

# A practical implicit solvent potential for NMR structure calculation



Ye Tian<sup>a,b</sup>, Charles D. Schwieters<sup>c</sup>, Stanley J. Opella<sup>b</sup>, Francesca M. Marassi<sup>a,\*</sup>

<sup>a</sup> Sanford-Burnham Medical Research Institute, 10901 North Torrey Pines Road, La Jolla, CA 92037, USA

<sup>b</sup> Department of Chemistry and Biochemistry, University of California San Diego, 9500 Gilman Drive, La Jolla, CA 92093-0307, USA

<sup>c</sup> Division of Computational Bioscience, Building 12A, Center for Information Technology, National Institutes of Health, Bethesda, MD 20892-5624, USA

## ARTICLE INFO

### Article history:

Received 22 January 2014

Revised 18 March 2014

Available online 2 April 2014

### Keywords:

Protein structure

Calculation

Implicit solvent

EEFx

XPLOR-NIH

NMR

## ABSTRACT

The benefits of protein structure refinement in water are well documented. However, performing structure refinement with explicit atomic representation of the solvent molecules is computationally expensive and impractical for NMR-restrained structure calculations that start from completely extended polypeptide templates. Here we describe a new implicit solvation potential, EEFx (Effective Energy Function for XPLOR-NIH), for NMR-restrained structure calculations of proteins in XPLOR-NIH. The key components of EEFx are an energy term for solvation energy that works together with other nonbonded energy functions, and a dedicated force field for conformational and nonbonded protein interaction parameters. The initial results obtained with EEFx show that significant improvements in structural quality can be obtained. EEFx is computationally efficient and can be used both to fold and refine structures. Overall, EEFx improves the quality of protein conformation and nonbonded atomic interactions. Moreover, such benefits are accompanied by enhanced structural precision and enhanced structural accuracy, reflected in improved agreement with the cross-validated dipolar coupling data. Finally, implementation of EEFx calculations is straightforward and computationally efficient. Overall, EEFx provides a useful method for the practical calculation of experimental protein structures in a physically realistic environment.

© 2014 Elsevier Inc. All rights reserved.

## 1. Introduction

The physical nature of a protein's physiological environment has long been recognized as an important factor governing molecular structure and biological function [1]. For example, the interactions of key amino acid residues with water and solvated molecules play key roles in protein activity and the highly anisotropic environment of the lipid bilayer membrane poses significant constraints on the structures and functions of membrane proteins.

Among the methods for three-dimensional molecular structure determination, the principal advantage of NMR spectroscopy is its ability to examine proteins in samples that are very close to their functional environments [2,3]. Yet, typically, even when NMR spectra are measured for soluble proteins in aqueous solutions or for membrane proteins in lipid bilayers, structure calculations are carried out with energy functions that do not include contributions from solvation energy and, instead, represent all the nonbonded electrostatic and van der Waals interactions by a single, purely repulsive term. Such a simplified treatment is very useful because it enables fast, restrained molecular dynamics (MD) calculations of high quality structures from fully extended polypeptide

templates by simulated annealing, to enhance sampling of conformational space and efficiently overcome the local-minimum problem [4,5]. However, this approach can also lead to structures with suboptimal quality parameters, such as poor packing, unsatisfied hydrogen bond donors or acceptors and unbalanced salt bridges.

Performing structure refinement in a full MD force field, with explicit atomic representation of the solvent molecules (water and/or lipid), is one way to improve structural quality and obtain information about a protein's interactions with its surroundings, as shown for structure refinement of both soluble proteins [6–11] and membrane proteins [12–14]. However, this approach is computationally expensive due to the large amount of time that must be devoted to calculating solvent–solvent interactions and, hence, is not practical for *ab initio* NMR structure calculations starting from completely extended polypeptide templates.

Methods in which solvent effects are treated implicitly have also been used in the refinement stages of NMR-restrained structure calculations. For example, refinement of initial structures using restrained MD calculations with generalized Born (GB) implicit solvent models have been shown to improve structural quality, particularly in cases where the experimental data are limited and the characterization of solvent effects is critical for identifying native fold [15–17]. However, GB methods have not been implemented as an integral parts of NMR structure calculation

\* Corresponding author.

E-mail address: [fmarassi@sbmri.org](mailto:fmarassi@sbmri.org) (F.M. Marassi).

rotocols from unfolded templates and remain too computationally intensive for routine calculations.

Various other models have been developed for the implicit treatment of solvent effects (reviewed in [18–22]). Of these, the effective energy function EEF1, developed by Lazaridis and Karplus [23], is particularly well suited to NMR applications for several reasons. EEF1 is based on the thermodynamic hypothesis that the native fold of a protein is the state of lowest free energy under physiological conditions and is determined by the amino acid sequence within the given solvent environment [1]. It contains terms for both intramolecular energy and solvation free energy and has been shown both to provide a realistic first approximation to the effective energy hypersurface of proteins and to work well for protein folding–unfolding studies [23–25]. Importantly, it enables fast calculation of energies and energy derivatives, and its extension to an implicit membrane model (IMM) [26] enables its application to membrane proteins, thus providing access to the majority of protein families found in nature.

Here we describe the development and implementation of an implicit solvation model based on EEF1 for NMR-restrained structure calculations in the program XPLOR-NIH [27,28]. The model is named EEFx (Effective Energy Function for XPLOR-NIH) to highlight its origins. The principal components of EEFx are a new nonbonded energy function for solvation free energy ( $E_{slv}$ ) and the related parameters that enable its implementation with the other XPLOR-NIH energy terms for NMR structure calculations. The XPLOR-NIH package is derived from XPLOR [29], which itself evolved from the CHARMM program [30,31]. XPLOR-NIH has many completely new features, designed to facilitate its applications and continuous development for NMR structure calculations. However, it also contains all of the original XPLOR functionality, including many energy functions derived from CHARMM. This is a significant factor facilitating the implementation of EEFx in XPLOR-NIH, since its EEF1 progenitor was originally developed for CHARMM and works in conjunction with CHARMM energy functions and force fields.

We show that EEFx yields significant improvements in structural quality, accuracy and precision for NMR-restrained calculations, from unfolded templates, of several proteins with sizes ranging from 60 to 260 amino acids. Structure calculations with EEFx are computationally efficient and can be easily implemented together with standard simulated annealing protocols. We anticipate that structure calculations using EEFx with extension to an implicit membrane environment potential (in progress) will be particularly useful for membrane proteins, where the highly anisotropic environment of the lipid bilayer membrane poses significant constraints on protein structure [32].

## 2. Description of EEFx

The XPLOR-NIH energy function ( $E_{TOTAL}$ ) can be grouped into three distinct classes [27–31]:

$$E_{TOTAL} = E_{EXP} + E_{KNOW} + E_{SYS} \quad (1)$$

where  $E_{EXP}$  contains experimental restraining energy terms derived from the NMR data,  $E_{KNOW}$  contains knowledge-based restraining terms and  $E_{SYS}$ , describing the energy of the molecular system, contains conformational and nonbonded energy terms. Many  $E_{EXP}$  and  $E_{KNOW}$  potentials have been developed for NMR-restrained structure calculations in XPLOR-NIH, including the widely used terms for distance, dipolar coupling and dihedral angle restraints [27,28], as well as various statistical torsion angle potentials [33,34].

In typical NMR structure calculations the conformational energy comprises terms for covalent bonds ( $E_{BON}$ ), covalent bond angles ( $E_{ANG}$ ) and improper dihedral angles ( $E_{IMP}$ ), and can include a

term for proper dihedral angles ( $E_{DIHE}$ ), although this is usually more effectively replaced by a statistical knowledge-based term. Furthermore, the nonbonded energy is described collectively by a single repulsive potential, implemented by turning on the *repel* option of the van der Waals energy function ( $E_{VDW-REPEL}$ ) and turning off the electrostatic energy function; this simplified term is used to prevent atomic overlap and can be scaled down to allow atoms to move through each other in the early stages of simulated annealing to accelerate the calculations [4,5].

By contrast, the nonbonded energy function of EEFx contains terms for three types of interactions: a Lennard–Jones van der Waals term ( $E_{VDW}$ ), to describe both repulsive and attractive forces; an electrostatic energy term ( $E_{ELEC}$ ), computed with the atomic charges specified in the topology file; and a new term for solvation free energy ( $E_{slv}$ ), introduced here to enable protein structure calculations in implicit solvent with EEFx.

$E_{slv}$  is an empirical function for the solvation free energy of a protein in water, parametrized with experimental solvation free energy data for small model molecules. It works together with the XPLOR functions for nonbonded energy ( $E_{VDW}$  and  $E_{ELEC}$ ) and conformational energy ( $E_{BON}$ ,  $E_{ANG}$ ,  $E_{IMP}$  and  $E_{DIHE}$ ), plus a new set of protein topology and parameters to generate the EEFx force field, such that:

$$\begin{aligned} E_{SYS} &= E_{EEFx} \\ &= (k_{BON}E_{BON}) + (k_{ANG}E_{ANG}) + (k_{IMP}E_{IMP}) + (k_{DIHE}E_{DIHE}) \\ &\quad + (k_{VDW}E_{VDW}) + (k_{ELEC}E_{ELEC}) + (k_{SLV}E_{SLV}) \end{aligned} \quad (2)$$

where  $E_{EEFx}$  is the effective energy of the solvated system and each energy term is scaled by its respective force constant  $k$ . The derivation of  $E_{slv}$  has been described [23] and its implementation in XPLOR-NIH is described below.

$E_{slv}$  is defined as the sum of the solvation free energy contributions from all  $i$  atomic groups in the protein, each described as the solvation free energy of group  $i$  in its fully solvated state minus the reduction in solvation due to the presence of surrounding groups  $j$ . The functional form of  $E_{slv}$ , expressed as an empirical energy function of the protein's atomic coordinates is given by:

$$\begin{aligned} E_{SLV} &= \sum_i \Delta G_i^{slv} \\ &= \sum_i \Delta G_i^{ref} - \sum_i \sum_{j \neq i} SW(r_{ij}) \left( \frac{2\Delta G_i^{free}}{4\pi\sqrt{\pi}\lambda_i} \right) \exp \left[ -\frac{(r_{ij} - R_i)^2}{\lambda_i^2} \right] \frac{V_j}{r_{ij}^2} \end{aligned} \quad (3)$$

where  $\Delta G_i^{ref}$  and  $\Delta G_i^{free}$  each represent the solvation free energy of atomic group  $i$  in its fully solvated state ( $\Delta G_i^{ref}$ ) and in its isolated, fully solvated state ( $\Delta G_i^{free}$ ),  $SW_{slv}$  is a switching function,  $r_{ij}$  is the distance between groups  $i$  and  $j$ ,  $R_i$  is the van der Waals radius of group  $i$ ,  $V_j$  is the volume of group  $j$ , and  $\lambda_i$  is the correlation length of the solvation free energy density centered at group  $i$ . In Eq. (3), the solvation free energy density of each group is modeled as a Gaussian function, exhibiting strong distance-dependence with its maximum magnitude centered at group  $i$  and decaying to zero away from it. The switching function  $SW_{slv}$  is similar to the XPLOR switching function used to control the atomic distances at which the nonbonded interaction terms for  $E_{VDW}$  and  $E_{ELEC}$  become effective [29].  $SW_{slv}$  has the form:

$$SW_{slv}(r_{ij}) = \begin{cases} 0 & \text{if } r_{ij} > r_{off} \\ \frac{\left( \frac{r_{ij}^2 - r_{off}^2}{r_{off}^2 - r_{on}^2} \right)^2 \left( \frac{r_{off}^2 + 2r_{ij}^2 - 3r_{on}^2}{r_{off}^2 - r_{on}^2} \right)}{1} & \text{if } r_{off} > r_{ij} > r_{on} \\ 1 & \text{if } r_{on} > r_{ij} \end{cases} \quad (4)$$

Eqs. (3) and (4), as well as equations for evaluating the analytical derivatives of  $\Delta G_i^{slv}$  needed to generate gradients of the pairwise

solvation free energies for all groups in the system, were coded in the C++ base framework of XPLOR-NIH, and new Python modules, *eefxPot* and *eefxPotTools*, were added to facilitate set up of the solvation energy term for EEFx calculations.

Evaluation of the solvation free energy using Eqs. (3) and (4) requires specification of the following parameters for each group  $i$  in the protein: atom type, atom radius  $R_i$ , atom volume  $V_i$ , atom correlation length  $\lambda_i$ , and values of  $\Delta G_i^{\text{ref}}$  and  $\Delta G_i^{\text{free}}$ . These parameters, plus related values of heat capacity and enthalpy, were encoded in the *eefxPotTools* module and are called by *eefxPot* during structure calculations. Their values (Table 1) were selected to be the same as those of Lazaridis and Karplus [23], with the values for hydrogen atoms taken to be zero. The values of  $\Delta G_i^{\text{ref}}$  were taken from the experimental data of Privalov and Makhatadze [35] while those of  $\Delta G_i^{\text{free}}$  are derived from  $\Delta G_i^{\text{ref}}$  as described [23]. The values of  $R_i$  correspond to CHARMM19 atomic radii and the volumes  $V_i$  are derived from  $R_i$ , as described [23]. The values of  $\lambda_i$  are set to the radius of the first hydration shell: 6.0 Å for NH3, NC2 and OC groups and 3.5 Å for all other groups.

Solvation free energy has a strong dependence on temperature. The CHARMM EEF1 model was found to reproduce thermodynamic parameters of protein folding-unfolding events very well and EEFx retains this thermodynamic functionality. The solvation free energy parameters encoded in *eefxPotTools* are values determined experimentally at 298.15 K [35]. For structure calculations based on NMR experiments performed at temperatures different than 298.15 K, *eefxPot* can perform temperature-dependent calibration of the free energy parameters in *eefxPotTools*, using values of heat capacity and enthalpy, also derived experimentally from model small molecules [35–38]. This is accomplished with the following functional forms of the Gibbs–Helmholtz equation, encoded in *eefxPot*:

$$\begin{aligned}\Delta G_i^{\text{ref}}(T) &= \Delta G_i^{\text{ref}}(T_0) - \Delta S_i^{\text{ref}}(T_0)(T - T_0) - \Delta C_{p,i} T \ln\left(\frac{T}{T_0}\right) \\ &\quad + \Delta C_{p,i}(T - T_0) \\ \Delta S_i^{\text{ref}}(T_0) &= \frac{\Delta H_i^{\text{ref}}(T_0) - \Delta G_i^{\text{ref}}(T_0)}{T_0}\end{aligned}\quad (5)$$

EEFx also requires dedicated protein topology and parameter files, containing the chemical information for specific residue and atom types and the various force constants for the conformational

and nonbonded energy terms. The original EEF1 model was designed to work with the CHARMM19 polar hydrogen force field, in which aliphatic hydrogens are treated implicitly by representing aliphatic groups as unified atoms with increased mass and specific van der Waals properties. Since NMR structure calculations require explicit inclusion of all hydrogen atoms, we generated a new set of parameter and topology files for specific use with EEFx. The new files, *proteinEEFx.par* and *proteinEEFx.top*, were derived from the CHARMM19 [30,31], PARALLHDG5.3 [6,8] and OPLS [39] force fields. Since PARALLHDG5.3 is itself derived from CHARMM and since PARALLHDG5.3 and OPLS are jointly effective for explicit water refinement in XPLOR-NIH [6,8], we reasoned that they would be good starting points for defining EEFx topology and parameters.

The *proteinEEFx* files were generated by making the following modifications to the amino acid parameters of PARALLHDG5.3: (i) the atom groupings were redefined to be those of CHARMM19 so as to be compatible with the solvation energy parameters in Table 1; (ii) for non-ionic residues, the partial atomic charges were replaced with those of CHARMM19 partial charges; (iii) for ionic residues (Arg, Lys, Asp, Glu and termini), the partial atomic charges were replaced with those of Lazaridis and Karplus [23], which were designed to obtain polar, albeit neutralized, residues that yield the proper stabilizing interactions for salt bridges. In addition, the force field also retains the full set of dihedral angle parameters defined in PARALLHDG5.3, thus enabling EEFx structure calculations to be performed with the XPLOR  $E_{\text{DIHE}}$  energy function instead of a statistical torsion angle potential, if desired.

Finally, to obtain effective EEFx calculations,  $E_{\text{slv}}$  is implemented together with the XPLOR van der Waals and electrostatics energy functions,  $E_{\text{VDW}}$  and  $E_{\text{ELEC}}$  [29], with their switching function effective between 7 Å and 9 Å; all nonbonded interactions beyond this value are neglected thus significantly reducing the computational cost. Dielectric screening, due to the influence of a protein's electrostatic properties on the density and distribution of the surrounding solvent molecules, is approximated by turning on the distance-dependent dielectric constant ( $\epsilon = r$ ) option of the  $E_{\text{ELEC}}$  function. This effect is neglected for nonpolar groups, where it is expected to be small.

### 3. Methods

#### 3.1. Structure calculations

All calculations were performed with XPLOR-NIH [27,28]. The new EEFx potential, *eefxPot*, is part of the XPLOR-NIH software suite (as of version 2.36), downloadable from the web (<http://nmr.cit.nih.gov/xplor-nih/>).

Free MD simulations were performed at 300 K, in Cartesian space, and implemented with four different models for nonbonded interactions (Table 2). The structures were downloaded from the protein data bank (PDB), energy minimized (500 steps of Powell minimization) and then subjected to 100 ps (and 1 ns in the case of SpAZ) of MD simulation, processed with normal atomic masses instead of the uniform mass setup that is routinely used in NMR structure calculation protocols. The trajectories were saved every 200 steps.

NMR-restrained structure calculations were performed using two conventional simulated annealing protocols [34]: the first for folding an initially extended conformation and the second for subsequent refinement of a folded model selected from the first folding protocol. Both protocols are based on the internal variable module [40] and share the same basic scheme comprising four stages: (i) torsion angle dynamics at high-temperature (3,500 K for folding, 3,000 K for refinement) for a time of 15 ps or 15,000

**Table 1**  
Parameters of  $E_{\text{slv}}$  used for EEFx calculations.<sup>a</sup>

Atom type	$V_i$	$\Delta G_i^{\text{ref}}$	$\Delta G_i^{\text{free}}$	$\Delta H_i^{\text{ref}}$	$\Delta C_{p,i}^{\text{ref}}$	$\lambda_i$	$R_i$
C	14.7	0.000	0.00	0.000	0.00	3.50	2.100
CR	8.3	−0.890	−1.40	2.220	6.90	3.50	2.100
CH1E	23.7	−0.187	−0.25	0.876	0.00	3.50	2.365
CH2E	22.4	0.372	0.52	−0.610	18.60	3.50	2.235
CH3E	30.0	1.089	1.50	−1.779	35.60	3.50	2.165
CR1E	18.4	0.057	0.08	−0.973	6.90	3.50	2.100
NH1	4.4	−5.950	−8.90	−9.059	−8.80	3.50	1.600
NR	4.4	−3.820	−4.00	−4.654	−8.80	3.50	1.600
NH2	11.2	−5.450	−7.80	−9.028	−7.00	3.50	1.600
NH3	11.2	−20.000	−20.00	−25.000	−18.00	6.00	1.600
NC2	11.2	−10.000	−10.00	−12.000	−7.00	6.00	1.600
N	0.0	−1.000	−1.55	−1.250	8.80	3.50	1.600
OH1	10.8	−5.920	−6.70	−9.264	−11.20	3.50	1.600
O	10.8	−5.330	−5.85	−5.787	−8.80	3.50	1.600
OC	10.8	−10.000	−10.00	−12.000	−9.40	6.00	1.600
S	14.7	−3.240	−4.10	−4.475	−39.90	3.50	1.890
SH1E	21.4	−2.050	−2.70	−4.475	−39.90	3.50	1.890
H	0.0	0.000	0.00	0.000	0.00	0.00	0.800

<sup>a</sup> Parameters are listed in *eefxPotTools*. Values of  $\Delta G_i^{\text{ref}}$ ,  $\Delta G_i^{\text{free}}$ ,  $\Delta H_i^{\text{ref}}$  and  $\Delta C_{p,i}^{\text{ref}}$  are from data of Privalov and Makhatadze [35–38]; they were determined experimentally at 298.15 K for small model molecules. Values of  $R_i$  correspond to CHARMM19 van der Waals radii. Values of  $V_i$  and  $\lambda_i$  were taken from the EEF1 model [23].

**Table 2**Energy functions, topology and parameters of the four different force fields used in the structure calculations.<sup>a</sup>

Model	$E_{\text{NONB}}^b$	Topology/parameters	Nonbonded parameters
REPEL	$E_{\text{VDW-REPEL}}$	protein.top/protein.par	$k_{\text{rep}} > 0$ , $C_{\text{rep}} > 0$
VDW	$E_{\text{VDW}}$	proteinEEF.top/proteinEEF.par	$k_{\text{rep}} = 0$ , group, vswitch, ctonnb = 7 Å, ctofnb = 9 Å
Vacuum	$E_{\text{VDW}} + E_{\text{ELEC}}$	proteinEEF.top/proteinEEF.par	$k_{\text{rep}} = 0$ , group, vswitch, ctonnb = 7 Å, ctofnb = 9 Å, switch, rdie
EEFx	$E_{\text{VDW}} + E_{\text{ELEC}} + E_{\text{solv}}$	proteinEEF.top/proteinEEF.par	$k_{\text{rep}} = 0$ , group, vswitch, ctonnb = 7 Å, ctofnb = 9 Å, switch, rdie, $r_{\text{on}} = 7$ Å, $r_{\text{off}} = 9$ Å

<sup>a</sup> All calculations were performed with  $nbxmod = 5$ , or  $nbxmod = 3$  for REPEL, to allow repulsions only between atoms separated by more than two covalent bonds. Calculations using the torsionDB potential were performed with  $nbxmod = 4$  to allow repulsions only between atoms separated by more than three covalent bonds.

<sup>b</sup>  $E_{\text{NONB}}$  is the XPLOR-NIH nonbonded energy function where  $E_{\text{VDW-REPEL}}$  is the simple repulsive form of the XPLOR van der Waals function,  $E_{\text{VDW}}$  is the switched Lennard-Jones form of the XPLOR van der Waals function and  $E_{\text{ELEC}}$  is the switched distance-dependent dielectric form of the XPLOR electrostatic function [29].  $E_{\text{solv}}$  and its switching function with parameters  $r_{\text{on}}$  and  $r_{\text{off}}$  are described in Eqs. (3) and (4).

timesteps; (ii) torsion angle dynamics with simulated annealing, where the temperature is reduced from the initial high temperature value to 25 K in steps of 12.5 K, for a time of 0.2 ps or 200 timesteps per temperature step (folding protocol), or a time of 0.63 ps or 630 timesteps per temperature step (refinement protocol); (iii) 500 steps of Powell torsion angle minimization; and (iv) 500 steps of Powell Cartesian minimization.

In the high temperature stage, experimental dihedral angle restraints and distance restraints were applied with respective force constants of  $k_{\text{CDIH}} = 10 \text{ kcal mol}^{-1} \text{ rad}^{-2}$  and  $k_{\text{DIST}} = 2 \text{ kcal mol}^{-1} \text{ Å}^{-2}$ . In the simulated annealing stage,  $k_{\text{CDIH}}$  was set to  $200 \text{ kcal mol}^{-1} \text{ rad}^{-2}$  and  $k_{\text{DIST}}$  was increased geometrically from 2 to  $30 \text{ kcal mol}^{-1} \text{ Å}^{-2}$ . In selected calculations, the torsionDB statistical torsion angle potential [34] was included with a force constant set to  $k_{\text{tDB}} = 0.02 \text{ kcal mol}^{-1} \text{ rad}^{-2}$  in the high temperature stage and ramped geometrically from 0.02 to  $2 \text{ kcal mol}^{-1} \text{ rad}^{-2}$  during simulated annealing. Atomic overlap was prevented by limiting allowed repulsions to those between atoms separated by three or more covalent bonds ( $nbxmod = 5$ ), except for calculations performed with torsionDB where allowed repulsions were limited to those between atoms separated by four or more covalent bonds ( $nbxmod = 4$ ).

Calculations were performed with either one of two different models for nonbonded interactions: REPEL or EEFx (Table 2). In the REPEL calculations, the simple repulsive van der Waals function [29] was used in conjunction with the default XPLOR-NIH protein topology and parameters (protein.top/par version 1.0). In the high temperature stage, only CA–CA atomic interactions were active, the van der Waals force constant was set to  $C_{\text{rep}} = 0.004 \text{ kcal mol}^{-1} \text{ Å}^{-4}$  and the van der Waals radius scale factor was set to  $k_{\text{rep}} = 1.2$ . In the simulated annealing stage, all atom–atom interactions were active,  $C_{\text{rep}}$  was ramped from 0.004 to  $4 \text{ kcal mol}^{-1} \text{ Å}^{-4}$ , and  $k_{\text{rep}}$  was ramped down from 0.9 to 0.8.

The EEFx calculations were performed with the proteinEEF topology and parameters. The van der Waals and electrostatic energy terms were implemented with a distance cutoff of 9 Å, a switching function for the Lennard-Jones potential between 7 and 9 Å, and distance-dependent dielectric. The new solvation energy term,  $E_{\text{solv}}$ , was implemented with a distance cutoff of 9 Å and the switching function for the potential between 7 and 9 Å. The force constants for van der Waals ( $k_{\text{VDW}}$ ), electrostatic ( $k_{\text{ELEC}}$ ) and solvation ( $k_{\text{solv}}$ ) energy terms were each set to 1.

During the folding protocol of the EEFx calculations, the 15 ps high temperature stage was further divided into two equal parts, the first performed as described for the REPEL calculations and the second performed with  $k_{\text{VDW}}$ ,  $k_{\text{ELEC}}$  and  $k_{\text{solv}}$  set to 0.1. This was done to prevent fatal atomic overlap in the early stages of calculations from extended templates. In the subsequent simulated annealing stage, the force constants were ramped geometrically from 0.1 to  $1 \text{ kcal mol}^{-1}$ . The values of  $k_{\text{VDW}}$ ,  $k_{\text{ELEC}}$  and  $k_{\text{solv}}$  were set to 1 throughout the refinement protocol of all EEFx calculations. For calculations of the largest protein EIN, the high

temperature stage of the folding protocol was performed with REPEL and without EEFx, while EEFx was used during the annealing stage. Explicit water refinement was implemented as described previously [7–10], using the *wrefine.py* script available in XPLOR-NIH.

### 3.2. Generation of partial NOE restraints

The complete data set of long-range distances (defined here as the 100% data set) included only distance restraints between atoms more than 4 residues apart in the protein sequence. Partial distance restraint data sets were generated by randomly selecting restraints from the 100% data set, to cover the percentage range from 1% to 100%. Five independent restraint sets of equal size were generated for each percentage value. Each set was used to fold 100 structures from extended templates, with the folding protocol described above, and the structure with lowest total energy was taken as input for the refinement protocol. Statistics were generated for the 10 structures with lowest energy from a total of 100 refined structures from each independent restraint set.

### 3.3. Structure analysis, validation and display

XPLOR-NIH was used to evaluate the precision and accuracy of the calculated structures, as well as to fit the experimental residual dipolar coupling (RDC) data to the calculated structures by singular value decomposition [41] and report the RMSD (root mean square deviation) measure of fit [42]. The backbone conformations, side-chain conformations and nonbonded atomic interactions of the calculated structures were assessed using the programs WHAT IF [43,44] and MolProbity [45–47]. Hydrogen bonds were computed in PyMol [48] using a script [49] with distance cutoff of 3.2 Å and angle cutoff of 55°. Structures were rendered with PyMol.

### 3.4. Generation of structural decoys

Blind structure predictions were performed using Rosetta [68], starting from the sequences of the proteins GB1 and BAF, excluding the PDB structural coordinates of the two proteins as well as their homologues from the structure prediction database. For each protein, 5,000 coarse-grained structural models were generated and then refined by all-atom relaxation, performed with the implicit aqueous environment protocol of Rosetta3.4. The refined all-atom structures were clustered according to their overall energy and their backbone CA atom RMSD to the lowest energy structure with a cutoff of 5 Å. For each protein, the most populated cluster encompassed more than 20% of the entire sampling space and contained 3,554 decoys for GB1 and 1,225 decoys for BAF. These decoys were subjected to two sets of Powell energy minimization (500 steps) in XPLOR-NIH, and then scored by calculating the total XPLOR-NIH energy with EEFx.



## 4. Results

### 4.1. Unrestrained molecular dynamics simulations

We first tested EEFx for its ability to model a physically realistic environment that can sustain stable, native protein structures. We performed unrestrained MD simulations of ten proteins at room temperature and examined the deviations of the resulting coordinates from those of the experimentally determined structures, taken to be representative of the native state. For each protein, we compared the results of 100 ps MD simulations performed with four different force fields (REPEL, VDW, vacuum, EEFx), each defined by a specific set of nonbonded energy function, topology and parameters for the system (Table 2). The proteins selected for analysis (Table 3) have sizes ranging from about 60 to 260 amino acids and a variety of structures, all determined by NMR spectroscopy, with coordinates and experimental restraints publicly available in the PDB. Three of these proteins (GB1, Eglin-c and Ubiquitin) were part of the set used for the initial development of CHARMM EEF1 [23] and, hence, provide useful benchmarks for direct comparison of EEFx performance.

A longer, 1 ns, MD simulation, performed for one of the proteins (SpAZ) demonstrates that simulations with EEFx are highly stable, as are those with VDW and vacuum (Fig. 1A). This behavior is similar to previous observations for CHARMM simulations with EEF1 [23]. Most of the changes in protein conformation occur within the first 30–40 ps of dynamics, indicating that 100 ps simulations are sufficiently long to compare the effects of the different force fields on protein structure. By contrast, as expected, the simple repulsive potential alone (REPEL) is unable to maintain the native structure in the absence of additional physical forces or experimental restraints, and the fold quickly and continuously unravels reaching RMSD values near 20 Å at 1 ns (Fig. 1A, black).

The data show that EEFx effectively maintains the native protein structure for the entire duration of dynamics (Fig. 1A, red). Furthermore, although simulations in vacuum (Fig. 1A, blue) or with VDW alone (Fig. 1A, gray) are also stable, EEFx yields a structure that is substantially closer to native (RMSD ~ 1 Å; Fig. 1B), while the structures from vacuum (RMSD > 2 Å; Fig. 1C) or VDW (RMSD > 3 Å; Fig. 1D) simulations both differ significantly from the native state. Notably, both the vacuum and EEFx simulations were performed with distance-dependent dielectric screening ( $\epsilon = r$ ), as described above. By contrast comparisons in the original development of EEF1 were made to the vacuum force field with fixed unity dielectric constant ( $\epsilon = 1$ ) because this is the accepted standard in the field and because it corresponds to an extreme of

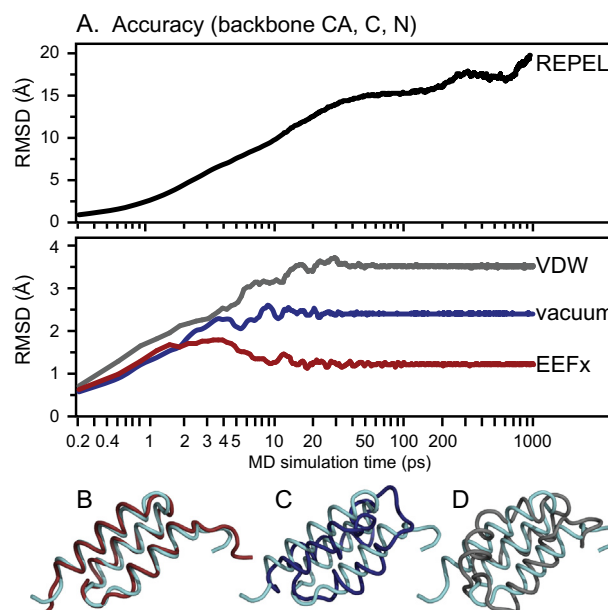
**Table 3**  
Proteins used for test structure calculations with EEFx.

Protein <sup>a</sup>	PDB	Length <sup>b</sup>	Residues <sup>c</sup>	Fold
GB1 [50,51]	3GB1	56	all	$\alpha\beta$
SpAZ [52]	1Q2N	58	6–55	$\alpha$
DDEF1-SH3 [53]	2RQT	61	all	$\beta$
Eglin-c [54]	1EGL	70	8–70	$\alpha\beta$
Ubi [55]	1D3Z	76	all	$\alpha\beta$
Din-I [56]	1GHH	81	1–71	$\alpha\beta$
BAF [57]	2EZX	89	all	$\alpha$
RNPK [58]	1KHM	89	12–84	$\alpha\beta$
IIBMt [59]	1VKR	125	12–107	$\alpha\beta$
ArfA-b [60]	2KSM	131	80–195	$\alpha\beta$
EIN [61,62]	1EZA	259	1–230	$\alpha\beta$

<sup>a</sup> GB1, protein G B1 domain; SpAZ, Staphylococcal protein A Z domain; Ubi, ubiquitin; DDEF1-SH3, human DDEF1 SH3 domain; Din, DNA-damage-inducible protein I; BAF, human barrier to autointegration factor; RNPK, nuclear ribonucleoprotein K KH domain; IIBMt, mannitol transporter enzyme II B domain; ArfA-B, *M. tuberculosis* ArfA B domain; EIN, enzyme I N-terminal domain.

<sup>b</sup> Full length of polypeptide.

<sup>c</sup> Residues used in calculations.



**Fig. 1.** Time dependence of unrestrained 1 ns MD simulations of SpAZ performed with EEFx (red), vacuum (blue), VDW (gray) or REPEL (black) force fields. (A) Structural accuracy reported as RMSD of backbone atoms (CA, C, N) to the experimental structure (PDB: 1Q2N). (B–D) Structures obtained from simulations with EEFx (B), vacuum (C) and VDW (D) superimposed on the experimentally-determined structure (cyan). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

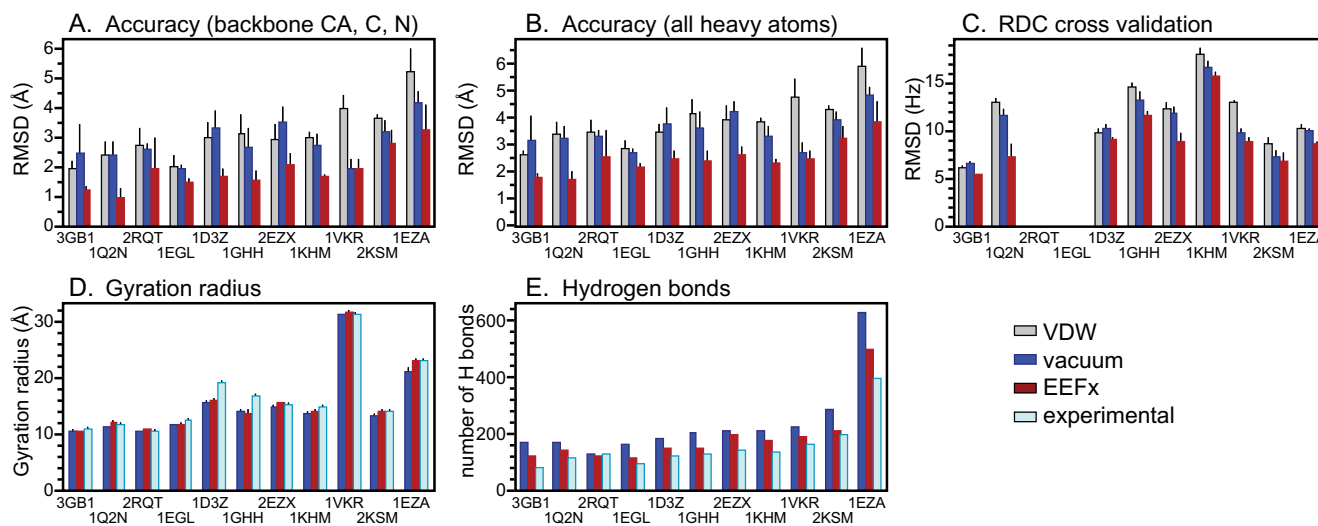
zero dielectric screening. Thus, even though the use of  $\epsilon = r$  and ionic sidechain neutralization (implemented in proteinEEFx.top) are expected to improve the results of vacuum simulations, EEFx still gives a significantly better result compared to vacuum.

The improved performance of EEFx compared to both vacuum and VDW, is also observed in the 100 ps simulations of the other test proteins (Fig. 2). Notable improvements in structural accuracy (Fig. 2A and B) are observed for all cases. Furthermore, for all proteins, simulations with EEFx yield significantly better cross validation with the experimental RDC data (Fig. 2C), providing independent evidence that EEFx produces close representations of the native structures. EEFx also relieves the molecular contraction that is observed in the vacuum simulations, as evidenced by the slightly higher gyration radii of the EEFx structures compared to vacuum (Fig. 2D). Contraction is a well-known effect of MD simulations performed in vacuum, where electrostatic interactions are amplified by the lack of solvent screening [63,64] and is readily visible in the 1 ns vacuum simulation of SpAZ (Fig. 1C).

Finally, simulations with EEFx yield a significant increase in the number of hydrogen bonds compared to the experimentally determined structures, further demonstrating that the model maintains stable secondary structural elements (Fig. 2E). For all test cases, the number of hydrogen bonds observed with EEFx is higher than the number observed for the experimental structures but lower than the number observed for vacuum simulations where the effects of electrostatics are not dampened by solvation screening, while VDW alone (not shown) significantly decreases the number of hydrogen bonds concomitant with structural distortion. We conclude that EEFx provides a physically realistic implicit solvent environment capable of supporting stable, native protein structures.

### 4.2. Recognition of native fold

We next tested EEFx for its ability to discriminate among native and unfolded protein states. Protein structure prediction has become a major tool in structural biology that can be used very



**Fig. 2.** Free MD simulations of native protein structures. (A and B) Accuracy to native structure reported as backbone atom (CA, C, N) and heavy atom RMSD. (C) Cross validation to experimental RDC data reported as RMSD. (D) Gyration radii of the proteins. (E) Number of hydrogen bonds observed for each simulation. Data are shown for simulations with EEFx (red), vacuum (blue) and VDW (gray) force fields, or for the experimentally-determined PDB structures (cyan). PDB codes correspond to protein names in Table 3. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

effectively to supplement sparse experimental restraints during protein structure determination by NMR [65–67]. Starting with the amino acid sequences of the proteins GB1 and BAF, we performed blind structure predictions using the Rosetta program, which is very successful at predicting three-dimensional structures of proteins from their amino acid sequences [68]. For each protein, we generated 5,000 coarse-grained decoys and then refined them by full-atom relaxation in Rosetta. The most populated clusters were subjected to energy minimization in XPLOR-NIH and then scored with the EEFx energy function. Minimization produces only very minor alterations (maximum 0.2 Å) of the decoys's original structures.

Rosetta represents proteins by their backbone heavy atoms plus CB atoms for the sidechains; its full-atom energy function is a hybrid of statistical, empirical and physically realistic terms, including: PDB-derived sidechain and backbone torsion angle potentials, orientation-dependent hydrogen bonds, short-range knowledge-based electrostatic energy, reference energies for the unfolded states of the twenty amino acids, solvation effects based on the EEF1 solvation free energy function, and Lennard–Jones nonbonded interactions. By contrast, EEFx does not include any statistical terms.

Analysis of the Rosetta and EEFx energy landscapes (Fig. 3A, B, G, H) and comparison of the lowest energy structures with the NMR structures of either GB1 (Fig. 3C–F) or BAF (Fig. 3I–L), show that both the Rosetta and EEFx energy functions effectively recognize the overall, native fold of the two proteins. It is remarkable that these decoys were generated *de novo* with no other input than the protein's amino acid sequence and the PDB, which now contains a sufficient number of structures to enable fragment-based structure predictions. For each protein, the ten lowest energy decoys selected by either Rosetta or EEFx have very similar precision (Fig. 3E, F, K, L), indicating the tendency of each energy function towards a specific structural ensemble. However, Rosetta and EEFx each select a different decoy based on lowest energy.

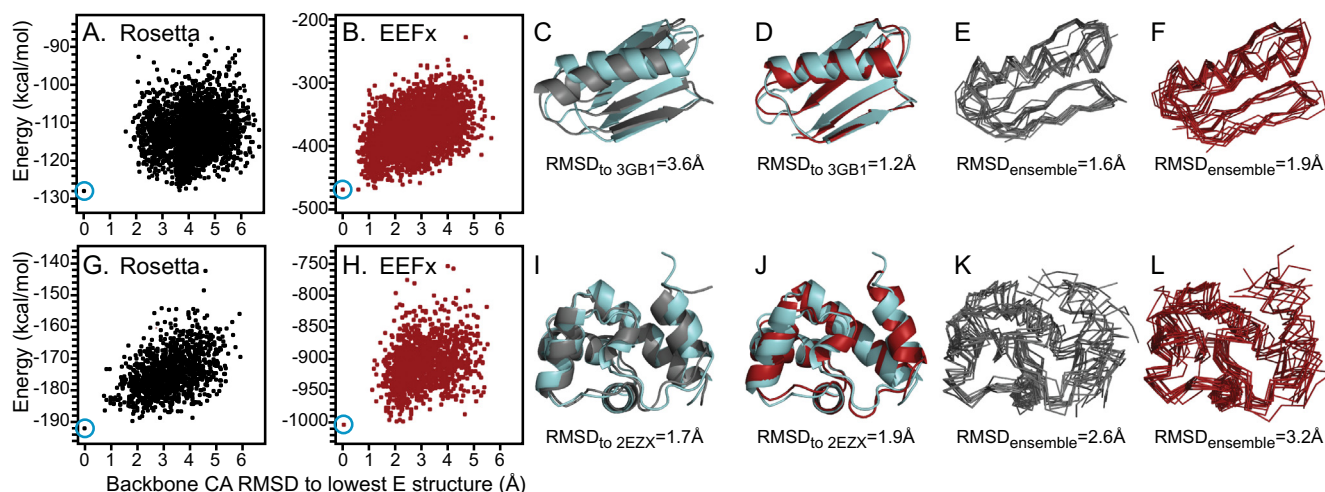
For GB1, the decoy with lowest Rosetta energy exhibits the overall features of the native fold, but is 3.6 Å RMSD away from the experimental structure (Fig. 3C). By contrast the decoy selected for lowest EEFx energy is very close to the experimental structure, with an RMSD of 1.2 Å (Fig. 3D). This difference in accuracy is also reflected in the shapes of the Rosetta and EEFx energy landscapes

of GB1: while the EEFx energy landscape has a marked funneling shape towards the native structure (Fig. 3B), the Rosetta landscape has much less pronounced funneling features (Fig. 3A). In the case of BAF, the decoys with lowest Rosetta and EEFx energies also correspond to the decoys with lowest RMSD relative to the experimental structure (Fig. 3I–L) and both have Rosetta and EEFx energy landscapes with marked funneling shape towards the native fold (Fig. 3G, H).

For both proteins, the EEFx energy landscapes have significantly (sixfold) greater energy dispersion. The EEFx energy bandwidth is 300 kcal mol<sup>−1</sup>, compared to the 50–60 kcal mol<sup>−1</sup> observed for Rosetta and, thus, provides a greater degree of discrimination among protein folds. Analysis of the decoy EEFx energies shows that electrostatic energy makes the most significant contribution to the overall value. On average, the ten decoys of GB1 with lowest Rosetta energy have 83 hydrogen bonds, while those with lowest EEFx energy have 95 hydrogen bonds. Similarly for BAF, the ten decoys with lowest Rosetta energy have, on average, 142 hydrogen bonds, while those selected by EEFx have 163. In the lowest EEFx energy decoys, additional hydrogen bonds are formed both among backbone and sidechain atoms and contribute to lowering the total energy of the system. We conclude that the new EEFx term gives results comparable to Rosetta over a wide range of conformations and may provide a wider dynamic range for discrimination of folded states.

#### 4.3. NMR-restrained protein structure calculations

EEFx was developed with the specific objective of providing a more physically realistic energy landscape for NMR-restrained structure calculations, without significantly sacrificing calculation speed and ease of implementation. To test its performance in this regard, we performed NMR-restrained calculations for six proteins in Table 3, using the experimental distance and dihedral angle restraints available in the PDB and retaining the RDC restraints only for cross validation. The calculations were started from extended templates, as is typically done in NMR structure determination, and performed with standard simulated annealing protocols, executed using either the simple repulsive function of the van der Waals energy term (REPEL) with the default XPLOR-NIH protein topology and parameters, or the EEFx energy function with



**Fig. 3.** Recognition of native protein fold. (A, B, G, H) Rosetta (black) and EEFx (red) energy landscapes of GB1 (A and B) and BAF (G and H). RMSD values are computed for CA atoms relative to the decoy with lowest energy (cyan circles). (C, D, I, J) Cartoon representations of the decoys of GB1 (C and D) and BAF (I and J) with lowest Rosetta energy (gray) or lowest EEFx energy (red) and superimposed experimental PDB structures (cyan). RMSDs represent structural accuracy, computed for CA atoms relative to the experimental PDB structures. (E, F, K, L) Ribbon representations of the 10 decoys of GB1 (E and F) and BAF (K and L) with lowest Rosetta energy (gray) or lowest EEFx energy (red). RMSDs represent precision of the structural ensembles, evaluated as average pairwise values for CA atoms. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

proteinEEFx topology and parameters, each with or without inclusion of the statistical torsion angle potential torsionDB [34].

The restraints used for structure calculations necessarily reflect heterogeneity both in the way they were measured and evaluated from the experimental data and also in their number relative to protein length. For example, the interpretation of NOE signals in terms of inter-atomic distances can vary substantially among research groups, and the number of long-range NOE restraints for each protein in Table 3 varies between 1.8 and 10.6 per residue. This situation reflects the typical range of variables associated with NMR structure calculations and hence provides a good test case for evaluating the performance of EEFx.

We first examined the ability of EEFx to produce folded protein structures with limited numbers of restraints. These tests were performed for two proteins, GB1 (Fig. 4) and ArfAB (Fig. S1), whose different sizes and distinct topologies make them excellent candidates for assessing the performance of EEFx. The NMR structure of GB1 is based on a large set of experimental restraints, including a complete set of NOEs, and is exceptionally well defined [50]. ArfAB is a larger polypeptide (131 residues) with an unusual fold whose structure was determined to very good precision [60].

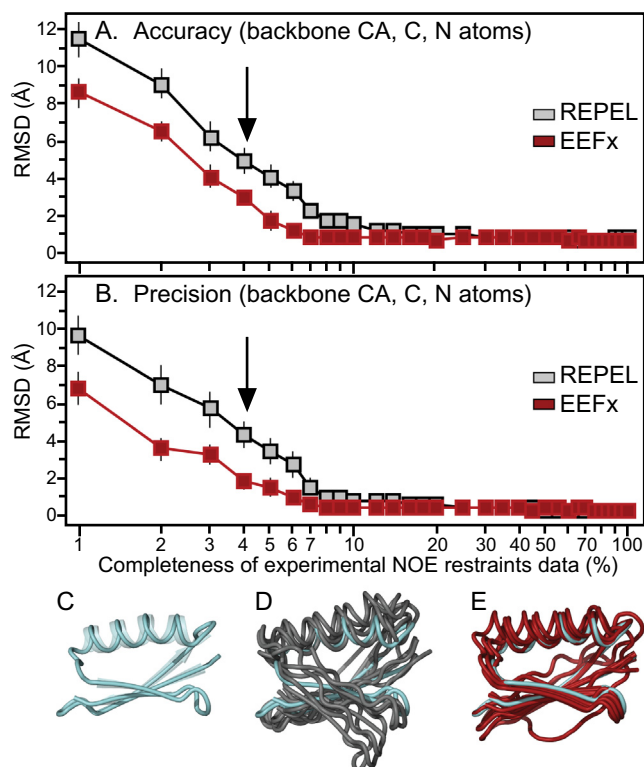
To examine the dependence of the calculations on the number of distance restraints, we used data sets with decreasing numbers of distances. Each set was generated from experimental hydrogen bonds and NOEs by first, removing all distances between sites separated by less than five residues in the protein sequence, and then, randomly eliminating long-range distances from the remaining restraints. Thus, each resulting data point reflects the average over five independent structure calculations performed with five unique sets of long-range distances of equal size. For each set, 100 structures were calculated and statistics were generated for the 10 structures with lowest total energy. This approach reduces the bias associated with the inherent information content of the distance restraints, a factor that also influences structural quality [69].

Both EEFx and REPEL are capable of determining the correct global folds of GB1 and ArfAB with as few as 0.2–0.4 long-range distances per residue. However, in this case of very limited restraints EEFx produces structures that are significantly closer to the native fold and more precise than those calculated with REPEL. For GB1 (Fig. 4), calculations performed with REPEL and a partial

data set containing only 4% of the long-range distance restraints ( $\sim 0.2$  restraints per residue) produce structures with an accuracy of 4.9 Å and a backbone precision of 4.4 Å. By contrast, structures calculated with EEFx, and the same partial data sets, have an accuracy of 2.9 Å and a precision of 1.8 Å. Similarly for ArfA-B (Fig. S1), calculations with EEFx using only 20% of the long-range distance data ( $\sim 0.5$  restraints per residue) yield structures with better accuracy (2.9 Å) and precision (2.3 Å) while structures calculated with REPEL have both lower accuracy (3.6 Å) and precision (3.0 Å). When all available long-range distances are used (5.4 per residue in GB1, and 2.5 per residue in ArfAB), structures calculated with EEFx and REPEL have similar accuracy but the EEFx structures have somewhat greater precision (Fig. 5). Similar trends are observed for the precision and accuracy determined for all heavy atoms.

To the extent that structural accuracy can be assessed relative to the actual native structure, the precision of NMR structures is typically higher than their accuracy [70]. The number of restraints available for calculation is the principal factor influencing the accuracy and precision of NMR structures, but the nature of the nonbonded energy function also plays an important role [70]. The principal effect of EEFx is to direct the calculation towards the native structure even in the absence of large numbers of restraints. The ability of EEFx to fold structures with limited numbers of distance restraints correlates with its ability to bury solvent accessible protein groups, form hydrogen bonds and optimize the radius of gyration. This is a significant advantage of EEFx, since modern methods for NMR structure determination are increasingly designed to shift the burden away from time-consuming measurements of multiple long-range distances and facilitate the determination of high-quality three-dimensional structures with very few or no distance restraints [65,66].

We next tested the performance of EEFx for the generation of high quality structures. Introduction of a new term in the target energy function can induce deterioration in the agreement between calculated structures and the other experimental and conformational energy terms. The data in Figs. 5 and 6, obtained for six proteins with increasing sizes and an assortment of structures, demonstrate that the improvements in precision and accuracy afforded by EEFx are not accompanied by significant costs to either conformational terms or terms associated with the NMR data – on the contrary.



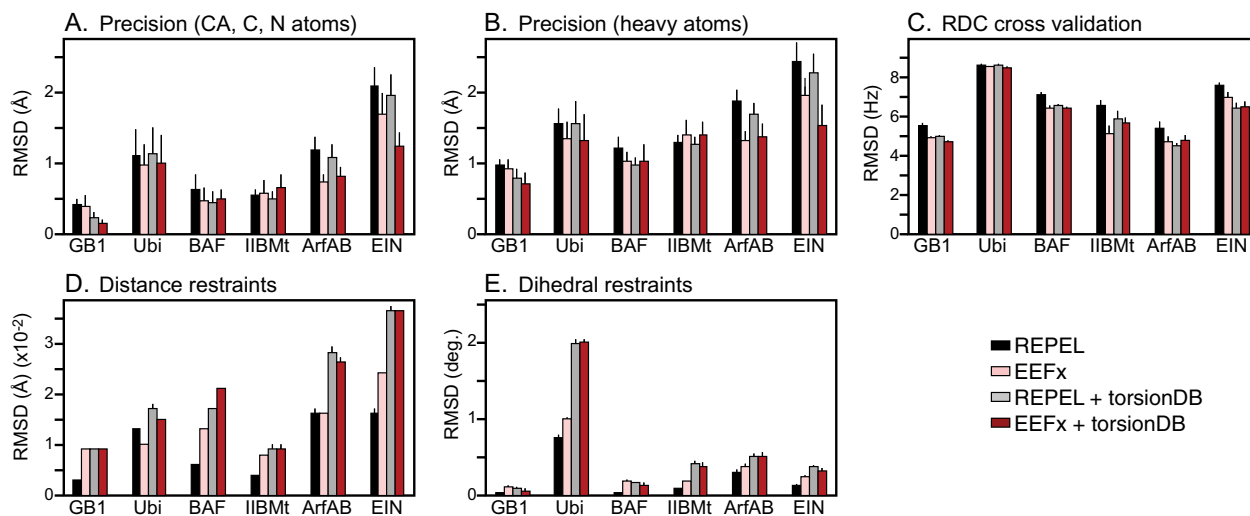
**Fig. 4.** Effect of EEfx on NMR-restrained structure calculations of GB1 with limited distance restraints. (A and B) Effect of the number of long-range (>4 residues apart) distance restraints on structural accuracy and precision. The total number of restraints was reduced by randomly eliminating distances from the full data set. Accuracy was evaluated as pairwise RMSD of backbone CA, C, N atoms relative to the experimental structure. Precision was evaluated as average pairwise RMSD of backbone CA, C, N atoms. (C–E) Cartoon representations of the native structure (cyan) and the ensembles of five lowest energy structures calculated with 4% of the NOE data, using REPEL (gray) or EEfx (red). Arrows indicate the 4% data points taken for structure illustration in panels C–E. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

In all cases, the precision of both backbone and heavy atom coordinates improves significantly for structures calculated with EEfx, with the sole exception of IIBMt, where the precision decreases slightly (Fig. 5A and B). Furthermore, calculations with EEfx produce similar or improved agreement with the experimental RDC data (Fig. 5C), which were purposely excluded from structure calculations. RDCs depend on the orientation of interatomic vectors relative to the external magnetic field, and their exclusion from structure calculation provides a useful independent test of structural accuracy [42]. All of the structures calculated and refined with EEfx have better or similar agreement with the RDCs reflecting an improvement in accuracy.

Finally, calculations with EEfx produce similar levels of agreement between the structures and experimental distance and dihedral angle restraints used in the calculations (Fig. 5D and E). Although, in some cases a slight deterioration is observed when EEfx is used, the combined use of EEfx with the statistical potential torsionDB [34] produces results with similar or better agreement than those obtained with torsionDB alone. In the case of ubiquitin, EEfx actually produces a slight improvement in the agreement between structure and distance restraints.

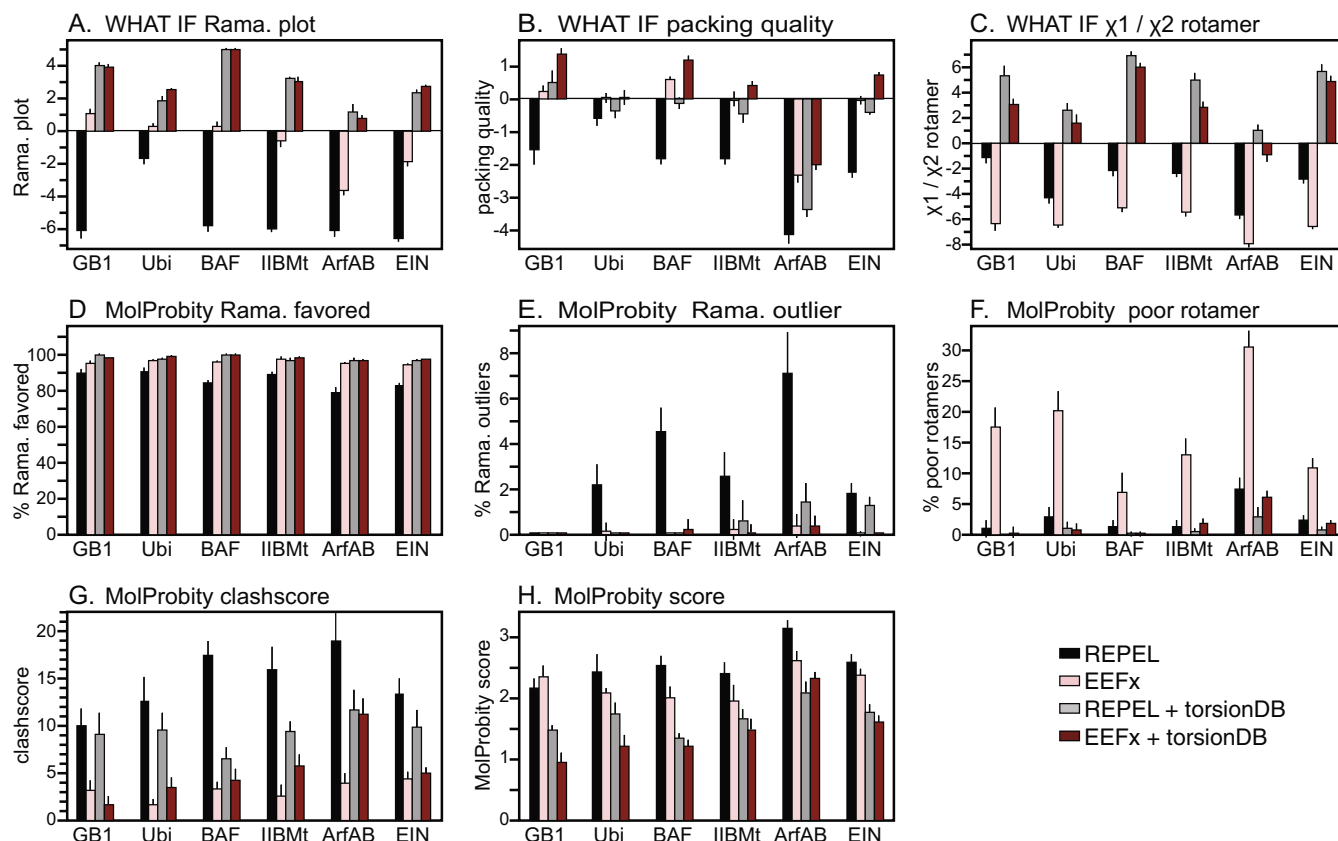
The structures calculated with EEfx also compare very favorably with those refined in explicit water, using the *wrefine.py* refinement protocol adapted from Refs. [7–10] and available in XPLOR-NIH. Correlations to the experimental data are similar for both EEfx and water-refined structures, while structural precision is somewhat better for EEfx (Fig. S2).

We further examined the quality of structures generated with EEfx with respect to WHAT IF [43] and MolProbity [46,47] validation metrics (Fig. 6). The results show that EEfx improves the quality of the backbone conformation in every case compared to results obtained with REPEL, regardless of whether torsionDB is included or not. Use of EEfx improves the WHAT IF Ramachandran plot appearance (Fig. 6A). Similarly, MolProbity indicates that EEfx causes the favored regions of the Ramachandran plot to become more populated (Fig. 6D) and the percent of Ramachandran outliers to drop significantly (Fig. 6E). With regards to sidechain conformation, both WHAT IF and MolProbity show that EEfx alone



**Fig. 5.** Structural statistics of NMR-restrained calculations performed with EEfx. (A and B) Structural precision evaluated as average pairwise RMSD of (A) backbone CA, C, N atoms and (B) all heavy atoms. (C) Agreement between structures and experimental RDC restraints excluded from structure calculations. (D, E) Agreement between structures and experimental distance and dihedral angle restraints used in the structure calculations. For each protein, the errors represent the mean  $\pm$  standard deviation evaluated for ensembles of 10 lowest energy structures. Bars represent data obtained in four ways: the standard simple repulsive XPLOR potential REPEL (black); EEfx (pink); REPEL plus torsionDB (gray); and EEfx plus torsionDB (red). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)





**Fig. 6.** Structural validation analyses of NMR-restrained calculations performed with EEfx. (A–C) WHAT IF validation statistics for (A) Ramachandran plot appearance; (B) protein packing quality; and (C)  $\chi_1/\chi_2$  torsion angles. (D–H) MolProbity validation statistics for (D) percent of residues in favored regions of the Ramachandran plot; (E) percent of residues in unfavored regions of the Ramachandran plot; (F) percent of residues with poor sidechain torsion angles; (G) clashscore; and (H) overall MolProbity score. For each protein, the errors represent the mean  $\pm$  standard deviation evaluated for ensembles of 10 lowest energy structures. Bars represent data obtained in four ways: the standard simple repulsive XPLOR potential REPEL (black); EEfx (pink); REPEL plus torsionDB (gray); and EEfx plus torsionDB (red). The MolProbity clashscore and MolProbity score are costs: the lower the better. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

results in worse  $\chi_1/\chi_2$  rotamer normality scores (Fig. 6C) and higher numbers of poor rotamers (Fig. 6F) for all proteins. This is expected for calculations performed without any dihedral angle potential term and is also observed for calculations performed with REPEL alone. However, these effects are readily corrected by the use of torsionDB [34], which was developed precisely for this purpose, or by inclusion of the XPLOR dihedral angle conformation energy term ( $E_{DIHE}$ ) in the calculations (Fig. S3), which is enabled by the more complete force field available in the proteinEEfx.top/par files that work with EEfx.

The validation results further show that EEfx improves the quality of protein conformation and nonbonded atomic interactions. The WHAT IF packing quality score (the atomic distributions around different molecular fragments) [71] and the MolProbity clashscore (the number of serious atomic overlaps per thousand atoms) [72] provide estimates of the quality of nonbonded atomic interactions or atomic packing. Notably, all structures generated with EEfx display marked improvements in both of these key metrics (Fig. 6B and G), even when compared with water-refined structures (Fig. S3). This is reflected in the overall MolProbity score [46,47] (the lower the better), which improves with EEfx for every protein tested (Fig. 6H). Generally, NMR structures tend to be somewhat less well packed and expanded relative to X-ray structures [73,74] and, often, the experimental NMR data are more consistent with high-resolution crystal structures than the corresponding NMR structures [75]. Here, the improved packing obtained with EEfx is also reflected in the improved agreement with the experimental RDC data.

Overall the best results are obtained when EEfx is used in conjunction with TorsionDB. However, the use of EEfx with the  $E_{DIHE}$  energy term also yields very favorable results (Fig. S3), thus providing a non-statistical, albeit empirical, alternative to the use of a statistical knowledge-based potential (torsionDB) for dihedral angles.

Finally, we report that calculations with EEfx are computationally efficient. For the proteins tested in this study, NMR-restrained calculations performed with EEfx were only 2.5 times longer in elapsed wall clock time than those with REPEL.

## 5. Conclusions

The benefits of protein structure refinement in water are well documented [6–11,15]. However, performing structure calculations with explicit atomic representation of the solvent molecules is computationally expensive and impractical for NMR-restrained structure determination. We conclude that the new EEfx potential described in this paper provides an effective energy function for the implicit solvation of proteins during NMR-restrained calculations.

The initial results show that EEfx outperforms the simple repulsive potential that is typically used in NMR structure calculations. The EEfx energy function effectively discriminates native from misfolded conformations and yields significant improvements in structural precision and accuracy, as well as conformational and nonbonded protein packing properties. Notably, EEfx can be used both to fold as well as refine NMR-restrained structures and improves the precision and accuracy of structure calculations

performed with limited numbers of experimental distance restraints. Finally, implementation of EEFx in XPLOR-NIH is straightforward and computationally efficient enabling structure calculations to be easily carried out on standard laboratory computers.

Additional studies on different proteins will be needed to fully explore the XPLOR-NIH EEFx energy landscape. However these initial results indicate that EEFx is a useful step forward towards the practical calculation of experimental protein structures in a physically realistic environment that closely resembles their native state.

## Acknowledgments

This research was supported by Grants from the National Institutes of Health (R01 GM100265; P01 AI074805 and R21 GM094727). It utilized the Resource for Molecular Imaging of Proteins at UCSD, supported by NIH Grant P41 EB002031). CDS was supported by funds from the NIH Intramural Research Program of The Center for Information Technology.

## Appendix A. Supplementary material

Supplementary data associated with this article can be found, in the online version, at <http://dx.doi.org/10.1016/j.jmr.2014.03.011>.

## References

- [1] C.B. Anfinsen, Principles that govern the folding of protein chains, *Science* 181 (1973) 223–230.
- [2] L. Banci, I. Bertini, C. Luchinat, M. Mori, NMR in structural proteomics and beyond, *Prog. Nucl. Magn. Reson. Spectrosc.* 56 (2010) 247–266.
- [3] H.X. Zhou, T.A. Cross, Influences of membrane mimetic environments on membrane protein structures, *Annu. Rev. Biophys.* 42 (2013) 361–392.
- [4] M. Nilges, A.M. Gronenborn, A.T. Brunger, G.M. Clore, Determination of three-dimensional structures of proteins by simulated annealing with interproton distance restraints. Application to crambin, potato carboxypeptidase inhibitor and barley serine proteinase inhibitor 2, *Protein Eng.* 2 (1988) 27–38.
- [5] G.M. Clore, A.M. Gronenborn, Determination of three-dimensional structures of proteins and nucleic acids in solution by nuclear magnetic resonance spectroscopy, *Crit. Rev. Biochem. Mol. Biol.* 24 (1989) 479–564.
- [6] J.P. Linge, M. Nilges, Influence of non-bonded parameters on the quality of NMR structures: a new force field for NMR structure calculation, *J. Biomol. NMR* 13 (1999) 51–59.
- [7] C.A. Spronk, J.P. Linge, C.W. Hilbers, G.W. Vuister, Improving the quality of protein structures derived by NMR spectroscopy, *J. Biomol. NMR* 22 (2002) 281–289.
- [8] J.P. Linge, M.A. Williams, C.A. Spronk, A.M. Bonvin, M. Nilges, Refinement of protein structures in explicit solvent, *Proteins* 50 (2003) 496–506.
- [9] S.B. Nabuurs, A.J. Nederveen, W. Vranken, J.F. Doreleijers, A.M. Bonvin, G.W. Vuister, G. Vriend, C.A. Spronk, DRESS: a database of Refined solution NMR structures, *Proteins* 55 (2004) 483–486.
- [10] A.J. Nederveen, J.F. Doreleijers, W. Vranken, Z. Miller, C.A. Spronk, S.B. Nabuurs, P. Guntert, M. Livny, J.L. Markley, M. Nilges, E.L. Ulrich, R. Kaptein, A.M. Bonvin, RECOORD: a recalculated coordinate database of 500+ proteins from the PDB using restraints from the BioMagResBank, *Proteins* 59 (2005) 662–672.
- [11] I. Bertini, D.A. Case, L. Ferella, A. Giachetti, A. Rosato, A Grid-enabled web portal for NMR structure refinement with AMBER, *Bioinformatics* 27 (2011) 2384–2390.
- [12] M. Sharma, M. Yi, H. Dong, H. Qin, E. Peterson, D.D. Busath, H.X. Zhou, T.A. Cross, Insight into the mechanism of the influenza A proton channel from a structure in a lipid bilayer, *Science* 330 (2010) 509–512.
- [13] X. Cheng, W. Im, NMR observable-based structure refinement of DAP12-NKG2C activating immunoreceptor complex in explicit membranes, *Biophys. J.* 102 (2012) L27–L29.
- [14] X. Cheng, S. Jo, F.M. Marassi, W. Im, NMR-based simulation studies of Pf1 coat protein in explicit membranes, *Biophys. J.* 105 (2013) 691–698.
- [15] B. Xia, V. Tsui, D.A. Case, H.J. Dyson, P.E. Wright, Comparison of protein solution structures refined by molecular dynamics simulation in vacuum, with a generalized Born model, and with explicit water, *J. Biomol. NMR* 22 (2002) 317–331.
- [16] J. Chen, W. Im, C.L. Brooks 3rd, Refinement of NMR structures using implicit solvent and advanced sampling techniques, *J. Am. Chem. Soc.* 126 (2004) 16038–16047.
- [17] J. Chen, H.S. Won, W. Im, H.J. Dyson, C.L. Brooks 3rd, Generation of native-like protein structures from limited NMR data, modern force fields and advanced conformational sampling, *J. Biomol. NMR* 31 (2005) 59–64.
- [18] B. Roux, T. Simonson, Implicit solvent models, *Biophys. Chem.* 78 (1999) 1–20.
- [19] M. Feig, C.L. Brooks 3rd, Recent advances in the development and application of implicit solvent models in biomolecule simulations, *Curr. Opin. Struct. Biol.* 14 (2004) 217–224.
- [20] N.A. Baker, Improving implicit solvent simulations: a Poisson-centric view, *Curr. Opin. Struct. Biol.* 15 (2005) 137–143.
- [21] J. Chen, C.L. Brooks 3rd, J. Khandogin, Recent advances in implicit solvent-based methods for biomolecular simulations, *Curr. Opin. Struct. Biol.* 18 (2008) 140–148.
- [22] D. Bashford, D.A. Case, Generalized born models of macromolecular solvation effects, *Annu. Rev. Phys. Chem.* 51 (2000) 129–152.
- [23] T. Lazaridis, M. Karplus, Effective energy function for proteins in solution, *Proteins* 35 (1999) 133–152.
- [24] T. Lazaridis, M. Karplus, “New view” of protein folding reconciled with the old through multiple unfolding simulations, *Science* 278 (1997) 1928–1931.
- [25] T. Lazaridis, M. Karplus, Discrimination of the native from misfolded protein models with an energy function including implicit solvation, *J. Mol. Biol.* 288 (1999) 477–487.
- [26] T. Lazaridis, Effective energy function for proteins in lipid membranes, *Proteins* 52 (2003) 176–192.
- [27] C.D. Schwieters, J.J. Kuszewski, N. Tjandra, G.M. Clore, The Xplor-NIH NMR molecular structure determination package, *J. Magn. Reson.* 160 (2003) 65–73.
- [28] C.D. Schwieters, J.J. Kuszewski, G. Marius Clore, Using Xplor, ÅNIH for NMR molecular structure determination, *Prog. Nucl. Magn. Reson. Spectrosc.* 48 (2006) 47–62.
- [29] A.T. Brünger, X-PLOR, Version 3.1: A System for X-ray Crystallography and NMR, Yale University Press, New Haven, 1992.
- [30] B.R. Brooks, R.E. Bruccoleri, B.D. Olafson, D.J. States, S. Swaminathan, M. Karplus, CHARMM: a program for macromolecular energy, minimization, and dynamics calculations, *J. Comput. Chem.* 4 (1983) 187–217.
- [31] B.R. Brooks, C.L. Brooks 3rd, A.D. Mackerell Jr., L. Nilsson, R.J. Petrella, B. Roux, Y. Won, G. Archontis, C. Bartels, S. Boresch, A. Caffisch, L. Caves, Q. Cui, A.R. Dinner, M. Feig, S. Fischer, J. Gao, M. Hodoscek, W. Im, K. Kucsera, T. Lazaridis, J. Ma, V. Ovchinnikov, E. Paci, R.W. Pastor, C.B. Post, J.Z. Pu, M. Schaefer, B. Tidor, R.M. Venable, H.L. Woodcock, X. Wu, W. Yang, D.M. York, M. Karplus, CHARMM: the biomolecular simulation program, *J. Comput. Chem.* 30 (2009) 1545–1614.
- [32] T.A. Cross, M. Sharma, M. Yi, H.X. Zhou, Influence of solubilizing environments on membrane protein structures, *Trends Biochem. Sci.* 36 (2011) 117–125.
- [33] J. Kuszewski, A.M. Gronenborn, G.M. Clore, Improving the quality of NMR and crystallographic protein structures by means of a conformational database potential derived from structure databases, *Protein Sci.* 5 (1996) 1067–1080.
- [34] G.A. Bermejo, G.M. Clore, C.D. Schwieters, Smooth statistical torsion angle potential derived from a large conformational database via adaptive kernel density estimation improves the quality of NMR protein structures, *Protein Sci.* 21 (2012) 1824–1836.
- [35] P.L. Privalov, G.I. Makhatazde, Contribution of hydration to protein folding thermodynamics. II. The entropy and Gibbs energy of hydration, *J. Mol. Biol.* 232 (1993) 660–679.
- [36] G.I. Makhatazde, P.L. Privalov, Contribution of hydration to protein folding thermodynamics. I. The enthalpy of hydration, *J. Mol. Biol.* 232 (1993) 639–659.
- [37] P.L. Privalov, G.I. Makhatazde, Contribution of hydration and non-covalent interactions to the heat capacity effect on protein unfolding, *J. Mol. Biol.* 224 (1992) 715–723.
- [38] P.L. Privalov, G.I. Makhatazde, Heat capacity of proteins. II. Partial molar heat capacity of the unfolded polypeptide chain of proteins: protein unfolding effects, *J. Mol. Biol.* 213 (1990) 385–391.
- [39] W.L. Jorgensen, J. Tirado-Rives, The OPLS [optimized potentials for liquid simulations] potential functions for proteins, energy minimizations for crystals of cyclic peptides and crambin, *J. Am. Chem. Soc.* 110 (1988) 1657–1666.
- [40] C.D. Schwieters, G.M. Clore, Internal coordinates for molecular dynamics and minimization in structure determination and refinement, *J. Magn. Reson.* 152 (2001) 288–302.
- [41] J.A. Losonczi, M. Andrec, M.W. Fischer, J.H. Prestegard, Order matrix analysis of residual dipolar couplings using singular value decomposition, *J. Magn. Reson.* 138 (1999) 334–342.
- [42] G.M. Clore, D.S. Garrett, R-factor, free R, and complete cross-validation for dipolar coupling refinement of NMR structures, *J. Am. Chem. Soc.* 121 (1999) 9008–9012.
- [43] G. Vriend, WHAT IF: a molecular modeling and drug design program, *J. Mol. Graph.* 8 (1990) 52–56.
- [44] J.F. Doreleijers, A.W. Sousa da Silva, E. Krieger, S.B. Nabuurs, C.A. Spronk, T.J. Stevens, W.F. Vranken, G. Vriend, G.W. Vuister, CING: an integrated residue-based structure validation program suite, *J. Biomol. NMR* 54 (2012) 267–283.
- [45] S.C. Lovell, I.W. Davis, W.B. Arendall 3rd, P.I. de Bakker, J.M. Word, M.G. Prisant, J.S. Richardson, D.C. Richardson, Structure validation by C $\alpha$  geometry: phi, psi and C $\beta$  deviation, *Proteins* 50 (2003) 437–450.
- [46] I.W. Davis, A. Leaver-Fay, V.B. Chen, J.N. Block, G.J. Kapral, X. Wang, L.W. Murray, W.B. Arendall 3rd, J. Snoeyink, J.S. Richardson, D.C. Richardson, MolProbity: all-atom contacts and structure validation for proteins and nucleic acids, *Nucleic Acids Res.* 35 (2007) W375–W383.
- [47] V.B. Chen, W.B. Arendall 3rd, J.J. Headd, D.A. Keedy, R.M. Immormino, G.J. Kapral, L.W. Murray, J.S. Richardson, D.C. Richardson, MolProbity: all-atom

- structure validation for macromolecular crystallography, *Acta Crystallogr. D Biol. Crystallogr.* 66 (2010) 12–21.
- [48] W.L. DeLano, PyMol, 2005. <<http://www.pymol.org>>.
- [49] R. Campbell, Personal Communication, 2013.
- [50] A.M. Gronenborn, D.R. Filpula, N.Z. Essig, A. Achari, M. Whitlow, P.T. Wingfield, G.M. Clore, A novel, highly stable fold of the immunoglobulin binding domain of streptococcal protein G, *Science* 253 (1991) 657–661.
- [51] J. Kuszewski, A.M. Gronenborn, G.M. Clore, Improving the packing and accuracy of NMR structures with a pseudopotential for the radius of gyration, *J. Am. Chem. Soc.* 121 (1999) 2337–2338.
- [52] D. Zheng, J.M. Aramini, G.T. Montelione, Validation of helical tilt angles in the solution NMR structure of the Z domain of Staphylococcal protein A by combined analysis of residual dipolar coupling and NOE data, *Protein Sci.* 13 (2004) 549–554.
- [53] S. Kaieda, C. Matsui, Y. Mimori-Kiyosue, T. Ikegami, Structural basis of the recognition of the SAMP motif of adenomatous polyposis coli by the Src-homology 3 domain, *Biochemistry* 49 (2010) 5143–5153.
- [54] S.G. Hyberts, M.S. Goldberg, T.F. Havel, G. Wagner, The solution structure of eglin c based on measurements of many NOEs and coupling constants and its comparison with X-ray structures, *Protein Sci.* 1 (1992) 736–751.
- [55] G. Cornilescu, J.L. Marquardt, M. Ottiger, A. Bax, Validation of protein structure from anisotropic carbonyl chemical shifts in a dilute liquid crystalline phase, *J. Am. Chem. Soc.* 120 (1998) 6836–6837.
- [56] B.E. Ramirez, O.N. Voloshin, R.D. Camerini-Otero, A. Bax, Solution structure of DinI provides insight into its mode of RecA inactivation, *Protein Sci.* 9 (2000) 2161–2169.
- [57] M. Cai, Y. Huang, R. Zheng, S.Q. Wei, R. Ghirlando, M.S. Lee, R. Craigie, A.M. Gronenborn, G.M. Clore, Solution structure of the cellular factor BAF responsible for protecting retroviral DNA from autointegration, *Nat. Struct. Biol.* 5 (1998) 903–909.
- [58] J.L. Baber, D. Libutti, D. Levens, N. Tjandra, High precision solution structure of the C-terminal KH domain of heterogeneous nuclear ribonucleoprotein K, a c-myc transcription factor, *J. Mol. Biol.* 289 (1999) 949–962.
- [59] P.M. Legler, M. Cai, A. Peterkofsky, G.M. Clore, Three-dimensional solution structure of the cytoplasmic B domain of the mannitol transporter Ilmannitol of the *Escherichia coli* phosphotransferase system, *J. Biol. Chem.* 279 (2004) 39115–39121.
- [60] P. Teriete, Y. Yao, A. Kolodzik, J. Yu, H. Song, M. Niederweis, F.M. Marassi, *Mycobacterium tuberculosis* Rv0899 adopts a mixed alpha/beta-structure and does not form a transmembrane beta-barrel, *Biochemistry* 49 (2010) 2768–2777.
- [61] D.S. Garrett, Y.J. Seok, A. Peterkofsky, A.M. Gronenborn, G.M. Clore, Solution structure of the 40,000 Mr phosphoryl transfer complex between the N-terminal domain of enzyme I and HPr, *Nat. Struct. Biol.* 6 (1999) 166–173.
- [62] D.S. Garrett, Y.J. Seok, D.I. Liao, A. Peterkofsky, A.M. Gronenborn, G.M. Clore, Solution structure of the 30 kDa N-terminal domain of enzyme I of the *Escherichia coli* phosphoenolpyruvate:sugar phosphotransferase system by multidimensional NMR, *Biochemistry* 36 (1997) 2517–2530.
- [63] M. Levitt, R. Sharon, Accurate simulation of protein dynamics in solution, *Proc. Natl. Acad. Sci. USA* 85 (1988) 7557–7561.
- [64] W.F. van Gunsteren, M. Karplus, Protein dynamics in solution and in a crystalline environment: a molecular dynamics study, *Biochemistry* 21 (1982) 2259–2274.
- [65] Y. Shen, O. Lange, F. Delaglio, P. Rossi, J.M. Aramini, G. Liu, A. Eletsky, Y. Wu, K.K. Singarapu, A. Lemak, A. Ignatchenko, C.H. Arrowsmith, T. Szyperski, G.T. Montelione, D. Baker, A. Bax, Consistent blind protein structure generation from NMR chemical shift data, *Proc. Natl. Acad. Sci. USA* 105 (2008) 4685–4690.
- [66] S. Raman, O.F. Lange, P. Rossi, M. Tyka, X. Wang, J. Aramini, G. Liu, T.A. Ramelot, A. Eletsky, T. Szyperski, M.A. Kennedy, J. Prestegard, G.T. Montelione, D. Baker, NMR structure determination for larger proteins using backbone-only data, *Science* 327 (2010) 1014–1018.
- [67] O.F. Lange, P. Rossi, N.G. Sgourakis, Y. Song, H.W. Lee, J.M. Aramini, A. Ertekin, R. Xiao, T.B. Acton, G.T. Montelione, D. Baker, Determination of solution structures of proteins up to 40 kDa using CS-Rosetta with sparse NMR data from deuterated samples, *Proc. Natl. Acad. Sci. USA* 109 (2012) 10873–10878.
- [68] R. Das, D. Baker, Macromolecular modeling with rosetta, *Annu. Rev. Biochem.* 77 (2008) 363–382.
- [69] S.B. Nabuurs, E. Krieger, C.A. Spronk, A.J. Nederveen, G. Vriend, G.W. Vuister, Definition of a new information-based per-residue quality parameter, *J. Biomol. NMR* 33 (2005) 123–134.
- [70] G.M. Clore, M.A. Robien, A.M. Gronenborn, Exploring the limits of precision and accuracy of protein structures determined by nuclear magnetic resonance spectroscopy, *J. Mol. Biol.* 231 (1993) 82–102.
- [71] G. Vriend, C. Sander, Quality control of protein models: directional atomic contact analysis, *J. Appl. Crystallogr.* 26 (1993) 47–60.
- [72] J.M. Word, S.C. Lovell, T.H. LaBean, H.C. Taylor, M.E. Zalis, B.K. Presley, J.S. Richardson, D.C. Richardson, Visualizing and quantifying molecular goodness-of-fit: small-probe contact dots with explicit hydrogen atoms, *J. Mol. Biol.* 285 (1999) 1711–1733.
- [73] A.M. Gronenborn, G.M. Clore, Structures of protein complexes by multidimensional heteronuclear magnetic resonance spectroscopy, *Crit. Rev. Biochem. Mol. Biol.* 30 (1995) 351–385.
- [74] R.A. Abagyan, M.M. Totrov, Contact area difference (CAD): a robust measure to evaluate accuracy of protein models, *J. Mol. Biol.* 268 (1997) 678–685.
- [75] G.M. Clore, A.M. Gronenborn, New methods of structure refinement for macromolecular structure determination by NMR, *Proc. Natl. Acad. Sci. USA* 95 (1998) 5891–5898.