



## Meeting Review

## Workshop — Predicting the Structure of Biological Molecules

Isaac Newton Institute, Centre for Mathematical Sciences, Cambridge University, UK, 26–30 April 2004

Damian Counsell\*

Rosalind Franklin Centre for Genomics Research, Wellcome Trust Genome Campus, Cambridge CB10 1SB, UK

\*Correspondence to:

Dr Damian Counsell, RFCGR,

Wellcome Trust Genome

Campus, Cambridge, CB10 1SB,

UK.

E-mail: d.counsell@rfcgr.mrc.ac.uk

### Abstract

This April, in Cambridge (UK), principal investigators from the Mathematical Biology Group of the Medical Research Council's National Institute of Medical Research organized a workshop in structural bioinformatics at the Centre for Mathematical Sciences. Bioinformatics researchers of several nationalities from labs around the country presented and discussed their computational work in biomolecular structure prediction and analysis, and in protein evolution. The meeting was intensive and lively and gave attendees an overview of the healthy state of protein bioinformatics in the UK. Copyright © 2004 John Wiley & Sons, Ltd.

**Keywords:** protein bioinformatics; structural genomics; protein structure prediction; molecular evolution; molecular simulation; structural bioinformatics

Received: 10 June 2004

Revised: 18 June 2004

Accepted: 21 June 2004

1 This workshop was organized by members of  
2 the Computational Biology Group at the UK  
3 Medical Research Council's National Institute for  
4 Medical Research (<http://mb2.nimr.mrc.ac.uk/>),  
5 Franca Fraternali, Richard Goldstein and Willie  
6 Taylor. It was a low-key affair, organized late, yet it  
7 was probably the best scientific meeting I have ever  
8 attended; I was interested in advance in the content  
9 of practically every session. Most of the seminars  
10 were well-prepared, clear, relevant and refreshingly  
11 concise. Even allowing for usually well-informed  
12 questions and interruptions, sessions rarely over-  
13 ran (or if they did, it didn't feel that way). Unfor-  
14 tunately, because I heard about the meeting only  
15 shortly before it took place, I was unable to attend  
16 every presentation in full. Although the speakers  
17 and attendees were of many nationalities, they are  
18 all currently working in the UK.

19 After Willie Taylor's introduction it was appro-  
20 priate that **Cyrus Chothia** (Laboratory of Molec-  
21 ular Biology, Cambridge, UK; [http://www.mrc-  
22 lmb.cam.ac.uk/genomes/Cyrus.html](http://www.mrc-lmb.cam.ac.uk/genomes/Cyrus.html)), one of the  
23 most prominent and longstanding researchers in  
24

25 the field of protein structure bioinformatics in  
26 Cambridge, should open proceedings. In his talk,  
27 'Structural constraints on protein mutations', he  
28 described his work with Rajkumar Sasidharan  
29 (<http://www.mrc-lmb.cam.ac.uk/genomes/sraj/>).  
30 It was good for the rest of the meeting that, regard-  
31 less of Chothia's standing, there was no reluctance  
32 to challenge his arguments and his contribution  
33 provoked the first of many stimulating discussions  
34 that took place both during and after presenta-  
35 tions.

36 When questioned, Chothia admitted to an inten-  
37 tional looseness with the term 'positive selection'  
38 in his description of the degree and type of residue  
39 type conservation in different locations in protein  
40 structures. He outlined how residue conservation  
41 varied with degree of site exposure and summa-  
42 rized the residue properties most likely to be shared  
43 in the same sites across homologues. Among the  
44 intriguing statistics he presented, Chothia noted  
45 that the normalized frequency of changes in sur-  
46 face residues was five to six times higher than core  
47 residues. The most 'selected' (conserved) residue  
48

1 positions were least likely to vary in their size  
2 first, followed by their physicochemical character.  
3 For the least 'selected' positions, the priorities were  
4 reversed.

5 In summary: average selectivity values for given  
6 sites in proteins are calculable, the frequency of  
7 variation can be explained in terms of the prop-  
8 erties and locations of the analysed sites, and the  
9 frequency with which residues vary at given sites  
10 had a medium correlation with the overall under-  
11 lying frequency of random mutations. Richard  
12 Goldstein asked if Sasidharan and Chothia's study  
13 showed that proteins tended towards robustness  
14 and Chothia admitted not. Willie Taylor asked  
15 about possible resemblances between the substitu-  
16 tion matrix derived from Chothia's structure-based  
17 alignments and the Dayhoff matrix. Chothia said  
18 that the two were similar.

19 **Juan Fernandez-Recio** (Crystallography and  
20 Biocomputing Unit, Department of Biochemistry,  
21 University of Cambridge; [http://www-cryst.bioc.  
22 cam.ac.uk/~juan/](http://www-cryst.bioc.cam.ac.uk/~juan/)) next presented 'Protein  
23 –protein docking by global energy minimization',  
24 work that he began in Ruben Abagyan's lab at the  
25 Scripps Institute [9]. He aims to find general meth-  
26 ods for predicting the structures of protein–protein  
27 complexes based solely on the structures of the  
28 members of those complexes. Useful because struc-  
29 tures of complexes are hard to determine, this  
30 has increased in importance as lower-resolution  
31 structure determination methods have become more  
32 powerful and generated more data.

33 While cheaper analyses treating docking part-  
34 ners as rigid bodies are easier to calculate, they  
35 produce unrealistic energy landscapes, unlikely to  
36 lead to even approximately correct solutions. Mod-  
37 els including fully flexible protein structure require  
38 the exploration of huge conformational spaces.  
39 Fernandez-Recio and co-workers seek a compro-  
40 mise: the first step is to treat structures as rigid  
41 bodies with 'soft' van der Waals' radii permitting  
42 atomic overlap; the second step is to permit flex-  
43 ibility elsewhere. Other efficiencies come through  
44 representing molecules in terms of their internal,  
45 rather than Cartesian, coordinates. This combina-  
46 tion resulted in one of the top 'blind' performers in  
47 the (CASP-like) CAPRI protein docking prediction  
48 competition (<http://capri.ebi.ac.uk/>). Unlike many  
49 reviews of bioinformatics methods by their devel-  
50 opers, Fernandez-Recio went on to give examples

51 of both the successful and unsuccessful applica- 52  
53 tions of his approach. He also discussed some of 54  
55 the other uses for the output of his docking sim- 56  
57 ulations — they can be used, for example, in the 58  
59 prediction of binding patches on proteins. 60

61 Everyone I spoke to was especially impressed 62  
63 with the volume and depth of analysis that had been 64  
65 performed by **Sanne Abeln** ([http://www.stats.ox.  
66 ac.uk/people/students.htm](http://www.stats.ox.ac.uk/people/students.htm)), still a first-year stu- 67  
68 dent in Charlotte Deane's bioinformatics group in 69  
70 the Statistics Department at Oxford University. In 71  
72 'Fold usage on genomes and protein structure evo- 73  
74 lution' she described her huge survey of protein 75  
76 structures across species. She compared the num- 77  
78 ber of distinct folds with genome size, examined 79  
80 the number of occurrences of folds, 'duplications' 81  
82 of folds, and families per fold and related them. 83  
84 She had asked what these data could say about 85  
86 the 'ages' of folds, evolutionary mechanisms and 87  
88 evolutionary relationships between folds. By tak- 89  
90 ing large sequence sets (150 + genomes from all 91  
92 kingdoms) and widely used bioinformatics tools 93  
94 (PSI-BLAST and SCOP), and applying them on 95  
96 a large scale, she not only made too many interest- 97  
98 ing observations to list here, but had already begun 99  
100 to devise plausible explanations for many of the 101  
102 phenomena she observed.

It seems that distributions of the popularity of 79  
80 folds are often described by power laws. Some 81  
82 folds at least appear to be missing in certain 83  
84 genomes. The data she collected for  $\alpha\beta$  proteins 85  
86 are different from folds in the other fold classes 87  
88 (similar comparisons against  $\alpha\beta$  proteins were 89  
90 made at several points over the course of the 91  
92 meeting). Abeln cautioned that it is very difficult 93  
94 to make phylogenetic trees from this kind of data 95  
96 since: 97

- There are no clear relations between the different 98  
99 measures of fold usage (i.e. occurrences of 100  
101 folds across genomes, duplications of folds on 102  
103 a genome, and families per fold). 104
- When a fold diverges to a new fold on one 105  
106 genome, occurrence and duplications are set 107  
108 back to one, and it is therefore difficult to obtain 109  
110 evolutionary relations between folds from these 111  
112 measures. 113

114 Interesting power law-based relations also emer- 115  
116 ged from their analyses of fold distributions across 117  
118 families and superfamilies. Just as there had been 119  
120 discussion of Chothia's use of the term 'positive 121  
122

1 selection', there was some debate over Abeln's  
2 allusions to 'old folds' in her discussion of the  
3 possible evolution of folds. The idea of 'trapped  
4 folds' having difficulty evolving was another theme  
5 which re-emerged later in the week, when Ben  
6 Blackburne described his hugely simplified *in silico*  
7 minimal proteins.

8 The second day was chaired by Richard Gold-  
9 stein and the first speaker, **Kenji Mizuguchi**  
10 (Department of Biochemistry, University of Cam-  
11 bridge, UK, [http://www-cryst.bioc.cam.ac.uk/~](http://www-cryst.bioc.cam.ac.uk/~kenji/)  
12 [kenji/](http://www-cryst.bioc.cam.ac.uk/~kenji/)), addressed 'Sequence-structure homology  
13 recognition'. Mizuguchi first clearly described the  
14 central problems of homology modelling: identify-  
15 ing the best structural templates against which to  
16 model the sequence of an unknown fold and find-  
17 ing the best alignment between that sequence and  
18 its target. He was classically biological in his use of  
19 terminology, distinguishing between the identifica-  
20 tion of analogous (corresponding, but not related)  
21 and truly homologous (corresponding and related)  
22 folds.

23 After an overview of existing methods for fold  
24 recognition and alignment he outlined FUGUE  
25 (<http://www-cryst.bioc.cam.ac.uk/fugue/>), a sys-  
26 tem he developed along with Jiye Shi and Tom  
27 Blundell [18]. FUGUE exploits structural data  
28 in the form of environment-specific substitution  
29 tables — 64 of them — and gap penalties. These  
30 are applied alongside modern sequence align-  
31 ment techniques and refined by testing to see  
32 how the environment definitions affect perfor-  
33 mance. Mizuguchi claimed 70–100 hits/day on  
34 the FUGUE Website and that the method out-  
35 performs other blind prediction servers in align-  
36 ment/assignment. Unfortunately, Mizuguchi's clear  
37 explanation of the problems and approach didn't  
38 leave him time to discuss other applications, but  
39 I look forward to reading about them elsewhere  
40 [19,20,22]. It was also satisfying during question-  
41 ing afterwards to hear him be sensibly dismissive of  
42 any attempt to attach statistical confidence values  
43 to FUGUE's output, given the absence of an under-  
44 lying mathematical model. For a wider view of the  
45 importance of fold recognition, he recommended  
46 his review in *Drug Discovery Today* [17].

47 **Franca Fraternali** ([http://mathbio.nimr.mrc.](http://mathbio.nimr.mrc.ac.uk/taylor/members/ffranca/)  
48 [ac.uk/taylor/members/ffranca/](http://mathbio.nimr.mrc.ac.uk/taylor/members/ffranca/)) was the first of  
49 the organizers to lead a seminar. She described the  
50 parametrization of a simple and easy-to-derive ana-  
51 lytical formula for taking account of solvent effects

52 in molecular dynamics simulations, using acces-  
53 sible surface areas. The method, parametrically  
54 optimised surfaces (POPS) [10], has already been  
55 integrated into GROMOS96, and demonstrated to  
56 be only about 30% slower than vacuum meth-  
57 ods — orders of magnitude cheaper than explicit  
58 water molecular dynamic simulations.

59 In order to obtain an energy term to add the sol-  
60 vent contribution to the force field, one needs to  
61 have solvation parameters that, multiplied with the  
62 surface terms, give the free energy of solvation.  
63 So far, theoreticians have used experimentally-  
64 obtained solvation energies of transfer of atoms  
65 from water to vapour. Fraternali sketched out a  
66 new approach to the calculation of these param-  
67 eters that makes use of explicit water simulations  
68 on a selected number of conformations of differ-  
69 ent peptides and proteins. From solute-restrained  
70 MD simulations of these conformers, calculated  
71 in explicit water, it is possible to obtain distribu-  
72 tions of the atomic forces exerted by that water and  
73 thereby parametrize the POPS forces accordingly.

74 For the second part of her talk, Fraternali con-  
75 centrated on more bioinformatic analysis of struc-  
76 tural data using POPS. The method has been  
77 parametrized in order to reproduce solvent acces-  
78 sibilities at atomic level (POPSA) and at the residue  
79 level (POPSR), based on a training set of about  
80 100 proteins of different sizes and topologies. The  
81 formula reproduces accessibilities calculated with  
82 the program NACCESS with less than 10% error.

83 Fraternali has shown how the formula proved  
84 useful in identifying protein-protein and pro-  
85 tein-RNA interactions in large macromolecular  
86 assemblies like the ribosome — even based on low  
87 resolution structures (C- $\alpha$  and P atoms only) like  
88 the 70S ribosome. Differences between the 30S as  
89 a separate subunit and as part of the 70S complex  
90 (with the 50S subunit) have been highlighted in  
91 this way. Because of the presence of the P-tRNA  
92 in the 70S ribosome, localized conformational rear-  
93 rangements occurring within the subunits, exposing  
94 Arg and Lys residues to negatively charged binding  
95 sites of P-tRNA, can be identified. POPSR can also  
96 be used to estimate the loss of free energy of sol-  
97 vation upon complex formation, particularly useful  
98 in designing new protein-RNA complexes and in  
99 suggesting more focused experimental work.

100 Like many of the most effective bioinformatic  
101 approaches, POPS is an approximation to make  
102

1 large-scale problems tractable. In this case, Frater-  
2 nali used it to tackle the problem of the large multi-  
3 component ribosome structures and to produce  
4 illuminating data. A POPS web server has been  
5 made available at [http://ibivu.cs.vu.nl/programs/  
6 popswww/](http://ibivu.cs.vu.nl/programs/popswww/) [6].

7 **Michele Vendruscolo's** ([http://www.ch.cam.ac.  
8 uk/CUCL/staff/mv.html](http://www.ch.cam.ac.uk/CUCL/staff/mv.html)) group, in Cambridge's  
9 Chemistry Department, studies non-native struc-  
10 tures of proteins and uses molecular dynamics  
11 to translate experimental measurements into struc-  
12 tures. Vendruscolo made the important point that  
13 we know far less about the cellular states of pro-  
14 teins than about their crystal states, as determined  
15 by X-ray crystallography. We urgently need to  
16 understand the forms proteins take when they form  
17 aggregates, intermediates, assemblies, or when they  
18 are the nuclei of misfolded forms.

19 Vendruscolo outlined his group's use of *rest-*  
20 *raind simulations* to investigate such problems.  
21 The approach generates an ensemble of struc-  
22 tures for study for which specific experimentally-  
23 measured restraints are satisfied. Various exper-  
24 imental techniques can be used to obtain the  
25 restraints. Vendruscolo outlined the technique with  
26 an example of three amino acids for which a dozen  
27 or so interactions and specific bonds had to be satis-  
28 fied. Once an experimental technique and a struc-  
29 tural interpretation of the derived data have been  
30 chosen, the model for the interactions emerges and  
31 a pseudo-energy function penalizes deviations from  
32 the experimentally derived restraints. Vendruscolo  
33 argued that these were essential because molecu-  
34 lar dynamics simulations cannot entirely replace  
35 experiments in structure determination problems.

36 He then detailed some specific case studies of  
37 published applications of the restrained simula-  
38 tion technique, beginning with a 2004 JACS paper  
39 [15] using data from site-directed spin-labelling  
40 of acyl co-enzyme A binding protein (ACBP)  
41 to investigate the residual structure present in  
42 the unfolded protein. Restraints were imposed  
43 on the average over a set of copies (replicas)  
44 of the molecule and the technique was imple-  
45 mented through 25 different non-interacting models  
46 of the molecule — multiple simulations increased  
47 the accuracy of the back-calculation of non-  
48 restrained values. Not all of several hundred pos-  
49 sible restraints are used in any given model, but  
50 those used have to be mutually consistent.  
51

52 Vendruscolo showed contact maps of the native  
53 and denatured states, maps of the average dis-  
54 tances between pairs of residues (these were, in  
55 fact, based on the probabilities of the interactions  
56 between pairs of residues). Although denatured  
57 ACBP molecules are highly heterogeneous, Ven-  
58 druscolo claimed that the sensitivity of the compu-  
59 tational technique allowed him and his co-workers  
60 to identify long-range conformational tendencies.

61 He also gave other example applications: the  
62 identification of rare (e.g. once a day) but large  
63 structural fluctuations from the native state [26],  
64 based on hydrogen exchange with solvent; the  
65 investigation of transition states too short-lived to  
66 be investigated properly experimentally; and the  
67 modelling of amyloid fibres using solid-state NMR-  
68 derived distance restraints.

69 **José Saldanha** ([http://mathbio.nimr.mrc.ac.  
70 uk/taylor/members/jsaldan/](http://mathbio.nimr.mrc.ac.uk/taylor/members/jsaldan/)) of Willie Taylor's  
71 lab then led us through a rich case history of the  
72 application of comparative modelling to the anal-  
73 ysis of a therapeutic target molecule. Although a  
74 useful technique, comparative modelling can be  
75 difficult to present scientifically because its applica-  
76 tion rarely makes a good 'story'. It is often a step  
77 in a larger process or a door to a wider biolog-  
78 ical question. Saldanha had worked in collaboration  
79 with Daruka Mahadevan, a consultant oncologist at  
80 the University of Arizona. Saldanha did bioinform-  
81 atics to analyse targets proposed by his collabor-  
82 ator; Mahadevan performed expression studies.

83 Saldanha first provided some background on  
84 prostate cancer, the second most common form of  
85 death in males, and on prostate-specific membrane  
86 antigen (PSMA), the main target for his investi-  
87 gations, giving reasons why it might well be a  
88 better marker for prostate disease than the widely-  
89 known prostate-specific antigen (PSA). PSMA is  
90 a 750 amino acid protein, implicated in many  
91 body functions — questions were later asked about  
92 the wisdom of choosing such a widely-used tar-  
93 get. Saldanha's choice rested on several bases:  
94 there are several isoforms of PSMA, and the form  
95 expressed in prostate cancer is distinct from the  
96 others; tumour endothelial cells express it, but not  
97 normal endothelial cells; and other researchers are  
98 targeting PSMA in prostate cancer. There is also  
99 good clinical evidence from early trials that PSMA  
100 can be manipulated specifically and safely.

101 Saldanha ran through the range of bioinformatics  
102

1 programs that were applied to the problem, includ- 52  
2 ing BLAST (sequence search), PSIPRED (sec- 53  
3 ondary structure prediction), THREADER (fold 54  
4 recognition), SAP (a structure-based sequence 55  
5 alignment program) and QUANTA (a commer- 56  
6 cial modelling suite). This process of bioinformatic 57  
7 characterization ran from determining its domain 58  
8 boundaries to alignment to structure prediction. 59  
9 It turned out that the transferring receptor was 60  
10 likely to be the best template. Although distantly 61  
11 related to PSMA, it has a similar domain structure. 62  
12 The two molecules may share similar properties 63  
13 of dimerization and a similar binding–recycling 64  
14 model. 65

15 Saldanha's model(s) proved consistent with mut- 66  
16 agenesis data and suggested an apical domain that 67  
17 might be involved in substrate binding. Docking of 68  
18 the natural dipeptide substrate, NAAG, hinted that 69  
19 the specificity pocket might be distinctive enough 70  
20 to help in the design of inhibitors, but a full 3-D 71  
21 structure is yet to be experimentally determined. 72

22 Workers in Janet Thornton's large group at 73  
23 the European Bioinformatics Institute (EBI) have 74  
24 been seeking to infer function from structural 75  
25 information for some time now. James Watson 76  
26 ([http://www.ebi.ac.uk/Information/Staff/person](http://www.ebi.ac.uk/Information/Staff/person_maint.php?person_id=345) 77  
27 [\\_maint.php?person\\_id=345](http://www.ebi.ac.uk/Information/Staff/person_maint.php?person_id=345)) outlined their efforts 78  
28 to obtain functional assignments within struc- 79  
29 tural genomics work, particularly in collaboration 80  
30 with the Midwest Center for Structural Genomics 81  
31 (<http://www.mcsg.anl.gov/>). 82

32 Watson pointed out that, when it works, func- 83  
33 tional assignment from three-dimensional structure 84  
34 is more appropriate to the identification of bio- 85  
35 chemical rather than biological function. Currently 86  
36 sequence methods are the most successful way to 87  
37 assign function, but structure-based methods can 88  
38 provide additional functional information. There 89  
39 are still plenty of occasions when no bioinformatic 90  
40 methods work and function can only be identified 91  
41 by direct experiment. 92

43 Watson described ProFunc, a bioinformatics 93  
44 pipeline combining a variety of methods [13]. 94  
45 The structural contributions come from match- 95  
46 ing homologous folds, a variety of 3-D template 96  
47 methods, binding site identification and structure 97  
48 motif (for example helix–turn–helix) conservation. 98  
49 Databases of 3-D templates describe enzyme active 99  
50 sites, ligand binding sites and DNA binding sites. 100  
51 Hits to these templates are ranked by comparing 101  
102

the surrounding environment of the match and cal- 52  
culating a similarity score. He also described the 53  
use of 'nests', small structural motifs involving 54  
protein backbones that are commonly found to sta- 55  
bilize some secondary structures and can also stabi- 56  
lize ligand binding. The structural alignments come 57  
from firstly centring on the 3-D template match 58  
(e.g. enzyme active site) then expanding the align- 59  
ment based on sections considered 'fittable' (within 60  
an RMSD cut-off) that consist of at least seven 61  
consecutive residues. 62

Sadly, I was only able to catch the end of David 63  
Burke's ([http://www-cryst.bioc.cam.ac.uk/~](http://www-cryst.bioc.cam.ac.uk/~dave/) 64  
[dave/](http://www-cryst.bioc.cam.ac.uk/~dave/)) presentation, 'Ab initio structure predic- 65  
tion' [2,4], and the subsequent discussion. When 66  
I arrived, Burke was addressing the question of 67  
how to filter tens of thousands of models of loops. 68  
Currently, van der Waals' overlap was the main cri- 69  
terion, but he suggested that molecular dynamics 70  
force fields, solvent accessibility and comparison 71  
with known structures could all be applied to win- 72  
now the output from modelling programs. Burke 73  
also summarized the questions that still concerned 74  
him — and concern many structural bioinformati- 75  
cians: 76

- Is it best to separate the selection of the models 77  
from the generation of models? 78
- Has the majority of the reasonable peptide con- 79  
formations in the protein universe been observed 80  
in the structures deposited in the PDB to date? 81
- How can distantly related molecules be mod- 82  
elled? 83

84 Many of us had heard Willie Taylor ([http://math](http://math.bio.nimr.mrc.ac.uk/taylor/members/wtaylor/) 85  
[bio.nimr.mrc.ac.uk/taylor/members/wtaylor/](http://math.bio.nimr.mrc.ac.uk/taylor/members/wtaylor/)) 86  
talk before, but he promised us that 'Folds, knots 87  
and tangles' would include both 'something old 88  
and something new' amongst a collection of meth- 89  
ods which, although apparently disconnected, all 90  
could contribute towards *ab initio* structure pre- 91  
diction. He began by describing the universe of 92  
non-redundant folds by type ( $\alpha$ ,  $\alpha\beta$  and  $\beta$ ) and 93  
pointed out that this division of foldspace, while 94  
superficially illuminating, says less about deep sim- 95  
ilarities between fold classes, than about how we 96  
look at proteins. 97

Now that Taylor and his co-workers are actively 98  
interested in model 'proteins' (i.e. non-biological 99  
structures devised *in silico*), he has found that they 100  
are difficult to classify by eye and they have used 101  
Ptitsyn and Finkelstein's concept of *structural* 102

1 layers to find a way to compare them without the  
2 perennial problems of using, say, RMS deviations  
3 between  $\alpha$  and  $\beta$  proteins.

4 Taylor's talks benefit from being supported by  
5 live demos of actual programs running on a  
6 Linux laptop, rather than static computer slides.  
7 He first used RasMol to show the cell matrices  
8 he plots from the distribution of his fold types  
9 along axes of complexity, and 'curl and stag-  
10 ger'. He has described this classification and its  
11 sub-classifications as a 'Periodic Table' of protein  
12 structures [25]. In his demonstration this repre-  
13 sentation was completely dynamic, with individual  
14 spheres being clickable to give the SAP represen-  
15 tation of each protein fold's superimposed struc-  
16 tures — colour-coded by their strength of mutual  
17 correspondence [23].

18 He now uses this scheme for the classification  
19 of model proteins. When asked about the RasMol  
20 renderings of such elements, Taylor pointed out  
21 that these projections represent the architecture of  
22 the protein, failing to discriminate, for example,  
23 between parallel and anti-parallel  $\beta$ -strands, but  
24 the full topology for each protein is recorded in  
25 a 'topology-string' and can be used if needed  
26 [11]. Taylor then moved on to questions of *ab*  
27 *initio* protein structure prediction and contrasted  
28 his whole-structure interests with the loop-focused  
29 work of David Burke, who had preceded him.

30 Taylor used a constrained random walk to gener-  
31 ate structures, along the way occasionally gener-  
32 ating secondary structure elements — sometimes  
33 domains. A random walk combined with a sys-  
34 tem for the generation of layers produces struc-  
35 tures which are more protein-like. Occasionally  
36 this approach results in the production of knots.  
37 This behaviour had to be suppressed with 'smooth-  
38 ing'. Some real proteins in the PDB could not be  
39 smoothed down to a line. It turned out that these  
40 special cases are *knotted*. This curious, almost-  
41 incidental discovery led to a publication in *Nature*  
42 [24].

43 Smoothing can be used to compare the complex-  
44 ity of proteins. According to the number of self-hits  
45 of smoothed proteins, TIM barrels are simpler than  
46 Rossmann folds, for example. It is possible to grow  
47 protein traces *in silico* through the building of local  
48 contacts and plot the ease of building a given fold  
49 making only local connections from a specified  
50 point in its structure.  
51

Finally, Taylor ran through some of the elements  
used in his *ab initio* folding experiments:

- Secondary structure predicted with PSI-PRED. 52
- Random walks generated with RAMBLE. 53
- Filtering performed using radius of gyration. 54
- Filtering for knots. 55
- Filtering for complexity. 56
- Folds scored (of the order of  $10^5$  in number) with  
CAO (Contact Accepted MutatiOn) [14]. 57
- POPS (the solvent accessibility algorithm des-  
cribed by Fraternali) and SPREK. 58

Alternative structures produced using his group's  
*ab initio* methods can be ranked in order by fold  
and clustered. He hoped to have a comprehensive  
system using these or similar techniques up and  
running in time for the next CASP meeting.

Another local speaker, **Vijayalakshmi Chelliah**  
of Cambridge University's Biochemistry Depart-  
ment, moved us on from protein structure determi-  
nation to protein function determination with her  
talk on 'The identification of interacting sites in  
protein families'. She started from the reasonable  
premise that critically important residues tend to  
be conserved by the members of protein fami-  
lies. She had used HOMSTRAD to generate 96  
environment-specific substitution tables for pro-  
tein residues and taken these as a background  
against which to detect important sites, those where  
residues are more conserved in families than would  
be expected from the tables.

The method is simple and logical:

- Make a structure-based alignment of family  
members. 59
- Compare the observed and expected substitution  
patterns. 60
- Measure the informational difference between  
the two. 61

The higher the score, the more distant the two  
distributions are. High-scoring positions identified  
in this way are those considered most likely to be  
functional. These scores can then be mapped onto  
structures to find high-scoring clusters. For this last  
stage, Chelliah used Kin3Dcont, part of the kin-  
contour program (<http://kinemage.biochem.duke.edu/index.php>)  
produced by the Richardson Lab at Duke University,  
North Carolina.

Chelliah was careful to ignore large gaps when  
making alignments and to restrict her analysis

1 to sequences with less than 80% mutual identity  
2 in order to minimize the noise from 'briefly'  
3 conserved, but functionally unimportant, residues.

4 In most of around 250 families the 'averaged  
5 out' active site predicted was between 0 Å and 9 Å  
6 from the true active site, but the method missed  
7 functional sites that were indirectly involved in  
8 the activity of proteins and sites that were buried.  
9 Along the way to these results she made some  
10 interesting observations:

- 11 • Critical residues that were also structurally  
12 important did not score as highly as might have  
13 been expected by this method.
- 14 • Even inaccessible residues turned out to be very  
15 highly conserved — Chelliah put this down to  
16 their being important to the structural integrity  
17 of active sites in the molecule.
- 18 • She felt that this might have been countered by  
19 looking for sites retained in both orthologues  
20 and paralogues and tested this by adding in  
21 phylogenetic information. As it turned out, the  
22 addition of close homologues generated more  
23 noise.  
24

25 She observed, as people often do with methods  
26 like this, that the predictions were best when  
27 residues were in truly equivalent positions within  
28 similar structures.  
29

30 Returning to structure prediction, 'Conforma-  
31 tional sampling for protein structure determination  
32 and prediction' was the title of **Mark DePristo's**  
33 talk. DePristo is another member of Cambridge's  
34 Biochemistry Department ([http://raven.bioc.cam.  
35 ac.uk/~mdepristo/](http://raven.bioc.cam.ac.uk/~mdepristo/)). He described a method devel-  
36 oped (and now used) to check protein models,  
37 but which turns out to have a range of useful  
38 structure-related applications. He introduced his  
39 hybrid approach by summarizing the problem in  
40 a series of simple figures. If the solution of a pro-  
41 tein structure is a global minimum on an energy  
42 (or other scoring function) landscape, then our aim  
43 should be to smooth out that landscape to avoid  
44 local minima and sample enough of it to find the  
45 true minimum. Since there is no definitive solu-  
46 tion, we must carefully choose heuristics. DePristo  
47 explained the advantages of molecular dynam-  
48 ics/simulated annealing approaches over conjugate  
49 gradient/steepest descent ones.

50 His framework for such investigations, RAP-  
51 PER, avoids optimizing a non-linear function.

52 Instead it chooses many starting points and applies  
53 local minimum-finding methods. Once a general  
54 class of structures has been specified, the poten-  
55 tial energies of those structures can be compared.  
56 Because small deviations from ideal geometry are  
57 allowed in the real world and flexibility comes at  
58 computational cost, RAPPER fixes many param-  
59 eters (bond lengths, angles) and samples residue-  
60 specific propensity tables and hand-curated con-  
61 formation libraries. The algorithm constructs rea-  
62 sonable 3-D models consistent with prior structural  
63 constraints and additional arbitrary ones, and pro-  
64 gresses from the N- to C-terminus of a structure,  
65 pruning additions in the wrong conformation.

RAPPER has been applied to loop modelling [2],  
66 (re)construction of native ensembles [7], compar-  
67 ative modelling, and crystallographic model gen-  
68 eration [8]. More details of the program and its  
69 variants are available from the RAPPER Website:  
70 <http://raven.bioc.cam.ac.uk/index.php>

71 **David Jones** (Department of Computer Sci-  
72 ence, Bioinformatics Unit, University College Lon-  
73 don, [http://www.cs.ucl.ac.uk/staff/D.Jones/index.  
74 html](http://www.cs.ucl.ac.uk/staff/D.Jones/index.html)) spoke on the 'Detection of native disorder  
75 in proteins'. To begin, he joked about the irony of  
76 his spending years trying to predict structure from  
77 sequence before trying to predict 'non-structure'  
78 from sequence. He also graciously credited  
79 Jon Ward (<http://www.cs.ucl.ac.uk/staff/J.Ward/>)  
80 with having done most of the work. After running  
81 through the basic assumptions of sequence-struc-  
82 ture interdependency, he discussed the various  
83 kinds of disordered proteins that were known to  
84 exist. Some proteins are partially or completely  
85 unfolded yet remain functional, and we assume that  
86 this is because their molecules form an ensemble of  
87 states, rather than a unique structure. These disor-  
88 dered states could be compact or extended molten  
89 globules or random coils and, interestingly, can  
90 fold fully on binding.  
91

92 Jones talked about the blurry line between true  
93 disorder and experimental uncertainty in deter-  
94 mining protein structures as well as the experimen-  
95 tal methods which can be used to detect disorder.  
96 He proposed functional classes of disordered  
97 regions in proteins: 'springs and linkers', modi-  
98 fication sites, regions important to the timing of  
99 complex assembly and molecular recognition sites.  
100 Functional importance is often assumed to corre-  
101 late with evolutionary conservation and the work  
102 on predicting disorder seems to produce results

1 consistent with this. He also outlined some pre-  
2 vious work to identify signals of disorder in pro-  
3 teins.

4 Ward and Jones had trained a support vector  
5 machine (SVM) on a non-redundant set of crystal  
6 structures and found that they could use it  
7 to identify 40% of disordered residues with a  
8 1% error rate. The performance was better for  
9 longer regions — over 30 amino acid residues  
10 in length — for which the detection fraction and  
11 error rates were 80% and 0.1%, respectively. The  
12 SVM was then applied across genomes and detec-  
13 tion rates compared with biological function (as  
14 assigned by gene ontology classifications) [27].  
15 He believed other workers' predictions of disorder  
16 in prokaryotic proteins were likely to be  
17 overestimates. In eukaryotes, molecules associ-  
18 ated with the actin cytoskeleton scored highly,  
19 while the bacteria-like environment of mitochon-  
20 dria seemed to contain few disordered protein  
21 components. There was also high correlation with  
22 DNA-transposition and development and mor-  
23 phogenesis. Molecular functions more likely to  
24 be associated with protein disorder predictions  
25 included transcription regulators, protein kinases  
26 and transcription factors. Metabolic and biosyn-  
27 thetic protein functions scored low. The disorder  
28 prediction server, DISOPRED, is available at  
29 <http://bioinf.cs.ucl.ac.uk/disopred/disopred.html>  
30

31 Another Chothia group member, **Martin Mad-**  
32 **era** (<http://stash.mrc-lmb.cam.ac.uk/mm238/>)  
33 talked about his work on 'Comparisons of sequence  
34 families' and his responsibility for the Chothia  
35 group's 'Superfamily' database at the LMB [16].  
36 This is a library of HMM models for all proteins  
37 of known 3-D structure. He recounted a history  
38 of protein sequence comparison methods, of the  
39 problems of characterizing more distantly related  
40 protein groupings, and he detailed more recent  
41 improvements in this resource. He gave a clear  
42 overview of pairwise vs. sequence profile vs. HMM  
43 methods and, having made the case for HMMs, he  
44 discussed the refinements implemented in Super-  
45 family, which relies on the segmentation of PDB  
46 structures into domains and the combination of  
47 multiple HMMs to represent its groupings. The  
48 domain-based analysis of Superfamily can now be  
49 used to compare whole genomes for their domain  
50 composition.

51 We moved from better models of real, stable,

52 folded proteins, to predictions of disordered pro-  
53 teins to completely imaginary proteins. **Benjamin**  
54 **Blackburne** ([http://slater.chem.nott.ac.uk/~](http://slater.chem.nott.ac.uk/~bpb/)  
55 **bpb/**), formerly of Jonathan Hirst's group at Not-  
56 tingham University and now a member of Richard  
57 Goldstein's group, talked about the properties of his  
58 phylogenies of minimalist proteins [3]. Blackburne  
59 had explored the relationships between hypotheti-  
60 cal 2-D proteins catalogued in the sort of protein  
61 database the inhabitants of 'Flatland' [1] might re-  
62 cognize. In Blackburne's planar protein universe,  
63 residues are of only two types, hydrophobic or  
64 hydrophilic. Proteins fold when strings of such  
65 residues arranged on a square or tetrahedral lat-  
66 tice of available points turn in on themselves in  
67 a plausible way. Folds that arrange those residues  
68 with the lowest energy are 'native'. A 'fit' pro-  
69 tein is one which has a pocket — i.e. two external  
70 residues around a hole that could be 'functional'.  
71

72 With so few degrees of freedom, all sequences  
73 of given short lengths and all structures derived  
74 from them can be known. The proteins can be  
75 arranged in families, where a family is a group  
76 in which all the possible relatives can be generated  
77 from another by mutation and yet still meet the  
78 rules for the formation of viable structures; the  
79 relationships between the model structures can be  
80 visualized in graphs, whose nodes are the structures  
81 and whose edges are point mutations between  
82 them. There are outliers, and some families are  
83 more weakly connected to related families than  
84 others. There are 'bottlenecks' where there are few  
85 evolutionary routes from one family to another.  
86 'Hubs' bridge multiple families. 'Funnels' form  
87 when the structures are arranged such that the  
88 nodes radiate out to variants of decreasing stability.  
89

90 Some phenomena can be compared in an illu-  
91 minating way with the evolution of real proteins.  
92 For example, in Blackburne's world neutral evolu-  
93 tion seems necessary for minimal proteins to  
94 reach functional states and longer chains offer more  
95 potential for such noisy change. Other characteris-  
96 tics of these artificial proteins are more problem-  
97 atic: their sequences are not directional and inser-  
98 tions and deletions cannot have the same meaning  
99 when there are so few residue types.

100 Of course, much of the subsequent discus-  
101 sion was about the relevance of such evolu-  
102 tionary landscapes to real proteins, whether the  
103 graphs had scale-free properties, other aspects of



1 real protein behaviour which ought to be mod- 52  
2 elled (Cyrus Chothia), and the correspondence 53  
3 between Blackburne's neutral mutation-tolerant 54  
4 proteins and Chothia's stable-to-mutation proteins 55  
5 (Willie Taylor). 56

6 **Richard Goldstein's** (<http://mathbio.nimr.mrc.ac.uk/goldstein/members/rgoldstein/>) talk, 57  
7 'Modelling molecular evolution', covered an area 58  
8 of growing interest, the effort to combine sequence 59  
9 and structural analysis to investigate the evolution 60  
10 and function of proteins. He described methods 61  
11 aimed at increasing our understanding of the struc- 62  
12 tural basis for variations in amino acid residue 63  
13 substitution rates, identifying functional sites and, 64  
14 in particular, for characterizing members of the 65  
15 large and pharmacologically important family of 66  
16 G protein-coupled receptors (GPCRs). 67

17 First, he highlighted a central flaw in compar- 68  
18 ative sequence analysis: most approaches are 69  
19 based on a model that assumes positions in 70  
20 sequences represent independent samplings from 71  
21 all possible sequences and ignore the phyloge- 72  
22 netic relationships between related proteins. He 73  
23 also reminded us — as molecular phylogeneticists 74  
24 often have to remind biochemists and molecu- 75  
25 lar biologists — that residues 'conserved' between 76  
26 closely related sequences are not as significant as 77  
27 investigators often believe. 78

28 Rather than ignore these problems or devise 79  
29 *ad hoc* fixes, Goldstein, Goldman and others have 80  
30 more recently attempted to model evolution explic- 81  
31 itly. To begin, Goldstein developed substitution 82  
32 matrices for different types of local structure, but 83  
33 has since devised a more general approach. Each 84  
34 protein can be divided up into zones, without mak- 85  
35 ing assumptions about which models apply where; 86  
36 the probability of any given location belonging 87  
37 to a particular site class is a parameter which is 88  
38 itself optimized by an expectation maximization 89  
39 algorithm. 90

40 Once a set of environment categories has emerg- 91  
41 ed, Goldstein and co-workers assign qualitative 92  
42 labels to them (e.g. 'hydrophilic'), and the *a pos-* 93  
43 *teriori* probabilities of each position belonging to 94  
44 class can be estimated. By applying this approach 95  
45 to large enough families of aligned sequences with 96  
46 structural information, he claimed, it is possible to 97  
47 identify locations where different types of selective 98  
48 pressure have been operating and obtain insights 99  
49 into the underlying basis of such selective pres- 100  
50 sure, e.g. how physicochemical properties such as 101  
51

size and hydrophobicity are differentially important 52  
in different classes of site. 53

This approach can be used to identify function- 54  
ally important locations — sites belonging to the 55  
slowest evolving rate classes — and different over- 56  
all probabilities that a position is involved in gen- 57  
eral function, stabilization, dimerization, packing, 58  
structure, or the extent to which a position con- 59  
strained [21]. 60

Goldstein then focused on the application of 61  
this general approach to the specific question of 62  
the GPCRs. Despite representing only 1% of the 63  
genome, they are estimated to be the target of 64  
almost half all drugs and only one signalling 65  
process does not involve a member of this family. 66  
Although only one known high-resolution structure 67  
is available, Goldstein's group worked with a 68  
dataset of about 200 GPCRs, and analysed them 69  
to produce patchworks of model assignments along 70  
the lengths of sequences. 71

Some properties of these molecules gave a strong 72  
signal. It is harder, for example, to identify the 73  
inner and outer surface of transmembrane (TM) 74  
helices, such as those in the 7-TM structure of 75  
the GPCRs, than it is to identify the inner and 76  
outer faces of 'normal' protein structure helices. 77  
Goldstein *et al.*'s site classes correlate with the 78  
'innerness' and 'outerness' of these helices. Also, 79  
a propensity to involvement in dimerization seems 80  
to correlate with slowly varying sites. 81

The European Bioinformatics Institute's **Hugh** 82  
**Shanahan** (<http://www.biochem.ucl.ac.uk/~shanahan/>) described more function-from-structure 83  
work, this time targeted at predicting DNA-binding 84  
proteins from 3-D motifs and electrostatic infor- 85  
mation. There is no shortage of important DNA- 86  
binding proteins and a huge and growing inter- 87  
est in the regulation of transcription. Shanahan 88  
quoted estimates of up to 7% of eukaryotic and 89  
3% of prokaryotic genes coding for DNA bind- 90  
ing proteins. Equally, structural genomics projects 91  
will generate many uncharacterized structures. 92  
Although he acknowledged the importance and 93  
utility of sequence-based approaches, he argued 94  
that function varies significantly as sequence iden- 95  
tity between unknown and known (template) pro- 96  
tein sequences falls below 40%. He pointed out 97  
that, although at least one neural net-based method 98  
exists for identifying DNA binding proteins, it 99  
100  
101  
102

1 has a high false-positive rate and requires high-  
2 resolution atomic data, and claimed that homology-  
3 based modelling produces lower false-positive  
4 scores.

5 Shanahan further contended that, of the four  
6 main known classes of structural motif:

- 7 • Helix–turn–helix.
- 8 • Helix–hairpin–helix.
- 9 • Helix–loop–helix.
- 10 • Zinc finger.

11 the middle two are more easily identified with  
12 Hidden Markov Model (HMM) methods; zinc fin-  
13 ger proteins are too structurally variable. Shanahan  
14 concentrated on the first, helix–turn–helix (H–T–  
15 H) structures. He began by summarizing the pro-  
16 cedure to identify structural templates:

- 17 • Search the literature for H-T-H motifs.
- 18 • Identify HMMs in Pfam or SMART.
- 19 • Identify structural templates from domains using  
20 the CATH super-structural family (the H-level of  
21 that database).
- 22 • Scan the Protein DataBank with templates.
- 23 • Add any new H–T–H DNA-binding proteins to  
24 the list.
- 25 • Repeat until no other structures are found.

26 The group obtained 90 non-redundant structures  
27 in the PDB and generated seven structural tem-  
28 plates to cover that set, applying an accessibility  
29 criterion. At first the results didn't seem much bet-  
30 ter than those obtained with HMMs: 0.5% false  
31 positives. Then they refined the method by inte-  
32 grating the potential over a region close to the  
33 accessible surface of motifs and tested this by using  
34 the electrostatic data to attempt to identify the bind-  
35 ing region in known DNA-binding proteins [12].

36 A method to detect DNA-binding sites on the  
37 surface of a protein structure is important for func-  
38 tional annotation. They analysed residue patches  
39 on the surface of DNA-binding proteins and pre-  
40 dicted DNA-binding sites using a single feature  
41 of these surface patches. They first surveyed sur-  
42 face patches and DNA-binding sites for accessi-  
43 bility, electrostatic potential, residue propensity,  
44 hydrophobicity and residue conservation. From  
45 this, they observed that the DNA-binding sites usu-  
46 ally fell in the top 10% of patches with the largest  
47 positive electrostatic scores. This knowledge led to  
48 their development of a prediction method in which  
49

50 patches of surface residues were selected such that  
51 they excluded residues with negative electrostatic  
52 scores.

53 They used this method to make predictions for a  
54 dataset of 56 non-homologous DNA-binding pro-  
55 teins and identified 68% of the dataset correctly.  
56 Using this data, they improved the false-positive  
57 score to 0.02%. Shanahan added that the hybrid  
58 method involves fewer parameters than sequence  
59 homology, does not require full electrostatic calcu-  
60 lations to be performed and that it might be possible  
61 to use data from homology models to provide a  
62 cross-check for HMM searches.

63 The final talk of the meeting rounded the  
64 event off perfectly. **Chris Calladine** ([http://www-  
65 civ.eng.cam.ac.uk/crc/crc-web.htm](http://www-civ.eng.cam.ac.uk/crc/crc-web.htm)), who retired  
66 only a couple of years ago from the Cambridge  
67 University Department of Structural Engineering,  
68 dazzled us with a multidisciplinary, multimedia  
69 presentation on the 'Mechanics of interfaces in  $\alpha$ -  
70 helical supercoils'. He used overheads, animation  
71 and a succession of cork-and-cardboard models to  
72 show how juxtaposed helices could abut in diverse  
73 ways, interlocking the 'knobs' of their respective  
74 sidechains. The knobs of one helix fit into the  
75 'holes' between the knobs of the other when they  
76 pack. For simple superhelices and four-helix bun-  
77 dles — as distinct from the helix-built cylinders  
78 Calladine later touched on [5] — there were three  
79 standard modes of knobs-into-holes packing, which  
80 he illustrated with overlaid interface figures pro-  
81 duced as overheads, as simple figures and as clev-  
82 erly constructed three-dimensional models.

83 One of the most pleasing things about structural  
84 bioinformatics is that its practitioners collaborate  
85 across specialisms to tackle difficult, interesting  
86 and messy problems out of both curiosity and  
87 necessity — not merely to meet the conditions of  
88 interdisciplinary funding programmes. Calladine's  
89 work exemplified this beautifully. He has worked  
90 in this area in collaboration with Charlie Laughton  
91 (molecular dynamics) at Nottingham University  
92 and Ben Luisi and Venkatesh Pratap (structural  
93 biology) at Cambridge. Pratap wrote software that  
94 finds  $\alpha$ -helices and their neighbours, identifies the  
95 local superhelical angle of their arrangement and  
96 categorizes those arrangements according to those  
97 angles. Pratap's animation of a bistable 'switch' in  
98 the packing of a right-handed, four-helix bundle  
99 of  $\alpha$ -helices in one of the three main classes  
100  
101  
102

1 of arrangement formed the finale of Calladine's  
2 presentation.

## References

1. Abbott EA. 1884. *Flatland: A Romance of Many Dimensions*. Shambhala: Boston, MA.
2. de Bakker PI, DePristo MA, Burke DF, Blundell TL. 2003. *Ab initio* construction of polypeptide fragments: accuracy of loop decoy discrimination by an all-atom statistical potential and the AMBER force field with the Generalized Born solvation model. *Proteins* **51**: 21–40.
3. Blackburne BP, Hirst JD. 2001. Evolution of functional model proteins. *J Chem Phys* **115**: 1935–1942.
4. Burke DF, Deane CM. 2001. Improved protein loop prediction from sequence alone. *Protein Eng* **14**: 473–478.
5. Calladine CR, Sharff A, Luisi BF. 2001. How to untwist an  $\alpha$ -helix: structural principles of an  $\alpha$ -helical barrel. *J Mol Biol* **305**: 603–618.
6. Cavallo L, Kleinjung J, Fraternali F. 2003. POPS: a fast algorithm for solvent accessible surface areas at atomic and residue level. *Nucleic Acids Res* **31**: 3364–3366.
7. DePristo MA, de Bakker PIW, Lovell SC, Blundell TL. 2002. *Ab initio* construction of polypeptide fragments: efficient generation of accurate, representative ensembles. *Proteins Struct Funct Genet* **51**: 41–55.
8. DePristo MA, de Bakker PIW, Lovell SC, Blundell TL. 2004. Heterogeneity and inaccuracy in protein structures solved by X-ray crystallography. *Structure* **12**: 831–838.
9. Fernandez-Recio J, Totrov M, Abagyan R. 2004. Identification of protein–protein interaction sites from docking energy landscapes. *J Mol Biol* **335**: 843–865.
10. Fraternali F, Cavallo L. 2002. Parameter optimized surfaces (POPS): analysis of key interactions and conformational changes in the ribosome. *Nucleic Acids Res* **30**: 2950–2960.
11. Johannissen LO, Taylor WR. 2004. Protein fold comparison by the alignment of topological strings. *Protein Eng* **16**: 949–955.
12. Jones S, Shanahan HP, Berman HM, Thornton JM. 2003. Using electrostatic potentials to predict DNA-binding sites on DNA-binding proteins. *Nucleic Acids Res* **31**: 7189–7198.
13. Laskowski RA, Watson JD, Thornton JM. 2003. From protein structure to biochemical function? *J Struct Funct Genomics* **4**: 167–177.
14. Lin K, Kleinjung J, Taylor WR, Heringa J. 2003. Testing homology with Contact Accepted mutatiOn (CAO): a contact-based Markov model of protein evolution. *Comput Biol Chem* **27**: 93–102.
15. Lindorff-Larsen K, Kristjansdottir S, Teilum K, *et al.* 2004. Determination of an ensemble of structures representing the denatured state of ACBP. *J Am Chem Soc* **126**: 3291–3299.
16. Madera M, Vogel C, Kummerfeld SK, Chothia C, Gough J. 2004. The SUPERFAMILY database in 2004: additions and improvements. *Nucleic Acids Res* **32**: (Database issue): D235–239.
17. Mizuguchi K. 2004. Fold recognition for drug discovery. *Drug Discovery Today: Targets* **3**: 18–23.
18. Shi J, Blundell TL, Mizuguchi K. 2001. FUGUE: sequence–structure homology recognition using environment-specific substitution tables and structure-dependent gap penalties. *J Mol Biol* **310**: 243–257.
19. Shirai H, Blundell TL, Mizuguchi K. 2001. A novel superfamily of enzymes that catalyze the modification of guanidino groups. *Trends Biochem Sci* **26**: 465–468.
20. Shirai H, Mizuguchi K. 2003. Prediction of the structure and function of AstA and AstB, the first two enzymes of the arginine succinyltransferase pathway of arginine catabolism. *FEBS Lett* **555**: 505–510.
21. Soyer OS, Dimmic MW, Neubig RR, Goldstein RA. 2003. Dimerization in aminergic G-protein-coupled receptors: application of a hidden-site class model of evolution. *Biochemistry* **42**: 14 522–14 531.
22. Stebbings LA, Mizuguchi K. 2004. HOMSTRAD: recent developments of the Homologous Protein Structure Alignment Database. *Nucleic Acids Res* **32**: (Database issue): D203–207.
23. Taylor WR. 2000. Protein structure comparison using SAP. *Methods Mol Biol* **416**: 657–660.
24. Taylor WR. 2000. A deeply knotted protein structure and how it might fold. *Nature* **406**: 916–919.
25. Taylor WR. 2002. A 'periodic table' for protein structures. *Nature* **416**: 657–660.
26. Vendruscolo M, Paci E, Dobson CM, Karplus M. 2003. Rare fluctuations of native proteins sampled during equilibrium hydrogen exchange. *J Am Chem Soc* **125**: 15 686–15 687.
27. Ward JJ, Sodhi JS, McGuffin LJ, Buxton BF, Jones DT. 2004. Prediction and functional analysis of native disorder in proteins from the three kingdoms of life. *J Mol Biol* **337**: 635–645.