} = (V, A)$ by removing some of the arcs incident with the feedback vertex set W in order to obtain an acyclic graph $G(V, A')$ with $A' \subset A$. Superficially, this suggests considering the *feedback arc set problem*, i.e., to find a maximal subset $A' \subset A$ such that $\tilde{G}' = (V, A')$ is acyclic. In the present context, however, we are interested in the set of noncollinear cliques W , while the feedback arcs $A \setminus A'$ do not seem to have an interpretation in terms of the phylogenetic footprints.

In both cases we compute the transitive closure of the acyclic graph by connecting two vertices with an arc from i to j if and only if there is a directed path from i to j . This can be achieved e.g. by Warshall's algorithm in $O(|U|^3)$ time.²³ The resulting graph \tilde{G}^* is again acyclic and represents a partial order.

The computation of a well-ordering of the vertex set of a directed acyclic graph \tilde{G} (such that arcs go only from vertices with smaller labels to vertices with larger labels) is known as *topological sorting* and can be solved in $O(|A|)$ time.^{24,25}

4. SORTING AS OPTIMIZATION PROBLEM

A completely different approach starts with the observation that the standard sorting problem can be reformulated as a combinatorial optimization problem. This is motivated by our requirement to obtain a sorted list of all footprints, including those that we previously identified as feedback vertex set W . Recall that the weight function $\omega(I)$ in eq 2 is just a heuristic approximation, hence one might want to investigate a collection of noncollinear footprints in detail (e.g. be explicitly considering the sequence alignments of all these footprints). A properly sorted footprint listing is of practical advantage.

Let $(X, <)$ be a finite ordered set, which, without losing generality, we can identify with the set $\{1, 2, \dots, |X|\}$ endowed with the standard order on \mathbb{N} . Furthermore let $U = (u_{ij})$ be a symmetric weight matrix with $u_{ij} = u_{ji} > 0$. We use $u_{ij} = (\omega(i) + \omega(j))/2$ in terms of the footprint weights (1). Write $\pi(i)$ for the X -element at the i th position in the sorted list and set

$$f_{ij}(\pi) = \begin{cases} u_{\pi(i)\pi(j)} & \text{if } i < j \text{ and } \pi(i) > \pi(j) \\ u_{\pi(i)\pi(j)} & \text{if } i > j \text{ and } \pi(i) < \pi(j) \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

The quantity $f_{ij}(\pi) = 0$ if $\pi(i)$ and $\pi(j)$ are correctly sorted with respect to each other in the particular ordering (permutation) π , while $f_{ij}(\pi) > 0$ if $\pi(i)$ and $\pi(j)$ are missorted. Thus the total cost of a particular ordering π is conveniently defined as

$$f(\pi) = \sum_{ij} f_{ij}(\pi) \quad (5)$$

which can be interpreted as the total weight of all conflicting pairs. It can be interpreted as a weighted variant of Kendall's τ parameter.²⁶

The problem of sorting an n -element set (given that the comparison of two elements can be evaluated in constant time) is solved in $O(n \log n)$ by standard algorithms such as *quick sort*, implemented e.g. in the C standard library function `qsort`. It is clear, therefore, that a heuristic approach based on minimizing f will not be computationally efficient. In contrast to classical sorting algorithms, however, we can generalize the combinatorial optimization approach as we shall see below. Let us first consider the properties of eq 4 for well-orders and partial orders, however.

Theorem 1. *If $(X, <)$ is well-ordered, then every adaptive walk that makes use of the canonical transpositions reaches the correct sorting.*

Proof. If π is not the correct sorting, then there is $m \in X$ such that $\pi(m) < \pi(m-1)$. It is convenient to define to $f_k(\pi) = \sum_i f_{ik}(\pi)$, the total weight of the conflicts of the k th object in the list. Let π' be the ordering obtained by exchanging m and $m-1$, i.e. $\pi'(i) = \pi(i)$ for $i \neq m, m-1$, $\pi'(m) = \pi(m-1)$ and $\pi'(m-1) = \pi(m)$. One easily verifies that $f_k(\pi) = f_k(\pi')$ for $k \neq m-1, m$. Tedious but straightforward computation shows that $f_{m-1}(\pi') + f_m(\pi') - f_{m-1}(\pi) - f_m(\pi) = -2u_{\pi(m-1)\pi(m)} < 0$, i.e. every canonical transposition that exchanges a missorted pair of adjacent objects strictly decreases the cost function f . Therefore there is no local minimum of f with the exception of the correct sorting (where a missorted object m does not exist by definition) (q.e.d.).

It follows immediately that adaptive walks using arbitrary transpositions and/or reversals are also guaranteed to find the correct sorting because these movesets contain the canonical transpositions.

In this formulation the sorting problem can be extended in an obvious way to an arbitrary relation \angle on X which need not be complete or even antisymmetric. Clearly, there is a perfect solution π satisfying $f(\pi) = 0$ if and only if (X, \angle) is a partially ordered set. A permutation thus satisfies $f(\pi) = 0$ if and only if π codes for a linear extension of \angle .

Theorem 2. *If (X, \angle) is a partially ordered set, then every adaptive walk that makes use of the (general) transpositions reaches a correct topological sorting.*

Proof. The argument in the proof above can be modified in the following form in the case of partial orders. We first observe that if π is not a linear extension of \angle then there are $m, m' \in X$ such that $\pi(m) \angle \pi(m')$, while $\pi(k)$ is incomparable with $\pi(m)$ and $\pi(m')$ for all $m < k < m'$. In the extreme case $m \neq m' - 1$. Now consider the permutation π' defined by $\pi'(m) = \pi(m')$, $\pi'(m') = \pi(m)$, and $\pi'(j) = \pi(j)$ for $j \neq m, m'$. Since the positions k between m and m' are incomparable with both m and m' we have $f_j(\pi) = f_j(\pi')$ for all $j < m, j > m'$, and $m < j < m'$. Furthermore $f_m(\pi') - f_m(\pi) = -u_{\pi(m)\pi(m')} < 0$ since these terms differ only by the contribution from comparing $\pi(m)$ and $\pi(m')$. By the same argument $f_{m'}(\pi') - f_{m'}(\pi) = -u_{\pi(m')\pi(m)} < 0$. Thus $f(\pi') < f(\pi)$ (q.e.d.).

Let us now turn to the general case where we are given a directed graph $\tilde{G} = (V, A)$ that is not acyclic in general. Let $U = (u_{ij})$ be a symmetric weight matrix satisfying $u_{ij} =$

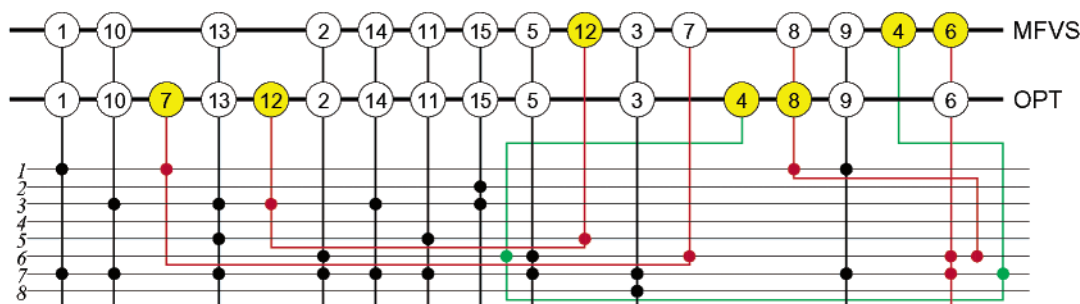


Figure 2. Alignment of two permutations. We use here the first 15 cliques from the comparison of eight *HoxA* clusters of various vertebrates discussed in ref 28. The two rows on the top display the alignment of the exact ordering derived from an exact solution of the minimum feedback vertex set problem and the ordering obtained by optimizing eq 6. The feedback vertices are indicated in gray. Below the positions of the underlying footprints on the 8 nucleotide sequences are shown.

$u_{ji} > 0$ whenever $(i, j) \in A$ or $(j, i) \in A$ is an arc in \bar{G} . Furthermore we define the cost function f as

$$f_{ij}(\pi) = \begin{cases} u_{\pi(i)\pi(j)} & \text{if } i < j \text{ and } (\pi(j), \pi(i)) \in A \\ u_{\pi(i)\pi(j)} & \text{if } i > j \text{ and } (\pi(i), \pi(j)) \in A \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

Equation 6 reduces to (4) provided $(x, y) \in A$ if and only if $x < y$ and $<$ is a partial order. We can hope to obtain a useful ordering π of the nodes of \bar{G} by minimizing f . We do not know of an exact solution to this problem. It would be easy to use adaptive walks or other simple local search algorithms such as simulated annealing,²⁷ however. Given the two theorems above, it appears that transpositions, and possibly also reversals of permutations, will be a good move set at least as long as the graph is nearly acyclic.

Given the permutation π that represents our best approximation to the true ordering of the vertices, we would also like to identify a minimum feedback vertex set. It is not clear whether this can be done exactly. A heuristic approximation proceeds by iteratively removing the vertex k with the largest total weight

$$g_k(\pi) = \sum_i \begin{cases} \omega(i) & \text{if } k < i \text{ and } (\pi(i), \pi(k)) \in A \\ \omega(i) & \text{if } k > i \text{ and } (\pi(k), \pi(i)) \in A \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

of all those objects with which it conflicts in the given ordering π .

More formally we have the following algorithm:

Input: $\bar{G} = (V, A)$, U , π
Output: W /* feedback vertex set */
 1: $W \leftarrow \emptyset$
 2: **while** $f(\pi) > 0$ **do**
 3: $\ell := \arg \max g_k(\pi)$.
 4: $W \leftarrow W \cup \{\pi(\ell)\}$.
 5: $u_{\pi(\ell)i}, u_{i\pi(\ell)} \leftarrow 0$ for all i .
 6: **end while**

Clearly, this procedure produces a feedback vertex set; its minimality is not guaranteed even if π minimizes f .

5. COMPARISON OF SORTED LISTS

To assess the quality of the orderings π obtained from the different approaches and variants described above we need a systematic way of comparing ordered lists. Since these

lists are represented as permutations, π' and π'' , of the same set n of objects, it seems natural to use distance measures $d(\sigma, \tau)$ on the symmetric group S_n . Natural metrics on S_n are associated with a length function that measures how much a group element is “different” from the identity element: $d(\sigma, \tau) = L(\sigma\tau^{-1})$.³⁰ For instance, the minimum number of transpositions that are necessary to generate π from the identity permutation satisfies $L_\tau(\pi) = n - \text{cyc}(\pi)$, where $\text{cyc}(\pi)$ is the number of cycles in the cycle representation of π . Length functions associated with sorting by canonical transpositions (in terms of so-called inversions) and reversals can also be computed.³¹

Another possibility, which allows a more convenient comparison of the two lists, is to *align* the two permutations π' and π'' of n elements such that the number D of insertions and deletions is minimized, Figure 2. This can be achieved by a simple dynamic programming scheme originally invented by Needleman and Wunsch³² for the alignment of two protein sequences. Starting from the initialization $D_{0i} = D_{i0} = i$ we have to compute

$$D_{ij} = \min \begin{cases} D_{i,j-1} + 1 \\ D_{i-1,j} + 1 \\ D_{i-1,j-1} \end{cases} \quad \text{whenever } \pi'(i) = \pi''(j) \quad (8)$$

The optimal number of insertions and deletions is then given by $D = D_{nm}$. The alignment of the orderings is then reconstructed from the matrix (D_{ij}) by a standard backtracking procedure. Alignments of more than two lists can be computed by means of the iterative procedure familiar from popular multiple sequence alignment programs such as **clustalw**.³³ We note in passing that such list alignments can be used as an alternative approach to the Top- k List Comparison problem discussed in ref 34.

Figure 3 shows an example of a manually sorted footprint list in comparison with the exact solution of the MFVS-based procedure and the optimization heuristic described above. These data are fairly noisy containing a significant number of noncollinear footprints. The resulting sortings and assignments of noncollinearities are rather different: the manual list differs by 16 and 21 indels from the exact and the optimization results, the two automatic sortings differ by 26 indels, i.e., the manually sorted list is much closer to the exact MFVS-based listing than to the result of the simple optimization procedure. This is not surprising since we cannot expect adaptive walks to produce good results in problems with a large number of order conflicts. More

Sorting			Exon	Footprints												
manual	MFVS	opt.		PmW	HsA	HsC	HsD	HfM	HfN	TrAa	TrAb	TrD	TrCa	Ci	TrBa	HsB
1	1	1	*	•	•	•	•	•	•	•	•	•	•	•	•	•
2	2	2		•	•	•	•	•	•	•	•	•	•	•	•	•
3	3	3		•	•	•	•	•	•	•	•	•	•	•	•	•
4	4	4		•	•	•	•	•	•	•	•	•	•	•	•	•
5	5	5		•	•	•	•	•	•	•	•	•	•	•	•	•
6	6	6		•	•	•	•	•	•	•	•	•	•	•	•	•
7	7	7		•	•	•	•	•	•	•	•	•	•	•	•	•
8	8	8		•	•	•	•	•	•	•	•	•	•	•	•	•
9	9	9		•	•	•	•	•	•	•	•	•	•	•	•	•
10	10	10		•	•	•	•	•	•	•	•	•	•	•	•	•
11	11	11	*	•	•	•	•	•	•	•	•	•	•	•	•	
12	12	12		•	•	•	•	•	•	•	•	•	•	•	•	•
13	13	13		•	•	•	•	•	•	•	•	•	•	•	•	•
14	14	14		•	•	•	•	•	•	•	•	•	•	•	•	•
15	15	15		•	•	•	•	•	•	•	•	•	•	•	•	•
16	16	16		•	•	•	•	•	•	•	•	•	•	•	•	•
17	17	17		•	•	•	•	•	•	•	•	•	•	•	•	•
18	18	18		•	•	•	•	•	•	•	•	•	•	•	•	•
19	19	19		•	•	•	•	•	•	•	•	•	•	•	•	•
20	20	20		•	•	•	•	•	•	•	•	•	•	•	•	•
21	21	21	*	•	•	•	•	•	•	•	•	•	•	•	•	
22	22	22		•	•	•	•	•	•	•	•	•	•	•	•	•
23	23	23		•	•	•	•	•	•	•	•	•	•	•	•	•
24	24	24		•	•	•	•	•	•	•	•	•	•	•	•	•
25	25	25		•	•	•	•	•	•	•	•	•	•	•	•	•
26	26	26		•	•	•	•	•	•	•	•	•	•	•	•	•
27	27	27		•	•	•	•	•	•	•	•	•	•	•	•	•
28	28	28		•	•	•	•	•	•	•	•	•	•	•	•	•
29	29	29		•	•	•	•	•	•	•	•	•	•	•	•	•
30	30	30		*	•	•	•	•	•	•	•	•	•	•	•	•
31	31	31	•		•	•	•	•	•	•	•	•	•	•	•	•
32	32	32	•		•	•	•	•	•	•	•	•	•	•	•	•
33	33	33	•		•	•	•	•	•	•	•	•	•	•	•	•
34	34	34	•		•	•	•	•	•	•	•	•	•	•	•	•
35	35	35	•		•	•	•	•	•	•	•	•	•	•	•	•
36	36	36	•		•	•	•	•	•	•	•	•	•	•	•	•
37	37	37	•		•	•	•	•	•	•	•	•	•	•	•	•
38	38	38	•		•	•	•	•	•	•	•	•	•	•	•	•
39	39	39	•		•	•	•	•	•	•	•	•	•	•	•	•
40	40	40	*	•	•	•	•	•	•	•	•	•	•	•	•	
41	41	41		•	•	•	•	•	•	•	•	•	•	•	•	•
42	42	42		•	•	•	•	•	•	•	•	•	•	•	•	•
43	43	43		•	•	•	•	•	•	•	•	•	•	•	•	•
44	44	44		•	•	•	•	•	•	•	•	•	•	•	•	•
45	45	45		•	•	•	•	•	•	•	•	•	•	•	•	•
46	46	46		•	•	•	•	•	•	•	•	•	•	•	•	•
47	47	47		•	•	•	•	•	•	•	•	•	•	•	•	•
48	48	48		•	•	•	•	•	•	•	•	•	•	•	•	•
49	49	49		•	•	•	•	•	•	•	•	•	•	•	•	•

Figure 3. Some footprint data sets are quite noisy and contain many noncollinear entries. The example shown here is the Hox-10 region from the lamprey *Petromyzon marinus* (Pm), human (Hs), the hornshark *Heterodontus francisci* (Hf), the pufferfish *Takifugu rubripes* (Tr), and the tunicate *Ciona intestinalis* (Ci). The first column gives the manual sorting based on the raw tracker output,²⁹ the second column is the sorting based on solving the MFVS problem exactly, and the third column was obtained using the optimization approach. The + signs indicate footprint cliques that are identified as noncollinear; one match was removed from two footprints in ref 29 after visual inspection, indicated by a = sign. The two exons of the Hox-10 gene itself are indicated by a *. Footprints are marked by •.

sophisticated heuristics such as simulated annealing³⁵ or genetic algorithms³⁶ will have to be used for such data sets. Most data sets that we have encountered so far, however, are much more well-behaved, containing a much smaller feedback set. In those cases the adaptive walk procedure works very well.

6. AN APPLICATION

For some applications a properly sorted list of phylogenetic footprints is not only a convenience but a necessary prerequisite for further data analysis. Fortunately, many data sets, in particular those with fewer input sequences, are much less noisy than the one shown in Figure 3 so that noncollinear cliques can be identified (almost) unambiguously. In this section we briefly consider an example of an application where unambiguity is important.

The genomic DNA of all organisms is permanently reorganized by mutation and recombination.³⁷ The rates of mutation and recombination are important characteristics of the dynamics of evolution. Variations of substitution (point-mutation) rates between genes and lineages, for instance,

need to be taken into account in phylogenetics.³⁸ Mutation rate variations are of interest in their own right e.g. as indicators of adaptive evolution.^{39,40} The distribution of SNPs (single nucleotide polymorphisms) throughout the genome depends strongly on the interplay between mutation and recombination;^{41,42} the understanding of these mechanisms is a prerequisite for the efficient usage of SNPs as disease markers.

The total size of a genome also varies over evolutionary time scales. These variations can be caused by large scale changes such as gene and chromosome duplications and by the accumulation of small-scale local insertions, deletions, and inversions. Insertion/deletion (indel) rates can of course be estimated from pairwise alignments of DNA sequences⁴³ as long as the organisms in question are sufficiently closely related so that reliable pairwise sequence alignments can be computed. In this section we show how the length variations of the nonconserved sequences between *adjacent* homologous footprints can be used to detect differences in the indel rates in a pair of sequences relative to an outgroup. Clearly, this requires a reliable method for excluding noncollinear

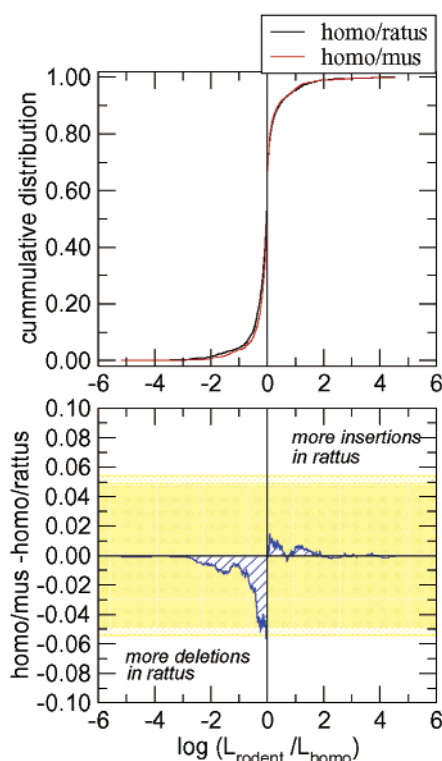


Figure 4. Comparison of length ratios of $N_1 = 1122$ human/rat and $N_2 = 1404$ human/mouse nonconserved sequence fragments located between homologous footprints in the vicinity of various immune genes. The maximum difference between the two cumulative probability density functions is $D^* = 0.05659$, compared to the threshold value $D_{95\%} = 0.05438$ for the Kolmogorov-Smirnov test with effective sample size $N_{\text{eff}} = N_1 N_2 / (N_1 + N_2) \approx 623.6$. The two distributions are therefore (just barely) significantly different.

(and hence nonhomologous) footprints and a proper ordering of the footprints so that adjacency can be determined.

Let l_j^h and l_j^k be the lengths of the intervening sequences between two adjacent footprints that are common to the input sequences h and k . Given a reference sequence, say human h , we are interested in the differences between the distributions of $\alpha_j^k = \log(l_j^k/l_j^h)$ for different species k , say mouse and rat. We assume that the total length of insertions and deletions will be approximately proportional to the length l_j^h of the piece of the reference sequence, hence we consider the length ratio rather than the length difference. (Depending on the mechanism that one envisages for the insertion and deletion processes a different scaling might be more appropriate, however.) We use the logarithm of the length ratio because it produces a distribution that is symmetric if there is no bias between insertions and deletions.

Recent analysis showed that chromosome 10 of the rat differs by several chromosome rearrangements from both the mouse and the human homologues,⁴⁴ see also ref 45. It is of interest whether this increased activity is associated with an elevated indel rate in the mouse. To address this issue we compare here the distribution of indels in mouse and rat relative to the corresponding human sequence. Our data set consists of immune genes with all intron sequences and with 5000nt of flanking sequence on each side.

Figure 4 summarizes the results. The effect of deletions and insertions appears to be slightly increased in the rat compared to the mouse. The difference is very small

however: the Kolmogorov-Smirnov test just barely yields a 95% confidence for a significant difference of the distribution of rat/man and mouse/man sequence length ratios.

7. DISCUSSION

We have shown here that the seemingly trivial task of sorting a list of phylogenetic footprints properly by their location along the sequences in fact gives rise to a complex optimization problem. As a chemical application one might for instance consider the problem of sorting a list of samples of complex mixtures according to their compositions in the presence of missing data, i.e., when not all components are measured in all samples. In this case, simpler approaches such as lexicographic sorting cannot be applied.

Here we have described two approaches to solving this task. An exact algorithm that requires the solution of the NP-complete minimum feedback vertex set problem and a formulation as a combinatorial optimization problem that essentially generalizes a “landscape version” of bubble sort.

Recently an algorithm has been published that produced all (inclusion-wise) minimal feedback vertex sets with polynomial delay.⁴⁶ In principle at least, this approach could be used to identify ambiguities among those feedback vertex sets with nearly optimal weight ω .

The “Footprint Sorting Problem” can also be viewed abstractly as a multiobjective optimization problem. Given a directed graph $\vec{G}(V, A)$, a vertex weight function $\omega: V \rightarrow \mathbb{R}^+$, and an edge weight function $u: AR^+$ find a permutation $\pi: V \rightarrow V$ and a “feedback set” $W \subset V$ such that both the weight $\sum_{i \in W} \omega(i)$ of W and the weight total of the conflicts

$$u(V \setminus W) = \sum_{i,j \in V \setminus W} f_{ij}(\pi)$$

among the remaining vertices are minimized.

This rather general version of a topological sorting problem arises in many different contexts. For example we might want to sort the results from different database queries for the same topic. In this case u_{ij} is e.g. confidence or score-difference with which a particular database ranks the results i better than j and $\omega(i)$ measures how much information the result contains. The goal would be to rank the results as good as possible in accordance with rankings from the individual queries and to focus on the most detailed results.

ACKNOWLEDGMENT

Financial support by the German DFG Bioinformatics Initiative (C.F., S.J.P., P.F.S.) and the Allan Wilson Centre for Molecular Ecology and Evolution in New Zealand (W.H.) is gratefully acknowledged.

REFERENCES AND NOTES

- (1) Tagle, D. A.; Koop, B. F.; Goodman, M.; Slightom, J. L.; Hess, D. L.; Jones, R. T. Embryonic epsilon and gamma globin genes of a prosimian primate (*Galago crassicaudatus*). Nucleotide and amino acid sequences, developmental regulation and phylogenetic footprints. *J. Mol. Biol.* **1988**, *203*, 439–455.
- (2) Arnone, M. I.; Davidson, E. H. The hardwiring of development: Organization and function of genomic regulatory systems. *Development* **1997**, *124*, 1851–1864.
- (3) Fickett, J. W.; Wasserman, W. W. Discovery and modeling of transcriptional regulatory regions. *Curr. Opin. Biotech.* **2000**, *11*, 19–24.

- (4) Davidson, E. *Genomic Regulatory Systems*; Academic Press: San Diego, 2001.
- (5) Ludwig, M. Z. Functional and evolution of noncoding DNA. *Curr. Opin. Genet. Dev.* **2002**, *12*, 634–639.
- (6) Schwartz, S.; Zhang, Z.; Frazer, K. A.; Smit, A.; Riemer, C.; Bouck, J.; Gibbs, R.; Hardison, R.; Miller, W. PipMaker - A Web server for aligning two Genomic DNA sequences. *Genome Res.* **2000**, *4*, 577–586.
- (7) Blanchette, M.; Schwikowski, B.; Tompa, M. Algorithms for phylogenetic footprinting. *J. Comput. Biol.* **2002**, *9*, 211–223.
- (8) Prohaska, S.; Fried, C.; Flamm, C.; Wagner, G.; Stadler, P. F. Surveying phylogenetic footprints in large gene clusters: applications to Hox cluster duplications. *Mol. Phylog. Evol.* **2004**.
- (9) Prohaska, S. J.; Fried, C.; Amemiya, C. T.; Ruddle, F. H.; Wagner, G. P.; Stadler, P. F. The Shark HoxN cluster is homologous to the Human HoxD cluster. *J. Mol. Evol.* **2004**.
- (10) Jegga, A. G.; Sherwood, S. P.; Carman, J. W.; Pinski, A. T.; Phillips, J. L.; Pestian, J. P.; Aronow, B. J. Detection and visualization of compositionally similar *cis*-regulatory element clusters in orthologous and coordinately controlled genes. *Genome Res.* **2002**, *12*, 1408–1417.
- (11) Kent, W. J.; Baertsch, R.; Hinrichs, A.; Miller, W.; Haussler, D. Evolution's cauldron: Duplication, deletion, and rearrangement in the mouse and human genomes. *Proc. Natl. Acad. Sci. U.S.A.* **2003**, *100*, 1484–1489.
- (12) Lempel, A.; Cederbaum, I. Minimum feedback arc and vertex sets of a directed graph. *IEEE Trans. Circuit Theory* **1966**, *13*, 399–403.
- (13) Festa, P.; Pardalos, P.; Resende, M. G. C. In *Encyclopedia of Optimization*; Kluwer: 2001; Vol. 2, pp 94–106.
- (14) Preusse, G.; Ruskey, F. Generating linear extensions fast. *SIAM J. Comput.* **1994**, *23*, 373–386.
- (15) Gusfield, D. *Algorithms on Strings, Trees, and Sequences*; Cambridge University Press: Cambridge, UK, 1997.
- (16) Garey, M. R.; Johnson, D. S. *Computers and Intractability. A Guide to the Theory of NP-Completeness*; Freeman: San Francisco, 1979.
- (17) Shamir, A. A linear time algorithm for finding minimum cutsets in reduced graphs. *SIAM J. Comput.* **1979**, *8*, 654–655.
- (18) Yehuda, B.; Geiger, D.; Naor, J.; Roth, R. M. Approximation algorithms for the vertex feedback set problem with applications to constraint satisfaction and Bayesian inference. In *Proc. 5th Annual ACM-SIAM Symp. on Discrete Algorithms*, 1994; pp 344–354.
- (19) Ashar, P.; Malik, S. Implicit computation of minimum-cost feedback vertex sets for partial scan and other applications. In *Proc. 31st Annual ACM-SIAM Conf. (DAC)*; ACM Press: 1994; pp 77–80.
- (20) Festa, P.; Pardalos, P.; Resende, M. G. C. Algorithm 815: Fortran subroutines for computing approximate solutions of feedback set problems using GRASP. *ACM Trans. Math. Software* **2001**, *27*, 456–464.
- (21) Brunetta, L.; Maffioli, F.; Trubian, M. Solving the feedback vertex set problem on undirected graphs. *Discr. Appl. Math.* **2000**, *101*, 37–51.
- (22) Fujito, T. A unified approximation algorithm for node-deletion problems. *Discr. Appl. Math.* **1998**, *86*, 213–231.
- (23) Marshall, S. A theorem on Boolean matrices. *J. ACM* **1962**, *9*, 11–12.
- (24) Toda, S. On the complexity of topological sorting. *Inf. Process. Lett.* **1990**, *35*, 229–233.
- (25) Hagerup, T.; Maas, M. Generalized topological sorting in linear time. *Nord. J. Comput.* **1994**, *1*, 38–49.
- (26) Kendall, M. G. A new measure of rank correlation. *Biometrika* **1938**, *30*, 81–93.
- (27) Kirkpatrick, S.; Gelatt, C. D.; Vecchi, M. P. Optimization by simulated annealing. *Science* **1983**, *220*, 671–680.
- (28) Chiu, C.-H.; Dewar, K.; Wagner, G. P.; Takahashi, K.; Ruddle, F. H.; Ledge, C.; Bartsch, P.; Scemama, J.-L.; Stellwag, E.; Fried, C.; Prohaska, S. J.; Stadler, P. F.; Amemiya, C. T. Bichir HoxA cluster sequence reveals surprising trends in ray-finned fish genomic evolution. *Genome Res.* **2004**, *14*, 2004.
- (29) Fried, C.; Prohaska, S. J.; Stadler, P. F. Independent Hox-cluster duplications in lampreys. *J. Exp. Zool., Mol. Dev. Evol.* **2003**, *299B*, 18–25.
- (30) Lyndon, R. C. Length functions in groups. *Math. Scand.* **1963**, *12*, 209–234.
- (31) Kececioğlu, J. D.; Sankoff, D. Exact and approximation algorithms for sorting by reversals, with application to genome rearrangement. *Algorithmica* **1995**, *13*, 180–210.
- (32) Needleman, S. B.; Wunsch, C. D. A general method applicable to the search for similarities in the amino acid sequences of two proteins. *J. Mol. Biol.* **1970**, *48*, 443–452.
- (33) Thompson, J. D.; Higgs, D. G.; Gibson, T. J. CLUSTALW: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position specific gap penalties, and weight matrix choice. *Nucl. Acids Res.* **1994**, *22*, 4673–4680.
- (34) Fagin, R.; Kumar, R.; Sivakumar, D. Comparing top-*k* lists. *SIAM J. Discr. Math.* **2003**, *17*, 134–160.
- (35) Aarts, E. H. L.; Korst, J. *Simulated Annealing and Boltzman Machines*; J. Wiley & Sons: New York, 1990.
- (36) Goldberg, D. E. *Genetic Algorithms in Search, Optimization and Machine Learning*; Addison-Wesley: Reading, MA, 1989.
- (37) Hardison, R. C.; Roskin, K. M.; Yang, S.; Diekhans, M.; Kent, W. J.; Weber, R.; Elnitski, L.; Li, J.; O'Connor, M.; Kolbe, D.; Schwartz, S.; Furey, T. S.; Whelan, S.; Goldman, N.; Smit, A.; Miller, W.; Chiaromonte, F.; Haussler, D. Covariation in frequencies of substitution, deletion, transposition, and recombination during Eutherian evolution. *Genome Res.* **2003**, *13*, 13–26.
- (38) Schadt, E. E.; Sinsheimer, J. S.; Lange, K. Applications of codon and rate variation models in molecular phylogeny. *Mol. Biol. Evol.* **2002**, *19*, 1550–1562.
- (39) Tajima, F. Simple methods for testing molecular clock hypothesis. *Genetics* **1993**, *135*, 599–607.
- (40) Huelsenbeck, J. P.; Nielsen, R. Variation in the pattern of nucleotide substitution across sites. *J. Mol. Evol.* **1999**, *48*, 86–93.
- (41) Nachman, M. W. N. Single nucleotide polymorphisms and recombination rate in humans. *Trends Genet.* **2001**, *17*, 481–485.
- (42) Hellmann, I.; Ebersberger, I.; Ptak, S. E.; Pääbo, S.; Przeworski, M. A neutral explanation for the correlation of diversity with recombination rates in humans. *Am. J. Hum. Genet.* **2003**, *72*, 1527–1535.
- (43) Gu, X.; Li, W. H. The size distribution of insertions and deletions in human and rodent pseudogenes suggests the logarithmic gap penalty for sequence alignment. *J. Mol. Evol.* **1995**, *40*, 464–473.
- (44) Behboudi, A.; Sjostrand, E.; Gomez-Fabre, P.; Sjöling, A.; Taib, Z.; Klinga-Levan, K.; Stahl, F.; Levan, G. Evolutionary aspects of the genomic organization of rat chromosome 10. *Cytogenet. Genome Res.* **2002**, *96*, 52–59.
- (45) Millwood, I. Y.; Bihoreau, M. T.; Gauguier, D.; Hyne, G.; Levy, E. R.; Kreutz, R.; Lathrop, G. M.; Monaco, A. P. A gene-based genetic linkage and comparative map of the rat X chromosome. *Genomics* **1997**, *40*, 253–261.
- (46) Schwikowski, B.; Speckenmeyer, E. On enumerating all minimal solutions of feedback problems. *Discr. Appl. Math.* **2002**, *117*, 253–265.

CI030411+