

# Development of a force field for conditional optimization of protein structures

Sjors H. W. Scheres and Piet  
Gros\*

Department of Crystal and Structural Chemistry,  
Bijvoet Center for Biomolecular Research,  
Utrecht University, Padualaan 8,  
3584 CH Utrecht, The Netherlands

Correspondence e-mail: p.gros@chem.uu.nl

Received 23 October 2002  
Accepted 11 December 2002

Conditional optimization allows the incorporation of extensive geometrical information in protein structure refinement, without the requirement of an explicit chemical assignment of the individual atoms. Here, a mean-force potential for the conditional optimization of protein structures is presented that expresses knowledge of common protein conformations in terms of interatomic distances, torsion angles and numbers of neighbouring atoms. Information is included for protein fragments up to several residues long in  $\alpha$ -helical,  $\beta$ -strand and loop conformations, comprising the main chain and side chains up to the  $\gamma$  position in three distinct rotamers. Using this parameter set, conditional optimization of three small protein structures against 2.0 Å observed diffraction data shows a large radius of convergence, validating the presented force field and illustrating the feasibility of the approach. The generally applicable force field allows the development of novel phase-improvement procedures using the conditional optimization technique.

## 1. Introduction

During the standard crystallographic diffraction experiment, information about the phases of the observed reflections is lost. In order to obtain a molecular model describing the crystal content, this information must be regained. In protein crystallography, this process has typically been divided into well separated steps: phase determination by experimental methods or molecular replacement, phase extension by density modification and iterative cycles of model building and refinement. Nowadays, it is realised that these steps are coupled more tightly than previously thought (Lamzin *et al.*, 2000) and programs have been developed that link these steps in an automated way. For example, the (*RE*-)*SOLVE* package (Terwilliger & Berendzen, 1999) links structure solution, density modification and model building, and the *ARP/wARP* program (Perrakis *et al.*, 1999) links density modification, model building and refinement. Owing to the typically low observation-to-parameter ratio in protein crystallography, the incorporation of additional information in this process is critical. We have presented a method, called conditional optimization, in which extensive prior stereochemical information may be formulated in terms of loose atoms (Scheres & Gros, 2001). With initial simplified test calculations, we showed that a structure can be obtained by this approach using 2.0 Å diffraction data without any prior phase information. Thus, in principle the entire process from phasing to refinement can be expressed in a single step. However, these tests were performed with calculated diffraction data from a highly simplified structure of four polyalanine  $\alpha$ -helices, which

can be described by a very limited parameter set defining the expected geometries. Here, we present a parameter set for conditional optimization of the far more complex structures that are protein molecules.

In the conditional formalism, we express geometrical knowledge by the definition of interaction functions, termed conditions. These conditions depend on expected numbers of neighbouring atoms, interatomic distances and torsion angles within protein molecules. Conditions are continuous functions ranging from zero to one and show similarities with the knowledge-based interaction functions defined by Sippl (1995). Conformations of protein fragments up to several residues long are described by joint conditions, which are products of conditions describing a set of geometrical features of a protein fragment. In principle, (joint) conditions could be defined for all possible conformations in protein molecules, but this would require a vast amount of interaction functions exceeding available computing power. Therefore, we have defined conditions describing the most common conformations observed in the Protein Data Bank (Berman *et al.*, 2002) for main-chain atoms and side-chain atoms up to the  $\gamma$  position.

With the defined parameter set, we show that a large radius of convergence can be obtained for conditional optimization of three small protein structures against 2.0 Å observed diffraction data.

## 2. Mean-force potential for protein structures

### 2.1. Brief review of the conditional formalism

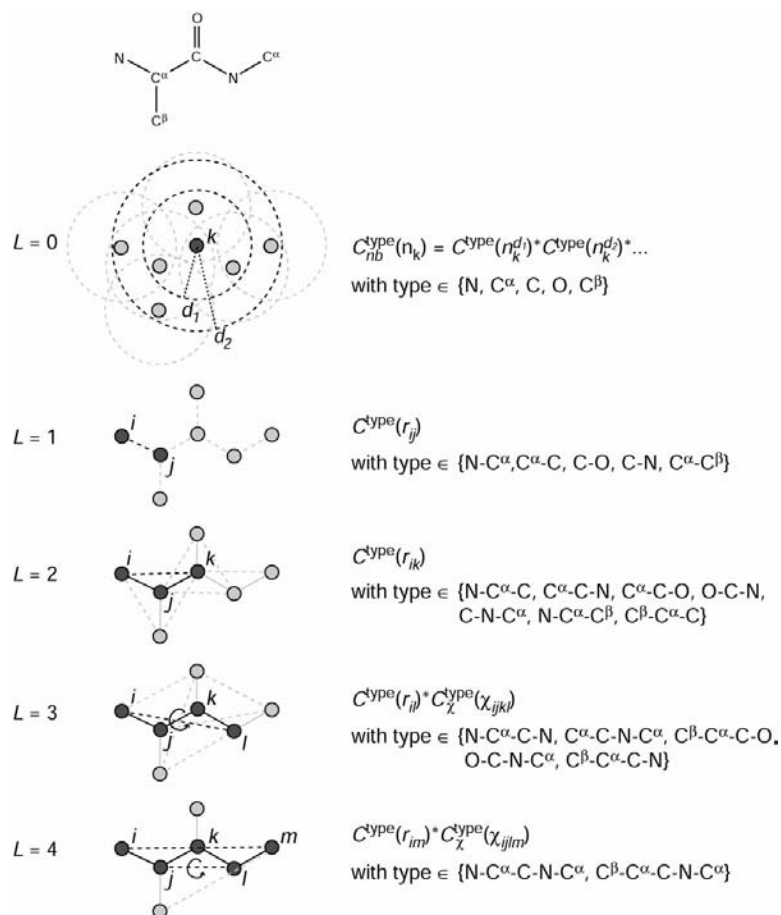
In conditional optimization, we express prior knowledge about protein structures without explicitly assigning chemical identities to the atoms. Instead, we take all possible assignments into account by using an  $N$ -particle approach. We define conditions  $C = [0, 1]$ , which are continuous interaction functions based on optimal values for the interatomic distances, torsion angles and numbers of neighbouring atoms in protein structures. We describe protein structures as a collection of linear elements (of length  $L$ ) which are non-branched sequences of  $L + 1$  atoms. Fig. 1 shows a common fragment present in protein structures and a schematic representation of the conditions that describe a linear element  $N-C^\alpha-C-N-C^\alpha$ , which depend on the number of neighbouring atoms per atom, interatomic distances and torsion angles.

As discussed in more detail previously (Scheres & Gros, 2001), a linear element of length  $L$  is composed of a total of  $L(L + 1)/2$  linear (sub-)elements of length  $l \leq L$ . Multiplication of all conditions corresponding to these (sub-)elements, gives the so-called joint condition  $JC$ .

For a linear combination of  $L + 1$  atoms  $i, j, \dots, p$  and  $q$ , the joint condition  $JC_{ij\dots pq}^{\text{type}}$  describes to what extent the conformation of the atoms resembles a defined target conformation of a particular type of

linear element. A minor change was made to the conditional formalism as presented previously. Originally, joint conditions were defined as binomial multiplications of the individual conditions, according to the binary combination of all (sub-)elements. In the current implementation, the resulting higher powers of individual conditions in the expression of joint conditions have been removed and joint conditions are defined as the multiplication of all corresponding individual conditions. Consequently, the factor  $n$  in the calculation of derivatives [see formulae (5) and (6) in Scheres & Gros, 2001] reduces to one. In Fig. 1 the atoms  $i, j, k$  and  $l$  resemble a linear element of type  $N-C^\alpha-C-N$ . The corresponding joint condition for  $N-C^\alpha-C-N$  is determined by the individual conditions as given in (1),

$$JC_{ijkl}^{N-C^\alpha-C-N} = C_{nb}^N(n_i)C_{nb}^{C^\alpha}(n_j)C_{nb}^C(n_k)C_{nb}^N(n_l)C^{N-C^\alpha}(r_{ij}) \\ \times C^{C^\alpha-C}(r_{jk})C^{C-N}(r_{kl})C^{N-C^\alpha-C}(r_{ik})C^{C^\alpha-C-N}(r_{jl}) \\ \times C^{N-C^\alpha-C-N}(r_{il})C_\chi^{N-C^\alpha-C-N}(\chi_{ijkl}). \quad (1)$$



**Figure 1**

Schematic representation of the conditions defining a linear element  $N-C^\alpha-C-N-C^\alpha$ . Shown are a protein fragment containing this linear element and schematic representations of the conditions involved. These conditions depend on number of neighbouring atoms  $n$  (for two shells with radii  $d_1$  and  $d_2$ ), interatomic distances  $r$  and torsion angles  $\chi$ . The interactions present in this fragment are shown in dashed lines for  $L = [0, 4]$ . For convenience, bonds are shown in solid lines. For each layer  $L$  a single example is highlighted in black and the conditions applicable are given on the right-hand side.

**Table 1**

Number of pentapeptide and heptapeptide configurations used for defining the force-field parameters.

Shown are the number of heptapeptide configurations extracted in  $\alpha$ -helical ( $\alpha$ ) conformation, the number of pentapeptides in  $\beta$ -strand ( $\beta$ ) and loop ( $I^{AA}$ ,  $I^{AB}$ ,  $I^{BA}$  and  $I^{BB}$ ) conformations and the number of  $\chi_1$  rotamer conformations,  $g^-$ ,  $t$  or  $g^+$ , observed for the middle residues of the pentapeptides and heptapeptides.

Secondary structure	No. of peptides	$g^-(\chi_1)$	$t(\chi_1)$	$g^+(\chi_1)$
$\alpha$	1999	910	609	0
$\beta$	2332	883	721	285
$I^{AA}$	1590	515	190	312
$I^{AB}$	2180	924	199	457
$I^{BA}$	1822	541	613	191
$I^{BB}$	2562	930	546	440

The function  $JC_{ijkl}^{N-C^\alpha-C-N}$  will take on the value 1 when the configuration of atoms  $i, j, k$  and  $l$ , with numbers of neighbouring atoms  $n$ , interatomic distances  $r$  and dihedral angle  $\chi$ , matches all individual conditions. This implies that these atoms have adopted an  $N-C^\alpha-C-N$  conformation.

A protein structure can be described by the sum of its linear elements. Therefore, we define a least-squares target function that depends on the expected number of conformations present in the target structure,

$$E = \sum_{\text{type}} E^{\text{type}} = \sum_{\text{type}} w^{\text{type}} \left( TC^{\text{type}} - \sum_{ij\dots pq} JC_{ij\dots pq}^{\text{type}} \right)^2, \quad (2)$$

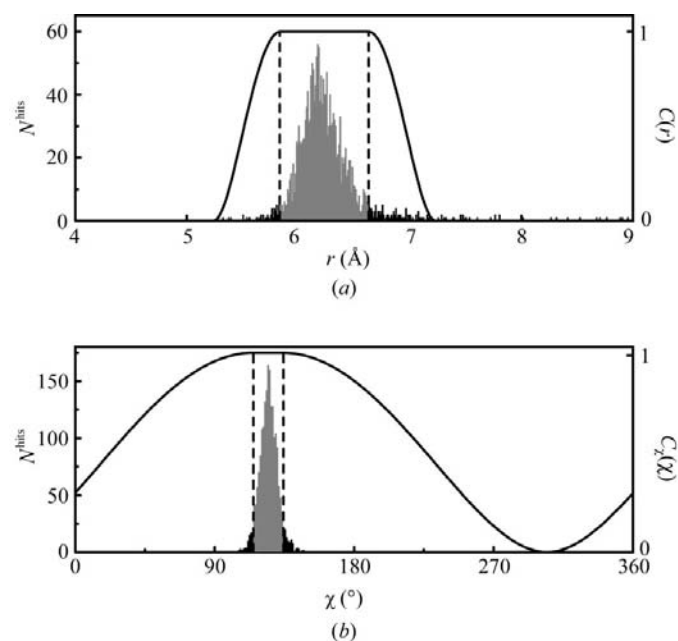
where  $TC^{\text{type}}$  is the expected sum of joint conditions for the types of linear elements in the target structure and  $w^{\text{type}}$  is a weighting factor. The first summation runs over all types of linear elements (of various lengths  $L^{\text{type}}$ ) that have been defined. The second summation runs over all possible combinations of  $L^{\text{type}} + 1$  atoms  $ij\dots pq$ . The minimum of this target function corresponds to a set of atoms with the expected number of linear elements in their expected types of conformations. Derivatives of this target function can be calculated with respect to all atomic coordinates. This allows the application of gradient-driven optimization techniques, which we term conditional optimization.

## 2.2. The general parameter set

For the mean-force potential, we defined 18 atom types ( $L = 0$ ) described by the expected numbers of neighbouring atoms within four neighbour shells of increasing radii (corresponding to typical distances for bonds, angles, torsion angles and Lennard–Jones interactions). We did not take into account glycine  $C^\alpha$  and proline N, which have deviating numbers of nearest neighbouring atoms. By combination of the atom types, we defined 26 bond types ( $L = 1$ ); these combine to form 43 types of angles ( $L = 2$ ). For longer fragments ( $L > 2$ ), separate conformations observed in  $\alpha$ -helices,  $\beta$ -strands and loops were defined. For  $\alpha$ -helices we defined conditions up to  $L = 12$ , for  $\beta$ -strands up to  $L = 9$  and for loops up to  $L = 7$ , taking into account the structural variability of the secondary-structure elements. For loops, separate conditions were defined for conformations corresponding to the A and B

regions of the Ramachandran plot, but conformations corresponding to the L region were not taken into account. For linear elements comprising two subsequent loop residues, separate conditions were defined for the possible combinations of  $\varphi/\psi$ -angle rotations AA, AB, BA and BB. For side-chain atoms up to the  $\gamma$  position we defined conditions according to the three preferred  $\chi_1$  rotamer conformations; in  $\alpha$ -helices only the two commonly observed  $\chi_1$  rotamers were defined. No distinction was made between the atoms at the  $\gamma$  position of different amino acids except for the cysteine  $S^\gamma$  atom; consequently,  $C^\gamma$  or  $O^\gamma$  atoms were treated equally. Side-chain atoms beyond the  $\gamma$  position were only defined up to  $L = 2$ , omitting information with respect to their rotamer conformations, which drastically reduced the number of possible combinations of defined conformations.

To determine minimum and maximum values for the condition parameters, distributions of observed numbers of neighbouring atoms, interatomic distances and torsion angles were calculated for the high-resolution protein structures in the *SCAN3D* database of the *WHATIF* program (Vriend *et al.*, 1994). The observed numbers of neighbouring atoms were calculated for 20 protein structures in this database, comprising a total of approximately 24 000 protein atoms. Observed interatomic distances and torsion angles were calculated from oligopeptides extracted using the *SCAN3D* structural annotation. Oligopeptides in a helical conformation were extracted as seven subsequent residues with an H (helix) assignment,  $\beta$ -strands were extracted as five subsequent residues with an S (strand) assignment and loops as five subsequent residues with a T (turn) or C (coil) assignment for the middle three residues. Backbone conformations with anno-



**Figure 2** (a) Observed distributions  $N_{\text{hits}}$  and defined conditions  $C$  for interatomic distance  $r$  and (b)  $N_{\text{hits}}$  and  $C_\chi$  for torsion angle  $\chi$  between the outermost atoms of a linear element comprising atoms  $C^\alpha(i)$  to  $C^\alpha(i + 4)$  in an  $\alpha$ -helical conformation. The minimum and maximum values for which  $C = 1$  are set to comprise 90% of the observed conformations (shown in grey).

**Table 2**

Number of conditions defined in the general parameter set for conditional optimization.

Conditions are defined for linear elements of different length ( $L$ ) and different chemical topologies. In addition, the defined conditions differentiate between distinct conformations of secondary-structure elements,  $\alpha$ ,  $\beta$ ,  $l^{AA}$ ,  $l^{AB}$ ,  $l^{BA}$  and  $l^{BB}$  and  $\chi_1$  rotamer conformations,  $g^-$ ,  $t$  and  $g^+$ .

Layer	No. of different topologies	Secondary-structure differentiation	$\chi_1$ rotamer differentiation	No. of conditions
$L = 0$	18	—	—	72
$L = 1$	26	—	—	26
$L = 2$	42	—	—	42
1	1	$\alpha, \beta, l$	—	3
$L = 3$	2	—	—	2
2	1	$\alpha, \beta, l$	—	6
4	1	$\alpha, \beta, l^A, l^B$	—	16
4	1	—	$g^-, t, g^+$	12
$L = 4$	4	$\alpha, \beta, l$	—	12
4	1	$\alpha, \beta, l^A, l^B$	—	16
3	1	$\alpha, \beta, l^A, l^B$	$g^-, t, g^+$	33
$L = 5$	9	$\alpha, \beta, l^A, l^B$	—	36
3	1	$\alpha, \beta, l^A, l^B$	$g^-, t, g^+$	33
$L = 6$	4	$\alpha, \beta, l^A, l^B$	—	16
4	1	$\alpha, \beta, l^A, l^B, l^{AB}, l^{BA}, l^{BB}$	—	24
4	1	$\alpha, \beta, l^A, l^B$	$g^-, t, g^+$	44
$L = 7$	4	$\alpha, \beta, l^A, l^B$	—	16
4	1	$\alpha, \beta, l^A, l^B, l^{AB}, l^{BA}, l^{BB}$	—	24
1	1	$\alpha, \beta, l^A, l^B$	$g^-, t, g^+$	11
2	1	$\alpha, \beta, l^A, l^B, l^{AB}, l^{BA}, l^{BB}$	$g^-, t, g^+$	34
$L = 8$	9	$\alpha, \beta$	—	18
3	1	$\alpha, \beta$	$g^-, t, g^+$	15
$L = 9$	8	$\alpha, \beta$	—	16
4	1	$\alpha, \beta$	$g^-, t, g^+$	20
$L = 10$	8	$\alpha$	—	8
3	1	$\alpha$	$g^-, t$	6
$L = 11$	9	$\alpha$	—	9
3	1	$\alpha$	$g^-, t$	6
$L = 12$	8	$\alpha$	—	8
4	1	$\alpha$	$g^-, t$	8
Total				592

† For helices, only conditions for  $\chi_1$  rotamers  $g^-$  and  $t$  were defined.

**Table 3**

Characteristics of the three test cases, human hyperplastic discs protein (PDB code 1i2t), erabutoxin (PDB code 3ebx) and turkey ovomucoid third domain (PDB code 1ds3).

PDB code	No. of atoms: protein/total	Secondary-structure content	Space group	$d_{\min}$ (Å)	No. of reflections† ( $d > 2$ Å)
1i2t	472/602	$\alpha$	$P2_12_12_1$	1.04	4662 (7)
3ebx	475/590	$\beta$ /loop	$P2_12_12_1$	1.4	3690 (0)
1ds3	378/426	$\alpha$ /loop	$P2_1$	1.65	2938 (13)

† The number of missing reflections is given in parentheses.

tated torsion angles  $-180 < \varphi < 0^\circ$  and  $-110^\circ < \psi < 50^\circ$  were termed A and conformations with  $-180 < \varphi < 0^\circ$  and  $50 < \psi < 180^\circ$  were termed B. Only  $\beta$ -strands with five subsequent residues in the B conformation were taken into account. For the middle residue of the extracted oligopeptides a distinction between the three  $\chi_1$  rotamers  $g^-$ ,  $t$  and  $g^+$  was made based on its value as annotated in the database:  $-120 < \chi_1 < 0^\circ$ ,  $120 < \chi_1 < 240^\circ$  and  $0 < \chi_1 < 120^\circ$ , respectively. Table 1 shows the total numbers of extracted oligopeptides in the different conformations that were used to determine the corresponding condition parameters.

For each condition type, the minimum and maximum values of the condition parameter were set so as to comprise 90% of the conformations as observed in the *SCAN3D* database. Histograms were made of the observed numbers of neighbouring atoms, interatomic distances or torsion angles. Bin widths were chosen such that the top of each histogram reached at least 50 hits, except for distributions with less than 200 hits, where the top should reach at least 20 hits. For each histogram, a frequency cutoff value was chosen such that 90% of all hits lie within the interval ranging from the first to the last bin for which the number of hits exceeds this cutoff value. For this interval region, condition  $C$  corresponds to 1. The widths of the slopes (see Fig. 2) were set to 0.05 Å at layer  $L = 1$  up to 0.75 Å at layer  $L = 12$  for distance conditions; widths of neighbour conditions were respectively set to 1.5, 4.8, 12.7 and 26.7 neighbouring atoms for the four shells with increasing radii; widths of torsion-angle conditions were set to  $(360^\circ - \chi_{\max} + \chi_{\min})/2$ , thus providing a continuous function for the entire range of torsion angles. As an example, Fig. 2 shows the histograms of observed distances and torsion angles and the resulting conditions for a linear element of  $C^\alpha(i)$  to  $C^\alpha(i + 4)$  in an  $\alpha$ -helical conformation.

The complete conditional parameter set that was obtained as described above has been deposited as supplementary material<sup>1</sup>. A summary of the numbers of all defined conditions is given in Table 2.

### 2.3. Protein-specific force fields

The force field parameters as defined in §2.2 represent geometric expectations of common conformations as observed in many protein structures. To define the expectations for a specific protein, a subset is extracted from this general parameter set specific for that particular protein. Based on the known amino-acid sequence and estimated fractions of  $\alpha$ -helical,  $\beta$ -strand and loop content, occurrences of all types of linear elements are determined and used to calculate expected sums of joint conditions  $TC^{\text{type}}$ . In this calculation, we also take into account contributions from reminiscent conformations that give non-zero values for  $JC^{\text{type}}$ . For differentiation of loops into A and B conformations and differentiation of  $\chi_1$  rotamers, the expected fractions are set to the observed relative occurrences of these conformations in the *SCAN3D* database. The target functions corresponding to these types are grouped,

$$E^{\text{group}} = \left\{ \sum_{\text{group}} g^{\text{type}} \left( TC^{\text{type}} - \sum_{ij\dots pq} JC_{ij\dots pq}^{\text{type}} \right) \right\}^2, \quad (3)$$

where  $g^{\text{type}}$  (with  $\sum_{\text{group}} g^{\text{type}} = 1$ ) corresponds to the relative occurrence of each group member and the summation runs over all types that are part of the group.

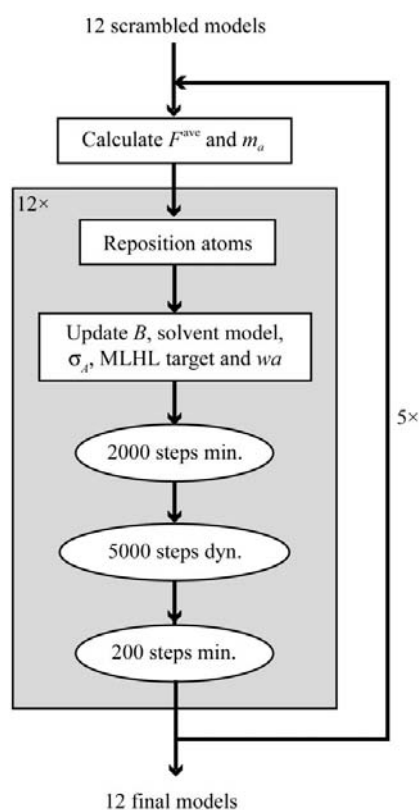
<sup>1</sup> Supplementary material has been deposited in the IUCr electronic archive (Reference: WD0002). Services for accessing this material are described at the back of the journal.

### 3. Experimental

Three small protein structures were selected for testing purposes: human hyperplastic discs protein (PDB code 1i2t), erabutoxin (PDB code 3ebx) and turkey ovomucoid third domain (PDB code 1ds3); see Table 3. These represent examples of all- $\alpha$ -helical, all- $\beta$ -sheet and mixed  $\alpha/\beta$ -fold, respectively. Published diffraction data sets were truncated at 2.0 Å resolution. All three data sets were nearly complete to this resolution limit. For the ovomucoid third domain test case, five of the lowest resolution reflections were marked as probable measurement errors and these reflections were removed from the reflection file. For these reflections an

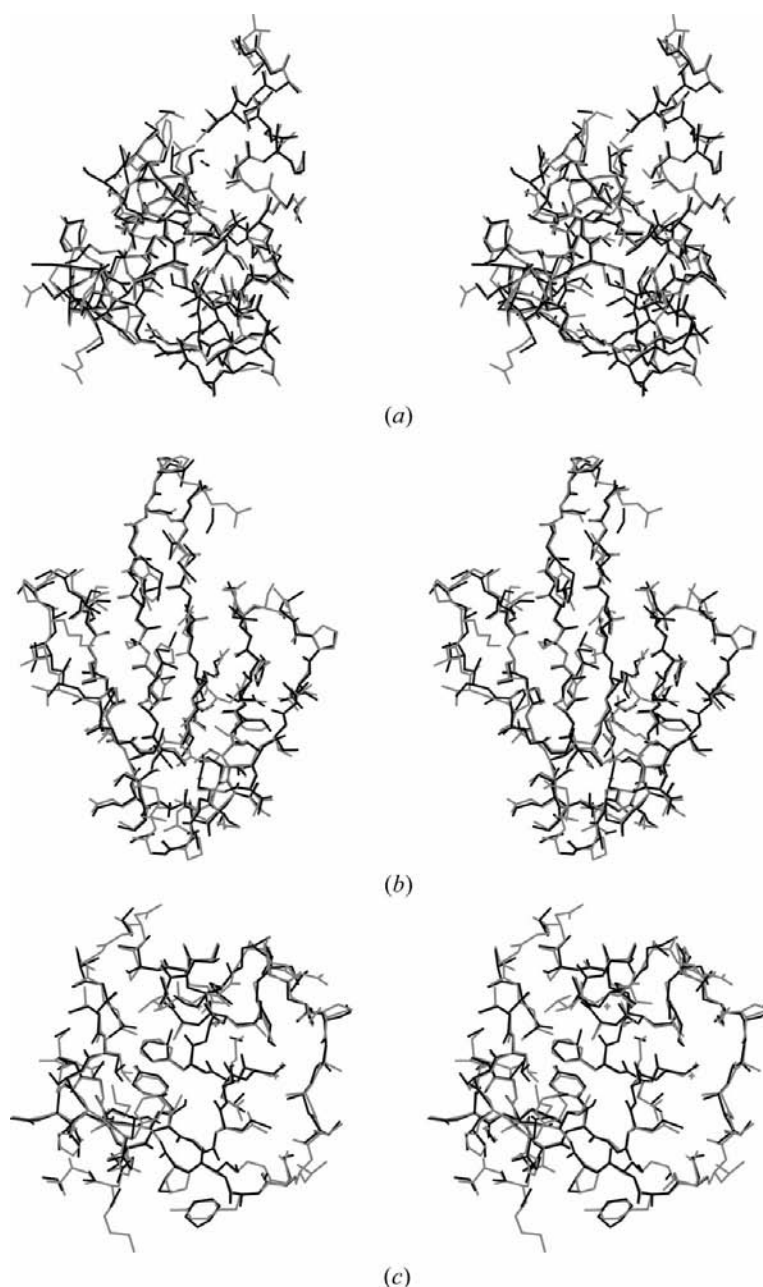
almost zero intensity was observed, while their calculated intensities were significantly higher.

Two aspects of conditional dynamics using the presented force field were tested: the stability of structures when starting with correct coordinates and the optimization behaviour for structures distant from the correct answer. To test the stability of structures corresponding to the correct answer, equilibrium runs were started from the deposited protein coordinates. These optimizations comprised 5000 steps of dynamics preceded and followed by 200 steps of minimization using conditional optimization implemented in *CNS* (Brünger *et al.*, 1998). We used the maximum-likelihood crystallographic



**Figure 3**

Refinement protocol for the optimization starting from 12 scrambled structures. Prior to every optimization cycle, average structure factor  $F^{\text{ave}}$  and figures of merit  $m_a$  were calculated from the 12 individual structure factor sets  $F^i$ . In every cycle a small amount of atoms was repositioned for each structure, based on its  $m_a|F^{\text{obs}}|\exp(i\varphi^{\text{ave}}) - D|F^{\text{calc}}|\exp(i\varphi^{\text{calc}})$  difference map. All atoms at density levels lower than  $-2.5\sigma$  and their neighbouring atoms (within 1.8 Å distance) with density lower than  $-1.5\sigma$  were selected for repositioning. These atoms were repositioned at the highest positive peaks of the difference map, with a minimum interatomic distance constraint of 1.2 Å and a triangulation constraint prohibiting the formation of a triangle of three bonded atoms. For each model, overall isotropic  $B$ -factor optimization, bulk-solvent correction, estimation of  $\sigma_A$  values based on figures of merit  $m_a$  and calculation of weight  $wa$  on the X-ray term of the target function (MLHL) were performed. Every optimization cycle comprised 5000 steps of dynamics calculations (dyn.) preceded and followed by 200 steps of energy minimization (min.) for each of the 12 structures.



**Figure 4**

Stereoviews of equilibrated structures in black superimposed on the target structures in grey of (a) the all- $\alpha$ -helical case 1i2t, (b) the all- $\beta$ -sheet (3ebx), and (c) the mixed  $\alpha/\beta$  (1ds3) test cases.

target function (MLF; Pannu & Read, 1996) with  $\sigma_A$  values estimated by Read's procedure (Read, 1986) based on 10% of free reflections (Brünger, 1993). Reflections for cross-

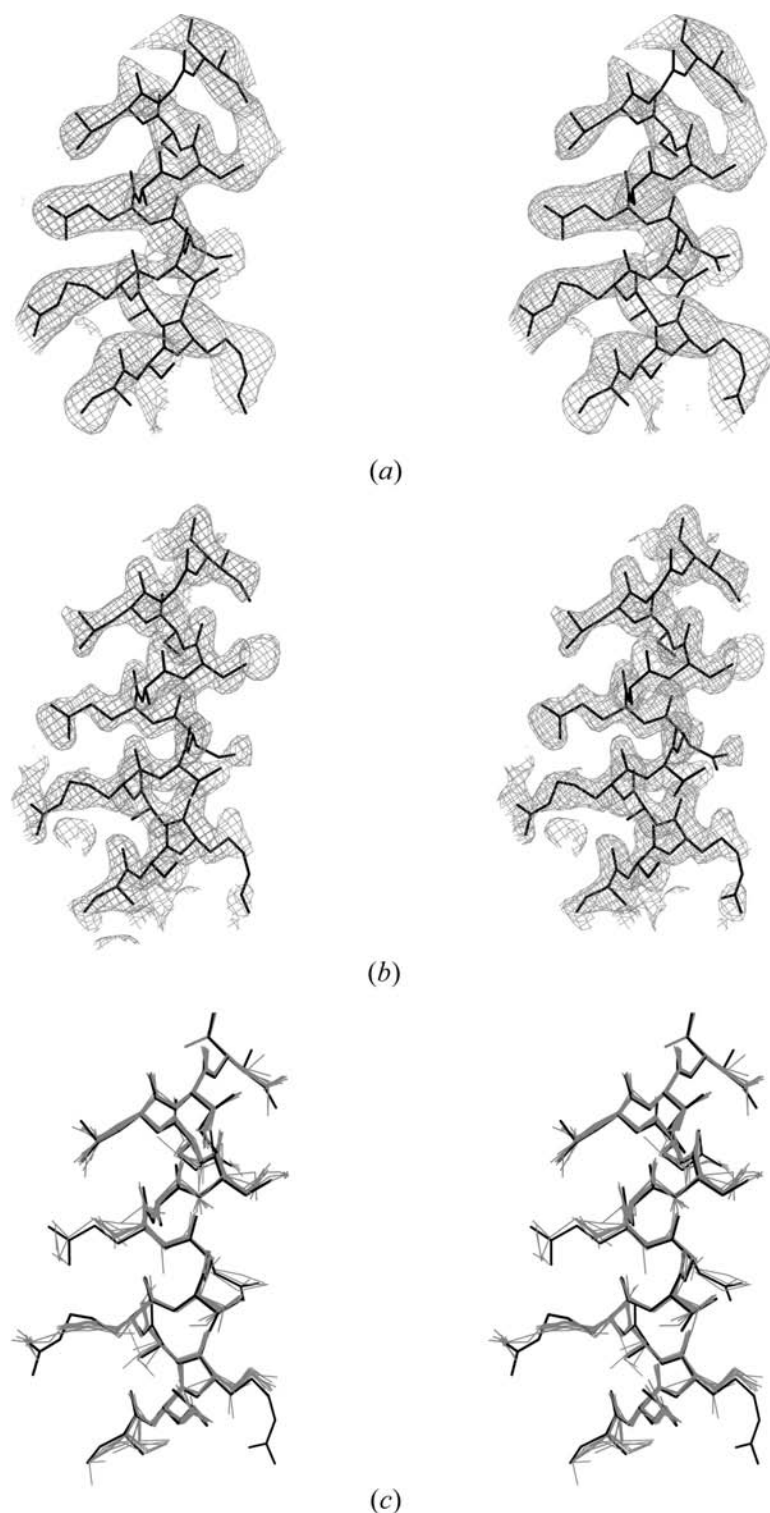
validation were selected randomly from reflections with a Bragg spacing  $d < 10 \text{ \AA}$ . To test the optimization behaviour for structures distant from the correct answer, optimization runs were performed starting from scrambled models with a root-mean-square (r.m.s.) coordinate error of  $1.5 \text{ \AA}$ . For each test case, 12 different starting models were generated by applying random coordinate shifts to all protein atoms. The scrambled models were refined according to the protocol as shown in Fig. 3. For each cycle of phase-restrained maximum-likelihood refinement (MLHL; Pannu *et al.*, 1998), target phases were obtained from the average structure factor  $F^{\text{ave}}$  of all 12 individual structure-factor sets  $F^i$ . (In the presented test cases, averaging the structure factors of the 12 starting models yielded phase errors of  $\sim 70^\circ$  for data to  $2 \text{ \AA}$  resolution. Phase errors of similar magnitude would result from a single model with an r.m.s. random coordinate error of  $\sim 1.1 \text{ \AA}$ .) Resolution-dependent figures of merit were calculated from the reflections in the test set as  $m'_a = \sum_{i=1}^N F^i / \sum_{i=1}^N |F^i|$  and extrapolated to  $N \rightarrow \infty$ :  $m_a = \{[N(m'_a)^2 - 1]/(N - 1)\}^{1/2}$ .  $\sigma_A$  estimates were calculated from these cross-validated figures of merit  $m_a$ , because the standard routine to estimate  $\sigma_A$  values gave spurious results for these structures with large errors and small numbers of reflections in the test set. Values for weights  $w_a$  on the X-ray restraint as determined with standard routines showed a strong variation over the 12 different structures. One common value for each cycle was determined by exploiting a relationship with the sum of the figure of merit over all reflections ( $w_a \propto 1/\sum m$ ), as observed during initial calculations with models of varying quality (results not shown).

For each test case equal atom labels, 'X', were given to all protein atoms and carbon scattering factors were assigned to all of them. Water and other non-protein atoms were not included in the calculations. Atomic  $B$  factors were assigned based on the number of neighbouring atoms as described previously (Scheres & Gros, 2001). Standard routines were used for scaling and bulk-solvent correction. To avoid negative atomic  $B$  factors after scaling, inverse scaling was applied to  $F^{\text{obs}}$  rather than scaling  $F^{\text{calc}}$ . Dynamics calculations were performed with a time step of  $0.2 \text{ fs}$  and the temperature was coupled to a bath of  $600 \text{ K}$ . All calculations were performed on four  $667 \text{ MHz}$  single-processor Compaq XP1000 workstations with at least  $1.2 \text{ Gb}$  of memory.

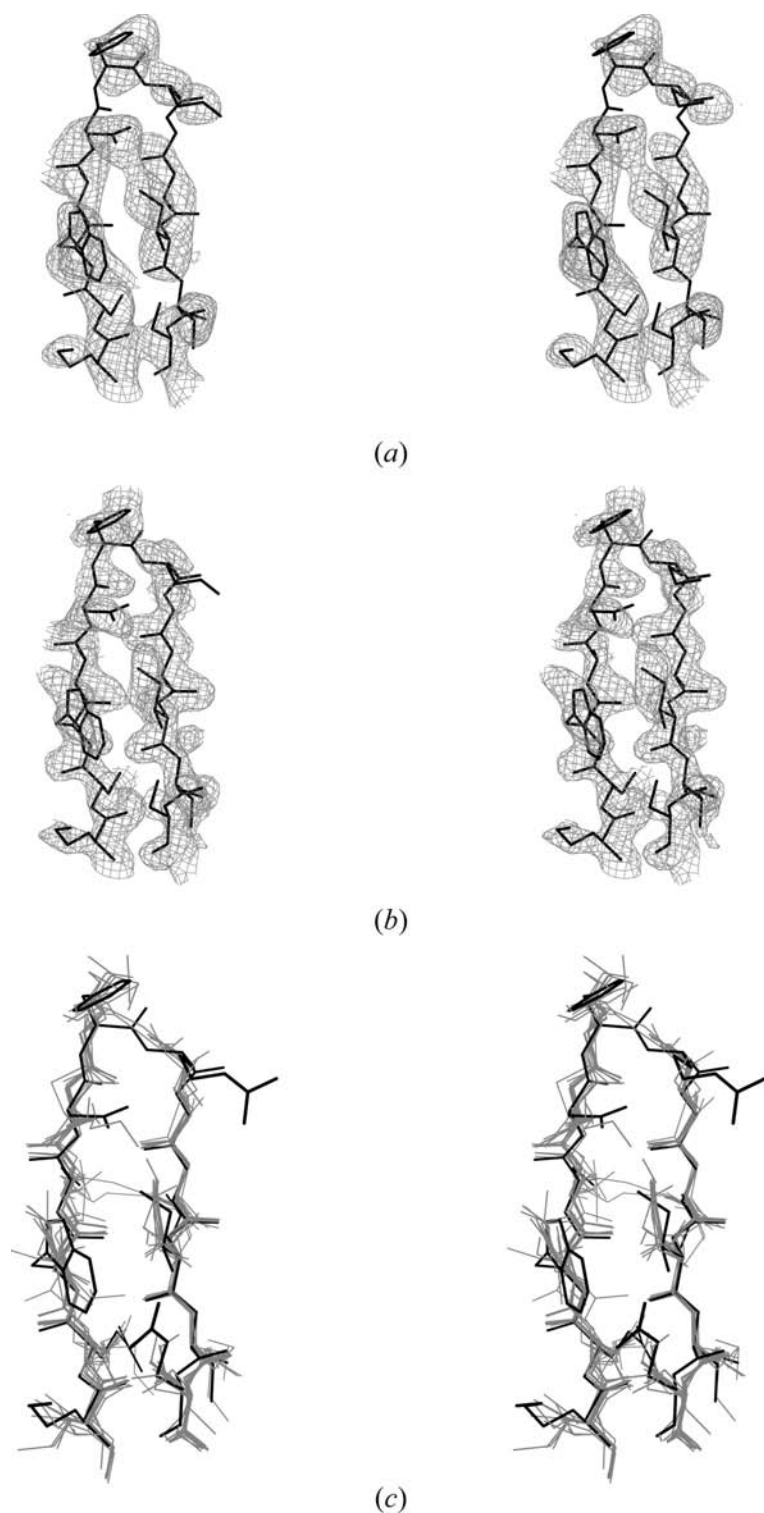
## 4. Results and discussion

### 4.1. Stability of correct structures

The first evaluation of the defined force field concerns the stability of correct protein structures in conditional optimization. Fig. 4 shows equilibrated



**Figure 5** Electron-density maps and models obtained by conditional optimization of 1i2t using 12 scrambled models. Shown are stereoviews of part of the  $m_a |F^{\text{obs}}| \exp(i\phi^{\text{ave}})$  electron-density maps obtained before (a) and after (b) optimization and the 12 final structures obtained in grey (c). The target structure of 1i2t is superimposed in black.


**Figure 6**

Electron-density maps and models obtained by conditional optimization of 3ebx using 12 scrambled models. Shown are stereoviews of part of the  $m_a|F^{\text{obs}}|\exp(i\varphi^{\text{ave}})$  electron-density maps obtained before (a) and after (b) optimization and the 12 final structures obtained in grey (c). The target structure of 3ebx is superimposed in black.

structures after dynamics calculations starting from the deposited coordinates of all three test cases. The mean phase errors of these structures increase from  $<20^\circ$  to  $\sim 30^\circ$  (see

**Table 4**

Results from equilibrium and optimization runs using the presented force field for conditional optimization.

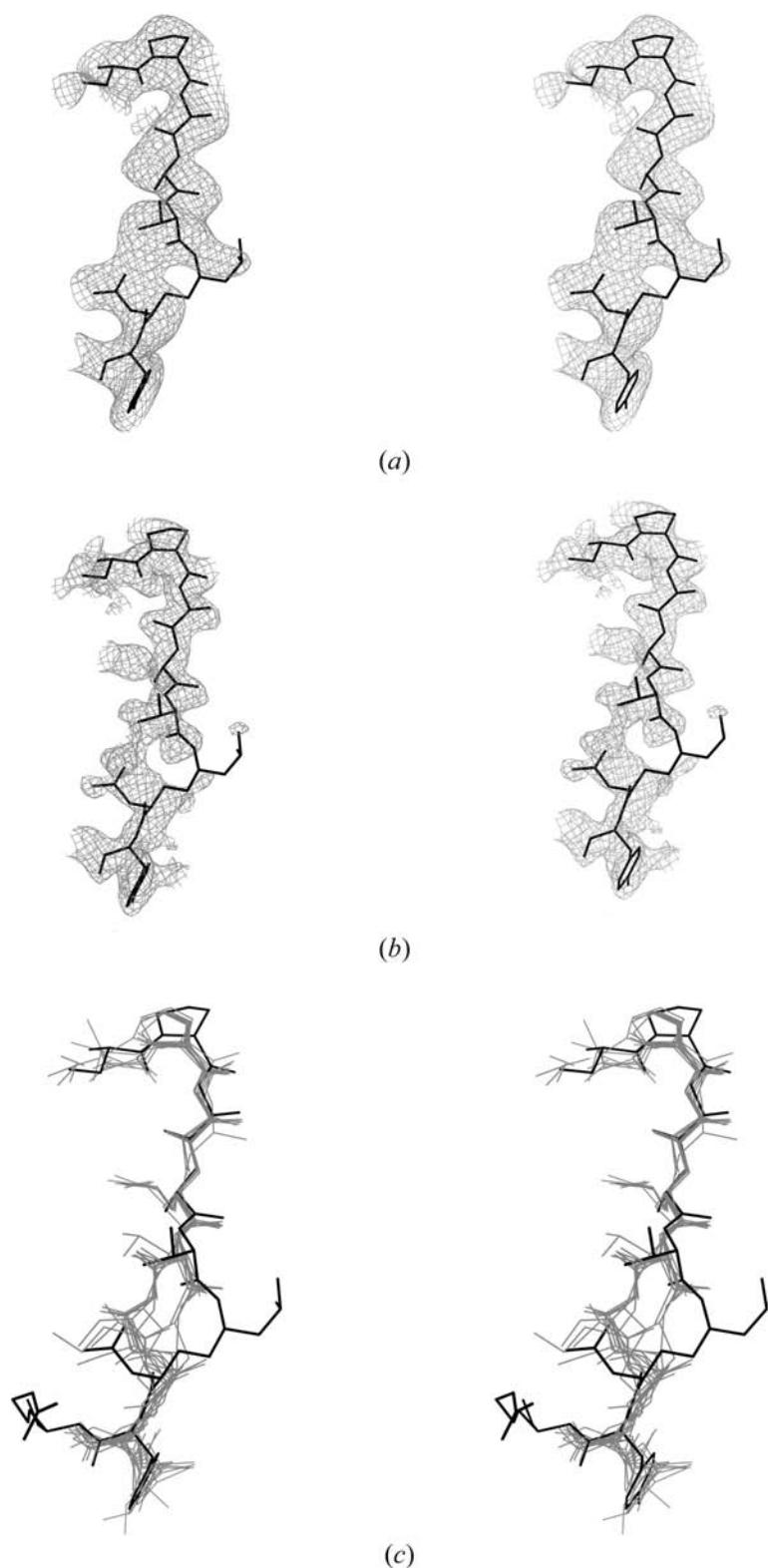
For each of the three test cases, human hyperplastic discs protein (PDB code 1i2t), erabutoxin (PDB code 3ebx) and turkey ovomucoid third domain (PDB code 1ds3), the secondary-structure content,  $\alpha$ -helix,  $\beta$ -sheet and loop, used to define the protein-specific force fields are given in percentages. Amplitude-weighted ( $|F^{\text{obs}}|$ ) mean phase errors before and after the equilibrium and optimization runs are given. Phase errors are calculated with respect to phases of the structures deposited in the PDB. In the case of optimization starting from 12 models, phase errors are given for the averaged structure factors. The CPU times are given that were required for each of the 12 models in these optimization runs.

Test case	Protein-specific force field	$\Delta\varphi$ ( $^\circ$ )				CPU time (h)
		Before equilibration	After equilibration	Before optimization	After optimization	
1i2t	100% $\alpha$	19	27	71	28	10
3ebx	100% $\beta$	19	32	70	45	12
1ds3	25% $\alpha$ , 25% $\beta$ , 50% loop	14	27	71	45	18

Table 4), but the corresponding electron-density maps are still easily interpretable. Errors that are introduced during these runs can be attributed to conformations for which no or limited conditions were defined. In the all- $\alpha$  case, a single main-chain break occurs in a turn next to a proline residue. The all- $\beta$  case shows three main-chain breaks that concern two residues with a conformation in the L region and one glycine in a conformation outside any of the three common regions of the Ramachandran plot. For the mixed  $\alpha/\beta$  case three main-chain breaks are also observed related to conformations outside the A and B region of the Ramachandran plot. For all three test cases side chains beyond the  $\gamma$  position are unstable and atoms at the  $\delta$ ,  $\epsilon$ ,  $\zeta$  and  $\eta$  positions of the side chains are displaced from their correct positions during equilibration. Since unstable parts in the protein structures coincide with conformations that were poorly or not defined, extension of the parameter set to describe these conformations may lead to better modelling of the target structure at the expense of requiring more computing power.

#### 4.2. Searching behaviour in optimization

A second requirement for the force field presented is favourable searching behaviour in the optimization of structures (far) away from the defined minimum. Optimization runs were performed for all three test cases, starting from 12 scrambled structures with coordinate errors of 1.5 Å r.m.s.d. Figs. 5, 6 and 7 show optimized structures and map improvements for the all- $\alpha$ , all- $\beta$  and mixed  $\alpha/\beta$  test cases, respectively. Corresponding phase improvements and CPU times required for these runs are given in Table 4. For the all- $\alpha$  test case, optimization converges readily towards the global



**Figure 7**  
Electron-density maps and models obtained by conditional optimization of 1ds3 using 12 scrambled models. Shown are stereoviews of part of the  $m_a|F^{\text{obs}}|\exp(i\varphi^{\text{ave}})$  electron-density maps obtained before (a) and after (b) optimization and the 12 final structures obtained in grey (c). The target structure of 1ds3 is superimposed in black.

minimum. Subsequent refinement cycles yield significant improvement of the electron-density map and phase infor-

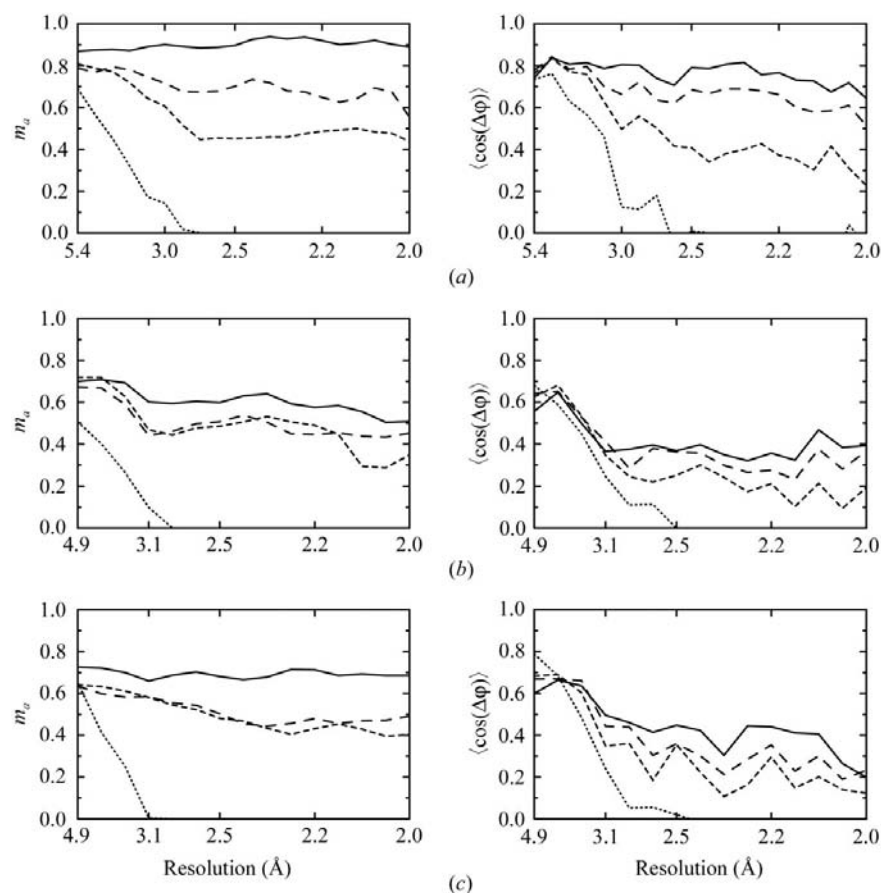
mation over the whole resolution range. Errors in the optimized structures coincide with conformations that were also unstable during equilibration. For the all- $\beta$  and mixed  $\alpha/\beta$  cases, optimization convergences less readily, but considerable phase improvement is still obtained. In addition to the parts unstable during equilibration, most of the errors in the optimized structures are observed in the loop regions and include missing and false main-chain connections. For some of the  $\beta$ -strands, we also observe inadvertent reversal of the chain direction.

The all- $\alpha$  helical test case performs significantly better in the conditional optimization than the all- $\beta$  sheet and mixed  $\alpha/\beta$  test cases. This difference may be attributed to various reasons: (i) for the all- $\alpha$  test case estimated figures of merit  $m_a$  are in good agreement with the mean cosine of the phase error, while for the all- $\beta$  and mixed  $\alpha/\beta$  test cases significant overestimation is observed (see Fig. 8); (ii) the all- $\alpha$  test case has a higher solvent content ( $\sim 50\%$ ) than the other two cases (both  $\sim 35\%$ ), resulting in a significantly larger number of reflections to 2.0 Å resolution (see Table 3); (iii) the information content of the used force field is higher for the all- $\alpha$  test case than for the all- $\beta$  and mixed  $\alpha/\beta$  test cases; (iv) the proteins from the all- $\beta$  and mixed  $\alpha/\beta$  test cases contain more conformations that are not accounted for in the used force fields.

## 5. Conclusions

We introduced a mean-force potential for conditional optimization of protein structures. The interaction functions in this force field describe protein fragments in  $\alpha$ -helical,  $\beta$ -strand and loop conformations of up to four, three and two residues long, respectively. Distinct interaction functions for the three preferred  $\chi_1$  rotamers describe corresponding geometries for side chains up to the  $\gamma$  position. Notably, we omitted glycine and proline residues, main-chain conformations involving the L region of the Ramachandran plot and torsion angles (and higher order information) for side-chain atoms beyond the  $\gamma$  position, owing to increasing computational costs. We tested the parameter set in conditional optimization of three small protein structures using 2.0 Å observed diffraction data. Dynamics runs starting from the deposited coordinates show that the definition of the global minimum is correct for the defined main-chain conformations and for side chains up to the  $\gamma$  position. Breaks are observed for main-chain conformations outside the A and B region and for side chains beyond the  $\gamma$  position that were undefined or poorly defined. A more precise definition of these conformations in the force field could improve the optimization behaviour. However, inclusion of the omitted elements would give rise to a large increase in the




**Figure 8**

Estimated figures of merit  $m_a$  and average cosine of the phase error,  $\langle \cos(\Delta\phi) \rangle$ , for the (a) all- $\alpha$  helical (1i2t), (b) all- $\beta$  sheet (3ebx) and (c) mixed  $\alpha/\beta$  (1ds3) test cases. Values are shown as calculated for the initial scrambled structures (dotted lines), structures after optimization cycle 1 (dashed lines) and cycle 2 (long-dashed lines) and for the final optimized structures (solid lines).

number of possible combinations, increasing the computational cost dramatically.

Optimization starting from 12 structures with 1.5 Å r.m.s.d. random coordinate shifts showed excellent convergence for the  $\alpha$ -helical hyperplastic discs protein. Considerable phase improvement was also obtained for the  $\beta$ -sheet protein erabutoxin and the ovomucoid third domain with mixed  $\alpha/\beta$  fold, but the optimized structures contain more errors: typically, chain reversals for  $\beta$ -strands and incorrect formation of loops. The applied multiple-model procedure proved crucial for these optimizations, since with the limited numbers of available test-set reflections standard procedures to estimate phase quality failed for starting models with such large errors. In contrast to the all- $\alpha$  helical case, significant overestimation of the phase quality was observed for the all- $\beta$  sheet and mixed  $\alpha/\beta$  test cases. This overestimation coincides with the more difficult convergence in the optimization runs of the all- $\beta$

and mixed  $\alpha/\beta$  test cases, which may indicate the importance of further improvement of this procedure.

Our results illustrate that a large radius of convergence may be obtained by conditional optimization of protein molecules with observed diffraction data to medium resolution. The coordinate errors of our starting models were generated in a completely random way and such favourable error distributions are hard to obtain when starting from a single electron-density map. In addition, we used truncated data, which also may have contributed favourably to the optimization behaviour. Still, the significant reduction in phase errors from  $\sim 70$  to  $45^\circ$  or better is promising. The generally applicable mean-force potential presented allows development of phase-improvement and automated model-building procedures using conditional optimization, as well as investigation of the efficacy of this approach in *ab initio* phasing of protein structures.

This work is supported by the Netherlands Organization for Scientific Research (NWO-CW: Jonge Chemici 99-564).

## References

- Berman, H. M., Battistuz, T., Bhat, T. N., Bluhm, W. F., Bourne, P. E., Burkhardt, K., Feng, Z., Gilliland, G. L., Iype, L., Jain, S., Fagan, P., Marvin, J., Padilla, D., Ravichandran, V., Schneider, B., Thanki, N., Weissig, H., Westbrook, J. D. & Zardecki, C. (2002). *Acta Cryst.* **D58**, 899–907.
- Brünger, A. T. (1993). *Acta Cryst.* **D49**, 24–36.
- Brünger, A. T., Adams, P. D., Clore, G. M., DeLano, W. L., Gros, P., Grosse-Kunstleve, R. W., Jiang, J.-S., Kuszewski, J., Nilges, M., Pannu, N. S., Read, R. J., Rice, L. M., Simonson, T. & Warren, G. L. (1998). *Acta Cryst.* **D54**, 905–921.
- Lamzin, V. S., Perrakis, A., Bricogne, G., Jiang, J., Swaminathan, S. & Sussman, J. L. (2000). *Acta Cryst.* **D56**, 1510–1511.
- Pannu, N. S., Murshudov, G. N., Dodson, E. J. & Read, R. J. (1998). *Acta Cryst.* **D54**, 1285–1294.
- Pannu, N. S. & Read, R. J. (1996). *Acta Cryst.* **A52**, 659–668.
- Perrakis, A., Morris, R. & Lamzin, V. S. (1999). *Nature Struct. Biol.* **6**, 458–463.
- Read, R. J. (1986). *Acta Cryst.* **A42**, 140–149.
- Scheres, S. H. W. & Gros, P. (2001). *Acta Cryst.* **D57**, 1820–1828.
- Sippl, M. J. (1995). *Curr. Opin. Struct. Biol.* **5**, 229–235.
- Terwilliger, T. C. & Berendzen, J. (1999). *Acta Cryst.* **D55**, 849–861.
- Vriend, G., Sander, C. & Stouten, P. F. (1994). *Protein Eng.* **7**, 23–29.