



# Quality assessment of protein NMR structures<sup>☆</sup>

Antonio Rosato<sup>1</sup>, Roberto Tejero<sup>2</sup> and Gaetano T Montelione<sup>3,4</sup>

Biomolecular NMR structures are now routinely used in biology, chemistry, and bioinformatics. Methods and metrics for assessing the accuracy and precision of protein NMR structures are beginning to be standardized across the biological NMR community. These include both knowledge-based assessment metrics, parameterized from the database of protein structures, and model versus data assessment metrics. On line servers are available that provide comprehensive protein structure quality assessment reports, and efforts are in progress by the world-wide Protein Data Bank (wwPDB) to develop a biomolecular NMR structure quality assessment pipeline as part of the structure deposition process. These quality assessment metrics and standards will aid NMR spectroscopists in determining more accurate structures, and increase the value and utility of these structures for the broad scientific community.

## Addresses

<sup>1</sup> Magnetic Resonance Center and Department of Chemistry, University of Florence, 50019 Sesto Fiorentino, Italy

<sup>2</sup> Departamento de Química Física, Universidad de Valencia, Avenida Dr. Moliner 50, 46100 Burjassot, Valencia, Spain

<sup>3</sup> Center for Advanced Biotechnology and Medicine, Department of Molecular Biology and Biochemistry, and Northeast Structural Genomic Consortium, Rutgers, The State University of New Jersey, Piscataway, NJ 08854, USA

<sup>4</sup> Department of Biochemistry and Molecular Biology, Robert Wood Johnson Medical School, Rutgers, The State University of New Jersey, Piscataway, NJ 08854, USA

Corresponding author: Montelione, Gaetano T ([guy@cabm.rutgers.edu](mailto:guy@cabm.rutgers.edu))

**Current Opinion in Structural Biology** 2013, **23**:715–724

This review comes from a themed issue on **Biophysical methods**

Edited by **Wah Chiu** and **Gerhard Wagner**

For a complete overview see the [Issue](#) and the [Editorial](#)

Available online 21st September 2013

0959-440X/\$ – see front matter, © 2013 The Authors. Elsevier Ltd. All rights reserved.

<http://dx.doi.org/10.1016/j.sbi.2013.08.005>

## Introduction

Nuclear Magnetic Resonance (NMR) spectroscopy, along with X-ray diffraction and cryo-electron microscopy, is one of the three major experimental techniques providing three-dimensional (3D) structures of biological

macromolecules. In addition to its unique role in characterizing biomolecular dynamics, NMR is routinely used for structure determinations of small (<20 kDa) proteins [1<sup>••</sup>,2,3<sup>••</sup>] and is beginning to be used more routinely for determining structures of larger (20–50 kDa) soluble and membrane proteins (e.g. Refs. [4<sup>•</sup>,5<sup>•</sup>,6,7]). NMR-derived structure models can be used interchangeably with models generated by X-ray crystallography in many biological applications. It is therefore natural that many quality assessment metrics are common between the two techniques. In addition, there is a portfolio of metrics that are specific to NMR, which take into account the distinctive features of the NMR data used in the structure determination process.

In this review, we outline some of the metrics in common use for protein NMR structure quality assessment. A large part of our review reflects recently published recommendations of the world-wide Protein Data Base (wwPDB) task forces on validation of biomolecular structures determined by X-ray crystallography [8<sup>••</sup>] and NMR methods [9<sup>••</sup>].

## Knowledge-based measures

Knowledge-based (KB) metrics describe how well the structure model conforms to expectations with respect to selected features that can be assessed by comparison with the extensive database of experimental structures. These include bond length and bond angle distributions, dihedral angle distributions, atomic packing, hydrogen bond geometries, and other geometric features. Ideal values or value distributions are derived from statistical analyses of high-resolution X-ray structures, and are generally consistent with basic principles of biophysical chemistry.

There has been some debate regarding the use of KB information derived from X-ray crystal structures of biomolecules in assessing solution NMR structures. There are often differences in the sample conditions used in determining NMR and X-ray structures, and particular conformations from the distribution present in solution may be selected by the requirements of the crystal lattice. Nevertheless, there is no cogent reason to adopt different KB parameter distributions for assessing solution or solid-state protein NMR structures with respect to those used to assess structures determined by X-ray crystallography. Problems indicated by KB assessments can be mapped onto the 3D structure to identify local hot spots of structural inaccuracy [10<sup>••</sup>,11<sup>••</sup>].

The most general protein model assessment tools look at residue pair-distribution functions (e.g. PROSA2 [12]) or

<sup>☆</sup> This is an open-access article distributed under the terms of the Creative Commons Attribution-NonCommercial-No Derivative Works License, which permits non-commercial use, distribution, and reproduction in any medium, provided the original author and source are credited.

distributions of hydrophobic and hydrophilic residues (e.g. Verify3D [13]) which are characteristic of native protein structures. These analyses are important first steps in structure validation. However, their value is primarily in identifying severely incorrect folds [14], which rarely result from NMR structure determinations done with high restraint densities. These scores may, however, be important when assessing structures determined from sparse restraint networks. Non-globular protein folds (e.g. coil-coil structures) may exhibit poor PROSA2 or Verify3D scores even when the models are accurate.

Dihedral angle distributions are the most prominent KB statistic used in assessing protein NMR structures. They are generally reported as *Z*-scores relative to distributions observed in high-resolution crystal structures, or across all structures that have been deposited in the PDB. Dihedral angle distributions are generally reported separately for amino-acid residue backbone and side chains. Backbone  $\phi$ ,  $\psi$  distributions are generally assessed based on compliance with the Ramachandran plot. Historically, the program ProCheck [15] has been used for this analysis, but recent work using a larger set of X-ray structures determined at high resolution suggests that more accurate assessments can be made using improved backbone  $\phi$ ,  $\psi$  distribution statistics [8<sup>\*\*</sup>,16]. The assessment of side chain dihedral angle distributions (also referred to as rotamer normality) can be more subtle. Protein side chains have been observed to largely adopt standard rotamer states (g<sup>-</sup>, t, g<sup>+</sup>) even when buried in the cores of protein crystal structures. While NMR data can in principle determine accurate side-chain conformations, surface side chains are often dynamically averaged complicating the interpretation of the corresponding NMR data.

Several tools have been used to assess core atom packing, including the Molprobit [11<sup>\*\*</sup>] program for assessment of overpacking, and both the Molprobit and Rosetta-Holes [17] programs for assessment of underpacking. Severe atomic overlaps are rarely an issue in NMR structures, unless there are errors in the restraints, because of the use of lower-distance bounds. High-energy contacts can occur due to simplified treatments of van der Waals and hydrogen-bonded interactions, but are generally relieved by energy refinement. Underpacking, however, may be a more general problem for NMR-derived structures than generally appreciated, and methods like RosettaHoles [17] should be a key component of protein NMR structure assessment.

The role of KB-driven energy refinement in determining protein NMR structures, such as KB potentials or fragment libraries [18<sup>\*</sup>,19], is also somewhat controversial. In our opinion and experience, appropriate protocols for energy refinement, including the use of KB potentials

and fragment libraries, can significantly improve the accuracy of NMR structures (see e.g. [1,3,20<sup>\*</sup>,21]), particularly for larger proteins determined with sparse restraint networks [4<sup>\*</sup>,5<sup>\*</sup>]. Approaches based on molecular dynamics can also often provide appreciable improvements in structure quality [22].

### Model versus data measures

Model versus data (MvD) metrics describe how well the NMR structure model matches experimental data. MvD quality assessments include data that have been used in the structure generation process and, where feasible, cross-validation using data that have *not* been used in structure generation calculations.

#### *NMR restraint analysis*

The most general form of MvD validation involves comparison of distances and dihedral angles in models with the corresponding experimental restraints. Table 1 provides a standard format used to report such restraint violations [9<sup>\*\*</sup>,20<sup>\*</sup>], although other formats are also common in the NMR literature. These metrics provide an overview of global (or average) restraint violation statistics, as well as information on the most significant outliers. Clusters of restraint violations in regions of the 3D structure may indicate errors in the local structure. Methods have been described to convert between restraint formats [20<sup>\*</sup>,23,24], allowing initial restraint lists generated for one structure generation program to be used with alternative structure generation programs. The NOE completeness score [20<sup>\*</sup>,25] is a useful metric of structural accuracy, assessing the fraction of short distances in the model structure that are consistent with restraint data set.

#### *NOESY data*

Restraint analysis has the significant shortcoming that the restraints are themselves interpretations of NOESY and other NMR data. Accordingly, NMR structure quality assessment should also include some metrics validating models against uninterpreted spectral data. In the case of NOESY data, several methods have been developed for back-calculating NOESY spectral data (e.g. Refs. [26–28]), although to date none of these has come into general use. A more rapid, though approximate, approach is to compare models with unassigned NOESY peak lists. The RPF program [29] quantifies the agreement between the NOESY peak list, chemical shift assignments, and NMR models by calculating recall, precision, and F measures, as well as a normalized F-measure (called a Discriminating Power DP score). The normalized DP score is highly correlated with the accuracy of the NMR model [30<sup>\*</sup>,31<sup>\*</sup>].

#### *RDC, scalar coupling, paramagnetic, and SAXS data*

MvD metrics could also include assessment of scalar coupling, residual dipolar coupling (RDC), chemical shift

Table 1

Summary of KB and MvS structural statistics for bacterial protein Alr2454. This analysis was generated by the PSVS protein NMR structure quality assessment server [10\*\*,14]. These results are adopted from Ref. [66]

		Alr2454 <sup>a</sup>
<b>Completeness of resonance assignments<sup>b</sup></b>		
Backbone (%)		99.4
Side chain (%)		98.3
Aromatic (%)		96.6
Stereospecific methyl (%)		100
<b>Conformationally restricting restraints<sup>c</sup></b>		
Distance restraints		
Total		2478
Intra-residue ( $i = j$ )		688
Sequential ( $ i - j  = 1$ )		619
Medium range ( $1 <  i - j  < 5$ )		462
Long range ( $ i - j  \geq 5$ )		709
Dihedral angle restraints		162
Hydrogen bond restraints		0
Disulfide bond restraints		0
No. of restraints per residue		25.5
No. of long range restraints per residue		6.8
<b>Residual restraint violations<sup>c</sup></b>		
Average no. of distance violations per structure:		
0.1–0.2 Å		8.75
0.2–0.5 Å		1.85
>0.5 Å		0
Largest distance violation (Å)		0.35
Average no. of dihedral angle violations per structure:		
1–10°		8.75
>10°		0
Largest dihedral angle violation (°)		3.8
<b>NOE completeness score</b>		0.692
<b>Model quality<sup>c</sup></b>		
RMSD backbone atoms (Å) <sup>d</sup>		0.6
RMSD heavy atoms (Å) <sup>d</sup>		0.9
RMSD bond lengths (Å)		0.018
RMSD bond angles (°)		1.1
MolProbity Ramachandran statistics <sup>c,d</sup>		
Most favored regions (%)		96.8
Allowed regions (%)		3.1
Disallowed regions (%)		0.1
Global quality scores (raw/Z-score) <sup>c</sup>		
Verify3D	0.40	–0.96
ProsaII	0.66	0.04
ProCheck (phi-psi) <sup>d</sup>	–0.15	–0.28
ProCheck (all) <sup>d</sup>	–0.03	–0.18
MolProbity clash score	12.51	–0.62
RPF scores <sup>e</sup>		
Recall/precision	0.976	0.934
F-measure/DP-score	0.955	0.817
Model contents		
Ordered residue range <sup>d</sup>		1–100
Total no. of residues		108
BMRB accession number:		17965
PDB ID		2LJW <sup>a</sup>

<sup>a</sup> Structural statistics computed for an ensemble of 20 structures.

<sup>b</sup> Computed using AVS software [67] from the expected number of assignable resonances, excluding: highly exchangeable protons (N-terminal and Lys amino groups, Arg guanido groups, hydroxyls of Ser, Thr, Tyr), carboxyl carbons of Asp and Glu, non-protonated aromatic carbons, and the C-terminal His<sub>6</sub> tag. Methyl protons are counted as a single assignable resonance.

<sup>c</sup> Calculated using PSVS 1.4 [10\*\*]. Ramachandran statistics were calculated by Molprobity [11\*\*]. Average distance violations were calculated using the sum over  $r^{-6}$  for degenerate protons and stereochemically distinct protons lacking stereospecific assignments.

<sup>d</sup> On the basis of 'well defined' residue ranges [S(phi) + S(psi) > 1.8].

anisotropy (CSA), paramagnetic resonance enhancement (PRE), paramagnetic pseudo-contact shift (PCS), solid-state dipolar coupling, and small angle X-ray or neutron scattering (SAXS or SANS) data. Several tools for validating structures against these data are available, including methods for validation of protein structures against RDC data [32–34], CSA data [35] and SAXS data [36]. Residual dipolar coupling data provide information about the orientations of internuclear bond vectors with respect to the molecular orientation tensor, and hence provide long-range information including, for example, the relative orientations of secondary structure elements, such as helix tilt angles, which are sometimes inaccurate in protein NMR structures. The most commonly used MvD metric for these data is the RDC Q-factor [37]. Aromatic RDCs also have been found to provide an important approach for validating accuracy in core structures of proteins [38\*].

#### Free R factors

In spite of extensive discussions and of various formulations that have been proposed [25,27–29,37,39], a cross-validation metric analogous to the free R-factor of X-ray crystallography has not yet been broadly adopted by the biological NMR community. This is attributable, at least in part, to the sparseness of NMR data, and the need to manually evaluate each individual peak in NOESY spectra when applying truly quantitative methods. The RDC free Q-factor [37] is potentially more accessible, as quantifying RDC data is generally more straightforward than the quantification of large numbers of NOESY peaks.

#### Chemical shift data

Chemical shift data have great potential for NMR structure validation. The most straightforward validation involves comparing experimental chemical shift data with values predicted from the 3D model structure [40–42,43\*,44\*\*,45,46,47\*]. Extensive chemical shift data must be obtained at the onset of any protein structure determination, and chemical shift data often used indirectly as restraints on dihedral angles [48\*\*]. Several recent papers have described important progress in calculating chemical shifts from molecular models using both empirical [49\*\*,50,51\*,52\*\*] and quantum chemical [44\*\*,47\*,53,54] approaches. However, systematic, large scale tests using chemical shifts as a standard metric for protein NMR structure validation are not yet available.

#### Which parts of the structure should be validated?

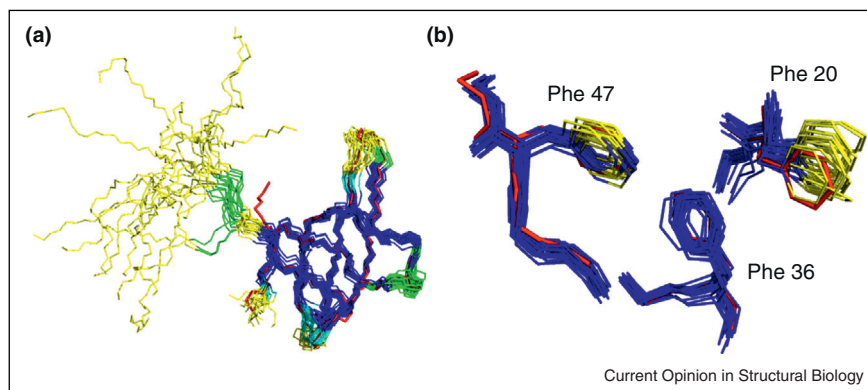
Protein NMR structures are generally represented by ensembles of conformers having the same level of agreement with the experimental data. Different regions of the structure are often converged to different degrees among conformers. In common practice, a distinction is made between well-defined and ill-defined (i.e. not-well-defined) regions [9\*\*,20\*]. In single-domain proteins,

the backbone or all-heavy-atom root mean square deviation of coordinates (RMSD) computed after superposition of all the well-defined regions is taken as a measure of structure precision. Similarly, the per-residue RMSD is taken as a measure of local precision. Alternative measures of local or global precision include dihedral angle order parameters (DAOP) [55,56\*,57\*\*] and distance variance matrix methods [57\*\*,58,59] (Snyder DA, Grullon J, Huang YJ, Tejero R, Montelione GT: **The expanded findcore method for identification of a core atom set for assessment of protein structure prediction.** *Proteins* 2013 (submitted for publication). This is an extension of the variance matrix method of Ref. [58] for annotating well-defined versus ill-defined atoms in an NMR ensemble. FindCore also provides criteria for identifying regions of the protein structure that are internally well-defined, but not well-defined with respect to one another, guiding the independent superimposition of these regions for RMSD calculations.) (Figure 1). Often, there is a good, albeit qualitative, correlation between local restraint density and local precision, implying that the ill-defined regions of the structure result from the experimental data providing insufficient information. This variability can be due to local protein dynamics (so that locally a single conformation in solution in fact does not exist) or to experimental factors limiting the information that can be extracted from spectra (e.g. extensive resonance degeneracy).

The precision of NMR structures indicated by the convergence across the ensemble of NMR conformers is operational, and does not provide a true representation of the Boltzmann distribution of conformations actually present in the NMR sample. Indeed, NMR structure ensembles do not even provide a statistically meaningful description of the true precision of coordinates given the experimental uncertainties in deriving distance restraints. For example, fast exchange may give rise to inconsistent restraints, which when simultaneously satisfied can pin the local conformational distribution into an unrealistically narrow range of conformations [60].

Because NMR experiments do not provide enough data to characterize them, the conformations observed in ill-defined regions largely result from of the combination of random initial conformations with the potential energy functions. Hence, these regions should not be included in global structure validation. However, even such ill-defined regions may be restrained to some degree by sparse experimental data, which may be biologically relevant. Ill-defined regions may also contribute to the back-calculation of NMR observables (e.g. NOESY peak lists). Hence, in our opinion all atoms for which experimental data are available (including only chemical shift data) should be included in the atomic coordinates that are deposited in the PDB, consistent with recommendations of the wwPDB NMR VTF [9\*\*].

Figure 1



Comparison of DOAP and variance distance matrix results for identifying well-defined atom sets. **(a)** NMR structure ensemble superimposition showing residues defined as 'well defined' by the PDBStat DAOP analysis [20\*] and those identified as 'well defined' by the FindCore [20\*,58] method, for protein SgR42 (PDB\_id 2jz2). Residues identified as 'well defined' by both methods are shown in dark blue, those identified as 'well defined' by variance matrix but not by DAOP in light blue, and those identified as 'well defined' by DAOP but not by variance matrix in green. Residues identified as 'ill defined' by both methods are shown in yellow. The backbone atoms of the corresponding X-ray crystal structure (PDB id 3c4s) are shown in red. **(b)** Expansion showing atom-specific 'well defined' (dark blue) and 'ill defined' (yellow) designations for the sidechains of residues Phe20, Phe36, and Phe47 in protein SgR42. This image demonstrates the value of atom-specific designations of well-defined regions of a protein NMR structure. Adopted from Ref. [20\*].

### Protein structure quality assessment servers

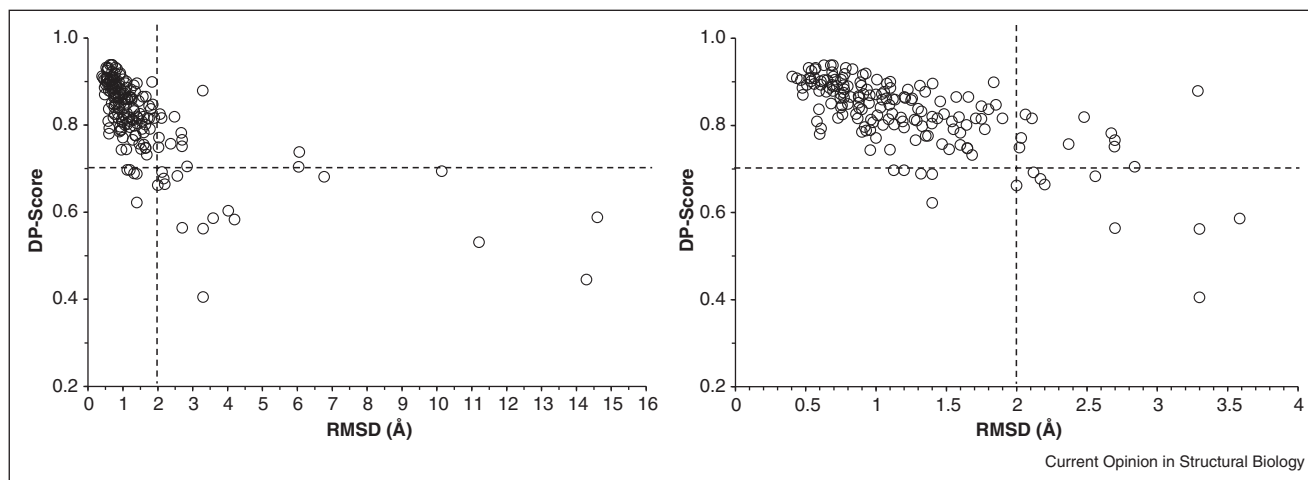
Recently, several on-line servers have become available which integrate both KB and MvD assessment metrics to provide a comprehensive NMR structure quality assessment report. Web-based tools described in Table 2 include the CiNG [61\*\*,62\*], Molprobit [11\*\*], PSVS [10\*\*,14], Vivaldi [63\*], and ResProx [43\*] servers. These, as well as other programs for NMR structure quality assessment, have been critically reviewed in a recent publication [57\*\*]. The wwPDB NMR VTF has also recommended

metrics and standards for biomolecular structure quality assessment [8\*\*,9\*\*], and these guidelines are being used in developing software pipelines that will generate standardized structure quality reports for all NMR and X-ray crystal structures submitted to the wwPDB.

### Quality measures versus structure accuracy

Satisfaction of KB metrics is a necessary but not sufficient criteria for validating the accuracy of protein NMR structures. This conclusion was illustrated by the recent

Figure 2



Correlation between backbone RMSD to the reference structure, which is taken as a measure of structural accuracy, and the DP-score. These results demonstrate the discriminating power of the DP score in distinguishing accurate from less-accurate protein NMR models. The horizontal and vertical lines indicate empirical thresholds for accurate structure models [30\*]. Data are from the 2010 and 2013 rounds of the CASD-NMR project [1,2], which have been merged. Left panel: full range; right panel: expansion of the 0–4.0 Å RMSD range.



**Table 2****Web accessible servers providing knowledge-based (KB) and/or model versus data (MvD) protein NMR structure quality assessment reports**

Server	url	Validation statistics provided
CING [62] v.1.0: Common Interface for NMR Structure Generation	<a href="https://nmr.le.ac.uk/icing">https://nmr.le.ac.uk/icing</a>	KB: Bond lengths and bond angles, backbone Ramachandran distributions, peptide omegas, packing, sidechain rotamer normality, disulfides, salt-bridges, chemical shift validations. Provides residue-specific ROC scores [68]. Uses DSSP [69,70], WHATIF [71], PROCHECK [15,72], and BMRB chemical shift validation MvD: DOAP analysis. Restraint violation analysis, redundancy and duplicated restraints, back calculation of chemical shift from structure. Uses SHIFTX [49**] and VASCO [73]
Molprobrity [11**]	<a href="http://molprobrity.biochem.duke.edu/">http://molprobrity.biochem.duke.edu/</a>	KB: Backbone Ramachandran distributions, peptide omegas, packing, H-bond satisfaction, sidechain rotamer normality MvD: RDC analysis with RDCvis [74]
PSVS [10**,14] v.1.5. Protein Structure Validation Software suite.	<a href="http://psvs.nesg.org/">http://psvs.nesg.org/</a>	KB: Bond lengths and bond angles, backbone Ramachandran distributions, peptide omegas, packing, sidechain rotamer normality, chemical shift validation and completeness using AVS [67]. Provides Z scores relative to high-resolution crystal structures for Prosall [12], Verify3D [13], Procheck_bb [15,72], Procheck_all [15,72], and MolProbrity [11**]. Also uses DSSP [69,70], PDBStat [20*], LACS [75], and PDB Clash score (deposit.rcsb.org/validate) MvD: DAOP or Variance Matrix analysis. Restraint violation analysis, redundancy and duplicated restraints, NOE completeness, NOE DP scores, RDC Q scores, and GLM-RMSD [31*] 'equivalent resolution' score. Provides mapping of Ramachandran outliers and RPF violations onto 3D structure. Also uses PDBStat [20*], FindCore [20*,58] and RPF [29,30*] software
ResProx [43*] Resolution by Proxy	<a href="http://www.resprox.ca/">http://www.resprox.ca/</a>	KB: Assesses 'equivalent resolution' based on 25 protein features including backbone Ramachandran distributions, peptide omegas, H-bond geometry, over and underpacking. Uses Vader [40], PROSESS [42], Molprobrity [11**], RosettaHoles [17], and GeNMR [41] software
Vivaldi [63*] v.1: Visualization and validation of biomacromolecular NMR structures from the PDB	<a href="http://www.ebi.ac.uk/pdbe/vivaldi/">http://www.ebi.ac.uk/pdbe/vivaldi/</a>	KB: Bond lengths and bond angles, peptide omegas, chemical shift validations, backbone Ramachandran distributions, CING [61**] ROC scores. Uses WHATIF [71], CING [61**], and VASCO [73] MvD: DAOP and Variance Matrix analysis. Restraint violation and RDC analysis. Uses NMRCORE [76]

Critical Assessment of Protein Structure Determination by NMR (CASD-NMR) [1,2], in which NMR data was publicly released, and the results of automated NOE assignment and structure generation methods were compared in a blind fashion with manually refined protein NMR structures. In many cases, models with excellent energetics were in fact quite different from the manually refined 'gold standard' structures. However, energetic considerations can identify inaccurate restraint data, and thus guide the experimentalist to more accurate interpretations of the raw NMR data in terms of experimental restraints [20\*,21].

The CASD-NMR experiment also compared correlations between several KB and MvD metrics with structural accuracy. In this analysis, the DP score was observed to provide the highest correlation between submitted models and the manually refined target structure (Figure 2). Although the DP score has significant shortcomings to the extent that it uses an interpreted

NOESY peak list rather than raw NOESY spectral data, these results demonstrate the importance of combining KB and MvD validation scores in structure quality assessment. For example, the DP score has recently been combined with several KB scores by linear regression analysis to provide an 'equivalent resolution' metric [31\*,64].

#### Areas of current research

The trend towards determining larger biomolecular structures using sparse NMR data and hybrid methods incorporating small angle X-ray scattering, chemical cross links, and other data, requires increasing use of KB information and potential energy functions in the structure determination process. Some solid-state NMR structure determinations may also rely on relatively sparse networks of experimental data. These methods provide less independent information for use in cross-validation, creating important challenges to the field in terms of robust structure quality assessment.

Quality assessment considerations for NMR structures of nucleic acids and carbohydrates, while generally similar to those for proteins, have not yet been extensively investigated. Methods also need to be developed for generating Boltzmann ensembles of conformers that best satisfy the experimental data (e.g. [37,65]), rather than fitting all the data to a single conformer multiple times, as is the common practice. Ensemble-averaging interpretations are particularly important in highly dynamic regions of a biomolecular structure, and new methods are needed for quality assessment of the ensembles of models proposed for such disordered regions.

## Conclusions

A high quality NMR-derived structure should meet minimal standards based on a wide range of KB and MvD validation assessment metrics. No single metric score is sufficient to validate a protein NMR structure. Users of biomolecular NMR structures also need to be informed about which parts of the structure are well-defined by the data, and which are less-reliably defined in terms of a unique structure. This information can be critical in interpreting structure–function relationships.

Structure determinations by NMR are often marginally underdetermined, and rely in either subtle or substantial ways on KB information, including at the very least standard values of bond lengths and bond angles. Although they are useful for biological studies, and in some cases quite accurate, the relatively low density of data constituting a typical protein NMR structure presents challenges to cross validation, using part of the data to generate the structure and another part of the data to validate the structure. In this regard, chemical shift data, which are available for most atoms reported in modern biomolecular NMR structures, hold promise for providing a robust and general approach for data-based protein structure quality assessments.

## Acknowledgements

We thank all of the scientists participating in the CASD-NMR Project [1\*\*,2] for generating the data summarized (Figure 2). Special thanks to J. Aramini, L. Ferella, Y.J. Huang, B. Mao, D. Snyder, and the members of the wwPDB NMR Validation Task Force for helpful scientific discussions. This work was supported by grants from the Protein Structure Initiative — Biology of the National Institutes of Health, Grant U54-GM094597 (to G.T.M), by the CONSOLIDER INGENIO CSD2010-00065 and Generalitat Valenciana PROMETEO 2011/008 (to R.T), and by the European Community FP7 e-Infrastructure ‘WeNMR’ project, Grant 261572 (to A.R.).

## References and recommended reading

Papers of particular interest, published within the period of review, have been highlighted as:

- of special interest
- of outstanding interest

1. Rosato A, Aramini JM, Arrowsmith C, Bagaria A, Baker D, Cavalli A, Doreleijers JF, Eletsky A, Giachetti A, Guerry P, Gutmanas A *et al.*: **Blind testing of routine, fully automated determination of protein structures from NMR data.** *Structure* 2012, **20**:227-236.

This paper reports the first results of the CASD-NMR initiative. In this collaborative effort, raw data were provided to eight teams developing software for automated NMR structure determination, while the corresponding 3D coordinates were not released. An independent team then compared the automatically generated structures with the reference structure, manually determined from the same data. All automated NOESY-based methods performed well on all targets.

2. Rosato A, Bagaria A, Baker D, Bardiaux B, Cavalli A, Doreleijers JF, Giachetti A, Guerry P, Guntert P, Herrmann T, Huang YJ *et al.*: **CASD-NMR: critical assessment of automated structure determination by NMR.** *Nat Methods* 2009, **6**:625-626.

3. Mao B, Guan R, Montelione GT: **Improved technologies now routinely provide protein NMR structures useful for molecular replacement.** *Structure* 2011, **19**:757-766.

This paper demonstrates that NMR structures determined and validated using modern methods provide coordinates which are sufficiently accurate to allow their routine use in determining phases of corresponding or homologous X-ray crystal structures by molecular replacement methods. In some cases, however, the accuracy of these NMR structures, as assessed by their phasing power, is improved by KB refinement with Rosetta.

4. Lange OF, Rossi P, Sgourakis NG, Song Y, Lee HW, Aramini JM, Ertekin A, Xiao R, Acton TB, Montelione GT, Baker D: **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* 2012, **109**:10873-10878.

This paper demonstrates the use of KB methods, including fragment libraries, in determining accurate protein structures by NMR, using perdeuterated, methyl protonated proteins, in the 20–40 kDa size range. Related results are presented in Ref. [5\*].

5. Raman S, Lange OF, Rossi P, Tyka M, Wang X, Aramini J, Liu G, Ramelet TA, Eletsky A, Szyperski T, Kennedy MA *et al.*: **NMR structure determination for larger proteins using backbone-only data.** *Science* 2010, **327**:1014-1018.

This paper demonstrates the use of KB methods, including fragment libraries, in determining accurate protein structures by NMR using perdeuterated proteins without methyl protonation in the 20–40 kDa size range. Related results are presented in Ref. [4\*].

6. Hiller S, Wagner G: **The role of solution NMR in the structure determinations of vDAC-1 and other membrane proteins.** *Curr Opin Struct Biol* 2009, **19**:396-401.

7. Hiller S, Garces RG, Malia TJ, Orekhov VY, Colombini M, Wagner G: **Solution structure of the integral human membrane protein vDAC-1 in detergent micelles.** *Science* 2008, **321**:1206-1210.

8. Read RJ, Adams PD, Arendall WB 3rd, Brunger AT, Emsley P, Joosten RP, Kleywegt GJ, Krissinel EB, Luthe T, Otwinowski Z, Perrakis A *et al.*: **A new generation of crystallographic validation tools for the protein data bank.** *Structure* 2011, **19**:1395-1412.

An overview of the recommendations of wwPDB X-ray Crystallography Validation Task Force, providing guidelines for both KB and MvD validation of biomolecular structures determined by X-ray crystallography. Recommendations are also outlined for the presentation of these results in structure validation reports to be prepared by the wwPDB.

9. Montelione GT, Nilges M, Bax A, Güntert P, Herrmann T, Richardson JS, Schwieters CD, Vranken WF, Vuister GW, Wishart DS, Berman HM *et al.*: **Recommendations of the wwPDB NMR validation task force.** *Structure* 2013, **21**:1563-1570.

An overview of consensus conclusions by the wwPDB NMR Validation Task Force providing recommendations for validation of protein NMR structures in three phases: (i) validation methods ready to implement immediately as part of a consensus NMR structure validation pipeline, (ii) validation methods which work well, but have not yet been broadly adopted by the biomolecular NMR community, and (iii) key issues of biomolecular structure validation which require additional research.

10. Bhattacharya A, Tejero R, Montelione GT: **Evaluating protein structures determined by structural genomics consortia.** *Proteins* 2007, **66**:778-795.

This paper describes the PSVS server, a popular tool for NMR structure validation. PSVS provides an extensive analysis of both KB and MvD validation metrics for protein NMR structures, and a comprehensive structure validation report. Although the original paper was published in 2007, the server is updated regularly by the introduction of new validation metrics. See also Ref. [14] exploring the use of PSVS for structure quality assessment of homology models.

11. Chen VB, Arendall WB 3rd, Headd JJ, Keedy DA, Immormino RM, Kapral GJ, Murray LW, Richardson JS, Richardson DC: **Molprobity: all-atom structure validation for macromolecular crystallography.** *Acta Crystallogr. Sect. D, Biol. Crystallogr.* 2010, **66(Pt 1)**:12-21.

Molprobity software is a key component of most protein NMR structure validation servers. It provides comprehensive KB validation of protein structures, and includes the RDCvis module (Ref. [75]) for visual validation of protein NMR structures using RDC data.

12. Sippl MJ: **Recognition of errors in three-dimensional structures of proteins.** *Proteins* 1993, **17**:355-362.
13. Luthy R, Bowie JU, Eisenberg D: **Assessment of protein models with three-dimensional profiles.** *Nature* 1992, **356**:83-85.
14. Bhattacharya A, Wunderlich Z, Monleon D, Tejero R, Montelione GT: **Assessing model accuracy using the homology modeling automatically software.** *Proteins* 2008, **70**:105-118.
15. Laskowski RA, Rullmann JA, MacArthur MW, Kaptein R, Thornton JM: **Aqua and procheck-NMR: programs for checking the quality of protein structures solved by NMR.** *J Biomol NMR* 1996, **8**:477-486.
16. Arendall WB 3rd, Tempel W, Richardson JS, Zhou W, Wang S, Davis IW, Liu ZJ, Rose JP, Carson WM, Luo M, Richardson DC: **A test of enhancing model accuracy in high-throughput crystallography.** *J Struct Funct Genomics* 2005, **6**:1-11.
17. Sheffler W, Baker D: **Rosettaholes: rapid assessment of protein core packing for structure prediction, refinement, design, and validation.** *Protein Sci Publ Protein Soc* 2009, **18**:229-239.
18. van der Schot G, Zhang Z, Vernon R, Shen Y, Vranken WF, Baker D, Bonvin AM, Lange OF: **Improving 3D structure prediction from chemical shift data.** *J Biomol NMR* 2013, **57**:27-35.

This work demonstrates a new and improved fragment picker and the iterative sampling algorithm RASREC to provide significant improvements in convergence and accuracy of the CS-Rosetta method (Ref. [19]). They also introduce improved criteria for assessing the accuracy of the models produced by CS-Rosetta.

19. Shen Y, Lange O, Delaglio F, Rossi P, Aramini JM, Liu G, Eletsky A, Wu Y, Singarapu KK, Lemak A, Ignatchenko A et al.: **Consistent blind protein structure generation from NMR chemical shift data.** *Proc Natl Acad Sci U S A* 2008, **105**:4685-4690.
20. Tejero R, Snyder D, Mao B, Aramini JM, Montelione GT: **PDBStat: a universal restraint converter and restraint analysis software package for protein NMR.** *J Biomol NMR* 2013, **56**:337-351.
- PDBStat software provides a robust tool for MvD validations, including extensive restraint analysis and conversion utilities. The software can also be used to convert an initial restraint list into restraint lists for other software (e.g. CYANA restraints to Rosetta restraints), allowing validation by consensus structure generation with multiple software packages. PDBStat also provides MvD validation against RDC data.
21. Ramelot TA, Raman S, Kuzin AP, Xiao R, Ma LC, Acton TB, Hunt JF, Montelione GT, Baker D, Kennedy MA: **Improving NMR protein structure quality by Rosetta refinement: a molecular replacement study.** *Proteins* 2009, **75**:147-167.
22. Bertini I, Case DA, Ferella L, Giachetti A, Rosato A: **A grid-enabled web portal for NMR structure refinement with amber.** *Bioinformatics* 2011, **27**:2384-2390.
23. Nederveen AJ, Doreleijers JF, Vranken W, Miller Z, Spronk CA, Nabuurs SB, Guntert P, Livny M, Markley JL, Nilges M, Ulrich EL et al.: **Recoord: a recalculated coordinate database of 500+ proteins from the PDB using restraints from the biomagresbank.** *Proteins* 2005, **59**:662-672.
24. Vranken WF, Boucher W, Stevens TJ, Fogh RH, Pajon A, Llinas M, Ulrich EL, Markley JL, Ionides J, Laue ED: **The CCPN data model for NMR spectroscopy: development of a software pipeline.** *Proteins* 2005, **59**:687-696.
25. Doreleijers JF, Raves ML, Rullmann T, Kaptein R: **Completeness of nodes in protein structure: a statistical analysis of NMR.** *J Biomol NMR* 1999, **14**:123-132.
26. Borgias BA, James TL: **Two-dimensional nuclear overhauser effect: complete relaxation matrix analysis.** *Methods Enzymol* 1989, **176**:169-183.

27. Gronwald W, Kirchhofer R, Gorler A, Kremer W, Ganslmeier B, Neidig KP, Kalbitzer HR: **Rfac, a program for automated NMR R-factor estimation.** *J Biomol NMR* 2000, **17**:137-151.

28. Zhu L, Dyson HJ, Wright PE: **A noesy-hsqc simulation program, spirit.** *J Biomol NMR* 1998, **11**:17-29.
29. Huang YJ, Powers R, Montelione GT: **Protein NMR recall, precision, and F measure scores (RPF scores): structure quality assessment measures based on information retrieval statistics.** *J Am Chem Soc* 2005, **127**:1665-1674.
30. Huang YJ, Rosato A, Singh G, Montelione GT: **Rpf: a quality assessment tool for protein NMR structures.** *Nucleic Acids Res* 2012, **40(Web Server issue)**:W542-W546.

This article describes a web server for the calculation of RPF scores for NMR structures. This MvD measure exploits unassigned NOESY data. One of the output measures, the DP-score, is a useful predictor of the accuracy of the query structure (see also Ref. [29]).

31. Bagaria A, Jaravine V, Huang YJ, Montelione GT, Guntert P: **Protein structure validation by generalized linear model root-mean-square deviation prediction.** *Protein Sci Publ Protein Soc* 2012, **21**:229-238.

This article describes a predictor of the quality of protein structures that integrates a variety of different measures into a single 'equivalent resolution' score. The authors demonstrate that this score provides an indication of structural accuracy, provided as the estimated RMSD from the true structure.

32. Clore GM, Robien MA, Gronenborn AM: **Exploring the limits of precision and accuracy of protein structures determined by nuclear magnetic resonance spectroscopy.** *J Mol Biol* 1993, **231**:82-102.
33. Valafar H, Prestegard JH: **Redcat: a residual dipolar coupling analysis tool.** *J Magn Reson* 2004, **167**:228-241.
34. Bryson M, Tian F, Prestegard JH, Valafar H: **Redcraft: a tool for simultaneous characterization of protein backbone structure and motion from rdc data.** *J Magn Reson* 2008, **191**:322-334.
35. Cornilescu G, Marquardt JL, Ottiger M, Bax A: **Validation of protein structure from anisotropic carbonyl chemical shifts in a dilute liquid crystalline phase.** *J Am Chem Soc* 1998, **120**:6836-6837.
36. Grishaev A, Wu J, Trewella J, Bax A: **Refinement of multidomain protein structures by combination of solution small-angle X-ray scattering and NMR data.** *J Am Chem Soc* 2005, **127**:16621-16628.
37. Clore GM, Schwieters CD: **How much backbone motion in ubiquitin is required to account for dipolar coupling data measured in multiple alignment media as assessed by independent cross-validation?** *J Am Chem Soc* 2004, **126**:2923-2938.

38. Sathyamoorthy B, Singarapu KK, Garcia AE, Szyperski T: **Protein conformational space populated in solution probed with aromatic residual dipolar  $^{13}\text{C}$ - $^1\text{H}$  couplings.** *Chembiochem* 2013, **14**:684-688.

An innovative approach demonstrating the correlation of aromatic Tyr and Phe  $^{13}\text{C}$ - $^1\text{H}$  RDC measurements with protein structure accuracy. These data are relatively easy to measure and provide important information useful for either structure determination, refinement, or validation.

39. Brunger AT, Clore GM, Gronenborn AM, Saffrich R, Nilges M: **Assessing the quality of solution nuclear magnetic resonance structures by complete cross-validation.** *Science* 1993, **261**:328-331.
40. Willard L, Ranjan A, Zhang H, Monzavi H, Boyko RF, Sykes BD, Wishart DS: **Vadar: a web server for quantitative evaluation of protein structure quality.** *Nucleic Acids Res* 2003, **31**:3316-3319.
41. Berjanskii M, Tang P, Liang J, Cruz JA, Zhou J, Zhou Y, Bassett E, MacDonell C, Lu P, Lin G, Wishart DS: **GeNMR: a web server for rapid NMR-based protein structure determination.** *Nucleic Acids Res* 2009, **37(Web Server issue)**:W670-W677.
42. Berjanskii M, Liang Y, Zhou J, Tang P, Stothard P, Zhou Y, Cruz J, MacDonell C, Lin G, Lu P, Wishart DS: **Process: a protein structure evaluation suite and server.** *Nucleic Acids Res* 2010, **38(Web Server issue)**:W633-W640.



43. Berjanskii M, Zhou J, Liang Y, Lin G, Wishart DS: **Resolution-by-proxy: a simple measure for assessing and comparing the overall quality of NMR protein structures.** *J Biomol NMR* 2012, **53**:167-180.

ResProx calculates equivalent X-ray crystal structure resolution of a protein NMR structure from 25 measurable protein features. It can be used to identify under-restrained, poorly refined or inaccurate NMR structures.

44. Vila JA, Aramini JM, Rossi P, Kuzin A, Su M, Seetharaman J, Xiao R, Tong L, Montelione GT, Scheraga HA: **Quantum chemical  $^{13}\text{C}$ (alpha) chemical shift calculations for protein NMR structure determination, refinement, and validation.** *Proc Natl Acad Sci USA* 2008, **105**:14389-14394.

Quantum chemistry calculations by the program CheShift are used to predict  $^{13}\text{C}_\alpha$  chemical shifts for various residue types, flanking sequences, and backbone  $\phi$ ,  $\psi$  values. The resulting web-accessible database provides a rapid approach for protein backbone structure quality assessment.

45. Martin OA, Villegas ME, Vila JA, Scheraga HA: **Analysis of  $^{13}\text{C}$ alpha and  $^{13}\text{C}$ beta chemical shifts of cysteine and cystine residues in proteins: a quantum chemical approach.** *J Biomol NMR* 2010, **46**:217-225.

46. Vila JA, Serrano P, Wuthrich K, Scheraga HA: **Sequential nearest-neighbor effects on computed  $^{13}\text{C}$ alpha chemical shifts.** *J Biomol NMR* 2010, **48**:23-30.

47. Martin OA, Vila JA, Scheraga HA: **Cheshift-2: graphic validation of protein structures.** *Bioinformatics* 2012, **28**:1538-1539.

CheShift-2 provides a graphical user interface to identify possible flaws in an NMR structure based on quantum mechanical computation of  $^{13}\text{C}_\alpha$  shifts. The shifts are calculated as averages over all members of the ensemble, and their average compared to the experimental data.

48. Shen Y, Bax A: **Sparta+: a modest improvement in empirical NMR chemical shift prediction by means of an artificial neural network.** *J Biomol NMR* 2010, **48**:13-22.

Together with ShiftX2, this is one of the most commonly used and reliable methods for NMR chemical shift predictions.

49. Han B, Liu Y, Ginzinger SW, Wishart DS: **Shiftx2: significantly improved protein chemical shift prediction.** *J Biomol NMR* 2011, **50**:43-57.

Together with Sparta+, this is one of the most commonly used and reliable methods for NMR chemical shift predictions. See also Ref. [50].

50. Wang B, Wang Y, Wishart DS: **A probabilistic approach for validating protein NMR chemical shift assignments.** *J Biomol NMR* 2010, **47**:85-99.

51. Sahakyan AB, Vranken WF, Cavalli A, Vendruscolo M: **Using side-chain aromatic proton chemical shifts for a quantitative analysis of protein structures.** *Angew Chem Int Ed Engl* 2011, **50**:9620-9623.

The authors demonstrate the use of aromatic chemical shift data in distinguishing correct structural models from incorrect ones. They also explore its use to detect differences caused by cofactor or ligand binding, or by sequence alterations between structures.

52. Sahakyan AB, Vranken WF, Cavalli A, Vendruscolo M: **Structure-based prediction of methyl chemical shifts in proteins.** *J Biomol NMR* 2011, **50**:331-346.

This paper describes the CH3Shift method of performing structure-based predictions of methyl chemical shifts. The method achieves high accuracy in the predictions of  $^1\text{H}$  chemical shifts for Ala, Thr, Val, Leu and Ile methyl groups. This approach is useful both for assignment of methyl groups from 3D protein structure, and for validation of 3D protein structures, particularly those determined using perdeuterated, methyl protonated samples.

53. Vila JA, Sue SC, Fraser JS, Scheraga HA, Dyson HJ: **Cheshift-2 resolves a local inconsistency between two X-ray crystal structures.** *J Biomol NMR* 2012, **54**:193-198.

54. Vila JA, Arnautova YA, Vorobjev Y, Scheraga HA: **Assessing the fractions of tautomeric forms of the imidazole ring of histidine in proteins as a function of pH.** *Proc Natl Acad Sci USA* 2011, **108**:5602-5607.

55. Hyberts SG, Goldberg MS, Havel TF, Wagner G: **The solution structure of eglin c based on measurements of many nodes and coupling constants and its comparison with X-ray structures.** *Protein Sci* 1992, **1**:736-751.

56. Kirchner DK, Guntert P: **Objective identification of residue ranges for the superposition of protein structures.** *BMC Bioinformatics* 2011, **12**:170.

The CYRANGE method for the identification of well-ordered regions in NMR structures combines different aspects of approaches based on both dihedral angle order parameters and distance matrices, and integrates them with new ideas. Overall, the aim is to define residue ranges that contain few gaps and cover most of the sequence. This method has been adopted by the wwPDB NMR as a standard for defining well-defined residue ranges for quality assessment.

57. Vuister GW, Fogh RH, Hendrickx PM, Doreleijers JF, Gutmanas A: **An overview of tools for the validation of protein NMR structures.** *J Biomol NMR* 2013.

This review reports a comparison of the application of nine different validation tools on two example NMR structures. Interestingly, the same property as analyzed by different tools is not always consistently assessed as being good or bad.

58. Snyder DA, Montelione GT: **Clustering algorithms for identifying core atom sets and for assessing the precision of protein structure ensembles.** *Proteins* 2005, **59**:673-686.

59. Andrec M, Snyder DA, Zhou Z, Young J, Montelione GT, Levy RM: **A large data set comparison of protein structures determined by crystallography and NMR: statistical test for structural differences and the effect of crystal packing.** *Proteins* 2007, **69**:449-465.

60. Tejero R, Bassolino-Klimas D, Bruccoleri RE, Montelione GT: **Simulated annealing with restrained molecular dynamics using congen: energy refinement of the NMR solution structures of epidermal and type-alpha transforming growth factors.** *Protein Sci* 1996, **5**:578-592.

61. Doreleijers JF, Sousa da Silva AW, Krieger E, Nabuurs SB, Spronk CA, Stevens TJ, Vranken WF, Vriend G, Vuister GW: **CiNG: an integrated residue-based structure validation program suite.** *J Biomol NMR* 2012, **54**:267-283.

CiNG is a powerful web-based tool for NMR structure validation. CiNG provides an extensive analysis of both KB and MvD parameters, and provide residue-specific quality scores, ROC scores, which are useful for protein structure validation and refinement.

62. Doreleijers JF, Vranken WF, Schulte C, Markley JL, Ulrich EL, Vriend G, Vuister GW: **NRG-CiNG: integrated validation reports of remediated experimental biomolecular NMR data and coordinates in wwPDB.** *Nucleic Acids Res* 2012, **40**(Database issue):519-524.

NRG-CiNG contains validation reports by CiNG for all NMR structures in the PDB, but on remediated, that is, corrected for a variety of possible errors or inconsistencies, NMR restraints.

63. Hendrickx PM, Gutmanas A, Kleywegt GJ: **Vivaldi: visualization and validation of biomacromolecular NMR structures from the PDB.** *Proteins* 2013, **81**:583-591.

Vivaldi provides a web-based graphical interface to interactively visualize and analyze NMR structures and the corresponding restraints that is directly embedded in the PDB. Vivaldi also permits the analysis of RDC data.

64. Bagaria A, Jaravine V, Guntert P: **Estimating structure quality trends in the protein data bank by equivalent resolution.** *Comput Biol Chem* 2013, **46C**:8-15.

65. Lindorff-Larsen K, Best RB, DePristo MA, Dobson CM, Vendruscolo M: **Simultaneous determination of protein structure and dynamics.** *Nature* 2005, **433**:128-132.

66. Aramini JM, Petrey D, Lee DY, Janjua H, Xiao R, Acton TB, Everett JK, Montelione GT: **Solution NMR structure of alr2454 from nostoc sp Pcc 7120, the first structural representative of pfam domain family pf11267.** *J Struct Funct Genomics* 2012, **13**:171-176.

67. Moseley HN, Sahota G, Montelione GT: **Assignment validation software suite for the evaluation and presentation of protein resonance assignment data.** *J Biomol NMR* 2004, **28**:341-355.

68. Doreleijers JF, Sousa da Silva AW, Krieger E, Nabuurs SB, Spronk CA, Stevens TJ, Vranken WF, Vriend G, Vuister GW: **CiNG: An integrated residue-based structure validation program suite.** *J Biomol NMR* 2012.

69. Kabsch W, Sander C: **Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features.** *Biopolymers* 1983, **22**:2577-2637.
70. Joosten RP, te Beek TA, Krieger E, Hekkelman ML, Hooft RW, Schneider R, Sander C, Vriend G: **A series of pdb related databases for everyday needs.** *Nucleic Acids Res* 2011, **39(Database issue)**:D411-D419.
71. Vriend G: **What if: a molecular modeling and drug design program.** *J Mol Graph* 1990, **8**:52-56 29.
72. Laskowski RA, MacArthur MW, Moss DS, Thornton JM: **Procheck: a program to check the stereochemical quality of protein structures.** *J Appl Crystallogr* 1993, **26**:283-291.
73. Rieping W, Vranken WF: **Validation of archived chemical shifts through atomic coordinates.** *Proteins* 2010, **78**:2482-2489.
74. Block JN, Zielinski DJ, Chen VB, Davis IW, Vinson EC, Brady R, Richardso JS, Richardson DC: **Kinimurse: macromolecular vr for NMR ensembles.** *Source Code Biol Med* 2009, **4**:3.
75. Wang L, Eghbalnia HR, Bahrami A, Markley JL: **Linear analysis of carbon-13 chemical shift differences and its application to the detection and correction of errors in referencing and spin system identifications.** *J Biomol NMR* 2005, **32**:13-22.
76. Kelley LA, Gardner SP, Sutcliffe MJ: **An automated approach for defining core atoms and domains in an ensemble of NMR-derived protein structures.** *Protein Eng* 1997, **10**:737-741.