



HAL
open science

Restraint Validation of Biomolecular Structures Determined by NMR in the Protein Data Bank

Kumaran Baskaran, Eliza Ploskon, Roberto Tejero, Masashi Yokochi,
Deborah Harrus, Yuhe Liang, Ezra Peisach, Irina Persikova, Theresa Ramelot,
Monica Sekharan, et al.

► **To cite this version:**

Kumaran Baskaran, Eliza Ploskon, Roberto Tejero, Masashi Yokochi, Deborah Harrus, et al.. Restraint Validation of Biomolecular Structures Determined by NMR in the Protein Data Bank. 2023. hal-04594578v1

HAL Id: hal-04594578

<https://hal.science/hal-04594578v1>

Preprint submitted on 27 Oct 2023 (v1), last revised 30 May 2024 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Restraint Validation of Biomolecular Structures Determined by NMR in the Protein Data Bank

Kumaran Baskaran^{1,‡}, Eliza Ploskon^{2,‡}, Roberto Tejero^{3,‡}, Masashi Yokochi^{4,5,‡}, Deborah Harrus⁶, Yuhe Liang⁷, Ezra Peisach⁷, Irina Persikova⁷, Theresa A. Ramelot⁸, Monica Sekharan⁷, James Tolchard⁶, John D. Westbrook^{7,*}, Benjamin Bardiaux⁹, Charles D. Schwieters¹⁰, Ardan Patwardnhan¹¹, Sameer Venankar⁶, Stephen K. Burley^{7,12,13,14,15}, Genji Kurisu^{4,5}, Jeff C. Hoch¹, Gaetano T. Montelione^{8,14,*}, Geerten W. Vuister^{2,*}, Jasmine Y. Young^{7,*}

‡These authors contributed equally

*Deceased

*Corresponding authors: Kumaran Baskaran, Gaetano T. Montelione, Geerten W. Vuister, Jasmine Y. Young

Lead Contact: Kumaran Baskaran

¹Biological Magnetic Resonance Data Bank, Department of Molecular Biology and Biophysics, UConn Health, Farmington, CT 06030, USA

² Department of Molecular and Cell Biology, Leicester Institute of Structural and Chemical Biology, University of Leicester, Leicester LE1 7RH, United Kingdom.

³ Departamento de Química Física, Universidad de Valencia, Dr. Moliner, 50 46100-Burjassot, Valencia, Spain.

⁴ Protein Data Bank Japan, Institute for Protein Research, Osaka University, Suita, Osaka 565-0871, Japan.

⁵ Protein Data Bank Japan, Protein Research Foundation, Minoh, Osaka 562-8686, Japan.

⁶ Protein Data Bank in Europe, EMBL-EBI, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, United Kingdom.

⁷ Research Collaboratory for Structural Bioinformatics Protein Data Bank, Institute for Quantitative Biomedicine, Rutgers, The State University of New Jersey, Piscataway, NJ 08854, USA.

⁸ Dept of Chemistry and Chemical Biology, Center for Biotechnology and Interdisciplinary Sciences, Rensselaer Polytechnic Institute, Troy, New York, 12180 USA.

⁹ Department of Structural Biology and Chemistry, Institut Pasteur, Université Paris Cité, CNRS UMR3528, 75015 Paris, France.

¹⁰ Computational Biomolecular Magnetic Resonance Core, National Institute of Diabetes and Digestive and Kidney Diseases, Bethesda, MD 20892, USA.

¹¹ The Electron Microscopy Data Bank, EMBL-EBI, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, United Kingdom.

¹² Research Collaboratory for Structural Bioinformatics Protein Data Bank, San Diego Supercomputer Center, University of California, La Jolla, California, USA.

¹³ Institute for Quantitative Biomedicine, Rutgers, The State University of New Jersey, Piscataway, NJ 08854, USA.

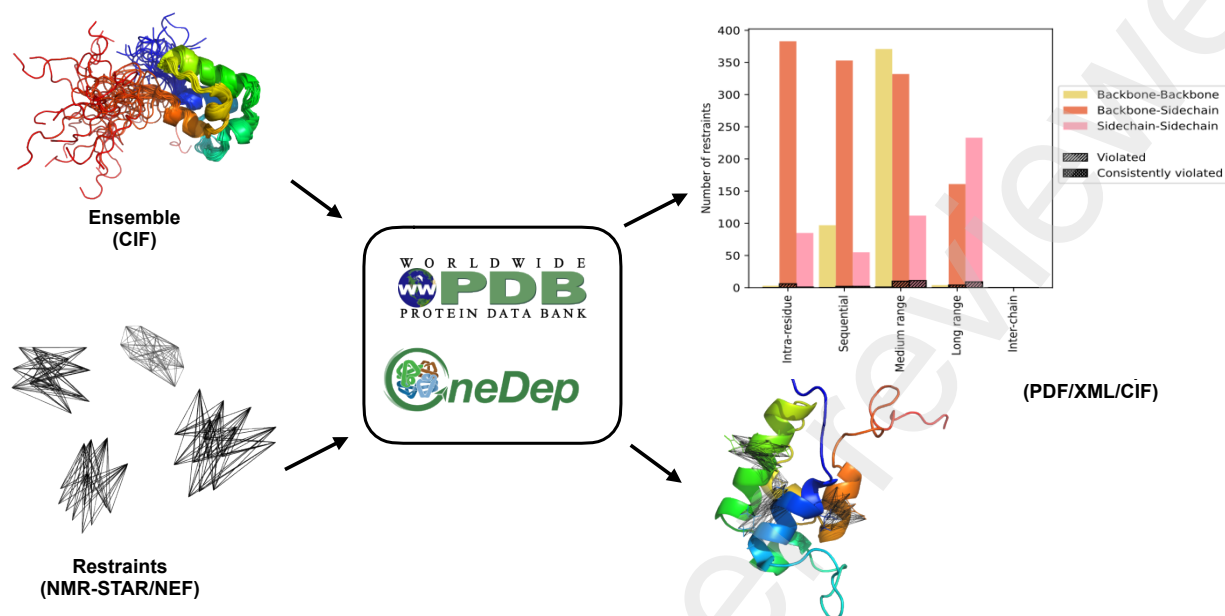
¹⁴ Cancer Institute of New Jersey, Rutgers, The State University of New Jersey, New Brunswick, NJ 08901, USA.

¹⁵ Department of Chemistry and Chemical Biology, Rutgers, The State University of New Jersey, Piscataway, NJ 08854, USA.

Key words: NMR, Validation, Protein Data Bank, biomolecular structures, data standard, distance restraints, NEF

Kumaran Baskaran orcid: 0000-0002-0418-0286 email: baskaran@uchc.edu
Eliza Ploskon orcid: 0009-0006-1863-4345 email: eapa2@leicester.ac.uk
Roberto Tejero orcid: 0000-0003-2504-5988 email: roberto.tejero@uv.es
Masashi Yokochi orcid: 0000-0002-3253-7449 email: yokochi@protein.osaka-u.ac.jp
Deborah Harrus orcid: 0000-0002-7651-672X email: dharrus@ebi.ac.uk
Yuhe Liang orcid: 0000-0002-0574-2041 email: yuhe.liang@rcsb.org
Ezra Peisach orcid: 0000-0002-7905-6327 email: ezra.peisach@rcsb.org
Irina Persikova orcid: 0000-0001-9544-8390 email: irina.persikova@rcsb.org
Theresa A. Ramelot orcid: 0000-0002-0335-1573 email: ramelt2@rpi.edu
Monica Sekharan orcid: 0000-0002-2694-7003 email: monica.sekharan@rcsb.org
James Tolchard orcid: 0000-0002-5779-4935 email: tolchard@ebi.ac.uk
John D. Westbrook orcid: 0000-0002-6686-5475 email: john.westbrook@rcsb.org
Benjamin Bardiaux orcid: 0000-0003-4014-9195 email: benjamin.bardiaux@pasteur.fr
Charles Schwieters orcid: 0000-0002-4216-4658 email: charles.schwieters@nih.gov
Ardan Patwardnhan orcid: 0000-0001-7663-9028 email: ardan@ebi.ac.uk
Sameer Velnankar orcid: 0000-0002-8439-5964 email: sameer@ebi.ac.uk
Stephen K. Burley orcid: 0000-0002-2487-9713 email: stephen.burley@rcsb.org
Genji Kurisu orcid: 0000-0002-5354-0807 email: gkurisu@protein.osaka-u.ac.jp
Jeffrey C. Hoch orcid: 0000-0002-9230-2019 email: hoch@uchc.edu
Gaetano T. Montelione orcid: 0000-0002-9440-3059 email: monteg3@rpi.edu
Geerten W. Vuister orcid: 0000-0001-6172-5097 email: gv29@leicester.ac.uk
Jasmine Y. Young orcid: 0000-0001-8896-6878 email: jasmine.young@rcsb.org

Graphical Abstract



Highlights

- PDB Structure Validation Report expanded to include Restraint Analysis
- NMR Exchange Format (NEF) and NMR-STAR for distance restraint representation
- Standard distance and dihedral restraint formats for model vs. restraint assessment
- Standardized restraint formats provide interoperability between modeling programs

Summary

Biomolecular structure analysis from experimental NMR generally relies on empirical (residual dipolar couplings, scalar couplings, paramagnetic relaxation enhancements, nuclear Overhauser effects) data and derived geometrical restraints (distance, dihedral angle). A challenge for the structural biology community has been a lack of standards for representing these restraints, preventing establishment of uniform methods of model-vs-data structure validation and limiting interoperability between restraint-based structure modeling programs. The NMR exchange (NEF) and NMR-STAR formats provide a standardized approach for representing commonly used NMR restraints. Using these restraint formats, a standardized validation system for assessing structural models of biopolymers against restraints has been developed and implemented in the wwPDB OneDep data harvesting system. The resulting wwPDB Restraint Violation Report provides a model vs. data assessment of biomolecule structures determined using distance and dihedral restraints, with extensions to other restraint types currently being implemented. These tools are useful for assessing NMR models, as well as for assessing biomolecular structure predictions based on distance restraints.

INTRODUCTION

Structure determination using NMR

Nuclear Magnetic Resonance (NMR) spectroscopy is a versatile experimental technique used not only for structure determination but also to probe conformational dynamics and interactions of biomolecules. Solution NMR emerged as a method for atomic-level structure determination in the mid-1980's when the three-dimensional structures of small proteins [e.g., proteinase inhibitor IIA (Williamson et al., 1985), Lac repressor headpiece (Kaptein et al., 1985), alpha1-purothionine (Clare et al., 1986), metallothionine (Wagner et al., 1987b), bovine pancreatic trypsin inhibitor (Wagner et al., 1987a), epidermal growth factor (Cooke et al., 1987; Montelione et al., 1987), and others (Kline et al., 1988)] as well as nucleic acids and protein / nucleic acid complex structures (Boelens et al., 1987; Gochin and James, 1990; Chuprina et al., 1993; Qian et al., 1994; Knegt et al., 1995) were first modeled using distance and dihedral angle restraints derived from Nuclear Overhauser Effect (NOE) and three-bond scalar coupling data. More recently, solution NMR methods providing chemical shift based dihedral angle restraints (Cheung et al., 2010; Shen and Bax, 2010; 2015), relative bond vector orientations determined from residual dipolar coupling (RDC) measurements (Cornilescu et al., 1998; Clare and Garrett, 1999; Losonczi et al., 1999; Lipsitz and Tjandra, 2004; Chen and Tjandra, 2012), solid-state methods using interatomic distance estimates based on dipolar coupling interactions (Lipsitz and Tjandra, 2004; Chen and Tjandra, 2012), and both solution and solid state NMR methods using paramagnetic NMR data (Sengupta et al., 2012; Trindade et al., 2021; Parigi et al., 2022) have also been developed, providing atomic resolution structures of biomolecules. Accordingly, NMR-derived biomolecular structures are generally modeled using primarily estimates of conformationally-averaged interatomic distances and dihedral angles between atoms or groups

of atoms in the form of distance and dihedral angle restraints. These restraints are provided as input to restrained molecular dynamics or other structural modeling programs, which incorporate empirical covalent bond geometry and conformational energy force fields, and output an ensemble of atomic-resolution models, the so-called NMR ensemble, which fit the experimental restraints.

The Protein Data Bank (PDB) is in its 52nd year of continuous operation. Established in 1971 as the first open-access digital data resource in biology (Protein_Data_Bank, 1971), it currently houses ~200,000 experimentally-determined 3D structures of proteins and nucleic acids (DNA and RNA) and their complexes with one another and with small-molecule ligands (*e.g.*, enzyme cofactors, inhibitors, peptides, and drugs). Since 2003, the PDB archive has been jointly managed by the Worldwide Protein Data Bank (wwPDB, <https://wwpdb.org>) partnership (Berman *et al.*, 2003; wwPDB_consortium, 2019). Full members of the wwPDB include three founding members - the US-funded RCSB Protein Data Bank (RCSB PDB) (Berman *et al.*, 2000; Burley *et al.*, 2021), the Protein Data Bank in Europe (PDBe) (Armstrong *et al.*, 2020) and Protein Data Bank Japan (PDBj) (Bekker *et al.*, 2022) - plus the Electron Microscopy Data Bank (EMDB) (Lawson *et al.*, 2016) and the Biological Magnetic Resonance Bank (BMRB) (Markley *et al.*, 2008; Ulrich *et al.*, 2008; Hoch *et al.*, 2023). Protein Data Bank China (PDBc) was recently admitted to the wwPDB as an Associate Member. The wwPDB is committed to managing PDB data for users around the world at no charge for data deposition or egress, nor limitations on data usage. All PDB data are made available by wwPDB partners under the most permissive Creative Commons CC0 license (<https://creativecommons.org/publicdomain/zero/1.0/>). The PDB has been a leader in adopting the emblematic principles of responsible data stewardship in the modern era: FAIR (Findability, Accessibility, Interoperability, and Reusability) (Wilkinson *et al.*, 2016) and FACT (Fairness, Accuracy, Confidentiality, and Transparency) (van der Aalst *et al.*, 2017).

Biomolecular structures provided by solution and solid-state NMR methods continue to make important contributions in biological research for both proteins and nucleic acids. Most recently, these methods have had unique and expanding impact on studies of artificially-designed proteins (Koga *et al.*, 2012; Lin *et al.*, 2015; Jacobs *et al.*, 2016; Anishchenko *et al.*, 2021; Koga *et al.*, 2021), multiple conformational states of biomolecules in dynamic equilibrium (Cicero *et al.*, 1995; Bertini *et al.*, 2004; Harish *et al.*, 2017; Gibbs *et al.*, 2018; Bhardwaj *et al.*, 2022), and for structural analysis of transient conformations that can be characterized using NMR data (Anthis and Clore, 2015; Alderson and Kay, 2020). Objective assessment and validation of these biomolecular NMR structure models remains an important activity of the wwPDB (Berman *et al.*, 2003; Montelione *et al.*, 2013; ww, 2019).

Biomolecular structure validation includes two general classes of assessment, knowledge-based validation, in which the model(s) are assessed in light of what is known about biomolecular structure from the existing database of experimental structures, and model vs. primary data validation, in which consistency is assessed between the structural model(s) and experimental data obtained for the subject biomolecule (Montelione *et al.*, 2013; Rosato *et al.*, 2013). The latter is crucially dependent upon the description of the measured quantities (*e.g.*,

NOEs, RDCs) as a function of the atomic coordinates. Accurate experimental models should score well across the multiple metrics available for these two assessment categories of (Bhattacharya et al., 2007; Rosato *et al.*, 2013). In such assessment methods, it is also important to estimate and consider the uncertainty of the model. This is generally done by comparing models generated from multiple runs of the model generation software in order to identify regions that are consistently modeled (*i.e.*, the “well-defined” regions of the model) and those that are not consistently modeled from the available data and methods (*i.e.*, the “not-well-defined” regions of the model) (Hyberts et al., 1992; Snyder and Montelione, 2005; Kirchner and Güntert, 2011; Montelione *et al.*, 2013; Tejero et al., 2013; Snyder et al., 2014). More rigorously, the precision of the model could be estimated by the propagation of experimental uncertainties using Bayesian methods (Rieping et al., 2005), but so far, this has been done in only a small number of biomolecular structure studies. Consensus recommendations for tools useful for knowledge-based validation and conventions for defining “well-defined regions” for biomolecular structures have been provided by the wwPDB NMR Structure Validation Task Force (Montelione *et al.*, 2013) and implemented in the wwPDB NMR Validation Report (Gore et al., 2017) using standardized knowledge-based validation methods (Tejero *et al.*, 2013; Vuister et al., 2014), with the understanding that NMR structure model validation methods continue to evolve and improve.

While the wwPDB NMR Validation Task Force also provided some recommendations for model-*versus*-data validation, including guidance for reporting how well models fit to experimentally-derived distance restraints, several factors have limited the implementation of these recommendations into the wwPDB NMR Validation Report. In the case of distance restraints analysis, the most challenging issue is that different structure generation software tools utilize NMR-derived distance restraints in different ways and formats. While tools have been developed to convert between some NMR restraint formats (Vranken et al., 2005; Tejero *et al.*, 2013; CCPN, 2023), and some large-scale remediation efforts have been performed at the Biological Magnetic Resonance Data Bank (Nederveen et al., 2005), in some cases it is challenging to represent distance-restraint information used by one structure generation program accurately in the restraint functions of a different program, without the active involvement of the software developers in this process.

These challenges have been addressed, at least in part, through the development of the NMR Exchange Format (NEF, (Gutmanas et al., 2015)), designed to provide reliable interoperability between NMR software programs and structure generation programs in particular. Its design was strengthened by involving software developers in creating and testing all aspects, including the NMR restraint representations. Here, we report an accurate two-way interconversion between the NEF restraint format and the NMR-STAR restraint format, the NMR data archive format of the wwPDB (Berman *et al.*, 2003; ww, 2019). We use the latter to implement a new restraint validation component of the wwPDB report, which builds on distance and dihedral-angle restraint representations in the NMR-STAR format, generating both a human-readable report in PDF format and machine-readable format in CIF (Westbrook et al., 2022). It also provides an XML representation of the data for further computer analysis. These innovations have also been validated against the stand-alone software PDBStat (Tejero *et al.*, 2013) for

restraint format interconversion. It has also been validated against the CcpNmr Analysis version-3 program suite (Skinner et al., 2016), as well against NEF implementations in the NMR structure calculation programs Xplor-NIH (Schwieters et al., 2006) and ARIA (Rieping et al., 2007), with implementation in the program CYANA/CANDID (Guntert and Buchner, 2015) close to completion as well. Together, these tools allow for straightforward generation and validation of experimental NMR-derived structure models against distance and dihedral-angle restraints using restraints in either NEF or NMR-STAR format.

In this paper, we describe the development of a comprehensive NMR model *versus* distance and dihedral-angle restraint data validation report for the wwPDB. Future extensions to other NMR restraint types (e.g., RDCs) can be readily implemented within the same framework. It is anticipated that this model vs. data NMR Restraint Validation Report, together with the already available knowledge-based validation report, will provide a more comprehensive and objective assessment of the reliability of biomolecular structures determined by NMR methods.

RESULTS

Once the resonance assignment process is sufficiently complete, distance restraints can be extracted from various NOE Spectroscopy (NOESY) and paramagnetic NMR (e.g., paramagnetic relaxation enhancement PRE) experiments. The primary observable in any NMR experiment is the chemical shift. Other types of information are derived by interpreting the intensity or the volume of the resonance peaks in the spectra. Not every resonance in NMR is well resolved or can be unambiguously assigned to a specific atom in the molecule under investigation. For example, the three protons of a methyl group give rise to a single resonance, which results in a single restraint involving all three atoms. Methylene protons, certain aromatic resonances, and side-chain amide protons of glutamine and asparagine may also lead to ambiguous assignment if they are not assigned stereo-specifically or individually. The ambiguity may sometimes be resolved by the overwhelming evidence resulting from the first round of structural modeling allowing the ambiguity to be determined (Guntert et al., 1991; Herrmann et al., 2002; Huang et al., 2006; Guntert and Buchner, 2015). Conformational exchange may also lead to multiple peaks, which can only be assigned unambiguously with the application of additional NMR experiments.

Restraint validation

Checking the validity of a given restraint between two atoms on a given model is not as trivial as it might appear, as several complicating aspects have to be taken into account for a proper analysis. Typically (although not in all programs), the distance restraint is modeled using a cost function corresponding to a square well potential and defined by the lower limit and upper limit of the distance between the two atoms, and if the measured distance in the model falls between these bounds, then the restraint is not violated, otherwise, it is violated. For methods that do not use bounds (e.g., log-harmonic potential as implemented in ARIA (Nilges et al., 2008)), the respective software is expected to specify how the strength of potentials needs to be translated into equivalent lower and upper bounds.

The complications for a meaningful restraint violation analysis arise from several factors. First, the overlap and ambiguity, as described above, need to be accounted for. Often, a so-called “ r^{-6} sum” over all the distances contributing to the restraint can be used (Hyberts *et al.*, 1992; Nilges, 1995; Bassolino-Klimas *et al.*, 1996). Second, all the restraints derived from the NMR experiments result from the spatial and temporal averages of the underlying dynamics exhibited by the biomolecule, and not all restraints may be fully satisfied at all instances of time. Molecular dynamics simulations, such as those employed to calculate structures on the basis of NMR data, ensure that on average, all restraints are satisfied for a maximum amount of time by minimizing the total energy of the system. A small fraction of restraints may be violated at every instant, but the set of restraints violated in each instant is different so that no restraint is consistently violated. Third, NMR structure calculations typically result in an ensemble comprised of multiple conformers. This conformational diversity results from both the experimental uncertainty and potentially to actual conformational variability in the sample. Unless the data are explicitly fit to multiple conformations or a conformational ensemble, the NMR-VTF has recommended that the NMR ensemble should be analyzed in terms of either well-defined or not-well-defined regions (Montelione *et al.*, 2013), reflecting the fact, by definition, that the biomolecular conformation in not-well-defined regions is not reliably modeled. Unusual dihedral angle values and steric clashes in the latter regions are not considered to be significant. Additionally, the extent that dynamical information can be faithfully represented in a fixed-size ensemble is an unanswered question.

Highly similar considerations apply to the validation of all types of NMR-derived restraints. Whereas assignment ambiguity is generally not pertinent for dihedral-angle restraints, conformational averaging is a similarly complicating factor. RDC restraints are also conformationally averaged, albeit on different timescales compared to the NOE-derived distance restraint. Additionally, analysis of RDC satisfaction requires establishing the alignment tensor (Losonczi *et al.*, 1999), with its own associated uncertainties.

Types of restraints and their validation

The distance restraint between two atoms is a derived quantity, sometimes originating from more than one spectrum. In a conventional NMR structure determination workflow, the chemical shift assignments are first derived from one or more (heteronuclear) through-bond type experiments, which are then used to assign the peaks of through-space NOESY-type spectra. The volumes or intensities of these assigned NOESY peaks are then used to estimate the distance between pairs of atoms. It is not uncommon to have identical chemical shifts within the spectral resolution for more than one atom. As mentioned earlier, the three protons of the methyl group usually have degenerate chemical shifts, and chemical shifts of different methyl groups from the same or different residues may also overlap. Such overlapping chemical shifts lead to ambiguous chemical shift assignment, which results in ambiguous restraints, in addition to the more easily interpretable unambiguous restraints. Below, we discuss the typical cases a distance restraint validation procedure needs to accommodate.

Type 1: Unambiguous distance restraints.

Unambiguous distance restraints are derived from well-resolved NOESY peaks between two atoms. The chemical shifts of the atoms involved in this type of restraint are non-degenerate and unambiguously assigned. Let r_{ij} be the distance between atom i and j in the molecular model (Figure 1(a)), $d_{min}(i,j)$ is the lower bound and $d_{max}(i,j)$ is the upper bound of the distance restraint. If $d_{min}(i,j) \leq r_{ij} \leq d_{max}(i,j)$ then the restraint is not violated, otherwise the violation is calculated as the lowest of the $|r_{ij} - d_{min}(i,j)|$ or $|r_{ij} - d_{max}(i,j)|$.

Type 2 Ambiguous restraints involving resonances with degenerate chemical shifts

Degenerate chemical shifts (e.g., those of magnetically equivalent methyl protons) give rise to ambiguous distance restraints. The NEF standard provides for the wild-card “%” identifier, e.g., HB%, to allow for ambiguous restraints involving such degenerate chemical shifts. To a first approximation, NOESY peak between the degenerate resonance, such as the methyl group, and another atom will have NOE contributions from each of the contributing protons [Figure 1(b) & (c)] which varies inversely to the distance. As an approximation, an effective distance (r_{eff}) can be calculated using the r^6 sum of the pairwise distances between the atoms contributing to the NOE peak (Nilges, 1995), and the potential violation of the restraint is assessed using the resulting r_{eff} . If $d_{min}(i,j) \leq r_{eff} \leq d_{max}(i,j)$ then the restraint is not violated, otherwise the violation can be calculated as $\min(|r_{eff} - d_{min}(i,j)|, |r_{eff} - d_{max}(i,j)|)$, where r_{eff} is given by

$$r_{eff} = (\sum r_{ij}^{-6})^{-\frac{1}{6}} \quad \text{Eqn 1.}$$

The r^6 sum distance restraint, r_{eff} , has the feature that it is dominated by the shortest distance to the set of ambiguously assignable (degenerate) atoms.

Type 3: Restraints between atoms involving resonances of stereospecifically- or individually-assignable atoms

Prochiral methylene protons, individually-assignable amide NH_2 groups (of for example, asparagine and glutamine), and aromatic ring protons (for example from phenylalanine and tyrosine side chains), may or may not have degenerate chemical shifts. If they are non-degenerate, they can potentially be stereo-specifically or individually assigned. Isopropyl methyl groups (of, for example valine and leucine) are also prochiral, and generally require stereo-specific assignments. However, unless individual or stereospecific assignments are established specific experimental or computational methods, these groups of resonances are also treated using ambiguous restraints. This ambiguity can be addressed by effectively treating the two separate resonances as though they are degenerate, and summing the volumes (or intensities) to create an ambiguous r^6 sum restraint (Nilges, 1995; Tejero *et al.*, 2013). NEF includes a robust standard to handle these situations, linking the NMR resonance assignment intimately with the calculated structures in the NMR ensemble (see Methods).

Case 3.1: Degenerate chemical shifts (HB% case, r^{-6} sum)

If a group of protons, such as those of methylene groups, have degenerate chemical shifts then they can be treated like Type 2 restraints described above, and the r^{-6} sum method (Eqn. 1) is used to define the restraint (Figure 1(d))

Case 3.2: Nondegenerate and stereo-specifically assigned

If the chemical shifts of a group of protons are non-degenerate and if they are stereo-specifically (or individually) assigned, then these restraints are treated as either a Type 1 or Type 2 restraint depending on whether the other atom is a single atom or group of atoms, respectively (Figure 1(e)).

Case 3.3: Nondegenerate and ambiguously assigned

If the chemical shifts of a group of protons are non-degenerate and if they are not stereo-specifically assigned, the situation is much more complex as the structure calculation algorithms employ different approaches in dealing with this issue. Thus, crucial information needs to be captured and adequately handled in the restraint validation protocols, which traditionally has presented a serious problem. The issue is best illustrated with an example. Let us assume that the chemical shifts of a group of protons, e.g., two methylene protons, are non-degenerate but cannot be assigned stereo-specifically. Following the NEF standard, these protons should be labeled with the “x” and “y” identifiers i.e., HBx and HBy for a pair of methylene protons attached to CB. Assume that a set of NOESY peak derived distance restraints were observed for HBx to atoms B1, B2, B3, B4, C1 and C2 and another set of restraints observed for HBy to atoms A1, A2, A3, A4, C1 and C2 (Figure 1(f)). The example indicates that in addition to the set of common restraints, i.e., to C1 and C2, there are also restraints exclusive to either HBx or HBy. It is a-priori undefined whether the stereospecific HB2 atom in the molecular structure maps onto HBx, and the HB3 atom to HBy, or vice versa. In this case, the different sets of NOEs to the two distinct resonances suggest that this issue could potentially be solved as part of the structure calculation protocol (Guntert *et al.*, 1991).

One straightforward approach in dealing with this issue is to collapse the two sets of restraints to HBx and HBy into one set involving an ambiguous HB%, with some form of treatment of the restraint limits, e.g., on the basis of the originating peak intensities or by taking either the shortest (most restricting) or longest (least restricting) limit. Simply put, the two stereospecifically distinct resonances are treated as degenerate and a r^{-6} sum restraint (Eqn. 1) is created (Nilges, 1995; Tejero *et al.*, 2013). In this approach, a Case 3.3 HBx/HBy restraint has been converted to a Case 3.1 ambiguous restraint.

A second approach, commonly also implemented by structure calculation programs like ARIA, Xplor-NIH or Cyana, is the concept of “floating chirality” (Folmer *et al.*, 1997). In the course of the structure calculation, the program adopts the most favorable mapping for all pairs at any

time, thus minimizing the resulting restraint energy. At some point during the calculation, typically in light of sufficient consistency, the mapping can be fixed. However, such structure-based stereospecific assignments cannot usually be obtained for all pairs of prochiral or individually-assignable atoms. In the past, structure calculation programs like ARIA/CYANA (Brünger et al., 1998) have provided such mapping information for subsets of restraints in the form of the so-called “float-files” or “stereo.aco” files, respectively. Unfortunately, this information has mostly been lost during the deposition process using data supplied in a program-specific format and this approach leads to inconsistency between atoms defined in the restraints and the model files. Extensive efforts to re-capture these mappings required great efforts and were only partially successful (Doreleijers et al., 2012). Chemical shift predication from model structures could offer an path to resolving ambiguous stere-specific assignments (Weiss and Hoch, 1987).

By documenting the actual restraints used in the structure calculation in NEF format and reporting structural data in PDBx/mmCIF format, this issue is solved. Using NEF, the stereospecific mapping, *i.e.*, HB2 onto HBx and HB3 on to HBy or vice versa, will be documented for each atomic position using the ambiguity tag, [atom site.pdbx atom ambiguity](#), and the restraint is validated accordingly (Figure 1(g)(h)). Note that this mapping could differ for each model of the structural ensemble. In the absence of such a documented mapping using a NEF/mmCIF pair, and in order to avoid the introduction of an erroneous restraint validation assessment, the restraint validation algorithm interprets both HBx and HBy as HB%, effectively treating them as Case 3.1 degenerate restraints, as described above. This procedure both assures that violations are not wrongly reported, and appropriate restraint and restraint-violation counts are maintained.

The usage of the NEF/mmCIF pair and ambiguity tag has the added advantage of solving a long-standing problem involving restraints to slowly rotating phenylalanine or tyrosine aromatic side chain protons, *i.e.*, in case of non-degenerate HD1/HD2 and HE1/HE2 chemical shifts. Structurally, the HD1/HD2 and HE1/HE2 designation, as well as the CD1/CD2 and CE1/CE2 designation, is determined by a the dihedral χ_1 angle and a small rotation beyond 180° will structurally swap CD1/HD1/CE1/HE1 with CD2/HD2/CE2/HE2. This, however, could have detrimental consequences for restraints formulated in terms of HD1/HD2/HE1/HE2 as large errors would be introduced by such a swap. Fortunately, this can be easily resolved using the HDx/HDy/HEx/HEy NEF nomenclature and ambiguity tag mapping in the mmCIF, as the mapping effectively provides the correct atomic coordinates to be used for the restraint validation. Note that restraints formulated in terms of the wild-card HD% or HE% atoms are always evaluated correctly.

Resolution of Case 3 non-degenerate ambiguously assigned restraints depends on the generation of consistent pairs of NEF / mmCIF NMR and structural data. Given the involvement of the software development community in the NEF project, we are confident that the common NMR structure calculation programs are poised to create such pairs, thereby assuring the correct data interpretation and a correct restraint validation.

Dihedral Angle Restraints

A dihedral angle is defined as the angle between half-planes defined by two sets of three atoms, having two atoms in common. Customarily, the common atoms are bonded, and each bonded to the other defining atoms. In proteins, the backbone dihedral angles ϕ and Ψ , defined by the backbone atoms C-N-CA-C and N-CA-C-N, respectively, provide crucial information to describe the main-chain geometry. The Ramachandran plot, which visualizes the distribution of these two backbone dihedral angles on a ϕ - Ψ graph, is widely used by various structure validation program suites (Laskowski et al., 1993; Lovell et al., 2003; Chen et al., 2010) to assess the quality of the structure.

For proteins, dihedral-angle restraints can be derived using the backbone chemical shift data (Cheung *et al.*, 2010; Shen and Bax, 2010; 2015). In the NEF format, any dihedral-angle restraint (z) is defined using its four relevant atoms, a target, upper- and lower-bound values. The sign of these angles indicates whether the angle is measured counterclockwise or clockwise. The convenient use of positive and negative signs in representing angles makes the upper and lower bounds an arbitrary choice. Without the target value, it is hard to tell which side of the angular region between upper and lower bound is the allowed region. Figure 2 illustrates how different choices of target values render either the acute or the obtuse angle as allowed regions.

For example, let $\phi_{measured}$ be the measured backbone dihedral angle in a given model. If both ϕ_{target} and $\phi_{measured}$ are in the angular region between ϕ_{min} and ϕ_{max} then the dihedral restraint is not violated, otherwise, it is violated, and the violation $\phi_{violation}$ is calculated as

$$\phi_{violation} = \min(|\phi_{min} - \phi_{measured}|, |\phi_{max} - \phi_{measured}|)$$

Dihedral angle restraints may also be ambiguous; *i.e.*, they may define multiple, discontinuous regions of the $\phi - \psi$ map. The target value could possibly also be assigned to more than one value, in order to define multiple conformations that are consistent with the data. Ambiguous dihedral restraints are defined in NEF as a set of restraints using the combination of `_nef_dihedral_restraint.restraint_id` and `_nef_dihedral_restraint.restraint_combination_id`. The ambiguous restraint is considered violated only if all of the possible restraints in the set are violated.

wwPDB OneDep deposition and validation

The global wwPDB OneDep tool (Young et al., 2017) supports deposition, validation (Gore *et al.*, 2017; Feng et al., 2021), and biocuration (Young et al., 2018) for macromolecular structures determined by macromolecular crystallography (MX), 3D electron microscopy (3DEM), and NMR since its launch in 2014. Recently the OneDep has been enhanced to further support NMR data generated by community software in either NEF (V1.1) (Gutmanas *et al.*, 2015) or NMR-STAR (V3.2) (Ulrich et al., 2019) format. The OneDep deposition interface allows authors to upload a single combined NMR data file that includes required chemical shift and restraint

data and optional peak list data in either NEF or NMR-STAR format while (presently) continuing to support native file formats from community software (e.g., CYANA, CNS, and Xplor-NIH). The latter option, however, will be phased out in the near future, in consultation with the NMR community, as the many problems associated with reliable interpretation are now addressed by using the NEF / CIF pair of NMR and structural data (vide infra). We encourage software developers to generate biomolecule structure deposition data in NMR-STAR / NEF formats for NMR data and PDBx / mmCIF format for atomic coordinate files for proper validation.

Assigned chemical shifts and the experimental restraint data are mandatory for deposition to the PDB of a biomolecular structure solved using NMR spectroscopy. Depositors are also highly encouraged to provide NOESY peak lists, as well as other relevant NMR information, such as RDC data, as part of their NEF file. The validation workflow converts the uploaded NEF data into NMR-STAR format, the archival format of NMR data in the PDB and BMRB Core Archives. The validation report is generated using the NMR data in NMR-STAR format and coordinate data in PDBx/mmCIF format.

The NMR-STAR/NEF NMR data file should contain the following mandatory data, encoded as so-called blocks/save frames in accordance with the respective defined data formats:

1. Sequence information
2. Assigned chemical shift data
3. Restraint data (various types)

Depositors are encouraged to also provide as much metadata as possible in their NMR-STAR/NEF file using the tags defined by their respective data dictionaries. However, some necessary metadata will be collected through the wwPDB deposition user interface and added to the NMR-STAR or NEF data file. Upon file upload, OneDep provides the following diagnostics:

1. Identifies the file type as an NMR unified data
2. Validates the NMR data, including checking NMR data content and providing data diagnostics such as identifying unusual chemical shift data values i.e, chemical shift values outside of the expected range. The various checks provide warnings for depositors to review (Table 1)
3. Cross-checks the sequence between the atomic coordinate and the chemical shift files, and provides a sequence alignment for depositors to review

Once the uploaded NMR data file has passed the file check, some metadata are automatically parsed at the deposition interface for authors to review, and is also captured in the atomic coordinate file to reference the corresponding NMR data. If a unified NMR data file is uploaded, it is then passed to the wwPDB validation package for generation of the wwPDB validation report (Gore *et al.*, 2017; Feng *et al.*, 2021), which includes restraint validation as outlined here. The resulting preliminary validation report is provided at the deposition interface for the author's review and correction of their data as needed. Subsequently, the official wwPDB Validation Report is generated by the wwPDB biocurators after data processing and is sent back to authors for the journal manuscript review process. The wwPDB Validation Reports are provided

in PDF, PDBx/mmCIF, and XML formats. At the time of release for PDB entries, the wwPDB validation reports are generated for public distribution at https://ftp.wwpdb.org/pub/pdb/validation_reports and the unified NMR data are made available at https://ftp.wwpdb.org/pub/pdb/data/structures/divided/nmr_data/ in both NEF and NMR-STAR formats.

Restraint violation analysis in wwPDB Validation Reports

The distance and dihedral-angle restraint analysis is presented in the wwPDB Validation Report in sections 8, 9 and 10. Section 8 describes the overall summary of the deposited restraints of all categories and restraint violations in different bins. For the distance restraints, a grouped data classification is both widespread in the NMR community and was recommended by the wwPDB NMR Validation Task Force (Montelione *et al.*, 2013). It provides a simple and convenient overview of the available data and their agreement with the molecular structure;

The conformationally restricting restraints are counted in categories of intra-residue, sequential, medium range ($|i-j| > 1$ and $|i-j| < 5$), long range ($|i-j| \geq 5$), and inter-chain restraints and listed (Table 2) along with the number of hydrogen bonds and disulfide bond restraints. Table 2 shows a summary of restraints data in the PDB ID 7m5t (Anishchenko *et al.*, 2021) as an example. Duplicate and redundant restraints are excluded from the statistics. Distance restraint values that do not restrict the conformations of the intervening dihedral angles are defined and these non-conformationally-restricting restraints are excluded from the restraint validation analysis. If an atom involved in a restraint has no corresponding atom in the coordinate file, it is counted as an unmapped restraint.

All conformationally-restricting restraints are validated against structural results from each model in the NMR ensemble. If the measured distance between a pair of atoms in a given model lies between the upper and the lower bound of the corresponding distance restraint as described above, then the restraint is not violated. If the measured distance in a model lies outside the boundaries defined by the restraint, then the absolute difference between the measured value and the nearest boundary is reported as the violation value. The results are reported, binned into small, medium, and large violation categories based on the magnitude of the violation values. In each bin, the average number of violations per model is calculated by dividing the total number of violations in each bin by the size of the ensemble. The maximum value of the violation in each bin is also reported. Table 3 lists distance violations per bin in the PDB ID 7m5t (Anishchenko *et al.*, 2021) as an example. If dihedral-angle restraints were included, similar overall and violation statistics are also provided for these. Violations less than 0.1 Å for distance restraints and less than 1° for angle restraints, which may have come from round-off error, are excluded from the statistics.

Sections 9 and 10 in the Validation Report provide a detailed analysis of distance and dihedral-angle restraints. Both sections have similar subsections and contents for their respective restraints categories and hence are discussed here together. Sections 9.1 and 10.1 describe the summary of violations in different restraint categories. For each category, the table provides

the total number of restraints, the percentage with respect to the total, the number of violated restraints, and the percentage with respect to that particular category and with respect to the total number of restraints. Restraints that are violated in at least one model are counted as violated and restraints that are violated in all the models are counted as consistently violated. The information in the table is also provided as a bar chart which gives a straightforward overview of any consistent violations. The example for PDB ID 7m5t (Figure 3) shows a typical pattern, with only a few violated restraints and no consistently violated ones.

The next subsections in the report, sections 9.2 and 10.2, provide the violation statistics for each model. The number of violations in each model and the mean, median, standard deviation and maximum values are listed in a table and are also presented as a bar chart in the report. Figure 4 shows the per-model bar chart for PDB ID 7m5t (Koepnick et al., 2019) The total number of violations for each model (~6) is low. The distribution of violations as indicated by the blue bars and indicators for median and mean (~0.17 Å) distance restraint violation is also low and highly similar for all models, suggesting that no model represents an outlier. Models 13, 14, and 20, however, appear slightly better for the agreement of local conformations with the experimental data as these models show no violations in the intra-residue and sequential categories.

The distance and dihedral angle violation statistics for the ensemble are presented in sections 9.3 and 10.3 of the report respectively. The table in these sections lists the number of violations for a given fraction of the ensemble. The number of restraints violated in all models, *i.e.*, the consistently violated restraints, are also listed. The bar chart (Figure 5) shows the violation statistics for the ensemble of PDB ID 7m5t. The figure shows that most restraints are only violated in 5% of the models or fewer, suggesting the absence of systematic violations. Together with the small magnitude of the observed violations (Figures 3-4), this indicates good agreement of these experimental data with the structural models.

Histograms of each restraint's average violation and the most violated restraints for the ensemble are given in sections 9.4 and 10.4 of the validation report, for distance and dihedral-angle restraints, respectively. Similarly, sections 9.5 and 10.5 of the report provides lists of all distance and dihedral angle violations in each model in the ensemble. It also provides the histogram of the magnitude of these violations (not shown).

DISCUSSION

In this paper, we describe a comprehensive set of new biomolecule distance- and dihedral-angle restraint validation tools for the wwPDB OneDep validation pipeline, generating both human and computer readable reports. Although examples provided herein pertain exclusively to proteins, the same restraint validation methods can be used for distance restraint validation of other biomolecules modeled from NMR-based restraints, including nucleic acids and carbohydrates.

As most biomolecular NMR structures are determined from distance (including both NOE and PRE) and dihedral angle restraints, it is necessary to provide a standardized validation of models against these restraints. The current implementation of wwPDB NMR Validation Software provides these tools, using upper and lower bound distance restraints, dihedral angle restraints, and chemical shift data in either NMR-STAR ver3 or NEF v1.1 data formats. In the course of this work, these features of the wwPDB NMR Validation Software have also been incorporated into the C/C++ program PDBStat ver5.23

(https://github.rpi.edu/RPIBioinformatics/PDBStat_public) (Tejero *et al.*, 2013), allowing cross-validation between PDBStat and wwPDB restraint conversation and of the implementation of rules and processes outlined in this paper. These tools are also useful in providing a standardized restraint validation protocol that can be applied across the PDB archive and used in providing standardized structure validation reports to support publication of NMR-derived biomolecular models. In conjunction with this validation software, the CcpNmr Analysis version-3 program suite (Skinner *et al.*, 2016) has also been made fully compatible with generating the required NEF input data and accepting the output of the validation pipeline for further inspection and analysis. In addition, programs for NMR structure generation, such as Xplor-NIH (Schwieters *et al.*, 2006), ARIA (Rieping *et al.*, 2007), and others under development, have been updated to accept both NEF and NMR-STAR formatted input data, as well as generating a consistent pair of NEF- and mmCIF formatted result files for data exchange and for OneDep deposition.

Results generated using the OneDep validation pipeline, both for structure-based validation and restraint validation, are available as PDF files for human inspection and interpretation, as well as in XML and PDBx/mmCIF formats for further processing by other software programs. For example, we used CcpNmr AnalysisStructure to import and process the XML-file for PDB entries 2png and 1pqx. We used the restraint violations to generate a per-model / per-residue metric and color-coded the structural ensembles. Fig. 6a shows the result for PDB ID 2png (<http://doi.org/10.2210/pdb2PNG/pdb>), essentially confirming the notion of very few and only incidental violations, as also evident from the data presented in Fig. 3-5. In contrast, PDB ID 1pqx (<http://doi.org/10.2210/pdb1PQX/pdb>) clearly shows localized hot-spots of substantial restraint violations, warranting further inspection (Fig. 6b). Fortunately, the NEF format also provides so-called linkage information, *i.e.*, between restraints and originating peaks, thus allowing for re-examination of the originating spectral data, *e.g.*, in CcpNmr AnalysisStructure.

The wwPDB is committed to improving data quality by making validation reports available to the public. The best effort to standardize existing restraints and chemical shifts into single NEF and NMR-Star formats is underway. This remediation effort will include PDB archive-wide re-generation of wwPDB validation reports with restraint validation which will enable archival statistical assessments for outlier detection in the near future.

Although the wwPDB restraint validation software is a significant advance, other valuable model-vs-data tools are also available in the NMR spectroscopists tool chest. It is also possible to validate models against a NOE completeness metric, assessing the percentage of restraints predicted by the model that are included in the restraint list (Doreleijers *et al.*, 1999). As

geometrical restraints are derived empirical NMR data, such as spectra or peak lists, various tools have also been described for validation of models against NOESY peak list and chemical shift data (Huang et al., 2005; Huang et al., 2012), or even directly against spectra (Thomas et al., 1991; Gorler and Kalbitzer, 1997; Ried et al., 2004). Models can also be validated by comparing metrics of flexibility based on chemical shift with models of flexibility derived from the structure models (Fowler et al., 2020; Fowler and Williamson, 2022), back calculation of chemical shifts from molecular models (Neal et al., 2003; Vila et al., 2008; Shen and Bax, 2010), or by back calculation of residual dipolar coupling data from models (Cornilescu *et al.*, 1998; Clore and Garrett, 1999; Losonczi *et al.*, 1999). Each of these methods has strengths and weaknesses (Rosato *et al.*, 2013). Several of these model-vs-data methods, including the NOE completeness score (Doreleijers *et al.*, 1999), the RPF-DP score for assessing models against NOESY peak lists (Huang *et al.*, 2005; Huang *et al.*, 2012), and the RDC Q factor (Cornilescu *et al.*, 1998; Clore and Garrett, 1999; Losonczi *et al.*, 1999) are available as servers and are also implemented in the software package PDBStat ver 5.xx (Tejero *et al.*, 2013). While use of one or more of these model-vs-data structure quality assessment methods is strongly recommended for depositors of NMR-based structural models to the wwPDB, these model-vs-data validation methods are not yet adopted by the wwPDB because there is not yet sufficient community consensus on their general applicability in the context of a global biomolecular structure archive.

Another important area of methods development involves representation of multiple conformational states of proteins within a single PDB entry. NMR structures are generally represented by a collection of models, representing the consistency and uncertainty of atomic positions in the structural model. Each of these models is consistent with all of the available NMR data. However, in some cases the NMR data should be more accurately interpreted in terms of multiple biomolecule conformations in dynamic equilibrium, *i.e.*, two or more conformations present in the same sample. In the past, these multiple conformational states have been represented in various ways in PDB depositions. Future expansions of the wwPDB NMR Structure Validation software will need to account for multiple chemical shift data, multiple restraint data, and multiple atomic coordinate sets that result from multiple conformational state modeling. Fortunately, NEF and NMR-STAR are inherently flexible and extensible to allow for an implementation of a standard for these situations.

While the wwPDB sites currently accept NMR-derived structures in legacy PDB format, since June 30, 2019, macromolecular crystallographic structures are accepted only in the PDB exchange macromolecular Crystallographic Information File (PDBx/mmCIF) format (Adams et al., 2019). The OneDep NMR Structure Validation software also uses PDBx/mmCIF as input format for atomic coordinates and either NMR-STAR or NEF formats for NMR restraint files. As the requirement for providing atomic coordinates for NMR-derived structures in PDBx/mmCIF and NMR data in NMR-STAR / NEF format is anticipated in the near future, it is important that the community begin the process of adopting these formats and conventions. Both NMR-STAR and NEF also support NOESY peak list and RDC data formats, anticipating support of these data types into the validation package in the future. While validation against NOESY peak lists and RDC data are anticipated for future expansions of the wwPDB NMR Structure Validation

software, additional consensus of the broader NMR community will be needed before standardizing these validation metrics.

The development and implementation of NEF is actively supported by NMR software developers. A recent round-robin NEF testing exercise, which included the wwPDB consortium implementing the current validation pipeline, yielded important knowledge on practical implementation details. A more detailed account of this exercise, including a detailed description of the NEF data format, will be presented elsewhere.

In conclusion, we presented the rationale for model-vs-data restraint validation by the wwPDB, together with a summary of validation tools for NMR distance and dihedral restraints, as implemented in the wwPDB validation pipeline and recommended by the wwPDB NMR-VTF committee (Montelione *et al.*, 2013). These tools will allow for a more comprehensive and therefore better assessment of the quality of biomolecular NMR structures and thereby benefit all users of the PDB biomolecular structure archive.

Acknowledgments

We thank structural biologists worldwide who have contributed structures to the PDB and members of the wwPDB NMR Validation Task Force and NEF Working Group for setting data standards and validation recommendations. We also thank all the staff members of the wwPDB partners for their support and feedback and wwPDB biocurators for testing and feedback on the wwPDB validation software. RCSB PDB is funded by the National Science Foundation (DBI-1832184), the U.S. Department of Energy (DE-SC0019749), and the National Cancer Institute, National Institute of Allergy and Infectious Diseases, and the National Institute of General Medical Sciences of the NIH under grant R01GM133198. Protein Data Bank in Europe is supported by the European Molecular Biology Laboratory-European Bioinformatics Institute and the Wellcome Trust (104948). Protein Data Bank Japan is supported by the Database Integration Coordination Program (JPMJND2205) from the department of National Bioscience Database Center (NBDC)-JST (Japan Science and Technology Agency), the Platform Project for Supporting in Drug Discovery and Life Science Research from AMED (22ama121001), and the joint usage program of Institute for Protein Research, Osaka University. BMRB is supported by the US National Institute of General Medical Sciences under grant R01GM109046. R.T., T.A.R., and G.T.M. are supported by US National Institute of General Medical Sciences 1R35-GM141818. E.P., G.W.V. and co-workers are supported by UKRI-MRC partnership grants MR/L000555/1 and MR/P00038X/1. RCSB PDB core operations are jointly funded by the National Science Foundation (DBI-1832184, PI: S.K.B.), the US Department of Energy (DE-SC0019749, PI: S.K.B.), and the National Cancer Institute, the National Institute of Allergy and Infectious Diseases, and the National Institute of General Medical Sciences of the National Institutes of Health (R01GM133198, PI: S.K. Burley). Other funding awards to RCSB PDB by the NSF and to PDBe by the UK Biotechnology and Biological Research Council are jointly supporting development of a Next Generation PDB archive (DBI-2019297, PI: S.K.B.; BB/V004247/1, PI: S.V.) and new Mol* features (DBI-2129634, PI: S.K.B.; BB/W017970/1, PI:

S.V.). CDS is supported by the Intramural Program of the National Institute of Diabetes and Digestive and Kidney Diseases. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

Author contributions

All authors contributed to the development of the restraint validation standards outlined in this manuscript. The NEF data standards and dictionary were developed by E.A.P., G.W.V., J.W. E.P. maintained the contacts with the NEF working group developers. The corresponding PDBx/mmCIF dictionary was developed by J.W. The wwPDB restraints validation software was developed by K.B. and tested against the PDBStat software by R.T. and E.A.P. The deposition and data checking of NEF and NMR-Star data were developed by M.Y. The testing and feedback of wwPDB validation software and NEF deposition within OneDep were performed by M.C., D.H., I.P., Y.L., M.S., E.P. and J.T. wwPDB project management on the development of validation software and NEF data deposition was provided by J.Y.Y. experimental data for structure validation were provided by T.A.R. The wwPDB validation software package is maintained by wwPDB partners headed by S.K.B., S.V., G.K., A.P. and J.C.H. G.T.M. and G.W.V. provided the scientific steering to the project. The article was written by K.B., G.T.M., G.W.V., and J.Y.Y., with contributions from all authors.

Declaration of interests

The authors declare no competing interests. G.T.M. is a founder of Nexomics Biosciences Inc., which though not a conflict of interest with respect to this work is a required disclosure.

Figures and Legends

Fig. 1. (a) Type1: Distance restraint between atoms i and j , (b) Type 2: Distance restraint between an atom and a group of atoms, (c) Type 2: Distance restraint between two groups of atoms, (d) Type 3: Distance restraint between non-stereo specifically assigned atoms with degenerate chemical shifts and group of atoms, (e) Type3: Distance restraint between stereo specifically assigned atoms with non-degenerate chemical shifts and group of atoms, (f) Type3: Distance restraint between non-stereo specifically assigned atoms with non-degenerate chemical shifts and group of atoms, (g) & (h) possible assignments for (f).

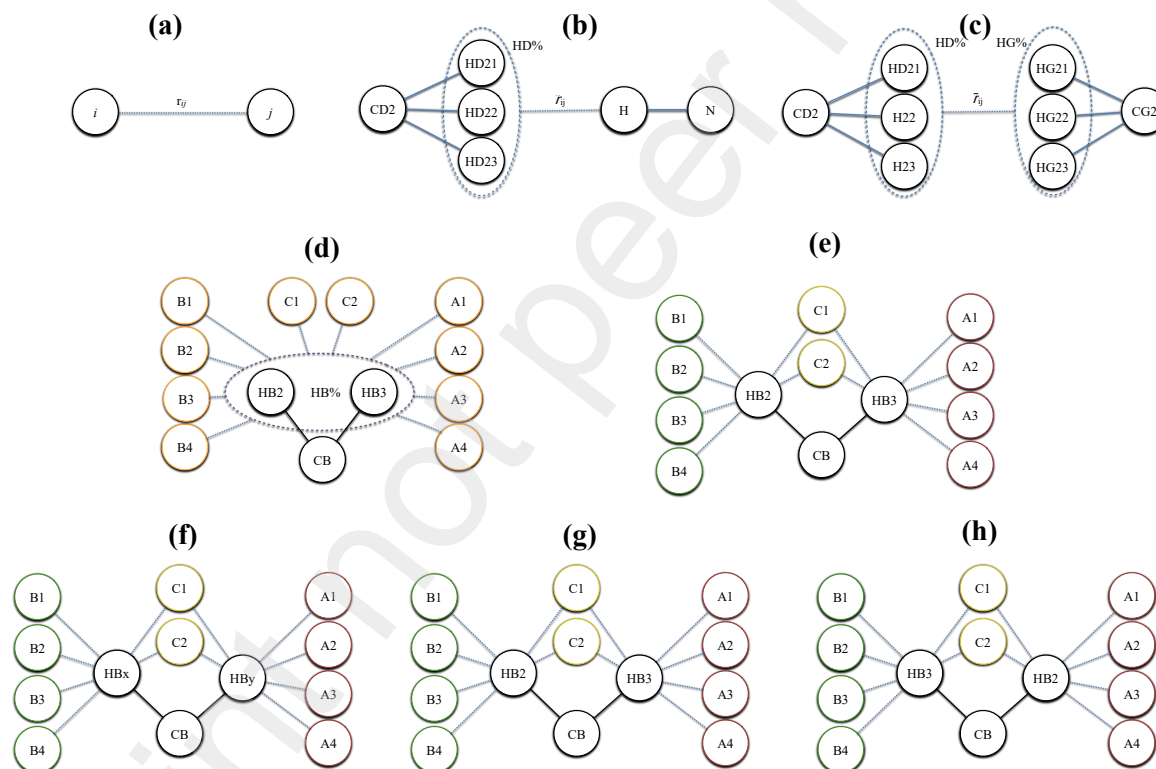


Fig. 2. (a) Angle restraint without target value (b) & (c) two possible target values for a given set of minimum and maximum, which makes either the counterclockwise (b) or the clockwise(c) angular region as allowed region.

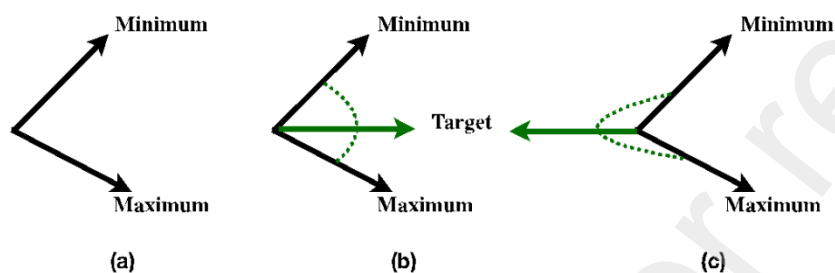


Fig. 3. Bar graph distribution of (a) distance and (b) dihedral angle restraints of PDB ID 7m5t. Violated and consistently violated portions are shown in different hash patterns.

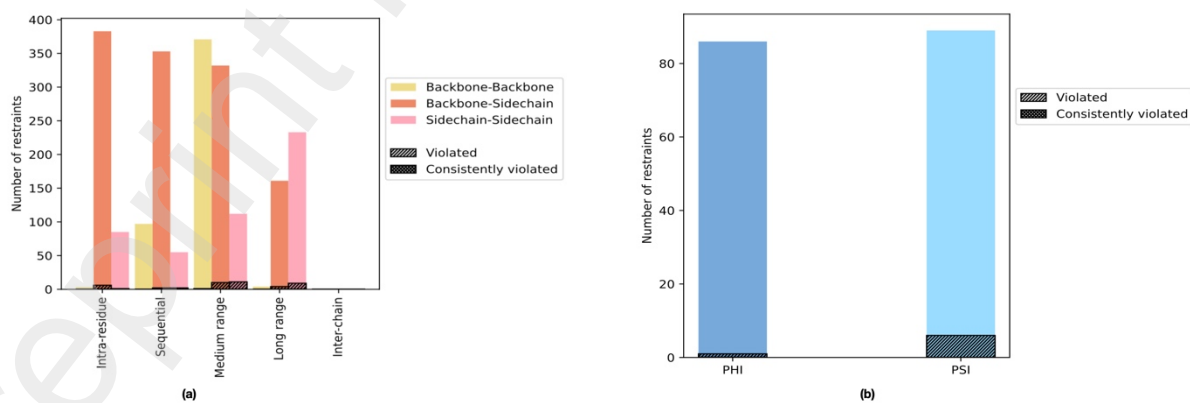


Fig. 4. Per-model distance (a) and dihedral (b) violation statistics of PDB ID 7m5t. The mean (dot), median (x) and the standard deviation (error bar) of the violation are shown in blue with respect to the y axis on the right.

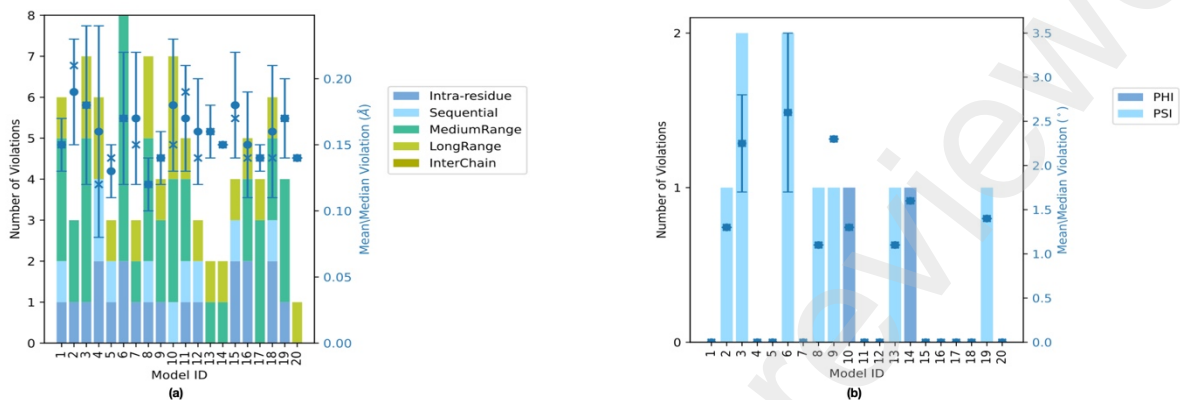


Fig. 5. Number of distance (a) and dihedral angle (b) violations *versus* the size of the ensemble for the PDB ID 7m5t.

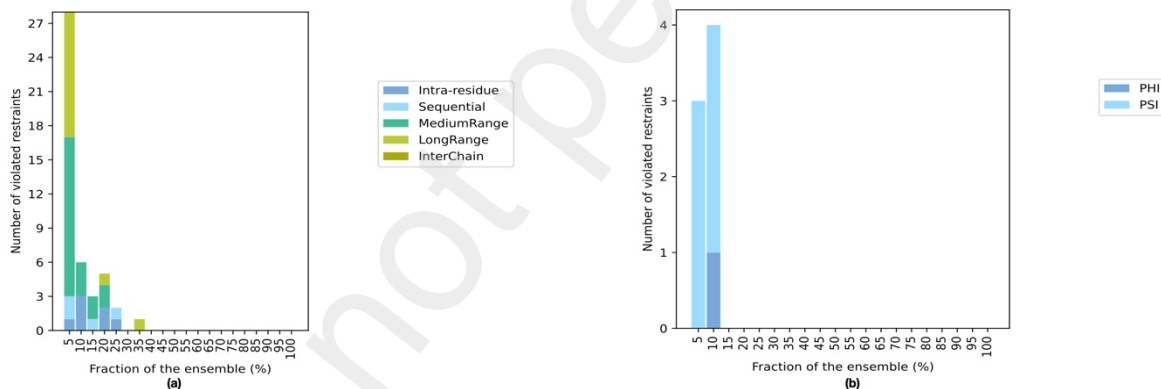


Fig. 6. Consistent distance restraint violations mapped on the structural ensemble of (a) PDB ID 2png and (b) PDB ID 1pqx. Per model, per-residue restraint violations were calculated by adding violations $> 0.3 \text{ \AA}$ observed in $>50\%$ of the models for all restraints involving the residue and mapped onto a blue-yellow color ramp for 0 to 5 violated restraints per model. Violations were calculated using validation reports generated by the wwPDB validation system (validate.wwpdb.org) with coordinates in PDBx/mmCIF format and NEF formatted restraints (available at the NEF GitHub repository).

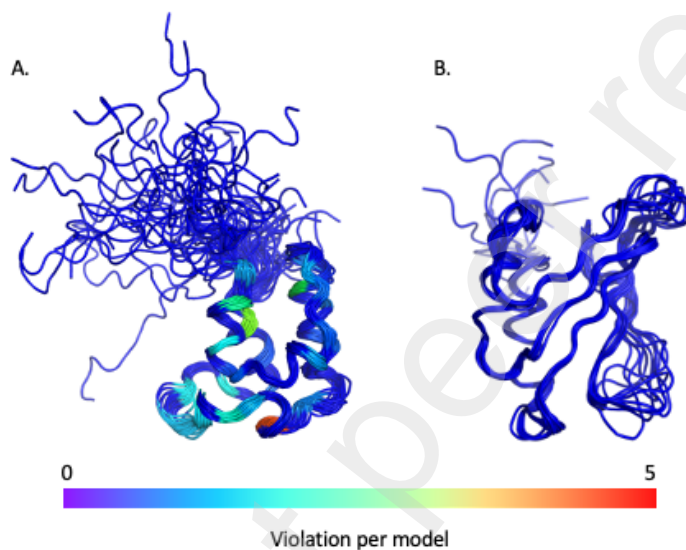


Table 1. NMR data checking at OneDep deposition

Types of data checking	Description of checking
File format	<ul style="list-style-type: none"> ● Check the file format whether it is unified NEF, or NMR-STAR or native format. File upload is blocked if the upload file type does not match to the selected file type. ● Check the presence of mandatory polymer sequences, chemical shifts, and restraints if a unified NEF or NMR-Star file is uploaded.
Polymer sequence	Cross-check author-provided sequences with assigned chemical shifts for the consistency within a NMR data file.
Nomenclature and atom assignments	<ul style="list-style-type: none"> ● Check atom naming in the NEF file and standardize nomenclature according to the Chemical Component Dictionary. ● Check the observed atoms are present in the chemical shifts.
PDBx/mmCIF dictionary compliant	Check data against PDBx/mmCIF dictionary for mandatory data, data type, and data boundaries (soft and hard limits). Provide a warning message if the value is outside the soft limit or an error message (blocked) if the value is outside the hard limit.
Cross-check consistency among NMR data types	<ul style="list-style-type: none"> ● Check between chemical shifts and restraints. Ambiguous methyl groups of distance restraints must have corresponding assigned chemical shifts ● Check between chemical shifts and assigned peak lists. Significant differences between chemical shift and spectral position are reported as errors.
Cross-check between atomic coordinates and experimental data	<ul style="list-style-type: none"> ● Polymer sequences in the atomic coordinates must be present in the data file. ● All atoms per residue must be present in the coordinates.

	<ul style="list-style-type: none"> • Check and validate bond distance for protonation state and disulfide bond
Ensemble check	Check that ensemble models are superimposed
Anomalous chemical shifts	<ul style="list-style-type: none"> • Check the value against archival statistical distribution and provide warning for unusual values • Check all methyl protons within a methyl group have identical chemical shift value unless different occupancies are provided

Table 2. Conformationally-restricting restraints using PDB ID 7m5t as an example.

Description	Value
Total distance restraints	2189
Intra-residue ($ i-j =0$)	471
Sequential ($ i-j =1$)	505
Medium range ($ i-j >1$ and $ i-j <5$)	675
Long range ($ i-j _5$)	398
Inter-chain	0
Hydrogen bond restraints	140
Disulfide bond restraints	0
Total dihedral-angle restraints	175
Number of unmapped restraints	0
Number of restraints per residue	23.6
Number of long range restraints per residue	4.0

Table 3. Average number of distance violations per model using PDB ID 7m5t as an example

Bins	Average number of violations per model	Max (Ångstrom)
1.0 - 0.2 Å (Small)	3.6	0.2
0.2 - 0.5 Å (Medium)	0.9	0.32
0.5 Å (Large)	None	None

STAR Methods

Key resources table

Resource availability

Lead Contact

Further information and requests for resources should be directed to and will be fulfilled by the Lead Contact, Kumaran Baskaran (baskaran@uchc.edu).

Materials availability

This study did not generate new unique reagents.

Data and code availability

[wwPDB](https://validate.wwpdb.org) validation tools are publicly accessible. The wwPDB anonymous validation server is provided at <https://validate.wwpdb.org> and the wwPDB validation API is accessible at <http://www.wwpdb.org/validation/onedep-validation-web-service-interface>. wwPDB validation report for each PDB ID is provided for users to download at PDB archive, https://ftp.wwpdb.org/pub/pdb/validation_reports/. These validation reports are also accessible at wwPDB website via PDB DOI links, e.g., https://www.wwpdb.org/pdb?id=pdb_00007m5t can be accessed via PDB DOI link: [10.2210/pdb7M5T/pdb](https://doi.org/10.2210/pdb7M5T/pdb).

The NEF standard and related code is available at <https://github.com/NMRExchangeFormat/NEF/>. The corresponding PDBx/mmCIF dictionary is accessible at https://mmcif.wwpdb.org/dictionaries/mmcif_nef.dic/Index/.

PDBStat is available under an open source license at <https://github.rpi.edu/RPIBioinformatics>

Experimental model and subject details

There was no model used.

Method details

The NMR exchange format (NEF) (Gutmanas *et al.*, 2015)

<https://github.com/NMRExchangeFormat/NEF/>) presents a community supported standard for the interchange of NMR data between different software programs. The format is based upon the STAR syntax (ref) and defines so-called saveframes, *i.e.*, self-contained blocks of data, for sequence, chemical shifts, resonance peaks in NMR spectra, dihedral-, distance-, and RDC-restraints, as well as relevant metadata and a linkage table connecting restraints and peaks. Importantly, the format is inherently extendable through so-called namespace specific tags, both with respect to the data contained in each saveframe or as complete additional saveframes.

The NEF defines a nomenclature convention for the twenty common protein amino acids and the eight RNA/DNA oligonucleotides in their common appearance in NMR spectra (*i.e.*, appropriately protonated at pH 7.0). This nomenclature follows the IUPAC convention, with extensions to accommodate NMR specific situations that follow from degenerate resonances and stereo-specificity, *e.g.*, for methylene protons and VAL, LEU methyl groups. Key aspects of this extension is the existence of a wild-card indicator (“%”, *e.g.*, as in HB%) and indicators for non-degenerate, but non stereo-specifically assigned resonances (“x” and “y”, as in HBx and HBy). Together with the presence of a specific atom-based tag (atom_site.pdbx_atom_ambiguity) in the structural mmCif file, this allows for an unambiguous and exact mapping of the NMR restraint onto the molecular structure. A detailed description of the NEF will be presented elsewhere.

Quantification and statistical analysis

No statistical analysis was performed.

References

- Adams, P.D., Afonine, P.V., Baskaran, K., Berman, H.M., Berrisford, J., Bricogne, G., Brown, D.G., Burley, S.K., Chen, M., Feng, Z., et al. (2019). Announcing mandatory submission of PDBx/mmCIF format files for crystallographic depositions to the Protein Data Bank (PDB). *Acta Crystallogr D Struct Biol* 75, 451-454. 10.1107/S2059798319004522.
- Alderson, T.R., and Kay, L.E. (2020). Unveiling invisible protein states with NMR spectroscopy. *Curr Opin Struct Biol* 60, 39-49. 10.1016/j.sbi.2019.10.008.

- Anishchenko, I., Pellock, S.J., Chidyausiku, T.M., Ramelot, T.A., Ovchinnikov, S., Hao, J., Bafna, K., Norn, C., Kang, A., Bera, A.K., et al. (2021). De novo protein design by deep network hallucination. *Nature* *600*, 547-552. 10.1038/s41586-021-04184-w.
- Anthis, N.J., and Clore, G.M. (2015). Visualizing transient dark states by NMR spectroscopy. *Q Rev Biophys* *48*, 35-116. 10.1017/S0033583514000122.
- Armstrong, D.R., Berrisford, J.M., Conroy, M.J., Gutmanas, A., Anyango, S., Choudhary, P., Clark, A.R., Dana, J.M., Deshpande, M., Dunlop, R., et al. (2020). PDBe: improved findability of macromolecular structure data in the PDB. *Nucleic Acids Res* *48*, D335-D343. 10.1093/nar/gkz990.
- Bassolino-Klimas, D., Tejero, R., Krystek, S.R., Metzler, W.J., Montelione, G.T., and Bruccoleri, R.E. (1996). Simulated annealing with restrained molecular dynamics using a flexible restraint potential: theory and evaluation with simulated NMR constraints. *Protein Sci* *5*, 593-603. 10.1002/pro.5560050404.
- Bekker, G.J., Yokochi, M., Suzuki, H., Ikegawa, Y., Iwata, T., Kudou, T., Yura, K., Fujiwara, T., Kawabata, T., and Kurisu, G. (2022). Protein Data Bank Japan: Celebrating our 20th anniversary during a global pandemic as the Asian hub of three dimensional macromolecular structural data. *Protein Sci* *31*, 173-186. 10.1002/pro.4211.
- Berman, H., Henrick, K., and Nakamura, H. (2003). Announcing the worldwide Protein Data Bank. *Nat Struct Biol* *10*, 980. 10.1038/nsb1203-980.
- Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N., and Bourne, P.E. (2000). The Protein Data Bank. *Nucleic Acids Res* *28*, 235-242. 10.1093/nar/28.1.235.
- Bertini, I., Del Bianco, C., Gelis, I., Katsaros, N., Luchinat, C., Parigi, G., Peana, M., Provenzani, A., and Zoroddu, M.A. (2004). Experimentally exploring the conformational space sampled by domain reorientation in calmodulin. *Proc Natl Acad Sci U S A* *101*, 6841-6846. 10.1073/pnas.0308641101.
- Bhardwaj, G., O'Connor, J., Rettie, S., Huang, Y.H., Ramelot, T.A., Mulligan, V.K., Alpkilic, G.G., Palmer, J., Bera, A.K., Bick, M.J., et al. (2022). Accurate de novo design of membrane-traversing macrocycles. *Cell* *185*, 3520-3532 e3526. 10.1016/j.cell.2022.07.019.
- Bhattacharya, A., Tejero, R., and Montelione, G.T. (2007). Evaluating protein structures determined by structural genomics consortia. *Proteins* *66*, 778-795. 10.1002/prot.21165.
- Boelens, R., Scheek, R.M., van Boom, J.H., and Kaptein, R. (1987). Complex of lac repressor headpiece with a 14 base-pair lac operator fragment studied by two-dimensional nuclear magnetic resonance. *J Mol Biol* *193*, 213-216. 10.1016/0022-2836(87)90638-3.
- Brünger, A.T., Adams, P.D., Clore, G.M., DeLano, W.L., Gros, P., Grosse-Kunstleve, R.W., Jiang, J.S., Kuszewski, J., Nilges, M., Pannu, N.S., et al. (1998). Crystallography & NMR system: A new software suite for macromolecular structure determination. *Acta Crystallogr D Biol Crystallogr* *54*, 905-921. doi: 10.1107/s0907444998003254.
- Burley, S.K., Bhikadiya, C., Bi, C., Bittrich, S., Chen, L., Crichlow, G.V., Christie, C.H., Dalenberg, K., Di Costanzo, L., Duarte, J.M., et al. (2021). RCSB Protein Data Bank: powerful new tools for exploring 3D structures of biological macromolecules for basic and applied research and education in fundamental biology, biomedicine, biotechnology, bioengineering and energy sciences. *Nucleic Acids Res* *49*, D437-D451. 10.1093/nar/gkaa1038.
- CCPN (2023). CCPN Converter. <https://doi.org/10.1002/prot.20449>.
- Chen, K., and Tjandra, N. (2012). The use of residual dipolar coupling in studying proteins by NMR. *Top Curr Chem* *326*, 47-67. 10.1007/128_2011_215.
- Chen, V.B., Arendall, W.B., 3rd, Headd, J.J., Keedy, D.A., Immormino, R.M., Kapral, G.J., Murray, L.W., Richardson, J.S., and Richardson, D.C. (2010). MolProbity: all-atom

- structure validation for macromolecular crystallography. *Acta crystallographica. Section D, Biological crystallography* 66, 12-21. 10.1107/S0907444909042073.
- Cheung, M.S., Maguire, M.L., Stevens, T.J., and Broadhurst, R.W. (2010). DANGLE: A Bayesian inferential method for predicting protein backbone dihedral angles and secondary structure. *J Magn Reson* 202, 223-233. 10.1016/j.jmr.2009.11.008.
- Chuprina, V.P., Rullmann, J.A., Lamerichs, R.M., van Boom, J.H., Boelens, R., and Kaptein, R. (1993). Structure of the complex of lac repressor headpiece and an 11 base-pair half-operator determined by nuclear magnetic resonance spectroscopy and restrained molecular dynamics. *J Mol Biol* 234, 446-462. 10.1006/jmbi.1993.1598.
- Cicero, D.O., Barbato, G., and Bazzo, R. (1995). NMR analysis of molecular flexibility in solution: A new method for the study of complex distributions of rapidly exchanging conformations. Application to a 13-Residue peptide with an 8-residue loop. *J. Am. Chem. Soc.* 117, 1027 - 1033.
- Clore, G.M., and Garrett, D.S. (1999). R-factor, free R, and complete cross-validation for dipolar coupling refinement of NMR structures. *J. Amer. Chem. Soc.* 121, 9008-9012.
- Clore, G.M., Nilges, M., Sukumaran, D.K., Brünger, A.T., Karplus, M., and Gronenborn, A.M. (1986). The three-dimensional structure of alpha1-purothionin in solution: combined use of nuclear magnetic resonance, distance geometry and restrained molecular dynamics. *EMBO J* 5, 2729-2735. 10.1002/j.1460-2075.1986.tb04557.x.
- Cooke, R.M., Wilkinson, A.J., Baron, M., Pastore, A., Tappin, M.J., Campbell, I.D., Gregory, H., and Sheard, B. (1987). The solution structure of human epidermal growth factor. *Nature* 327, 339-341. 10.1038/327339a0.
- Cornilescu, G., Marquardt, J.L., Ottiger, M., and Bax, A. (1998). Validation of protein structure from anisotropic carbonyl chemical shifts in a dilute liquid crystalline phase. *J. Amer. Chem. Soc.* 120, 6836-6837.
- Doreleijers, J.F., Raves, M.L., Rullmann, T., and Kaptein, R. (1999). Completeness of NOEs in protein structure: a statistical analysis of NMR. *Journal of biomolecular NMR* 14, 123-132. 10.1023/a:1008335423527.
- Doreleijers, J.F., Vranken, W.F., Schulte, C., Markley, J.L., Ulrich, E.L., Vriend, G., and Vuister, G.W. (2012). NRG-CING: integrated validation reports of remediated experimental biomolecular NMR data and coordinates in wwPDB. *Nucleic Acids Res* 40, D519-524. 10.1093/nar/gkr1134.
- Feng, Z., Westbrook, J.D., Sala, R., Smart, O.S., Bricogne, G., Matsubara, M., Yamada, I., Tsuchiya, S., Aoki-Kinoshita, K.F., Hoch, J.C., et al. (2021). Enhanced validation of small-molecule ligands and carbohydrates in the Protein Data Bank. *Structure* 29, 393-400 e391. 10.1016/j.str.2021.02.004.
- Folmer, R.H., Hilbers, C.W., Konings, R.N., and Nilges, M. (1997). Floating stereospecific assignment revisited: application to an 18 kDa protein and comparison with J-coupling data. *Journal of biomolecular NMR* 9, 245-258. 10.1023/a:1018670623695.
- Fowler, N.J., Sljoka, A., and Williamson, M.P. (2020). A method for validating the accuracy of NMR protein structures. *Nat Commun* 11, 6321. 10.1038/s41467-020-20177-1.
- Fowler, N.J., and Williamson, M.P. (2022). The accuracy of protein structures in solution determined by AlphaFold and NMR. *bioRxiv*, 2022.2001.2018.476751. 10.1101/2022.01.18.476751.
- Gibbs, A.C., Steele, R., Liu, G., Tounge, B.A., and Montelione, G.T. (2018). Inhibitor Bound Dengue NS2B-NS3pro Reveals Multiple Dynamic Binding Modes. *Biochemistry* 57, 1591-1602. 10.1021/acs.biochem.7b01127.
- Gochin, M., and James, T.L. (1990). Solution structure studies of d(AC)4.d(GT)4 via restrained molecular dynamics simulations with NMR constraints derived from two-dimensional NOE and double-quantum-filtered COSY experiments. *Biochemistry* 29, 11172-11180. 10.1021/bi00503a004.

- Gore, S., Sanz Garcia, E., Hendrickx, P.M.S., Gutmanas, A., Westbrook, J.D., Yang, H., Feng, Z., Baskaran, K., Berrisford, J.M., Hudson, B.P., et al. (2017). Validation of Structures in the Protein Data Bank. *Structure* 25, 1916-1927. 10.1016/j.str.2017.10.009.
- Gorler, A., and Kalbitzer, H.R. (1997). Relax, a flexible program for the back calculation of NOESY spectra based on complete-relaxation-matrix formalism. *J Magn Reson* 124, 177-188. 10.1006/jmre.1996.1033.
- Guntert, P., Braun, W., and Wuthrich, K. (1991). Efficient computation of three-dimensional protein structures in solution from nuclear magnetic resonance data using the program DIANA and the supporting programs CALIBA, HABAS and GLOMSA. *J Mol Biol* 217, 517-530. 10.1016/0022-2836(91)90754-t.
- Guntert, P., and Buchner, L. (2015). Combined automated NOE assignment and structure calculation with CYANA. *Journal of biomolecular NMR* 62, 453-471. 10.1007/s10858-015-9924-9.
- Gutmanas, A., Adams, P.D., Bardiaux, B., Berman, H.M., Case, D.A., Fogh, R.H., Guntert, P., Hendrickx, P.M., Herrmann, T., Kleywegt, G.J., et al. (2015). NMR Exchange Format: a unified and open standard for representation of NMR restraint data. *Nat Struct Mol Biol* 22, 433-434. 10.1038/nsmb.3041.
- Harish, B., Swapna, G.V., Kornhaber, G.J., Montelione, G.T., and Carey, J. (2017). Multiple helical conformations of the helix-turn-helix region revealed by NOE-restrained MD simulations of tryptophan aporepressor, TrpR. *Proteins* 85, 731-740. 10.1002/prot.25252.
- Herrmann, T., Guntert, P., and Wuthrich, K. (2002). Protein NMR structure determination with automated NOE assignment using the new software CANDID and the torsion angle dynamics algorithm DYANA. *J Mol Biol* 319, 209-227. 10.1016/s0022-2836(02)00241-3.
- Hoch, J.C., Baskaran, K., Burr, H., Chin, J., Eghbalnia, H.R., Fujiwara, T., Gryk, M.R., Iwata, T., Kojima, C., Kurisu, G., et al. (2023). Biological Magnetic Resonance Data Bank. *Nucleic Acids Res* 51, D368-D376. 10.1093/nar/gkac1050.
- Huang, Y.J., Powers, R., and Montelione, G.T. (2005). Protein NMR recall, precision, and F-measure scores (RPF scores): structure quality assessment measures based on information retrieval statistics. *J Am Chem Soc* 127, 1665-1674. 10.1021/ja047109h.
- Huang, Y.J., Rosato, A., Singh, G., and Montelione, G.T. (2012). RPF: a quality assessment tool for protein NMR structures. *Nucleic Acids Research* 40, W542-546. 10.1093/nar/gks373.
- Huang, Y.J., Tejero, R., Powers, R., and Montelione, G.T. (2006). A topology-constrained distance network algorithm for protein structure determination from NOESY data. *Proteins* 62, 587-603. 10.1002/prot.20820.
- Hyberts, S.G., Goldberg, M.S., Havel, T.F., and Wagner, G. (1992). 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, 736-751. 10.1002/pro.5560010606.
- Jacobs, T.M., Williams, B., Williams, T., Xu, X., Eletsy, A., Federizon, J.F., Szyperski, T., and Kuhlman, B. (2016). Design of structurally distinct proteins using strategies inspired by evolution. *Science* 352, 687-690. 10.1126/science.aad8036.
- Kaptein, R., Zuiderweg, E.R., Scheek, R.M., Boelens, R., and van Gunsteren, W.F. (1985). A protein structure from nuclear magnetic resonance data. lac repressor headpiece. *J Mol Biol* 182, 179-182. 10.1016/0022-2836(85)90036-1.
- Kirchner, D.K., and Guntert, P. (2011). Objective identification of residue ranges for the superposition of protein structures. *BMC Bioinformatics* 12, 170. 10.1186/1471-2105-12-170.
- Kline, A.D., Braun, W., and Wuthrich, K. (1988). Determination of the complete three-dimensional structure of the alpha-amylase inhibitor tendamistat in aqueous solution by

- nuclear magnetic resonance and distance geometry. *J Mol Biol* 204, 675-724. 10.1016/0022-2836(88)90364-6.
- Knegtel, R.M., Fogh, R.H., Otleben, G., Ruterjans, H., Dumoulin, P., Schnarr, M., Boelens, R., and Kaptein, R. (1995). A model for the LexA repressor DNA complex. *Proteins* 21, 226-236. 10.1002/prot.340210305.
- Koepnick, B., Flatten, J., Husain, T., Ford, A., Silva, D.A., Bick, M.J., Bauer, A., Liu, G., Ishida, Y., Boykov, A., et al. (2019). De novo protein design by citizen scientists. *Nature* 570, 390-394. 10.1038/s41586-019-1274-4.
- Koga, N., Koga, R., Liu, G., Castellanos, J., Montelione, G.T., and Baker, D. (2021). Role of backbone strain in de novo design of complex alpha/beta protein structures. *Nat Commun* 12, 3921. 10.1038/s41467-021-24050-7.
- Koga, N., Tatsumi-Koga, R., Liu, G., Xiao, R., Acton, T.B., Montelione, G.T., and Baker, D. (2012). Principles for designing ideal protein structures. *Nature* 491, 222-227. 10.1038/nature11600.
- Laskowski, R.A., MacArthur, M.W., Moss, D.S., and Thornton, J.M. (1993). PROCHECK - a Program to Check the Stereochemical Quality of Protein Structures. *J Appl Crystallogr* 26, 283-291.
- Lawson, C.L., Patwardhan, A., Baker, M.L., Hryc, C., Garcia, E.S., Hudson, B.P., Lagerstedt, I., Ludtke, S.J., Pintilie, G., Sala, R., et al. (2016). EMDatabank unified data resource for 3DEM. *Nucleic Acids Res* 44, D396-403. 10.1093/nar/gkv1126.
- Lin, Y.R., Koga, N., Tatsumi-Koga, R., Liu, G., Clouser, A.F., Montelione, G.T., and Baker, D. (2015). Control over overall shape and size in de novo designed proteins. *Proc Natl Acad Sci U S A* 112, E5478-5485. 10.1073/pnas.1509508112.
- Lipsitz, R.S., and Tjandra, N. (2004). Residual dipolar couplings in NMR structure analysis. *Annu Rev Biophys Biomol Struct* 33, 387-413. 10.1146/annurev.biophys.33.110502.140306.
- Losonczi, J.A., Andrec, M., Fischer, M.W., and Prestegard, J.H. (1999). Order matrix analysis of residual dipolar couplings using singular value decomposition. *J Magn Reson* 138, 334-342. 10.1006/jmre.1999.1754.
- Lovell, S.C., Davis, I.W., Arendall, W.B., 3rd, de Bakker, P.I., Word, J.M., Prisant, M.G., Richardson, J.S., and Richardson, D.C. (2003). Structure validation by Calpha geometry: phi,psi and Cbeta deviation. *Proteins* 50, 437-450. 10.1002/prot.10286.
- Markley, J.L., Ulrich, E.L., Berman, H.M., Henrick, K., Nakamura, H., and Akutsu, H. (2008). BioMagResBank (BMRB) as a partner in the Worldwide Protein Data Bank (wwPDB): new policies affecting biomolecular NMR depositions. *Journal of biomolecular NMR* 40, 153-155. 10.1007/s10858-008-9221-y.
- Montelione, G.T., Nilges, M., Bax, A., Guntert, P., Herrmann, T., Richardson, J.S., Schwieters, C.D., Vranken, W.F., Vuister, G.W., Wishart, D.S., et al. (2013). Recommendations of the wwPDB NMR Validation Task Force. *Structure* 21, 1563-1570. 10.1016/j.str.2013.07.021.
- Montelione, G.T., Wuthrich, K., Nice, E.C., Burgess, A.W., and Scheraga, H.A. (1987). Solution structure of murine epidermal growth factor: determination of the polypeptide backbone chain-fold by nuclear magnetic resonance and distance geometry. *Proc Natl Acad Sci U S A* 84, 5226-5230. 10.1073/pnas.84.15.5226.
- Neal, S., Nip, A.M., Zhang, H., and Wishart, D.S. (2003). Rapid and accurate calculation of protein ¹H, ¹³C and ¹⁵N chemical shifts. *Journal of biomolecular NMR* 26, 215-240. 10.1023/a:1023812930288.
- Nederveen, A.J., Doreleijers, J.F., Vranken, W., Miller, Z., Spronk, C.A., Nabuurs, S.B., Guntert, P., Livny, M., Markley, J.L., Nilges, M., et al. (2005). RECOORD: a recalculated coordinate database of 500+ proteins from the PDB using restraints from the BioMagResBank. *Proteins* 59, 662-672. 10.1002/prot.20408.

- Nilges, M. (1995). Calculation of protein structures with ambiguous distance restraints. Automated assignment of ambiguous NOE crosspeaks and disulphide connectivities. *J Mol Biol* 245, 645-660. 10.1006/jmbi.1994.0053.
- Nilges, M., Bernard, A., Bardiaux, B., Malliavin, T., Habeck, M., and Rieping, W. (2008). Accurate NMR structures through minimization of an extended hybrid energy. *Structure* 16, 1305-1312. 10.1016/j.str.2008.07.008.
- Parigi, G., Ravera, E., and Luchinat, C. (2022). Paramagnetic effects in NMR for protein structures and ensembles: Studies of metalloproteins. *Curr Opin Struct Biol* 74, 102386. 10.1016/j.sbi.2022.102386.
- Protein_Data_Bank (1971). Crystallography: Protein data bank. *Nature* 233, 223.
- Qian, Y.Q., Resendez-Perez, D., Gehring, W.J., and Wuthrich, K. (1994). The des(1-6)antennapedia homeodomain: comparison of the NMR solution structure and the DNA-binding affinity with the intact Antennapedia homeodomain. *Proc Natl Acad Sci U S A* 91, 4091-4095. 10.1073/pnas.91.9.4091.
- Ried, A., Gronwald, W., Trenner, J.M., Brunner, K., Neidig, K.P., and Kalbitzer, H.R. (2004). Improved simulation of NOESY spectra by RELAX-JT2 including effects of J-coupling, transverse relaxation and chemical shift anisotropy. *Journal of biomolecular NMR* 30, 121-131. 10.1023/b:jnmr.0000048945.88968.af.
- Rieping, W., Habeck, M., Bardiaux, B., Bernard, A., Malliavin, T.E., and Nilges, M. (2007). ARIA2: automated NOE assignment and data integration in NMR structure calculation. *Bioinformatics* 23, 381-382. 10.1093/bioinformatics/btl589.
- Rieping, W., Habeck, M., and Nilges, M. (2005). Inferential structure determination. *Science* 309, 303-306. 10.1126/science.1110428.
- Rosato, A., Tejero, R., and Montelione, G.T. (2013). Quality assessment of protein NMR structures. *Curr Opin Struct Biol* 23, 715-724. 10.1016/j.sbi.2013.08.005.
- Schwieters, C.D., Kuszewski, J.J., and Clore, G.M. (2006). Using Xplor-NIH for NMR molecular structure determination. *Prog Nucl Mag Res Sp* 48, 47-62. Doi 10.1016/J.Pnmrs.2005.10.001.
- Sengupta, I., Nadaud, P.S., Helmus, J.J., Schwieters, C.D., and Jaroniec, C.P. (2012). Protein fold determined by paramagnetic magic-angle spinning solid-state NMR spectroscopy. *Nat Chem* 4, 410-417. 10.1038/nchem.1299.
- Shen, Y., and Bax, A. (2010). SPARTA+: a modest improvement in empirical NMR chemical shift prediction by means of an artificial neural network. *Journal of biomolecular NMR* 48, 13-22. 10.1007/s10858-010-9433-9.
- Shen, Y., and Bax, A. (2015). Protein structural information derived from NMR chemical shift with the neural network program TALOS-N. *Methods in molecular biology* 1260, 17-32. 10.1007/978-1-4939-2239-0_2.
- Skinner, S.P., Fogh, R.H., Boucher, W., Ragan, T.J., Mureddu, L.G., and Vuister, G.W. (2016). CcpNmr AnalysisAssign: a flexible platform for integrated NMR analysis. *Journal of biomolecular NMR* 66, 111-124. 10.1007/s10858-016-0060-y.
- Snyder, D.A., Grullon, J., Huang, Y.J., Tejero, R., and Montelione, G.T. (2014). The expanded FindCore method for identification of a core atom set for assessment of protein structure prediction. *Proteins* 82 Suppl 2, 219-230. 10.1002/prot.24490.
- Snyder, D.A., and Montelione, G.T. (2005). Clustering algorithms for identifying core atom sets and for assessing the precision of protein structure ensembles. *Proteins* 59, 673-686. 10.1002/prot.20402.
- Tejero, R., Snyder, D., Mao, B., Aramini, J.M., and Montelione, G.T. (2013). PDBStat: a universal restraint converter and restraint analysis software package for protein NMR. *Journal of biomolecular NMR* 56, 337-351. 10.1007/s10858-013-9753-7.
- Thomas, P.D., Basus, V.J., and James, T.L. (1991). Protein solution structure determination using distances from two-dimensional nuclear Overhauser effect experiments: effect of

- approximations on the accuracy of derived structures. *Proc Natl Acad Sci U S A* **88**, 1237-1241. 10.1073/pnas.88.4.1237.
- Trindade, I.B., Invernici, M., Cantini, F., Louro, R.O., and Piccioli, M. (2021). PRE-driven protein NMR structures: an alternative approach in highly paramagnetic systems. *FEBS J* **288**, 3010-3023. 10.1111/febs.15615.
- Ulrich, E.L., Akutsu, H., Doreleijers, J.F., Harano, Y., Ioannidis, Y.E., Lin, J., Livny, M., Mading, S., Maziuk, D., Miller, Z., et al. (2008). BioMagResBank. *Nucleic Acids Res* **36**, D402-408. 10.1093/nar/gkm957.
- Ulrich, E.L., Baskaran, K., Dashti, H., Ioannidis, Y.E., Livny, M., Romero, P.R., Maziuk, D., Wedell, J.R., Yao, H., Eghbalnia, H.R., et al. (2019). NMR-STAR: comprehensive ontology for representing, archiving and exchanging data from nuclear magnetic resonance spectroscopic experiments. *Journal of biomolecular NMR* **73**, 5-9. 10.1007/s10858-018-0220-3.
- van der Aalst, W.M.P., Bichler, M., and Heinzl, A. (2017). Responsible data science. *Bus Inf Syst Eng.* **59**, 311–313.
- Vila, J.A., Aramini, J.M., Rossi, P., Kuzin, A., Su, M., Seetharaman, J., Xiao, R., Tong, L., Montelione, G.T., and Scheraga, H.A. (2008). Quantum chemical ¹³C(alpha) chemical shift calculations for protein NMR structure determination, refinement, and validation. *Proc Natl Acad Sci U S A* **105**, 14389-14394. 10.1073/pnas.0807105105.
- Vranken, W.F., Boucher, W., Stevens, T.J., Fogh, R.H., Pajon, A., Llinas, M., Ulrich, E.L., Markley, J.L., Ionides, J., and Laue, E.D. (2005). The CCPN data model for NMR spectroscopy: development of a software pipeline. *Proteins* **59**, 687-696. 10.1002/prot.20449.
- Vuister, G.W., Fogh, R.H., Hendrickx, P.M., Doreleijers, J.F., and Gutmanas, A. (2014). An overview of tools for the validation of protein NMR structures. *Journal of biomolecular NMR* **58**, 259-285. 10.1007/s10858-013-9750-x.
- Wagner, G., Braun, W., Havel, T.F., Schaumann, T., Go, N., and Wuthrich, K. (1987a). Protein structures in solution by nuclear magnetic resonance and distance geometry. The polypeptide fold of the basic pancreatic trypsin inhibitor determined using two different algorithms, DISGEO and DISMAN. *J Mol Biol* **196**, 611-639. 10.1016/0022-2836(87)90037-4.
- Wagner, G., Frey, M.H., Neuhaus, D., Worgotter, E., Braun, W., Vasak, M., Kagi, J.H., and Wuthrich, K. (1987b). Spatial structure of rabbit liver metallothionein-2 in solution by NMR. *Experientia Suppl* **52**, 149-157. 10.1007/978-3-0348-6784-9_8.
- Weiss, M.A., and Hoch, J.C. (1987). Interpretation of ring-current shifts in proteins: Application to phage λ repressor. *J Magn Reso* **72**, 324-333.
- Westbrook, J.D., Young, J.Y., Shao, C., Feng, Z., Guranovic, V., Lawson, C.L., Vallat, B., Adams, P.D., Berrisford, J.M., Bricogne, G., et al. (2022). PDBx/mmCIF Ecosystem: Foundational Semantic Tools for Structural Biology. *J Mol Biol* **434**, 167599. 10.1016/j.jmb.2022.167599.
- Wilkinson, M.D., Dumontier, M., Aalbersberg, I.J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.W., da Silva Santos, L.B., Bourne, P.E., et al. (2016). The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* **3**, 160018. 10.1038/sdata.2016.18.
- Williamson, M.P., Havel, T.F., and Wuthrich, K. (1985). Solution conformation of proteinase inhibitor IIA from bull seminal plasma by ¹H nuclear magnetic resonance and distance geometry. *J Mol Biol* **182**, 295-315. 10.1016/0022-2836(85)90347-x.
- ww, P.D.B.c. (2019). Protein Data Bank: the single global archive for 3D macromolecular structure data. *Nucleic Acids Res* **47**, D520-D528. 10.1093/nar/gky949.
- wwPDB_consortium (2019). Protein data bank: The single global archive for 3d macromolecular structure data. *Nucleic Acids Res* **47(D1)**, D520–D528.

- Young, J.Y., Westbrook, J.D., Feng, Z., Peisach, E., Persikova, I., Sala, R., Sen, S., Berrisford, J.M., Swaminathan, G.J., Oldfield, T.J., et al. (2018). Worldwide Protein Data Bank biocuration supporting open access to high-quality 3D structural biology data. *Database (Oxford)* 2018. 10.1093/database/bay002.
- Young, J.Y., Westbrook, J.D., Feng, Z., Sala, R., Peisach, E., Oldfield, T.J., Sen, S., Gutmanas, A., Armstrong, D.R., Berrisford, J.M., et al. (2017). OneDep: Unified wwPDB System for Deposition, Biocuration, and Validation of Macromolecular Structures in the PDB Archive. *Structure* 25, 536-545. 10.1016/j.str.2017.01.004.