

tsize2

On the Design of Optimization Criteria for Multiple Sequence Alignment^{*}

Dannie Durand¹ and Martin Farach-Colton²

¹ Department of Biological Sciences, Carnegie Mellon University, Pittsburgh, PA 15213, USA

² Department of Computer Science, Rutgers University, Piscataway, NJ 08855, USA

Abstract. Multiple sequence alignment (MSA) is important in functional, structural and evolutionary studies of sequence data. While MSA construction has traditionally been an interactive process, the rapid growth of genetic sequence data has engendered a need for automated sequence analysis without human intervention. This requires more accurate methods based on rigorous mathematical models that reflect sequence biology in a realistic way. Focusing on MSA as an optimization problem, we examine the problem of unifying mathematical tractability with biological accuracy in cost function design. In particular, we consider tree alignment, which is often viewed as the most “biological” of the rigorous approaches to MSA. We point out several important pitfalls in current optimization approaches to MSA and identify characteristics for good cost function design. Design issues specific to approximation algorithms are also addressed. We hope these ideas will lead to future research on a biologically realistic and mathematically rigorous approach to MSA.

1 Introduction

Multiple sequence alignment (MSA) is a formal method for identifying shared features in a set of related nucleic acid or protein sequences, such as the sequences of a gene of interest (e.g., myoglobin) in several different species or the sequences of a family of related genes (e.g., myoglobin, the α -globins and the β -globins) in a single species. The goal is to align the sequences in such a way that biological relationships between sequence elements are revealed. MSA's are used to build evolutionary trees, extract information about the function and structure of proteins, identify patterns for data base searching and motif recognition and to design probes and primers for laboratory experiments. In some of these applications, such as evolutionary tree reconstruction, the MSA will be used as input to another program.

To state the MSA problem formally, we represent each biopolymer, P_1, P_2, \dots, P_k , as a string of symbols chosen from an alphabet, Σ , where $\Sigma = \{A, C, G, T\}$ for nucleic acid sequences and Σ contains the twenty amino

^{*} To appear in *Biological Evolution and Statistical Physics*, M. Lässig and A. Valeriani, Eds, Springer Verlag, 2002, pp. 22–36

acids for protein sequences. Then an MSA is a set of sequences $\hat{P}_1, \hat{P}_2, \dots, \hat{P}_k$ drawn from $\Sigma \cup \{-\}$, where $-$ is called a *blank*, such that P_i is the string derived by removing all blanks from \hat{P}_i . Additionally, each \hat{P}_i should be of the same length. We will call $\hat{P}_1[j], \dots, \hat{P}_k[j]$ the *jth column* of the MSA. The standard interpretation of blanks is that they represent mutations in the form of insertions or deletions, and therefore they are sometimes called *indels*. However, the sequences P_1, \dots, P_n need not be evolutionarily related, a fact which we will explore further below. As an example, an MSA of four amino acid sequence fragments, taken from human, rabbit, pig and chicken albumin, is shown in Fig. 1. Two indels have been inserted after the central valine in the first three sequences to reveal several conserved features in the alignment such as the KSE motif starting at column 31.

```

HUMAN   MKWVTFISLL FLFSSAYSRG V--FRRDA-H KSEVAHRFKD LGEENFKALV
RABBIT  MKWVTFISLL FLFSSAYSRG V--FRREA-H KSEIAHRFND VGEEHFIGLV
PIG     --WVTFISLL FLFSSAYSRG V--FRRDT-Y KSEIAHRFKD LGEQYFKGLV
CHICK   MKWVTLISFI FLFSSATSRN LQRFARDAEH KSEIAHRYND LKEETFKAVA

```

Fig. 1. A multiple alignment of four partial albumen sequences

Columns in an MSA should share a target biological feature. The feature sought depends on the application of the alignment. For example, if the multiple alignment is used to illustrate evolutionary relationships, then the residues in each column should have a shared evolutionary history. If the multiple alignment is used to determine structure or function, then the residues in each column should have a shared structural or functional role. It is exactly this kind of feature extraction that is the goal of MSA computation, the underlying assumption of which is that the sequences of a molecule contain enough information to extract the feature.

Numerous programs are available to construct MSA's, many of them based on heuristic approximations. A traditional approach has been to use such software to generate an initial multiple alignment and then adjust that alignment "by hand". Under this paradigm, highly accurate methods are not needed. However, as sequence data bases grow in the era of whole genome sequencing, the need to generate large numbers of MSA's automatically is also growing. Methods that can construct reliable MSA's without human oversight are required. Such methods must depend on formal models that accurately reflect the biology of sequences.

A number of formal, mathematical models of MSA have been proposed and a wealth of papers elaborating on these models has appeared in the theoretical computational molecular biology literature in the last ten to fifteen years. These approaches generally formulate MSA as an *optimization problem*, by defining a *cost function* over the set of all possible MSA's and then

seeking the MSA which optimizes¹ the cost. A multiple sequence alignment may be viewed as a hypothesis concerning the true relationship of a set of biopolymers, be that relationship evolutionary, structural or functional. A cost function is a method for evaluating the quality of such hypotheses and may be used to compare alternate alignments, whether generated by exact algorithms or by heuristics. The advantage of this optimization approach is that the cost function makes explicit the assumptions upon which the optimization is based. Because they are explicit, these assumptions are open to scrutiny and falsification.

The success of this approach depends on the assumption that it is possible to design a cost function such that the MSA with optimal cost is also the best explanation of the biological relationships between the sequences. In the face of genome scale data, research that attempts to unify biological and mathematical notions of accurate alignment is needed. The goal of the current paper is to examine the optimization approaches currently in use and identify problems that must be addressed on the path to robust, reliable MSA for automatic analysis of genomic data. In Section 2, we examine several popular cost functions that have received widespread attention in the literature, including sum-of-pairs, star and tree alignment. We focus primarily on tree alignment, the most biologically motivated of the three. We explore mathematical aspects of cost function design in Section 3. In particular, design issues related to approximation algorithms are discussed. Important biological considerations that must be addressed by any sound optimization criterion are elucidated in Section 4. Here we also specifically address the biological reasonableness of tree alignment for MSA computation. Biological validation of multiple sequence alignments is discussed in Section 5. We conclude that while optimization and approximation offers the best hope for well founded MSA computation, the current crop of optimization criteria are seriously wanting. We suggest some avenues of attack in correcting the situation.

2 Optimization Criteria

The desired output of an MSA computation is the column relationships of the alignment. Thus, it is natural to focus on optimization criteria which associate a score with each column, and most criteria in the literature take this approach. We will not give a comprehensive survey of multiple sequence alignment algorithms here. Such surveys can be found in [1] and in the introduction to [2]. Statistical approaches, such as hidden Markov models to MSA are discussed in [3]. These approaches are used primarily for local multiple alignments. An experimental comparison of multiple sequence algorithms is described in [4]. The column score expresses how well matched the residues

¹ We will interchangeably minimize and maximize the function in our examples, as needed.

in the column are. The score of the MSA is then the sum of the scores of its columns. We seek a set of \hat{P} maximizing

$$D(A) = \sum_j d(\hat{P}_1[j], \dots, \hat{P}_k[j]),$$

for column scoring function $d(\cdot)$. The choice of the column scoring function, $d(\hat{P}_1[j], \dots, \hat{P}_k[j])$, should reflect the application of the alignment. Residues that maximize $d(\hat{P}_1[j], \dots, \hat{P}_k[j])$ should belong together. Note that in our formulation, the scores of columns are completely independent. However, in some formulations, the cost of a gap might depend on the length of the gap in many ways. For example, for so-called affine gap penalties, the cost of a gap of length l is $\alpha + \beta l$ for some constants α and β . This separation of the gap penalty into α , a gap initiation cost, and β , a gap extension cost, reflects the assumption that whole pieces of sequence can be deleted in single events. Such gap functions typically are used for pairwise alignments and for multiple alignments that are heuristically constructed from pairwise alignments, and they will not be considered further here.

Three cost functions commonly used to evaluate each column are sum-of-pairs (SP), tree alignment (TA) and star alignment (SA), shown graphically in Fig. 2. The sum-of-pairs cost [5,6] of a column $\hat{P}_1[j], \dots, \hat{P}_k[j]$ is the sum of the costs of all unordered pairs in the column

$$d_{SP}(\hat{P}_1[j], \dots, \hat{P}_k[j]) = \sum_{p < q} \delta(\hat{P}_p[j], \hat{P}_q[j]),$$

for some binary cost function $\delta(x, y)$. This definition is mathematically natural but not biologically intuitive.

Tree alignment [7,8] is based on the assumption that the residues in the columns of the multiple sequence alignment share an evolutionary history and that this history can be expressed as a single tree for all columns. Under this model, a column is scored by computing the cost of the underlying tree. The score of the tree expresses the strength of our belief, under some model of evolution, that these residues are related in the manner described by this tree.

In order to use this approach, several issues must be resolved. First, a tree topology is needed. In general, the underlying tree is not known. In fact, multiple sequence alignments are generally used to estimate evolutionary trees and not vice versa. Second, it is generally assumed that every column in the alignment has the same underlying tree topology. As we shall see below, this may not always be the case. Third, in order to compute the branch costs of the tree, we need to know the ancestral sequences associated with the internal nodes. Fourth, given a tree topology with ancestral sequences at the nodes, a cost must be associated with each branch of the tree. This implies an underlying model of evolutionary change. The appropriate model will vary with the data.

Let us consider how to infer the ancestral sequences and compute the branch costs when the topology is known. The k extant sequences are associated with the k leaves of the tree. Sequences for the internal nodes are selected so that the tree is the best estimate of the true tree under some model of evolutionary change. The model in use is *maximum parsimony*, in which it is assumed that the true tree required a minimum number of evolutionary steps. Under this model, the cost of an edge (X_i, X_j) in the tree is the minimum number of mutations required to transform sequence X_i into sequence X_j . That is, internal nodes are selected such that the sum of edge costs, i.e. the total number of mutations required along the branches of the tree, is minimized. This assumption that evolution is parsimonious, is a subject of much debate.

The above approach assumes that the topology of the evolutionary tree is known, e.g. from morphological data. If the tree topology is not known, then the tree which yields the best score must be found. Since the number of trees with k leaves grows exponentially with k , this approach rapidly becomes prohibitively expensive if exhaustive enumeration is used. In general, no fast algorithm is known.

A variant on tree alignment is star alignment (SA), in which it is assumed that the underlying tree is a star. Restricting the topology makes this approach much more tractable, and it has been used as an underlying structure in many algorithms and analyses (e.g., [9,10]).

Altschul and Lipman [9] pointed out that SP, TA and SA costs for the same column can be very different, as shown in Fig. 2. SP overcounts mu-

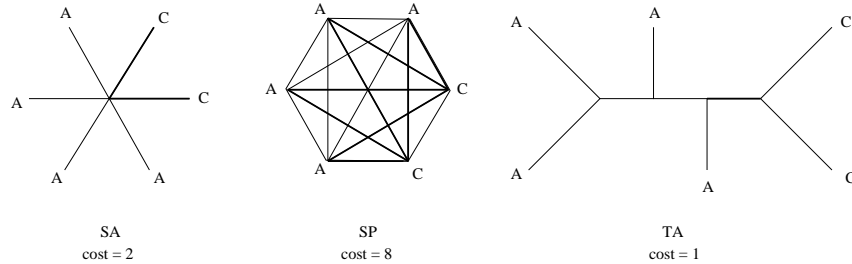


Fig. 2. Comparison of alignment costs for a multiple sequence alignment column, (A,A,A,A,C,C), under three different cost functions. Edges representing mutations are shown in bold

tations in a column because it considers that a separate mutation occurred between each pair in the column. In some sense, it is a “least parsimonious,” or profligate, model of mutation. In contrast, TA has been traditionally used with the maximum parsimony criterion. For most data sets the true answer is somewhere in between.

Efficiency considerations The MSA problem has been shown to belong to the class of *NP-complete problems* [11,12], for both sum-of-pairs and tree alignment. Problems in this class, which could be described colloquially as “officially hard problems,” have the property that the only known way to obtain an optimal solution to the problem is to generate all possible solutions and compare them. (See, for example, [13], for a discussion of NP-completeness.) Since the number of possible MSA’s increases exponentially with the number of sequences, the time required to find the optimal solution grows exponentially with the number of sequences as well. While no one has been able to prove that it is necessary to examine every possible solution to an NP-complete problem to find the optimum, no one has been able to solve any NP-complete problem exactly without enumerating all solutions. Since any NP-complete problem can be converted into any other NP-complete problem, a fast method to solve just one of these problems would provide a fast solution to all of them. Proving or disproving the existence of fast solutions to NP-complete problems is considered one of the most important unsolved mathematical challenges of our time. Since this problem has been the subject of intense scrutiny, without definitive resolution, by the theoretical computer science community for the last thirty years, we should not hope for a fast, exact MSA algorithm anytime soon.

In the face of this harsh reality, there are three approaches to MSA construction. For small data sets, typically eight to ten sequences of length at most 500 elements, it is possible to enumerate all alignments and select the optimal one. An optimal alignment of k sequences of length at most N can be obtained using dynamic programming in $O(2^k N^k)$ evaluations² of the SP cost function using $O(N^k)$ space [15]. An exact tree alignment algorithm for a given tree topology has been presented by Sankoff [7]. This requires $O(M(2N)^k)$ steps, where M is the number of internal nodes.

For larger data sets, one may turn to one of the numerous heuristic approaches that generate MSA’s with no optimality guarantee. Many of these are based on the *progressive alignment* approach, in which an optimal pairwise alignment is computed for every pair of sequences in the data set. These pairwise alignments are then merged to obtain a multiple alignment. Progressive alignment was first introduced by Taylor [16] and is the approach upon which the widely used CLUSTALW program [17] is based. The various progressive alignment methods differ in the rules used to determine in what order the pairwise alignments should be merged. Since the time required to compute an optimal pairwise alignment is $O(N^2)$ and the number of pairwise alignments is $O(k^2)$, the time required to compute a progressive alignment grows quadratically with the size of the input in comparison with the exponential time required to compute an optimal MSA.

² O -notation is used to describe the asymptotic behavior of a program in terms of the size of its input. Colloquially, $O(f(n))$ refers to a function whose behavior is proportional to $f(n)$ for large n . For a full treatment of O -notation, see [14].

Typically, it is possible to analyze up to 5000 sequences using a progressive alignment heuristic. However, progressive alignment will not, in general, give an optimal solution and will generate different solutions depending on the order in which the pairwise alignments are merged. This is illustrated in the multiple alignment of the sequences ACTCCAT, AGTCCT and ACGTCAT in Fig. 3. Consider a sum-of-pairs cost function where the cost of an

<p>(1) ACTCCAT (2) AGTCC-T</p>	<p>(1,2) + (2,3) (1) A-CTCCAT (2) A-GTCC-T (3) ACGTCA-T</p>
<p>(2) A-GTCCT (3) ACGTCAT</p>	<p>(1,2) + (1,3) (1) AC-TCCAT (2) AG-TCC-T (3) ACGTC-AT</p>

Fig. 3. Multiple alignment of the sequences CTCCAT, AGTCCT and ACGTCAT. Pairwise alignments these sequences are shown on the left. Two different progressive multiple alignments appear on the right, showing that the final form of the multiple sequence alignment depends on the order in which the pairwise alignments are merged

indel is $i = 2$ and the cost of a substitution is $s = 3$. The costs of the three optimal pairwise alignments shown in Fig. 3 are $s + i = 5$, $s + i = 5$ and $2i = 4$, respectively. The progressive alignment algorithm requires that we select one pairwise alignment as the seed of the multiple sequence alignment. At each step of the procedure, a new pairwise alignment is selected that contains one sequence not yet included in the growing MSA. The new sequence is added to the MSA using the pairwise alignment as a guide. The juxtaposition of the elements of the sequences in partial multiple alignment from the previous iteration must remain fixed. These elements may not be rearranged later to improve the score as additional sequences are added to the MSA. This restriction is responsible for both the increase in speed and the loss in accuracy associated with the progressive alignment heuristic. Fig. 3 shows two MSA's, both using the alignment of Sequences 1 and 2 as the seed. In the first MSA, the third sequence is added by aligning it with Sequence 2. In the second, the third sequence is aligned with Sequence 1. As the figure shows, two different MSA's result. The SP cost for the first alignment is $(s + i) + (s + i) + (2s + 2i) = 4s + 4i = 20$, while the second alignment costs $(s + i) + (s + 3i) + (2i) = 2s + 5i = 16$.

As the above example shows, heuristic approaches will not in general yield an optimal solution, nor indicate whether the solution presented differs from the optimum and, if so, how. *Approximation algorithms* [18] offer a promis-

ing compromise between exact, but computationally intractable methods and the heuristics currently most often used in practice. An a -approximation is a solution to an optimization problem that is guaranteed to have a score that differs from the optimal score by a factor of at most a . Approximation algorithms for both SP and tree alignment include those in [2,10,19–21]. Notably, a polynomial time approximation scheme, yielding a $(1+\epsilon)$ -approximation for arbitrarily small ϵ , has been presented by Jiang, Lawler and Wang [22,23] for tree alignment on a fixed topology. Approximation algorithms are powerful because they rigorously quantify the price to be paid by using an approximate method. However, to be practically useful, we need a better understanding of this price in *biological* terms. The work cited above has not, for the most part, investigated the biological implications of the approximation factor, a .

3 Approximations and Convergence

Since any reasonable cost function for MSA will probably be NP-hard to optimize [11,12], it is natural for researchers to look for approximation algorithms which find near optimal solutions to the MSA problem. These solutions will not be meaningful unless they satisfy (at least) the following two properties.

Smoothness: The success of the optimization approach to MSA depends on finding a cost function such that the optimal solution with respect to that cost function is close enough to the true solution. The definition of “true” solution depends on the application (e.g., reconstructing phylogeny, determining structure or function), as will be discussed in the next section. So the cost function should reflect the application in some way. Unlike problems such as the Traveling Salesman Problem, where the cost function (total distance traveled) is directly related to the problem under consideration (finding the shortest tour), mathematical cost functions used in MSA often have little relationship to the solution sought.

For approximation, stricter constraints on the properties of MSA cost functions are required. Intuitively, an MSA that is close in score to the optimal MSA should also be similar biologically. We call this notion *smoothness*. Unfortunately, “biologically similar” is a vague notion. It would be useful to make this notion concrete, for approximation algorithms are only meaningful if we understand the behavior of the cost function in the neighborhood of the optimal solution. What mathematical properties does a cost function require in order to exhibit smoothness? Under tree alignment, for example, it is not generally true that MSA’s with similar scores have similar biological interpretation, since it is possible for many different trees to have the same parsimony score.

Consistency: Biological wisdom holds that MSA computation should get easier as more sequences are introduced, and harder as the sequences get longer.

Roos [24] and McClure, Vasi and Fitch [4] have described specific instances of this behavior in their research. By contrast, phylogeny computation is held to get harder as the number of species increases but easier as the length of the sequences describing the species increases. Our optimization criteria should capture this varied dependence of problem instance difficulty on the size of the data set if they have any hope of capturing the underlying biology.

Ideally, the maximum of our optimization criterion should coincide with the “true” or biological alignment. A weaker notion, and one inspired by the considerations above, is that as the number of sequences being considered increases, the maximum of our optimization criterion should *converge* on the true alignment. This idea of converging on the true answer as we observe more data is basically the notion of *consistency* [25,26]. For concreteness, we define an algorithm to be *consistent* with respect to some parameter k of data size (e.g., number of species for MSA or sequence length for phylogeny) if it converges to the true solution as k goes to infinity. Similarly, a cost function is *consistent* if an algorithm that outputs its maximum is consistent. Notice that for MSA we would not expect a cost function to be consistent with respect to sequence length.

Suppose we have a cost function which is consistent. Would an approximation algorithm for the function be consistent? We consider a couple of examples. Suppose

$$g(\phi, d_k) = \begin{cases} 2 & \text{if } \phi \text{ is the true alignment for data } d_k \text{ on } k \text{ sequences;} \\ 1 & \text{otherwise.} \end{cases}$$

Now trivially, g is a consistent cost function. But a $1/2$ -approximation algorithm for g would be inconsistent, because *any* alignment yields a score at least half the maximum, so a $1/2$ approximation algorithm need not converge on the true answer. On the other hand, suppose

$$h(\phi, d_k) = \begin{cases} k & \text{if } \phi \text{ is the true alignment for data } d_k \text{ on } k \text{ sequences;} \\ 1 & \text{otherwise.} \end{cases}$$

Then any constant factor approximation algorithm for h will be consistent. This is because a $1-\epsilon$ approximation algorithm for h will be forced to output the true alignment for $k > 1/(1-\epsilon)$.

So we see that consistency for approximation algorithms requires more than an understanding of the biology of MSA, but also an understanding of the entire solution space for the cost function of choice. In particular, we need to know what happens to the set of alignments that are within a factor of $1-\epsilon$ of the best alignment, and how this set changes with k . We will say that a cost function f is ϵ -consistent if the set of all solutions within $1-\epsilon$ of the optimal for f converges to the true solution as k goes to infinity. We say f is *strongly consistent* if f is ϵ -consistent for every constant ϵ . In our example above, function g is ϵ consistent, for any $\epsilon < 1/2$, while function h is strongly consistent.

For phylogeny, a strongly consistent cost function was proposed in [27]. It was possible to do so in the case of phylogeny because the stochastic process relating phylogenies to DNA sequences is relatively well understood. We must seek such an understanding of MSA's before we can propose meaningful optimization criteria for this problem.

4 Evaluating Tree Alignment as a Model for MSA

How well do existing cost functions reflect biological relationships in MSA's? There is an unstated assumption that tree alignment is the gold standard, the ideal cost function, for multiple sequence alignment because the residues in each column are thought to be related by an evolutionary tree. The implication is that tree alignment is not used only because it is intractable. In this section, we present some examples that demonstrate that a tree is not always the appropriate model for multiple sequence alignment. In considering whether tree alignment is appropriate for a given set of sequences, two issues must be addressed:

- Is a tree the correct model for describing the relationship between residues in each column? A tree may not be a suitable model because the relationship between residues is functional or structural rather than historical. Even if the relationship is historical, for some data sets, no single tree will describe all columns in the alignment.
- What is the correct mutational model for scoring the branches of the tree? Tree alignment has historically been based exclusively upon the parsimony criterion. Data that does not happen to be parsimonious can favor the wrong tree model. In addition, column-oriented optimization approaches to MSA usually assume that sequence positions are independent and identically distributed. In general, these assumptions do not hold for biological sequence data.

The residues in a column do not share a common ancestor: When alignment is used to study function or structure, residues in a column do not always share a common ancestor. The goal is to align residues that share the same role. Although functional or structural residues usually share an evolutionary history, sometimes functional or structural roles can migrate to neighboring residues.

A possible example of shifting function occurs in dihydrofolate reductase (DHFR) gene – an important chemotherapeutic target in treating cancer and various infectious diseases. In their studies on protozoan parasites, Roos and colleagues have sought to design drugs that inhibit a metabolic protein in the parasite without affecting the infected host [28,29,24]. This requires identifying regions of structural or functional importance that differ substantially between protozoan and human versions of the DHFR protein.

Early sequence alignments placed the malaria parasite DHFR residue Phe²²³ downstream of structurally conserved regions of the protein (within the linker region which joins DHFR to thymidylate synthase, forming a bifunctional protein in protozoa and plants)[30,31]. This result was puzzling because mutational studies had suggested that Phe²²³ plays an important role in drug sensitivity. Realignment using sequences from the related parasite *Toxoplasma gondii* indicated that Phe²²³ is more likely homologous to a portion of the β -sheet which comprises the enzyme backbone [24]. The *T. gondii* sequence thus provided additional information that suggested an alternative alignment. Other protozoan sequences in the alignment have substantial, and different, nucleotide biases³. The *T. gondii* gene, which has relatively equal nucleotide distribution links the other protozoan sequences, facilitating alignment. This is another example of the observation, also made by McClure, Vasi and Fitch [4], that MSA's are very sensitive to sequence choice.

Phe²²³ is thought to play an indirect role in enzyme activity, interacting with His³⁴, within the active site [32]. The residue that plays this stabilizing role may have changed over time: a residue in a different position may provide this stabilizing effect in certain taxa (e.g. kinetoplastid parasites such as *Leishmania* and *Trypanosoma*) [33]. In this scenario, a random mutation allows a previously inactive residue to take on the functional role played by Phe²²³. In Roos' alignment, the residues currently thought to provide this stabilizing effect appear in the same column, but these residues in this column may not all share a common ancestor.

The tree is not unique: Another case where a tree is not an appropriate model occurs when the residues in any particular column share a common ancestor but the columns themselves have different evolutionary histories. A tree may describe any given column, but the columns taken together cannot be modeled by a single tree.

One situation where this occurs is exon shuffling. Most vertebrate genes consist of coding regions (exons) separated by DNA segments that are not translated into proteins (introns). The discovery of this intron/exon structure led to the theory of exon shuffling, first proposed by Gilbert in 1978 [34]. This theory posits that exons represent functionally and/or structurally important subunits of proteins, that introns occur at the boundaries of these modules and that proteins share and reuse the modules that exons encode.

The first evidence that the same exons appear in more than one gene was found in the human low-density lipoprotein (LDL) receptor gene [35,36]. The LDL receptor gene was shown to share exons with genes for epidermal growth factor, blood clotting factor IX and complementation factor C9. Since then,

³ The statistical profile of the primary sequence of genes can vary substantially, resulting in variations in, for example, the percentage of *GC* nucleotides. Such differences tend to obscure similarities between related sequences.

many such “mosaic” genes have been discovered [37–40]. In aligning mosaic genes, a different tree may be needed for each exon. If the exon boundaries are known, each exon could be aligned separately. However for many sequences, the splice sites have not been determined. Detecting such boundaries would require that the alignment already be known.

Residues with different evolutionary histories within a single gene can also occur due to horizontal gene transfer, the transfer of genetic material between species. For example, sequences similar to the Fn3 module in fibronectin have been found in bacterial proteins [41]. Fibronectin is a protein found in animals. Since Fn3 is found in both bacterial and animal proteins, one would expect to find Fn3 modules throughout the tree of life. However, Fn3 sequences have not been found in simpler eukaryotes, plants or fungi [41], suggesting a direct transfer of genetic material between bacteria and animals. Thus, in an alignment of genes containing the Fn3 sequence, one would not expect the residues in the Fn3 module to share an evolutionary history with other residues in the alignment. More than one tree is needed to model the sequence and, as in the case of exon shuffling, it may not be possible to know where the module boundaries occur. Other examples of mixing of genetic material possibly requiring more than one tree include transposition and gene conversion.

The tree is not parsimonious: Sequence data are not generally parsimonious, especially between distantly related sequences. Multiple substitution (e.g., $A \rightarrow C \rightarrow T$), coincidental substitution (e.g., $A \rightarrow C$ vs. $A \rightarrow G$), parallel substitution (e.g., $A \rightarrow C$ vs. $A \rightarrow C$) and back substitution (e.g., $A \rightarrow C \rightarrow A$), can all obscure the evolutionary history of a sequence.

Convergent evolution results in a situation where the parsimony criterion will lead to the wrong tree model. Residues that appear to be closely related may simply be similar due selective pressure because they perform the same function. An example of this occurs in cows and colobine monkeys, species that independently evolved foregut fermentation [42]. In order to maximize the nutritional benefits of foregut fermentation, the enzyme lysozyme had to evolve in these species to function in the acidic, pepsin-rich environment of the stomach. In other species, lysozyme is only needed in the intestines. As a result, lysozymes in cows and colobine monkeys exhibit amino acid substitutions not found in other species, suggesting, wrongly, that the two species are closely related. This would result in the wrong tree for those residues.

Tetraloops in rRNA provide another example of convergent evolution [43,44]. Tetraloops are strings of six bases that form loops at the end of helices in rRNA structures. The two end bases bond, allowing the internal four bases to form a loop. Although there are 256 possible inner loop sequences, only a small number actually occur in nature. Tetraloop sequences will tend to appear to be closely related, even when they are not.

Sequence data is not i.i.d.: Structural constraints prevent sequence data from being independent. Structural integrity depends on interactions between non-adjacent residues in the sequence. For example, α -helices are characterized by a heptad repeat, so that there are chemical interactions between every seventh residue in helical regions. Similarly, in order to maintain the structure of the RNA molecule, distant residues bind to form a structural interaction. Compensatory mutations between distant residues that form structural bonds are selectively favored.

Sequence data are not identically distributed either. Structural constraints on protein sequences result in variations in selective pressure at different positions, depending on whether they are located in α -helices, β -sheets or random coils and whether they have a role in tertiary structure or biochemical function. This fact has been recognized and exploited by researchers who developed structure-specific substitution matrices for recognizing specific secondary structures and motifs [45,46].

Additional variations in selective pressure occur at the DNA level. In protein coding regions, substitutions can be non-synonymous (resulting in an amino acid substitution in the protein coded for) or synonymous (resulting in a different codon for the same amino acid). Due to differences in selective pressure, synonymous changes are seen with more frequency than non-synonymous changes. Originally, it was thought that sequence positions could be classified as replacement sites, synonymous sites and non-coding sites and that mutation rates within each class would be relatively constant. More recently, evidence has emerged that suggests that selective pressure can vary within each class, even within a single gene or intron [47,48]

5 Biological Validation of MSA

As new optimization approaches to MSA construction are proposed, how do we determine if the resulting alignments are biologically sound? Structural and functional information has been used to validate multiple sequence alignments. In a comprehensive study of twelve different multiple alignment programs, McClure, Vasi and Fitch [4] measured algorithm performance by computing a numerical score based on the ability to find known motifs in four different data sets, for which supporting evidence from structural or mutational studies was available. Barton and Sternberg [49] have also used structural alignments to validate sequence-based alignments.

Structural information has also been used to guide the computation of MSA's. Some of the earliest work in structure-based alignments was presented by Lesk and Chothia [50], who used superposition of secondary structures to align globin sequences. Today, there are several common approaches to structural alignment, as surveyed in [51]. First, one may associate a secondary structure type (α -helix, β -sheet or random coil) with each residue and then impose the additional constraint that only residues associated with the same type

of secondary structure may be aligned (for example, [52,24]). This type of structural alignment is often done manually. Second, structural alignments may be performed by minimizing the root mean square distance between the aligned α -carbons in the backbone (see [53,54], for an example). A third approach is to minimize the difference between the distance matrices of the two proteins, where a distance matrix represents the distance between every pair of α -carbons in the protein (e.g., [55,56]). Structural alignments can also be evaluated by comparing the contact maps of the aligned proteins [57]. A contact map is a matrix describing the interactions of the amino acid side chains within the protein.

The use of structural information to guide or validate alignments is only possible for data sets where structural information is available. While the fraction of biopolymers for which structural information is available is small, the majority of newly sequenced biopolymers turn out to be members of families for which some structure is known. Another consideration is whether a structural approach to MSA is useful in all cases. For example, as discussed in Section 4, residues that share a structural or functional role do not always share an evolutionary history.

6 Conclusion

Optimization is an extremely effective tool for attacking many computational problems. In the case of problems from computational biology, one must be careful that the optimization criterion used closely follows some specific biological model of how the data relate to the underlying structure sought, e.g. a phylogeny or an MSA. In this vein, tree alignment has been touted as a “biological” approach to the MSA problem. However, as argued in Section 4, the definition of tree alignment (i.e. a single tree under a parsimonious scoring scheme) renders it meaningless in many biological settings. Furthermore, we argue in Section 3 that approximation algorithms for tree and other alignment criteria do not exploit biologically intuitive notions of convergence.

None of this need be the case. We are not pointing out fundamental shortcomings of the optimization/approximation approach in computational biology. Indeed, we believe that this approach offers the best hope for a principled attack on MSA and other problems. However, the strength of this approach, most notably a transparent association between underlying assumptions and computational methods, is rendered meaningless if the underlying assumptions are not biological. We therefore call for a more focused search for a biologically grounded model for MSA computation. Our paper points out many pitfalls of any such search, and suggests new directions for improvement. It should be considered a call to action.

7 Acknowledgements

The authors wish to thank Sergei Agulnik, Maja Bucan, Dan Gusfield, Laura Landweber, Eugene Myers, Dalit Naor, R. Ravi, David Roos, Ilya Ruvinsky, Mona Singh and, especially, the late Chris Overton for helpful discussions. D.D. was supported by NSF Grants BIR-94-13215 A01 and BIR-94-12594. M.F-C. was supported by NSF Career Development Award CCR-95-01942, NSF Grant BIR-94-12594, an Alfred P. Sloan Research Fellowship and NATO Grant 96-0215.

References

1. S. Chan, A. Wong, D. Chiu: *Bulletin of Mathematical Biology* **54**, 563–598 (1992)
2. P. Pevsner: *Journal of Applied Mathematics* **52**, 1763–1779 (1992)
3. R. Durbin, S. Eddy, A. Krogh, G. Mitchison: *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids* (Cambridge University Press, 1999)
4. M. McClure, T. Vasi, W. Fitch: *Mol. Biol. Evol.* **11**, 571 – 592 (1994)
5. D.J. Bacon, W.F. Anderson: *Journal of Molecular Biology* **191**, 153–161 (1986)
6. M. Murata, J.S. Richardson, J.L. Sussman: *Proc.Natl.Acad.Sci. USA* **82**, 3073–3077 (1985)
7. D. Sankoff: *Journal of Applied Mathematics* **28**, 443–453 (1975)
8. D. Sankoff, R.J. Cedergren: “Simultaneous Comparison of Three or More Sequences Related by a Tree”, in *Timewarps, Edits and Macromolecules: The Theory and Practise of Sequence Comparison* (Addison-Wesley, Reading, MA, 1983), pp. 253–258
9. S. Altschul, D. Lipman: *Journal of Applied Mathematics* **49**(1), 197–209 (1989)
10. D. Gusfield: *Bulletin of Mathematical Biology* **55**, 141–154 (1993)
11. E. Sweedyk, T. Warnow: (1992), “Manuscript”
12. L. Wang, T. Jiang: *Journal of Computational Biology* **1**(4), 337–348 (1994)
13. M.R. Garey, D.S. Johnson: *Computers and Intractability: A Guide to the Theory of NP-Completeness* (W. H. Freeman and Company, 1979)
14. T.H. Cormen, C.E. Leiserson, R.L. Rivest: *Introduction to Algorithms* (MIT Press/McGraw-Hill, 1990)
15. J.C. Setubal, J. Meidanis: *Introduction to Computational Molecular Biology* (PWS Publishing Company, Boston, 1997)
16. W.R. Taylor: *CABIOS* **3**(2), 81–7 (1987)
17. J.D. Thompson, D.G. Higgins, T.J. Gibson: *NAR* **22**(22), 4673–80 (1994)
18. D. Hochbaum: *Approximation Algorithms for NP-hard Problems* (PWS Publishing Company, Boston, 1997)
19. V. Bafna, E.L. Lawler, P. Pevzner: “Approximation Algorithms for Multiple Sequence Alignment”, in *5th Ann. Symp. On Pattern Combinatorial Matching* Vol. 807 (1994), pp. 43–53
20. R. Ravi, J.D. Kececioglu: “Approximation algorithms for multiple sequence alignment under a fixed evolutionary tree”, in *6th Ann. Symp. On Pattern Combinatorial Matching, Springer Verlag Lecture notes in Computer Science* (1995)

21. L. Wang, D. Gusfield: “Improved Approximation Algorithms for Tree Alignment”, in *7th Ann. Symp. On Pattern Combinatorial Matching* Vol. 1075 (1996), pp. 220–33
22. J. Jiang, L. Wang, E. Lawler: *Algorithmica* **16**, 302–15 (1996)
23. T. Jiang, E. Lawler, L. Wang: “Aligning sequences via an evolutionary tree: Complexity and approximation”, in *Proceedings of the Symposium on the Theoretical Aspects of Computer Science* (1994), pp. 760–769
24. D. Roos: *J. Biol. Chem.* **268**, 6269–6280 (1993)
25. J. Cavender: *Mathematical Biosciences* **40**, 271–280 (1978)
26. J. Felsenstein: *Syst. Zool.* **22**, 240–249 (1978)
27. M. Farach, S. Kannan: “Efficient algorithms for inverting evolution”, in *Proceedings of the Symposium on the Theoretical Aspects of Computer Science* (1996)
28. R.G. Donald, D.S. Roos: *Proc.Natl.Acad.Sci. USA* **90**, 11 703–11 707 (1993)
29. R.G.K. Donald, D.S. Roos: *Molec. Biochem. Parasitol.* **63**, 243–253 (1994)
30. J. Hyde: *Pharmacol Ther* **48**(1), 45–59 (1990)
31. M. Tanaka, H.M. Gu, D.J. Bzik, W.B. Li, J.W. Inselburg: *Mol Biochem Parasitol* **39**, 127–134 (1990)
32. M. Reynolds, D. Carter, M. Schumacher, D.S. Roos: “Personal communication”
33. D.S. Roos: “Personal communication”
34. W. Gilbert: *Nature* **271**, 501 (1978)
35. T.C. Sudhof, J.L. Goldstein, M.S. Brown, D.W. Russell: *Science* **228**, 815–822 (1985)
36. T.C. Sudhof, D.W. Russell, J.L. Goldstein, M.S. Brown, R. Sanchez-Pescador, G.I. Bell: *Science* **228**, 893–895 (1985)
37. R.L. Dorit, W. Gilbert: *Curr Opin Genet Dev* **1**, 464–469 (1991)
38. M.D.A. *et al.*: *Science* **287**(5461), 2185–9 (2000)
39. I.H.G.S. Consortium: *Nature* **409**(682), 860–921 (2001)
40. J.C.V. *et al.*: *Science* **291**(5507), 1304–51 (2001)
41. P. Bork, R.F. Doolittle: *Proc.Natl.Acad.Sci. USA* **89**, 8990–8994 (1992)
42. C.B. Stewart, A.C. Wilson: *Cold Spring Harbor Symposium on Quantitative Biology* **52**, 891–899 (1987)
43. R. Gutell, N. Larsen, C. Woese: *Microbiological Reviews* **58**(1), 10–26 (1994)
44. C.R. Woese, S. Winker, R.R. Gutell: *Proc.Natl.Acad.Sci. USA* **87**, 8467–8471 (1990)
45. R. Luthy, A.D. McLachlan, D. Eisenberg: *Proteins* **10**, 229–239 (1991)
46. P. Mehta, J. Heringa, P. Argos: *Protein Science* **4**, 2517–2525 (1995)
47. M. Kreitman, R.R. Hudson: *Genetics* **127**, 565–582 (1991)
48. S.W. Schaeffer, C.F. Aquadro: *Genetics* **117**, 61–73 (1987)
49. G. Barton, M. Sternberg: *Protein Engineering* **1**, 89–94 (1987)
50. A. Lesk, C. Chothia: *Journal of Molecular Biology* **136**, 225–270 (1980)
51. A. Godzik: *Protein Science* **5**, 1325–1338 (1996)
52. A. Aevarsson: *Journal of Molecular Evolution* **41**, 1096 – 1104 (1995)
53. A. Valencia, M. Kjeldgaard, E.F. Pai, C. Sander: *Proc.Natl.Acad.Sci. USA* **88**, 5443–5447 (1991)
54. G. Vriend, C. Sander: *Proteins* **11**(1), 52–58 (1991)
55. L. Holm, C. Sander: *Journal of Molecular Biology* (1993)
56. S. Pascarella, P. Argos: *Protein Engineering* **5**, 121–37 (1992)
57. A. Godzik, J. Skolnick, A. Kolinski: *Protein Engineering* **6**(8), 801–10 (1993)