

Using a β -Contact Predictor to Guide Pairwise Sequence Alignments for Comparative Modelling

Filippo Ledda¹, Giuliano Armano¹ and Andrew C.R. Martin^{*1,2}

¹Dipartimento di Ingegneria Elettrica ed Elettronica, University of Cagliari, Piazza d'Armi, I-09123, Cagliari, Italy

²Institute of Structural and Molecular Biology, Division of Biosciences, University College London, Darwin Building, Gower Street, London WC1E 6BT, UK

Email: Filippo Ledda - filippo.ledda@diee.unica.it; Giuliano Armano - armano@diee.unica.it; Andrew Martin* - andrew@bioinf.org.uk;

*Corresponding author

Abstract

Background: With the exponential rise in the number of available protein sequences, prediction of protein tertiary structure has become one of the most important tasks in bioinformatics; 'comparative', or 'homology', modelling is able to provide accurate models, but sequence alignment is a critical task. A strong correlation holds between the RMS deviation of models and the occurrence of errors in the alignment.

Results: In order to correct such errors, we developed BCAlign, based on an optimization procedure taking into account the correctness of the assignments of β -contacts, together with a standard scoring system. A β -contact evaluator (BCEval), based on a mixture of neural networks, is used to evaluate the assignments. Overall, compared with Needleman and Wunsch pairwise alignment, BCAlign improved alignments by 11.3% ('fraction of correct substitutions', FCS), on a set of 743 alignments of domains not showing any homology with the data used to train the evaluator. Three-dimensional models obtained from the alignments with the same proteins show an average RMSD improvement of 7.1%. On average, BCAlign results are comparable with multiple alignments obtained with MUSCLE (BCAlign improves FCS by 1.4%; RMSD is worse by 2.6%), but resulted in 42% of models having RMSD below 3Å, compared with 36% of models generated from a Needleman and Wunsch alignment and just 35% of models from a MUSCLE alignment. By choosing the 20% best-scoring alignments accord-

ing to the evaluator, models obtained with BCAlign provide a considerable improvement in the RMSD of about 10% over MUSCLE resulting in 48% of models having RMSD below 3Å, compared with 41% from Needleman and Wunsch and 39% from MUSCLE for the same set.

Conclusions: The evaluation of β -contacts has proved to be a useful measure in improving alignments for comparative modelling. An automatic procedure, BCAlign, which uses an iterative search strategy, has been developed to exploit a novel scoring scheme. The method shows significant improvements in the models generated, particularly where it has high confidence in the alignments generated. The method has been made available as a web server at <http://iasc.diee.unica.it/bcserver/> with a REST-style interface also available.

Introduction

The difference between the number of protein sequences translated from sequences held in GenBank [1] and the number of protein structures held by the PDB (Protein DataBank) [2] is vast. Only recently have high throughput methods started to be put in place to solve protein structure. Comparative modelling [3] offers a way to bridge the gap between the

number of sequences and structures.

Comparative modelling generally relies on knowing the structure of a homologous protein and using that as a template to build the structure of a protein. Methods include 3D-JIGSAW [4], FAMS [5], ESyPred3D [6], RAPPER [7]. COMPOSER [8, 9] and the particularly popular SwissModel [10, 11] and MODELLER [12–14].

However, the limiting factor in all these methods is obtaining the correct alignment. This is the most important stage of comparative modelling [15, 16], but unfortunately, particularly at low sequence identity, it can be the most difficult to get right. The sequence alignment one wishes to achieve is the alignment that would be obtained by performing a structural alignment and reading off the resulting sequence alignment. Of course the structure of the target is not available so one must rely on a sequence alignment. While multiple alignment can help, the sequence alignment can often differ substantially from the structural alignment.

There are numerous methods for performing structural alignment which often differ in the precise details of their results (e.g. CE [17], SSAP [18], STRUCTAL [19], DALI [20], MATRAS [21], VAST [22], SSM [23]). Since there are many different ways to superimpose two or more protein structures, if the proteins are not identical (or at least extremely similar), then there can be no single optimal superposition [24]. For our purposes, we have chosen SSAP as the gold standard, ‘correct’ alignment.

The most extreme types of misalignment (Misleading Local Sequence Alignments, MLSAs) are areas where the sequence alignment for a region is very clear, yet it does not match the structure-derived alignment [25]. We define less extreme misalignments, where the sequence and structural alignments do not agree, as SSMA (‘Sequence-Structure Misalignments’). For example, Figure 1 shows the sequence and structural alignment of a region from 1igmH00 and 1ap2A00 (a human and mouse antibody heavy chain variable region respectively) where an SSMA can clearly be seen.

In their analysis of the CASP2 comparative modelling section, Martin *et al.* [15] showed that there was a relationship between the percentage of correctly aligned residues and the sequence identity (Figure 2 of their paper). We have reproduced that analysis using approximately 56,000 pairs of homologous protein domains from CATH [26, 27], each of which was aligned on the basis of structure us-

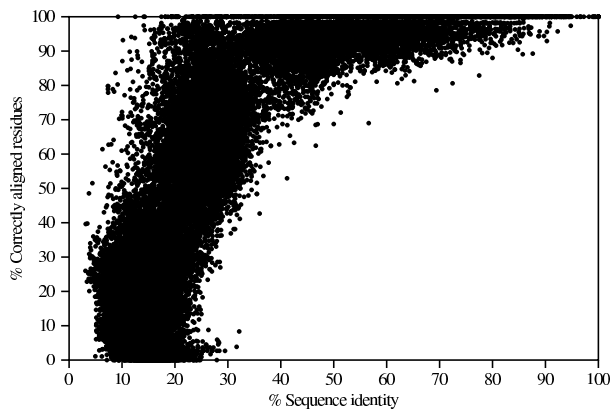


Figure 2: The relationship between the percentage correct sequence alignment and the percentage sequence identity. Each pair of NRep domains in each CATH homologous family has been structurally aligned by SSAP and sequence aligned using a Needleman and Wunsch global alignment. The structural alignment is taken as the correct alignment. Twelve outlying points have been removed after being identified as occurring owing to errors in the CATH database.

ing SSAP and on sequence using a Needleman and Wunsch sequence alignment [28]. Figure 2 clearly shows that if there is a high sequence identity between two sequences, then the sequence alignment is likely to match the structural alignment. However as sequence identity decreases, particularly below 30%, the accuracy of the alignment decreases and the sequence-based alignment can be completely different from the structural alignment. In this paper, we concentrate on improving the alignment in β -sheets and therefore hope to improve the models obtained.

Previous work by Lifson & Sander [29], Wouters & Curmi [30], Hutchinson *et al.* [31] and Fooks *et al.* [32] has shown clear residue pairing preferences between adjacent β -strands. With this in mind, we believe that some sequence mis-alignments can be detected and corrected by detecting errors in the assignment of β -contacts. Given a pair of β -strands (a ‘ β -pair’) assigned to a target from a template after initial sequence alignment, a measure of the likelihood of the register between the paired being formed in a real protein can be used as part of a scoring system of an alignment algorithm. Thus we developed BCEval, a β -contact evaluator based on a mixture

```

1ap2A00                               DIVMTQSPSSLTVTAGEKVTM
1igmH00 Sequence alignment            EVHLLSEGGNL-VQPGGSLRL
1igmH00 Structural alignment          EVHLLESG-GNLVQPGGSLRL
                                     ****

```

Figure 1: An example of an SSMA found between CATH domains 1igmH00 and 1ap2A00. The SSMA is indicated with asterisks.

of neural networks, able to predict whether a pair of β -strands is in the correct register. In addition, a pairwise sequence alignment method (BCAlign) has been developed able to take into account the β -contact evaluations. A search algorithm controlled by an iterative procedure had to be adopted to find the alignment instead of a classical dynamic programming technique such as Needleman and Wunsch. This is because the score of a substitution in the alignment will depend on the mutual register with another substitution along the sequence (because the register will affect the β -pairing), thus breaking the basic assumption of dynamic programming. In other words, while searching for the best alignment, the contacts of the parent template are assigned to the target; the scoring system then takes into account both of the assigned β -strands at the same time, so that the substitutions within a strand cannot be scored without taking into account the information about the neighbouring strand.

In this paper, we introduce both BCEval and BCalign. The accuracy of BCalign is assessed against (i) the standard Needleman and Wunsch pairwise sequence alignment, (ii) multiple alignments obtained with MUSCLE [33] and (iii) an equivalent of BCalign without the use of the evaluator (NoBCAlign). Additionally, the RMSD of models built using the different alignments is compared. The method has been made available as a web server at <http://iasc.diee.unica.it/bcserver/>.

Methods

When the homology modelling target and template sequences are aligned, the structural characteristics of the template are assigned to the target. Thus the secondary structure and the relative position within the structure (including interactions with other residues) are immediately known for the target sequence. A mis-alignment will lead to a wrong structural assignment. Thus we are able to examine contacts between residues in adjacent β -strands in

an attempt to detect mis-alignments using an evaluation of an assigned β -pair being correct based on machine learning (BCEval).

At first glance, including these evaluations in the scoring system of a typical dynamic programming algorithm seems straightforward. Unfortunately, the main dynamic programming assumption (that the optimal solution of the problems should depend on the optimal solution of its sub-problems) is broken. In order to overcome this limitation, we developed a technique which adopts a heuristic search algorithm (BCAlign).

Developing the β -Contact Evaluator (BCEval)

The evaluation of β -contacts can be tackled as a prediction problem, similar to contact map prediction. We must (i) define the training data, (ii) find a suitable representation of the input and output data and (iii) set up a proper architecture and learning algorithm(s). Methods were implemented using the GAME framework [34], written in Java 6.0.

Why not use a Generic Contact Map Predictor?

Generic contact predictors such as those by Cheng & Baldi [35] and Tegge *et al.* [36] have low accuracy owing to the difficulty of predicting all possible contacts occurring in a protein (including between α -helices). Even more specific predictors, specialized in β -contacts, report accuracies below 50% [37].

Fortunately, we already know which strands are in contact and we can concentrate on small shifts around a given position. Thus, we developed a new system specialized in recognizing a contact from the ‘shifted’ versions that could be identified from an alignment procedure.

Data Representation

The β -pairs must be represented in a fixed-length vector to obtain an input suitable for a neural net-

```

----AA-----AA-----BBBB-----CCCC-----CCCC-----BBBB-----
----12-----12-----1234-----12345-----54321-----4321-----
QSPVDIDTHTAKYDPSLKPLSVSYDQATSLRILNNGHAFNVEFDDSDQKAVLKGKGLDGT
CCCEEECCCEEECCCEEECCCEEECCCEEECCCEEECCCEEECCCEEECCCEEECCCEEECCCEEE

```

Figure 3: An example chain indicating the residues in contact. The letters in the first line indicate the β -strand pairs. The numbers in the second line indicate the residues in contact within the same pair. For example, the two residues labelled B1 form a contact.

work. The input vector must contain the residues of the two strands involved in the pairing and shifted versions of the same pair must be clearly recognizable.

Figure 3 shows an example in which the contacting residues belonging to different β -strands are indicated. While the length of the β -segments is variable, a fixed-length vector is needed for the data representation. A window of N residues would be perfectly suited to strands of length N , while information would be lost for pairings of longer strands and shorter strands would include residues not involved in contacts.

In addition, one must account of both parallel and anti-parallel strands. For instance, taking a window of four residues along the anti-parallel strands, B , in Figure 3, the encoding must indicate that the leucine at the first position in the first strand is in contact with the glycine in the last position of the second strand, not the valine in the first position. The different hydrogen-bonding patterns observed in parallel and anti-parallel sheets also result in different propensities in the contacts between residues, as shown by Hutchinson *et al.* [31] and Fooks *et al.* [32]. For these reasons, a ‘mixture of experts’ approach has been adopted: one expert only deals with strands of one type and length.

Profiles, obtained after three iterations of a PSI-BLAST search of the whole protein against *uniref90*¹, (inclusion threshold = 10^{-3} ; defaults for other parameters) were used to encode the residues in the window. A simple position-independent coding of the residues gave worse performance.

The Architecture

Figure 4 shows the architecture of BCEval. The ‘core evaluation module’ of BCEval consists of a mixture of 13 neural networks, each one specialized for

¹<http://www.ebi.ac.uk/uniref/>

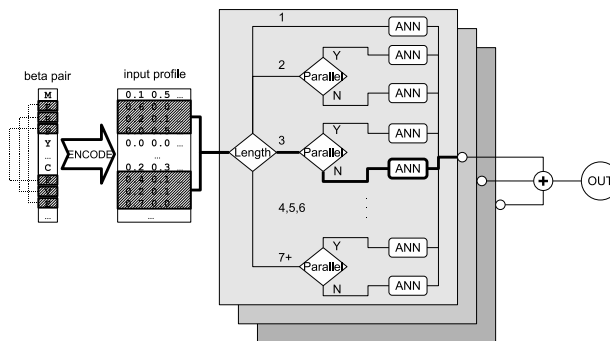


Figure 4: The BCEval architecture. The guarding functions ensure that only one neural network is activated at a time. The ‘parallel’ guard is able to distinguish between parallel and anti-parallel strand pairs, while the ‘length’ guard dispatches based on the length. In the example, an anti-parallel pair of length 3 is given so activating the path shown in bold. Three core units consisting of independently trained neural networks are averaged to obtain the final evaluation.

a specific length (1,2,3,4,5,6,7+) and type (parallel or anti-parallel) of β -pairing. The window length includes all (and only) the residues involved in each pairing, such that each neural network has a fixed-length vector as input, representing the residues involved in the contact. Simpler architectures with only one neural network and fixed input length (1, 2, 3) were tried first, but gave lower accuracy. The final output is obtained by averaging three core evaluation modules trained separately.

Training and Test Data Composition

A reference test set, TESTDOM, was built by selecting 10% of the total codes in the CATH database [26] at the homologue level and extracting the corresponding domains. A subset of the possible pairs of homologous domains in TESTDOM, mostly dis-

tant homologues (67% were below 30% sequence identity), was used to build a set of domain pairs. The resulting set, TESTALIGN, consists of 743 proteins, which have been used to test the alignment algorithms. In the same way, another set, TRAINALIGN, was obtained from the domains excluded from TESTDOM in order to train the parameters of the alignment algorithms. Finally, a set of protein chains, TRAINCH, consisting of protein chains from a dataset with identity $< 25\%$ ², selected in order not to include any chain containing the domains in TESTDOM, was used as a starting point to obtain contacts used in training BCEval; whole chains were used rather than domains in order to use DSSP [38] outputs from the EBI³ directly. Contacts included in the training of BCEval were obtained from TRAINCH DSSP files; these files include the position of the contacts between paired residues in β -strands. Negative examples (i.e. pairings not observed to be in a β -contact) were obtained by a synthetic sampling around the actual contacts. A Gaussian distribution ($\mu = 0, \sigma = \sqrt{5}$) around the positive examples was used to perform the sampling. The negative and positive samples were balanced, without taking into account the observed distribution. As seen in Figure 2, the extent of the observed shifts depends greatly on the sequence identity, making it hard to model the observed distribution correctly and, in any case, balanced inputs generally result in better learning. This partial synthetic sampling was preferred to sampling real data in order to obtain more, and more varied, samples. In practice, the negative data are randomly generated at each training iteration, so improving the diversity given to the training algorithm. All datasets are provided in Supplementary Material.

Training Technique and Parameter Setting

Each expert is a 3-layer feed-forward neural network, trained with a variant of back-propagation, with initial learning rate = 0.001 and momentum = 0.1. The learning rate is adjusted between iterations with an inverse-proportion law. The number of input neurons is $20N$, (where N is the size of the input window) and the number of hidden neurons is 75 for each neural network. A single output neuron indicates whether the given input is a contact or not. To

²http://bio-cluster.iis.sinica.edu.tw/~bioapp/hyprosp2/dataset_8297.txt

³<ftp://ftp.ebi.ac.uk/pub/databases/dssp/>

help the training algorithm avoid local minima, the training set was randomly shuffled at each iteration. Furthermore, each protein provides only a subset of its inputs, according to a random choice performed in accordance with a parameter, n . In particular, a random value k is generated in the range $[0, n - 1]$ and the inputs with index $k, k + n, k + 2n, \dots$ are provided to the learning algorithm.

To prevent the training process from stopping with a local oscillation of accuracy (evaluated on a validation set consisting of a 10% of TRAINCH, not used in the back-propagation process), weights are recorded when a minimum is encountered on the validation set, but the training continues until the error on the validation set increases for 10 consecutive iterations.

Developing the Pairwise Sequence Alignment (BCAlign)

The definition of an alignment algorithm includes two separate parts: (i) the cost function i.e. a scoring scheme used to evaluate an alignment; (ii) the alignment strategy, i.e. a strategy which gives the succession of substitutions, insertions and deletions which minimize the cost function. Here we describe a cost function which includes the evaluation of β -pairings made by BCEval and an alignment strategy suitable for use with the given cost function. Full details of both are provided in the Supplementary Material.

Defining the Cost Function

In brief, given a pairwise sequence alignment, A , between a template and a target sequence (the structure of the template being known), its cost,⁴ $c(A, S_{tpl})$, in BCalign consists of the sum of three main contributions (Equation 1).

$$c(A, S_{tpl}) = c_{nw}(A) + c_{\beta i}(A, S_{tpl}) + c_{bc}(A, S_{tpl}) \quad (1)$$

where S_{tpl} is the structure of the template sequence.

The component c_{nw} is the result of a classical similarity-based scoring scheme with affine gap penalties. The cost of the substitutions is obtained

⁴Note that scores, usually preferred in the scoring systems of sequence alignments, can also be viewed as the opposite of costs.

from a similarity scoring matrix M_s (e.g. BLO-SUM62), reversed so as to obtain a cost matrix $M_c = -(M_s - \max(M_s))$.

The term c_{β_i} in Equation 1 is related to the total number of gaps within the β -strands in the template. This is similar to the approach adopted in PIMA [39]. This component has been included in order to increase the number of β -pairs available for the evaluation: the larger costs help to avoid insertions and deletions inside β -strands, which rarely occur during evolution.

The last term in Equation 1, c_{bc} , results from the evaluation of β -pairs in the target sequence (assigned from the template, based on the alignment). This has two effects:

- to increase the cost of β -pairs that appear to be mistakenly assigned (i.e. shifted),
- to decrease the cost of β -pairs that appear to be assigned correctly.

The first of these changes the equilibrium of the alignment space, moving away from the solutions suggested by the other two terms that lead to wrong β -pair assignments. The second, although not directly improving the alignment, prevents drifts when correct assignments are found with the standard scoring scheme. The change of a pair of assignments may affect the solution in many different places within the alignment.

c_{bc} is composed of two elements summed over all β -pairs:

1. A term proportional to $-\tilde{p}(bp_{tgt})$, where $\tilde{p}(bp_{tgt})$ is the estimation of the probability of the β -pair bp_{tgt} being formed, given by BCEval (see Supplementary Material),
2. A term to stabilize the algorithm in the presence of wrong estimations which also takes account of the corresponding β -pair in the template (bp_{tpl}), for which the estimation error is known to be $1 - \tilde{p}(bp_{tpl})$. The requirement for this term derives from the assumption that, with the correct alignment, the errors in the estimation of p on the template and target are correlated. Therefore, if \tilde{p} for the template is significantly larger than for the target, we have a strong indicator of a probable error in the alignment (see Supplementary Material).

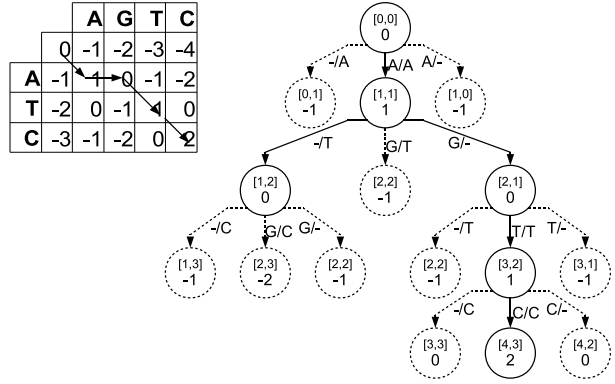


Figure 5: An example of alignment performed with a search algorithm. The search can be represented by a tree, in which the edge score is given by a simple scoring system (-1 for gaps, 1 match, -2 mismatch). Each circle represents a node, indicating the position in the two sequences and the path score. With a best-first search (i.e. the most promising nodes are opened first), the nodes shown with solid lines are expanded. In addition, nodes outside the solution path (in dashed lines) are explored, according to the local score. On the left, the corresponding Needleman and Wunsch matrix is indicated: note that the values in the Needleman and Wunsch matrix correspond to the scores of a node only when the best path to that node is followed.

Minimizing the Cost Function

Dynamic programming is generally used to minimize a cost function for sequence alignments and is the best choice when the optimal solution can be built incrementally by calculating the best solution for its sub-problems. With the proposed cost function, this assumption is broken, since the cost of a substitution is related to other substitutions along the sequences. A natural generalization of dynamic programming is represented by a search algorithm, which allows us to evaluate the path dynamically.

Using a global-search algorithm, the best alignment can be found by searching for the path in a tree which optimizes a score or cost function, leading to the end of the sequences. Figure 5 gives an example of a simple alignment performed with a best-first search strategy.

However, search algorithms may lead to an explosion in computational cost; in Figure 5, a blind (brute-force) search strategy is adopted, with the consequence that many nodes are expanded unnece-

essarily before finding the solution. The expected number of expanded nodes grows exponentially with the length of the path, which grows linearly with the length of the sequences. Consequently, to reduce the number of expanded nodes, heuristic search strategies, such as A* [40], can be adopted. A perfect heuristic (i.e. one which provides perfect estimates) for the components of the cost c_{nw} and $c_{\beta i}$ (Equation 1) can be obtained by adapting the approach used by dynamic programming algorithms. Hence, with only these two components, only the nodes in the optimal path are expanded by the A* algorithm, making this equivalent to a global dynamic programming approach. However, the component c_{bc} in Equation 1 cannot take advantage of any heuristic cost estimator and is relatively expensive to compute — a search algorithm computing this component dynamically would be computationally too expensive. Consequently, an iterative approach was adopted: after each iteration (the first being run without the component c_{bc}), the resulting β -pairs are collected and evaluated for use in the next iteration. The additional information is thus introduced step-by-step, permitting the algorithm, at each iteration, to escape from misleading pairings reached by following the other two components of the cost (c_{nw} and $c_{\beta i}$). The Iterative-Deepening A* (IDA*) [41] algorithm is used to perform the search. For further details, see Supplementary Material.

Evaluation Criteria

Two criteria were used to evaluate the results: (i) the fraction of correct substitutions (FCS) was measured by comparing the sequence alignment against a reference structural alignment obtained using SSAP [18], (ii) the RMSD of models generated from the alignments using MODELLER [12,14] in fully automatic mode with default parameters⁵. Fitting of models to the crystal structures was performed using the McLachlan algorithm [42] as implemented in the program ProFit (Martin, A.C.R., <http://www.bioinf.org.uk/software/profit/>).

For each experiment, we calculated the average RMSD of the models obtained as well as the percentage of ‘acceptable’ models, i.e. those with RMSD

⁵Only 637 models of the 743 of TESTALIGN were obtained from the alignments owing to problems in the automatic process which extracted the indexes for the domains from the PDB files. The problem is often caused by fragmented domains which include non-consecutive parts of sequence.

below 3Å, which is considered to be quite a strict criterion for distant homologues. The percentage of acceptable models is more indicative of the utility of the alignments than is the mean since the latter can easily be skewed by very bad models. In practical terms there is no difference between a ‘bad’ and a ‘very bad’ model.

An additional parameter, the ‘SSMA distance’ (SSMAD), defined as the mean distance of each residue from its correct position in the reference structural alignment, as used in our earlier work [15] was also tested, but was found to correlate less well with RMSD than the simpler FCS measure.

Results

BCEval

On average over a 7-fold cross validation on TRAINCH, BCEval achieved an accuracy of 0.785, precision of 0.771, recall of 0.811 and Matthews Correlation Coefficient of 0.571 (full results are shown in Supplementary Material, Table S1). In order to assess the use of BCEval in the evaluation of alignments, the correlation between the actual performance for a series of alignments and an evaluation metric from BCEval for that alignment was analyzed. This metric was the mean of the evaluations for the target protein: $\tilde{p}(bp_{tgt})$. Needleman and Wunsch alignments were also generated, scored with the BLOSUM45 matrix and gap opening/extension penalty of 13/1.

Figures 6 and 7 plot the BCEval metric against the fraction of correct substitutions (FCS) and the RMSD respectively. The existence of a considerable correlation between the scores and the alignment quality suggests that BCEval scores can be used effectively to choose the best alignment in a set and that β -pairs can be exploited to enhance pairwise sequence alignments.

BCAlign

Preliminary experiments were run on a subset of 500 domain pairs from TRAINALIGN to optimize parameters and the following were then used for all runs: $c_{go} = 22$, $c_{gx} = 9$, $c_{g\beta} = 6$, $\gamma_{abs} = 75$ and $\gamma_{rel} = 5$ (see Supplementary Material). Substitutions are scored using BLOSUM45; the algorithm is best suited to distant homologues since, for sequence alignments between close homologues, a standard se-

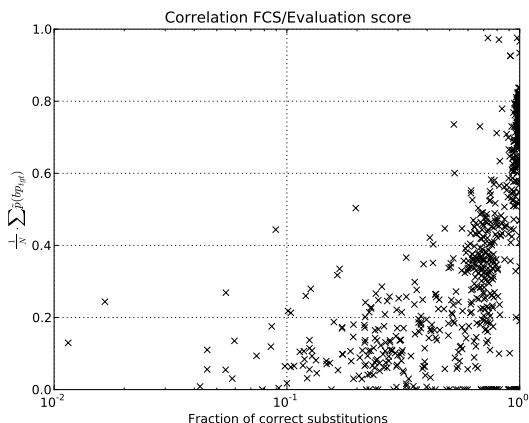


Figure 6: Plot of the β -contact predictor (BCEval) score *vs.* the ‘fraction of correct substitutions’ (FCS). Where BCEval scores zero, no β -pairs were assigned after the alignment because no contacts were present or because all were broken by gaps in the alignment.

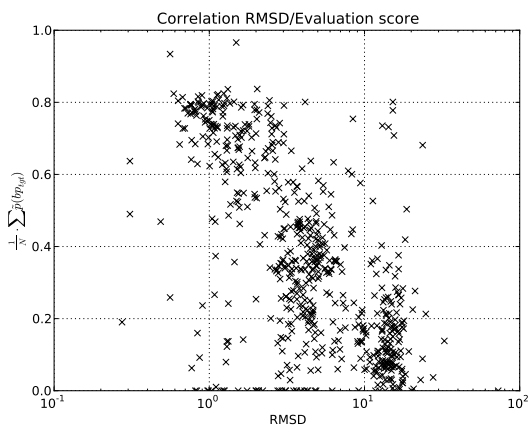


Figure 7: Plot of the β -contact (BCEval) score *vs.* RMSD (\AA) for three-dimensional comparative models generated using MODELLER.

quence alignment is usually sufficiently reliable, and very few SSMA are detected. The maximum number of iterations was set to 5 with a limit of one minute imposed on the search algorithm at each iteration using an Intel SU9600 CPU. The main code was written in Java and experiments were scripted using Python via the Jython 2.5 interpreter.

The performance of BCAlign was assessed in three comparisons: (i) with a Needleman and Wunsch alignment (scored using the BLOSUM45 matrix and gap opening/extension penalties 13/1, optimized as above), (ii) with multiple alignments obtained using MUSCLE [33], and (iii) with the same search technique, but without the use of the evaluator (‘NoBCAlign’) i.e. using optimized parameters as above, but setting $\gamma_{abs} = 0$ and $\gamma_{rel} = 0$. MUSCLE was run using standard parameters, including all the CATH homologous sequences contained in TESTDOM in the multiple alignments.

On the TESTALIGN dataset, BCAlign shows a relative improvement⁶ of 11.3% (0.628 *vs.* 0.703) in FCS compared with Needleman and Wunsch, 1.4% compared with MUSCLE and 6.2% compared with NoBCAlign. The RMSD improves by 7.14% (5.99 \AA *vs.* 6.43 \AA) compared with Needleman and Wunsch, and by 6.59% compared with NoBCAlign. However BCAlign performs slightly worse than MUSCLE (−2.6%) when assessed on RMSD, but see below. The large values of RMSD result from the fact that the majority of the alignments in the test set have sequence identity below 25%. In addition, as seen in Figure 7, a few models have extremely large RMSDs, skewing the mean value.

The percentage of acceptable models (i.e. with $\text{RMSD} < 3.0 \text{\AA}$) is probably a more useful measure of the success of an alignment method. In this experiment, this was 42% for BCAlign, 36% for Needleman and Wunsch, 35% for MUSCLE and 39% for NoBCAlign. Unexpectedly, multiple alignment using MUSCLE performed worst in this evaluation.

Better results are obtained by restricting comparisons to data for which we expect BCAlign to perform well, i.e. where a large number of β -pairings are present and the BCEval score improves. For structures with at least 8 β -pairs (58% of the alignments) the RMSD improvement is 1.22% over MUSCLE and 8.83% over NoBCAlign. The percentage of acceptable models improves to 48% for BCAlign, compared

⁶Relative improvements are calculated with $RI(a, b) = \frac{a-b}{(a+b)/2} \cdot 100$

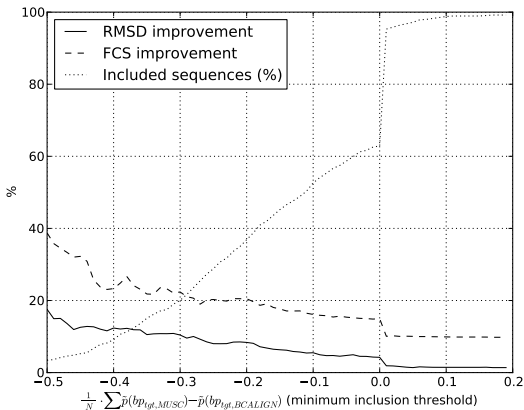


Figure 8: Average improvement in FCS and RMSD compared with MUSCLE at different inclusion thresholds. The threshold consists of the difference in the BCEval score between alignments obtained with BCAlign and MUSCLE. At each point in the plot, the alignments below the given threshold are included. The percentage of included alignments at each threshold is also shown.

with 41% for Needleman and Wunsch, 39% for MUSCLE and 44% for NoBCAlign, evaluating the same set of models.

In addition, the BCEval scores can be used to select those cases where BCEval makes confident predictions. Figures 8 and 9 show the average relative improvement in RMSD and FCS between BCAlign pairwise alignment and MUSCLE multiple alignment, when varying an inclusion threshold based on the improvement in the assignment of β -pairs when comparing MUSCLE and BCAlign alignments, as evaluated using BCEval. The graphs clearly show that, by using an inclusion threshold of less than -0.3 (thus including up to 20% of alignments), substantial improvements in FCS and RMSD can be obtained compared with other methods. For example, taking alignments with at least 8 β -pairings and an inclusion threshold of -0.3 (Figure 9), the percentage of proteins with RMSD lower than 3\AA is 39% for BCAlign, 21% for Needleman and Wunsch, 15% for MUSCLE and 26% for NoBCAlign.

Conclusions

Sequence alignment is the most critical task in comparative modelling: a strong correlation holds be-

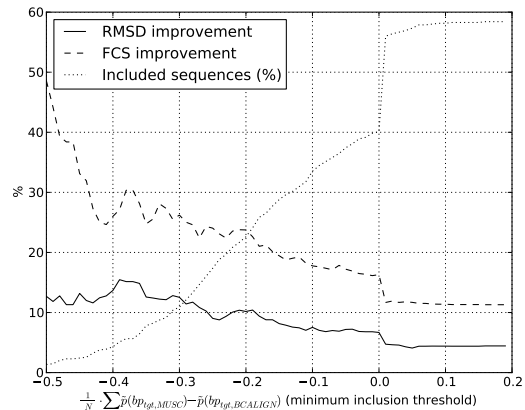


Figure 9: Average improvement in FCS and RMSD compared with MUSCLE at different inclusion thresholds, for proteins containing at least 8 β -pairs. The threshold consists of the difference in the BCEval score between alignments obtained with BCAlign and MUSCLE. At each point in the plot, the alignments below the given threshold are included. The percentage of included alignments at each threshold is also reported.

tween the RMS deviation of models and the occurrence of errors in the alignment. In order to improve alignments, we have exploited the likelihood of a given pairing between β -strands being correct. Since the location of β -strands is known for the template it can be assigned to the target sequence after the alignment. Our β -contact evaluator, BCEval, estimates the likelihood of assigned β -pairings occurring in real proteins by using a mixture of neural networks.

BCEval has then been exploited in a novel sequence alignment technique, BCAlign. We have presented a scoring system which combines a normal system based on a substitution matrix with BCEval. Since it is not possible to use standard dynamic programming with this scoring system, BCAlign resorts to a search algorithm, guided by an external loop to control the maximum run time.

Experiments confirm the validity of the approach: BCEval predictions show a considerable correlation with correct β -pair assignments and alignments obtained with BCAlign show that the evaluation of assigned β -pairs can be successfully exploited to enhance sequence alignments.

Overall, BCAlign showed a considerable improve-

ment compared with conventional pairwise Needleman and Wunsch alignment of 11.3% in FCS on a set of 743 alignments of domains not showing homology with the data used to train the evaluator. Three-dimensional models obtained from the alignments show an average RMSD improvement of 7.14% compared with standard Needleman and Wunsch sequence alignments. In addition, BCEval results are, on average, comparable with multiple alignments obtained with MUSCLE. However, Figures 8 and 9 show that choosing the 20% best-scoring alignments according to the evaluator, models obtained with BCAlign show a considerable improvement in the RMSD of about 10% over MUSCLE. The percentage of acceptable models shows an improvement of about 22% over MUSCLE when all proteins are considered and about 23% when only proteins containing at least 8 β -pairs are considered.

In conclusion, BCAlign appears to perform best when used in a mixed environment, in which different techniques compete while taking into account the scores assigned by BCEval. Restricting the use of BCAlign to those cases where BCEval makes the most confident predictions greatly increases its effectiveness. Even including the best 50% of the alignments shows BCAlign to be a good strategy (5% improvement over MUSCLE).

Finally, the implementation of the algorithms can probably be further improved. The computation is still not sufficiently efficient, frequently reaching the time limit for long sequences which, on average, will have more β -strands that can be exploited by the method and therefore are likely to show the best improvements. The search algorithm could be improved, particularly by enhancing the heuristic function to decrease the alternative paths that are explored. Alternatively, it may be possible to design a better control loop able to include the evaluations without overloading the search algorithm, or to use stochastic local search algorithms, including genetic algorithms.

Authors contributions

FL implemented the code, performed the experiments and drafted the paper; GA provided advice, discussions and access to resources, and enabled the collaboration; ACRM conceived the ideas, provided the data and completed the paper.

Acknowledgements

This work was supported by the Italian Ministry of Education – Investment funds for basic research, under the project ITALBIONET (Italian Network of Bioinformatics) and by the Visiting Professor Programme of the University of Cagliari.

References

1. Benson DA, Karsch-Mizrachi I, Lipman DJ, Ostell J, Rapp BA, Wheeler DL: **GenBank**. *Nuc. Ac. Res.* 2002, **30**:17–20.
2. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE: **The Protein Data Bank**. *Nuc. Ac. Res.* 2000, **28**:235–242.
3. Blundell TL, Carney D, Gardner S, Hayes F, Howlin B, Hubbard T, Overington J, Singh DA, Sibanda BL, Sutcliffe MJ: **Knowledge-based Protein Modelling and Design**. *Eur. J. Biochem.* 1988, **172**:513–520.
4. Bates PA, Sternberg MJ: **Model building by comparison at CASP3: using expert knowledge and computer automation**. *Proteins: Struct., Funct., Genet.* 1999, **37**:47–54.
5. Ogata K, Umeyama H: **An automatic homology modeling method consisting of database searches and simulated annealing**. *J. Mol. Graph.* 2000, **18**:305–306.
6. Lambert C, Léonard N, De Bolle X, Depiereux E: **ESyPred3D: Prediction of Proteins 3D Structures**. *Bioinformatics* 2002, **18**:1250–1256.
7. DePristo MA, De Bakker PIW, Shetty RP, Blundell TL: **Discrete Restraint-based Protein Modeling and the Calpha-trace Problem**. *Protein Sci* 2003, **12**:2032–2046.
8. Sutcliffe MJ, Haneef I, Carney D, Blundell TL: **Knowledge based modelling of homologous proteins. 1. Three-dimensional frameworks derived from simultaneous superposition of multiple structures**. *Protein Eng.* 1987, **1**:377–384.
9. Sutcliffe MJ, Hayes FRF, Blundell TL: **Knowledge based modelling of homologous proteins. 2. Rules for the conformations of substituted side chains**. *Protein Eng.* 1987, **1**:385–392.
10. Peitsch MC: **ProMod and Swiss-Model: Internet-based tools for automated comparative modelling**. *Biochem. Soc. Trans. (London)* 1996, **24**:274–279.
11. Arnold K, Bordoli L, Kopp J, Schwede T: **The SWISS-MODEL Workspace: A web-based environment for protein structure homology modelling**. *Bioinformatics* 2006, **22**:195–201.
12. Šali A, Blundell TL: **Comparative protein modelling by satisfaction of spatial restraints**. *J. Mol. Biol.* 1993, **234**:779–815.
13. Sánchez R, Šali A: **Evaluation of Comparative Protein Structure Modeling by MODELLER-3**. *Proteins: Struct., Funct., Genet.* 1997, **1**:50–58.
14. Fiser A, Do RK, Šali A: **Modeling of loops in protein structures**. *Protein Sci.* 2000, **9**:1753–1773.

15. Martin ACR, MacArthur MW, Thornton JM: **Assessment of comparative modeling in CASP2.** *Proteins: Struct., Funct., Genet.* 1997, **Suppl. 1**:14–28.
16. Cristobal S, Zemla A, Fischer D, Rychlewski L, Elofsson A: **A Study of Quality Measures for Protein Threading Models.** *Bioinformatics* 2001, **2**:5–5.
17. Shindyalov IN, Bourne PE: **Protein Structure Alignment by Incremental Combinatorial Extension (CE) of the Optimal path.** *Protein Eng.* 1998, **11**:739–747.
18. Taylor WR, Orengo CA: **Protein Structure Alignment.** *J. Mol. Biol.* 1989, **208**:1–22.
19. Subbiah S, Laurents DV, Levitt M: **Structural Similarity of DNA-binding Domains of Bacteriophage Repressors and the Globin core.** *Curr. Biol.* 1993, **3**:141–148.
20. Holm L, Sander C: **Protein Structure Comparison by Alignment of Distance Matrices.** *J. Mol. Biol.* 1993, **233**:123–138.
21. Kawabata T: **MATRAS: A Program for Protein 3D Structure Comparison.** *Nuc. Ac. Res.* 2003, **31**:3367–3369.
22. Gibrat JF, Madej T, Bryant SH: **Surprising Similarities in Structure Comparison.** *Curr. Opin. Struct. Biol.* 1996, **266**:540–553.
23. Krissinel E, Henrick K: **Secondary-structure Matching (SSM), a new tool for fast Protein Structure Alignment in Three Dimensions.** *Acta Crystallogr.* 2004, **60**:2256–2268.
24. Novotny M, Madsen D, Kleywegt GJ: **Evaluation of Protein fold Comparison Servers.** *Proteins: Struct., Funct., Genet.* 2004, **54**:260–270.
25. Saqi MAS, Russell RB, Sternberg MJE: **Misleading local sequence alignments: implications for comparative modelling.** *Protein Eng.* 1998, **11**:627–630.
26. Orengo CA, Michie AD, Jones S, Jones DT, Swindells MB, Thornton JM: **CATH—a Hierarchic Classification of Protein Domain Structures.** *Structure* 1997, **5**:1093–1108.
27. Pearl F, Todd A, Sillitoe I, Dibley M, Redfern O, Lewis T, Bennett C, Marsden R, Grant A, Lee D, Akpor A, Maibaum M, Harrison A, Dallman T, Reeves G, Diboun I, Addou S, Lise S, Johnston C, Sillero A, Thornton J, Orengo C: **The CATH Domain Structure Database and Related Resources Gene3D and DHS Provide Comprehensive Domain Family Information for Genome Analysis.** *Nuc. Ac. Res.* 2005, **33**:D247–D251.
28. Needleman SB, Wunsch CD: **A general method applicable to the search for similarities in the amino acid sequence of two proteins.** *J. Mol. Biol.* 1970, **48**:443–453.
29. Lifson S, Sander C: **Specific Recognition in the Tertiary Structure of Beta-sheets of Proteins.** *J. Mol. Biol.* 1980, **139**:627–639.
30. Wouters MA, Curmi PM: **An Analysis of side Chain Interactions and pair Correlations Within Antiparallel Beta-sheets: the Differences Between Backbone Hydrogen-bonded and Non-hydrogen-bonded Residue Pairs.** *Proteins* 1995, **22**:119–131.
31. Hutchinson EG, Sessions RB, Thornton JM, Woolfson DN: **Determinants of Strand Register in Antiparallel Beta-sheets of Proteins.** *Protein Sci* 1998, **7**:2287–2300.
32. Fooks HM, Martin ACR, Woolfson DN, Sessions RB, Hutchinson EG: **Amino acid Pairing Preferences in Parallel Beta-sheets in Proteins.** *J. Mol. Biol.* 2006, **356**:32–44.
33. Edgar RC: **MUSCLE: a Multiple Sequence Alignment Method with Reduced time and Space Complexity.** *BMC Bioinformatics* 2004, **5**:113–113.
34. Ledda F, Milanese L, Vargiu E: **GAME: A Generic Architecture based on Multiple Experts for Predicting Protein Structures.** *International Journal Communications of SIWN* 2008, **3**:107–112.
35. Cheng J, Baldi P: **Improved Residue Contact Prediction Using Support Vector Machines and a Large Feature set.** *BMC Bioinformatics* 2007, **8**:113–113.
36. Tegge AN, Wang Z, Eickholt J, Cheng J: **NNcon: Improved Protein Contact map Prediction Using 2D-recursive Neural Networks.** *Nucleic Acids Res* 2009, **37**:W515–W518.
37. Lippi M, Frasconi P: **Prediction of protein - residue contacts by Markov logic networks with grounding-specific weights.** *Bioinformatics* 2009, **25**(18):2326–2333, [<http://bioinformatics.oxfordjournals.org/content/25/18/2326.abstract>].
38. Kabsch W, Sander C: **Dictionary of Protein Secondary Structure: Pattern Recognition of Hydrogen-bonded and Geometrical Features.** *Biopolymers* 1983, **22**:2577–2637.
39. Smith RF, Smith TF: **Pattern-induced multi-sequence alignment (PIMA) algorithm employing secondary structure-dependent gap penalties for use in comparative protein modelling.** *Protein Eng.* 1992, **5**:35–41.
40. Hart PE, Nilsson NJ, Raphael B: **A Formal Basis for the Heuristic Determination of Minimum Cost Paths.** *Systems Science and Cybernetics, IEEE Transactions on* 1968, **4**(2):100–107.
41. Korf R: **Depth-First iterative-deepening: An optimal admissible tree search.** *Artificial Intelligence* 1985, **27**:97–109.
42. McLachlan A: **Rapid Comparison of Protein Structures.** *Acta Cryst* 1982, **A38**:871–873.