

*Deriving structural information from experimentally
measured data on biomolecules: a review*

W. F. van Gunsteren^{1,*}, J. R. Allison², X. Daura³, J. Dolenc¹, N. Hansen⁴, A.
E. Mark⁵, C. Oostenbrink⁶, V. H. Rusu¹, L. J. Smith⁷

¹Laboratory of Physical Chemistry, Swiss Federal Institute of Technology, ETH, Zürich, Switzerland.

²Centre for Theoretical Chemistry and Physics & Institute of Natural and Mathematical Sciences, Massey University, Auckland, New Zealand.

Biomolecular Interaction Centre, University of Canterbury, Christchurch, New Zealand.

Maurice Wilkinson Centre for Molecular Biodiscovery, New Zealand.

³Institute of Biotechnology and Biomedicine, Universitat Autònoma de Barcelona (UAB), 08193, Barcelona, Spain.

Catalan Institution for Research and Advanced Studies (ICREA), 08010, Barcelona, Spain

⁴Institute of Thermodynamics and Thermal Process Engineering, University of Stuttgart, Pfaffenwaldring 9, 70569, Stuttgart, Germany.

⁵School of Chemistry and Molecular Biosciences, University of Queensland, St. Lucia, QLD 4072, Australia

⁶Institute of Molecular Modeling and Simulation, University of Natural Resources and Life Sciences, Vienna, Austria.

⁷Department of Chemistry, University of Oxford, Inorganic Chemistry Laboratory, South Parks Road, Oxford OX1 3QR, United Kingdom.

Submitted to:

Angew. Chemie

Manuscript ID:

Date:

July 4, 2016

*Corresponding author. Email: wfvgn@ethz.ch

Abstract

1
2 During the past half century, the number and accuracy of experimental
3 techniques that can deliver values of observables for biomolecular systems
4 have been steadily increasing. The conversion of a measured value Q^{exp}
5 of an observable quantity Q into structural information is, however, a
6 task beset with theoretical and practical problems: (i) insufficient or in-
7 accurate values of Q^{exp} , (ii) inaccuracies in the function $Q(\vec{r})$ used to
8 relate the quantity Q to structure \vec{r} , (iii) how to account for the averag-
9 ing inherent in the measurement of Q^{exp} , (iv) how to handle the possible
10 multiple-valuedness of the inverse $\vec{r}(Q)$ of the function $Q(\vec{r})$, to mention
11 a few. These apply to a variety of observable quantities Q and measure-
12 ment techniques such as X-ray and neutron diffraction, small-angle and
13 wide-angle X-ray scattering, free-electron laser imaging, cryo-electron mi-
14 croscopy, nuclear magnetic resonance, electron paramagnetic resonance,
15 infrared and Raman spectroscopy, circular dichroism, Förster resonance
16 energy transfer, atomic force microscopy and ion-mobility mass spectrom-
17 etry. The process of deriving structural information from measured data
18 is reviewed with an eye to non-experts and newcomers in the field using
19 examples from the literature of the effect of the various choices and ap-
20 proximations involved in the process. A list of choices to be avoided is
21 provided.

22 1 Introduction

23 Over the past half century it has become increasingly feasible to obtain struc-
24 tural information on biomolecules at the atomic level of resolution due to the
25 development of a range of experimental techniques.

- 26 1. The analysis of X-ray diffraction patterns can be used to elucidate the
27 spatial structure of proteins.^[1-3]
- 28 2. Neutron diffraction measurements can be used to determine the positions
29 of hydrogen (i.e. deuterium) atoms,^[4] and provide information comple-
30 mentary to X-ray diffraction.
- 31 3. Small-angle X-ray scattering (SAXS) and wide-angle X-ray scattering (WAXS)
32 can be used to derive low resolution structural constraints for large biomolecules.^[5]
- 33 4. Free-electron X-ray lasers may be used for coherent imaging of biomolecules.^[6, 7]
- 34 5. Cryo-electron microscopy can provide structural information approaching
35 that of X-ray crystallography.^[8]
- 36 6. One- and higher-dimensional NMR techniques^[9] can be used to obtain
37 atom-atom distances within molecules, which in turn can be used to derive
38 spatial structures for biomolecules.^[10]
- 39 7. Electron Paramagnetic Resonance (EPR) measurements yield long-range
40 distance information between particular types of atoms.^[11]
- 41 8. Infrared (IR), Raman and fluorescence spectroscopy can be used to derive
42 local structural information on proteins.^[12, 13]

- 43 9. Circular Dichroism (CD) can provide information on secondary structure
44 content of a protein fold, but spectra possess a much lower information
45 density than can be obtained by NMR spectroscopic or X-ray diffraction
46 techniques.^[14]
- 47 10. Förster Resonance Energy Transfer (FRET) can be used to trace distances
48 between particular moieties attached to biomolecules.^[15–17]
- 49 11. Atomic Force Microscopy (AFM) can be used to obtain mechanical or
50 electrical response functions regarding particular systems.^[18]
- 51 12. Ion-mobility mass spectrometry (IM-MS) yields collision cross sections
52 that can be used to obtain size and shape information on biomolecules in
53 the gas phase.^[19]

54 A number of quantities Q , which are related to the mentioned measurement
55 techniques and for which values can be more or less directly measured, are
56 listed in Table 1.

57 In order to extend our understanding of biomolecular processes at the atomic
58 level of resolution the measured values of Q^{exp} of an observable quantity $Q(\vec{r}^N, \vec{p}^N)$
59 that depends on the Cartesian coordinates $\vec{r}^N \equiv (\vec{r}_1, \vec{r}_2, \dots, \vec{r}_N)$ and conjugate
60 momenta $\vec{p}^N \equiv (\vec{p}_1, \vec{p}_2, \dots, \vec{p}_N)$ of the N atoms in the molecular system, must
61 be converted into spatial, structural and dynamical information in terms of the
62 statistical-mechanical phase-space variables \vec{r}^N and \vec{p}^N of the system. The quan-
63 tities $Q(\vec{r}^N, \vec{p}^N)$ listed in Table 1 primarily depend on the atomic coordinates
64 \vec{r}^N . For this reason and for simplicity, the dependence of Q on the momenta
65 \vec{p}^N will not be discussed further.

66 If the molecules in the system move, the value $Q(\vec{r}^N(t))$ becomes depen-

Table 1: Some quantities Q and their dependence on configuration $\vec{r}^N = (\vec{r}_1, \vec{r}_2, \dots, \vec{r}_N)$ of N particles. θ denotes a torsional angle or an angle between a bond in a molecule and the direction of an applied magnetic field, and $\vec{\mu}$ an electronic transition dipole. $\hat{r} = \vec{r}/r$.

Quantity $Q(\vec{r}^N)$	Dependence of $Q(\vec{r}^N)$ on configuration \vec{r}^N of N particles
Structure factors (amplitudes)	$F_{hkl}(\vec{r}^N)$
NOE intensities	$I_{i,j}(\vec{r}^N)$
distances	r_{ij}^{-p} with $p = 3$ or 6
3J -coupling constants	$^3J_{i,j} = a \cos^2 \theta_{i,j} + b \cos \theta_{i,j} + c$
Residual dipolar couplings	$D_{i,j} = a \cos^2 \theta_{i,j} + b$
S^2 order parameters	$S_{i,j}^2 = \frac{1}{2} \left\{ 3 \sum_{\alpha=1}^3 \sum_{\beta=1}^3 \left\langle \frac{\hat{r}_{ij\alpha} \hat{r}_{ij\beta}}{r_{ij}^3} \right\rangle^2 - \left\langle \frac{1}{r_{ij}^3} \right\rangle^2 \right\}$
Chemical shielding	$\sigma(\vec{r}^N)$
FRET efficiencies	$E_{DA} = E(\vec{\mu}_D, \vec{\mu}_A, r_{DA})$
CD spectra	$I(\lambda; \vec{r}^N)$

67 dent on time. Depending on the type of measurement technique used, a mea-
68 sured value of Q^{exp} of an observable quantity Q corresponds to an average over
69 molecules in the system and over the time period of the measurement. Different
70 experimentally measured quantities will be sensitive to different time windows.
71 For example, NMR implicitly performs time averaging due to the finite dura-
72 tion of the various radio-frequency pulses and mixing times as well as the nature
73 of the specific phenomenon measured, e.g. magnetisation transfer. In a NMR
74 Nuclear Overhauser Effect (NOE) measurement the time window of sensitivity
75 depends on the rotational tumbling time and thus on the size and the solvent
76 environment of the molecule of interest. Such a time window of sensitivity may
77 in some cases be used to infer the time scale of particular molecular motions.

78 However, for most biomolecular systems studied the measured values Q^{exp} are
79 used to derive structural information on biomolecules that is independent of
80 time.

81 A value $Q(\vec{r}^N)$ for a particular configuration \vec{r}^N can be calculated using
82 theory that connects an observable Q to electronic and nuclear coordinates and
83 properties of molecules. For example, the intensity I_{hkl} of a diffracted X-ray
84 beam as function of the crystal lattice or Miller indices h , k and l can be
85 obtained as the square of the structure factor amplitude $F_{hkl}(\vec{r}^N)$ that results
86 from a spatial Fourier transform of the electron density of the system.^[20] NOE
87 intensity $I_{i,j}(\vec{r}^N)$ for a pair of nuclear spins belonging to atoms i and j can be
88 obtained by a relaxation matrix calculation.^[21] $^3J_{i,j}$ -couplings between atoms
89 i and j connected by three covalent bonds can be estimated using the Karplus
90 relation^[22] with its coefficients a , b and c commonly empirically determined
91 from chosen data sets of 3J -coupling values and molecular structures.^[23] The
92 chemical shielding $\sigma_i(\vec{r}^N)$ for an atom i is related to the local electron density
93 which can in principle be calculated quantum-chemically. However for proteins
94 empirical relations are currently commonly used.^[24] These are based on fitting
95 calculated chemical shifts to measured ones for a chosen set of molecules and
96 structures.

97 In order to obtain a value $Q(\vec{r}^N)$ from a molecular configuration \vec{r}^N , gener-
98 ally the following procedure is followed.^[25]

- 99 1. The atomic, electronic or nuclear degrees of freedom \vec{r} of the solute or sol-
100 vent essential to describe the relation $Q(\vec{r}^N)$ are determined. For example,
101 a $^3J_{i,j}$ -coupling between the spins of atoms i and j depends primarily on
102 the torsional angle $\theta_{i,j}$ defined by the four covalently connected atoms

103 $i - k - l - j$.

104 2. A relation $Q(\vec{r}^N)$ between Q and a particular configuration \vec{r} is postulated.

105 This may be based on quantum-chemical theory, empirical relations or a
106 combination of both. For example, the form of the Karplus relation $Q(\vec{r}) =$
107 ${}^3J(\theta)$ follows from quantum-chemical considerations, but its parameters
108 a , b and c are empirically determined. Such a relation or function $Q(\vec{r})$
109 generally involves a variety of assumptions and approximations, which
110 determine its accuracy.^[26]

111 3. An empirical interaction function $V^{\text{phys}}(\vec{r}^N)$ or force field that governs
112 the statistical-mechanical distribution of configurations \vec{r}^N is then used to
113 generate an ensemble of configurations \vec{r}^N with probability of occurrence
114 $P(\vec{r}^N)$. For example, at constant temperature T a Boltzmann-weighted
115 ensemble with configurational probability

$$P(\vec{r}^N) = \frac{\exp(-V^{\text{phys}}(\vec{r}^N)/(k_{\text{B}}T))}{\int \exp(-V^{\text{phys}}(\vec{r}^N)/(k_{\text{B}}T)) d\vec{r}^N}, \quad (1)$$

116 is normally required, where k_{B} denotes Boltzmann's constant.

117 4. Either equations of motion for the atoms in the system can be integrated
118 as function of time^[27] or an alternative method to generate an ensemble
119 of Boltzmann-weighted configurations can be chosen.

120 5. The ensemble average

$$\langle Q \rangle_{\text{sim}} \equiv \langle Q \rangle_{\vec{r}^N} \equiv \int Q(\vec{r}^N) P(\vec{r}^N) d\vec{r}^N, \quad (2)$$

121 is then calculated and compared to the measured value $Q^{\text{exp}} \equiv \langle Q \rangle_{\text{time,space}}$,

122 which is normally an average, over time and a set of molecules.

123 If $V^{\text{phys}}(\vec{r}^N)$ and $Q(\vec{r}^N)$ are correct, i.e. they correspond exactly to reality,
124 and the sampling of configurations \vec{r}^N is infinite, one would find,

$$\langle Q \rangle_{\text{sim}} = Q^{\text{exp}}, \quad (3)$$

125 that is there is an exact match between simulation and experiment. Assuming
126 Q^{exp} is correct further experiments would no longer be needed. However, in
127 practice these four conditions are rarely satisfied. This leads to the question how
128 to derive a structure or rather a set of structures, \vec{r}^N , from a set of measured
129 values Q^{exp} .

130 When trying to derive a structure \vec{r}^N from Q^{exp} values the following prob-
131 lems are generally encountered.

- 132 1. For most biomolecular systems the number of independent Q^{exp} values
133 available is much smaller than the number of degrees of freedom of the
134 system. This means the problem is underdetermined.
- 135 2. Q^{exp} values contain uncertainty or error.
- 136 3. It is not possible to fully account for averaging over time and space inherent
137 in the experimental measurement: inversion of the averaging operation in
138 Eq. 2 is impossible.
- 139 4. The function $Q(\vec{r})$ is not known or the accuracy of the function is uncer-
140 tain.
- 141 5. The inverse $\vec{r}(Q)$ of the function $Q(\vec{r})$ may not exist or if it does, it may
142 be multiple-valued, as in the case of the Karplus relation.

143 6. The sampling of configurations \vec{r}^N must be biased, i.e. guided towards
144 Q^{exp} . This is especially challenging if the inverse $\vec{r}(Q)$ of $Q(\vec{r})$ is multiple-
145 valued.

146 Over the past decades, this structure, i.e. ensemble determination problem
147 has been approached in a variety of ways. While some of these approaches have
148 a sound physical basis, others are highly speculative or ad-hoc.

149 Since for all but the smallest molecules in crystalline form the number of val-
150 ues of observable quantities in a biomolecular system is much smaller than the
151 number of degrees of freedom of the system, virtually all procedures to derive
152 structure \vec{r}^N from Q^{exp} values rely on an atomic model, i.e. a function specify-
153 ing likely structural parameters of a system. This normally takes the form of an
154 atomic interaction function $V^{\text{phys}}(\vec{r})$ that yields low-energy values for configura-
155 tions that are physically most reasonable. Such a function can be volume-based
156 (only atom-atom contact interactions, no long-range electrostatic or hydrogen
157 bonding terms), such as used in the structure refinement program X-PLOR.^[28]
158 The X-PLOR refinement function contains bond-length, bond-angle, torsional-
159 angle and atomic volume parameters, but no electrostatic terms. The quality
160 of the set or ensemble of structures derived from the set of Q^{exp} values depends
161 on (i) the number and accuracy of the Q^{exp} values; (ii) the accuracy of the
162 function $V^{\text{phys}}(\vec{r}^N)$ used in the structure determination procedure; and (iii) the
163 validity and accuracy of the procedure that was used to derive a structure \vec{r}^N
164 from the Q^{exp} values.

165 A common way to derive a structure from measured data is to represent the
166 latter by a biasing or restraining function $V^{Q,\text{restr}}(Q(\vec{r}^N); Q^\circ)$ that limits the
167 deviation of $Q(\vec{r}^N)$ from the measured target value $Q^\circ = Q^{\text{exp}}$. This function

168 $V^{Q,\text{restr}}$ is added to the physical interaction function V^{phys} used to describe the
169 system,

$$V(\vec{r}^N) = V^{\text{phys}}(\vec{r}^N) + V^{Q,\text{restr}}(Q(\vec{r}^N); Q^{\circ}). \quad (4)$$

170 This function is then minimised or used in a simulation. In the early days of
171 biomolecular structure determination, the function $V(\vec{r}^N)$ was minimised with
172 respect to a variation of the structure \vec{r}^N .^[29] Because molecular dynamics
173 (MD) simulation is able to overcome energy barriers in the order of $k_{\text{B}}T$, energy
174 minimization was replaced by MD simulation first in structure determination
175 based on NMR data (1983)^[30,31] and later in structure determination based on
176 X-ray diffraction data (1987).^[32,33] Currently, MD simulation of biomolecular
177 systems in which a biasing or restraining function $V^{Q,\text{restr}}$ is added to V^{phys} in
178 order to bias the sampling towards Q^{exp} values is the standard method to derive
179 structural information from measured data.^[34] The various procedures and
180 techniques that can be applied lead to the question how best to use experimental
181 data in biomolecular structure refinement. When using experimental data to
182 bias the sampling, the following issues should be considered.

- 183 1. Whether the experimental information corresponds to observed as opposed
184 to derived data.
- 185 2. Whether the number of experimental data Q^{exp} is sufficient.
- 186 3. The accuracy of the (physical) force field $V^{\text{phys}}(\vec{r}^N)$.
- 187 4. The accuracy and consistency of the experimental data Q^{exp} .
- 188 5. The accuracy of the function $Q(\vec{r})$ relating observable Q to structure \vec{r} .
- 189 6. The choice of biasing function $V^{Q,\text{restr}}(\vec{r})$ to guide $\langle Q \rangle_{\text{sim}}$ towards Q^{exp} .

- 190 7. The weighting of the biasing function $V^{Q,\text{restr}}(\vec{r})$ relative to that of the
191 force field energy $V^{\text{phys}}(\vec{r}^N)$ in $V(\vec{r}^N)$.
- 192 8. How to deal with averaging over time and space.
- 193 9. How to bias the sampling especially when the inverse $\vec{r}(Q)$ of the function
194 $Q(\vec{r})$ is multiple-valued.
- 195 10. How to ensure Boltzmann sampling or weighting of configurations \vec{r} .

196 In the present manuscript these ten issues are discussed using examples per-
197 taining to various observables Q . The examples were chosen from our own work
198 and mainly regard observables Q for which values can be determined through
199 NMR spectroscopy. The examples given mostly relate to small molecules and
200 peptides in order to avoid conformational sampling deficiencies that hamper the
201 convergence of averages in larger systems such as proteins when the conforma-
202 tional space that is to be sampled is too large. It should be noted, however,
203 that the issues they illustrate also relate to all other observable quantities and
204 measurement techniques mentioned above. The impact of choices made on the
205 structural interpretation of experimental data is evaluated, and a list of choices
206 to be avoided in structure refinement based on measured data is presented. Fi-
207 nally, the consequences of the different ways to derive biomolecular structure
208 from experimental data for the management of structural databases such as the
209 Protein Data Bank (PDB)^[35] are considered.

210 **2 Deriving structure from measured values of observables**

211 The modelling techniques described are based on statistical mechanics, in which
212 \vec{r}^N corresponds to a configuration. This applies to the solute and solvent com-

213 bined. By grouping together slightly different configurations of a solute molecule
214 that show particular geometric features, we refer to such a set of configurations
215 as a conformer. For example, all configurations of a solute that possess slightly
216 varying bond lengths and bond angles and torsional angles within a limited
217 range, e.g. 30° from a minimum-energy angle, may be considered a single con-
218 former. This is done when the differences between the configurations that belong
219 to this conformer are irrelevant for the phenomenon that is observed.

220 2.1 Use of observed as opposed to derived “experimental” 221 data

222 Quantities $Q(\vec{r}^N)$ that can be calculated from configurations \vec{r}^N can be either
223 observable quantities $Q^{\text{obs}}(\vec{r}^N)$ that are directly measurable in an experiment,
224 or non-observable quantities $Q^{\text{der}}(\vec{r}^N)$ that are derived from $\langle Q^{\text{obs}} \rangle_{\text{exp}}$ values
225 by applying a given procedure, f , based on various assumptions and approxi-
226 mations.^[36]

$$Q^{\text{der}} = f\left(\langle Q^{\text{obs}} \rangle_{\text{exp}}\right). \quad (5)$$

227 Peak location and intensity from X-ray diffraction or NMR spectroscopic ex-
228 periments are examples of observable quantities Q^{obs} , whereas molecular struc-
229 ture, torsional angles, or NMR order parameters are examples of derived quanti-
230 ties Q^{der} . The latter reflect to some extent the assumptions and approximations
231 associated with the procedure f used to convert $\langle Q^{\text{obs}} \rangle_{\text{exp}}$ into Q^{der} . Depend-
232 ing on the validity of the assumptions and the reliability of the approximations
233 involved, the obtained Q^{der} values may in reality carry little experimental infor-
234 mation. If so their use in the restraining function may lead to a flawed structural

235 interpretation of the $\langle Q^{\text{obs}} \rangle_{\text{exp}}$ values. In principle, only values Q° of observable
236 quantities Q^{obs} should be used in the restraining function.

237 An example of the effect of using a flawed procedure f to derive structural
238 restraints can be found in Ref. 37. NMR measurements of a β -octa-peptide
239 solvated in methanol yielded 40 NOE atom-atom distance bounds and 12 3J -
240 coupling constants.^[38] Using single-structure refinement and the volume-based
241 interaction function V^{phys} of the structure refinement program X-PLOR,^[28] 20
242 NMR model structures were generated. These suggested the peptide adopted
243 a hitherto unknown 2_8 - P -helical structure.^[38] MD simulation of this peptide
244 in methanol starting from an extended structure and without any experimental
245 restraints, viz. $V^{Q,\text{restr}} = 0$, resulted in a structural ensemble dominated by the
246 well-known 2.5_{12} - P -helix.^[37] Since the set of 20 NMR model structures and the
247 ensemble of MD generated trajectory structures showed virtually no overlap,^[37]
248 one might naively conclude that “the simulation did not match experiment”.
249 However, this would be incorrect. The 3J -couplings and NOE distances calcu-
250 lated from the MD trajectories showed an average deviation of 0.44 Hz and two
251 NOE bound violations of 0.05 nm. The 20 NMR model structures, which were
252 derived using the NOE bounds, showed an average deviation of 0.57 Hz and, as
253 expected, no NOE bound violations.

254 This example shows that derived quantities Q^{der} such as molecular structures
255 should not be relied on to validate simulations. In particular, they should not
256 be used in a restraining function $V^{Q,\text{restr}}$. For some non-observable, derived
257 quantities Q^{der} the assumptions and approximations involved in the procedure
258 f relating Q^{der} to $\langle Q^{\text{obs}} \rangle_{\text{exp}}$ are of less importance than in the case sketched
259 above. For example, in the absence of spin diffusion, NOE distance bounds

260 to r_{ij}^{-p} with $r_{ij} = |\vec{r}_i - \vec{r}_j|$ and $p = 3$ or 6 , may represent the NOE intensities
261 $I_{i,j}(\vec{r}^N)$ well. Another example are S^2 NMR order parameters, which are derived
262 from spectra based on well defined assumptions. The S^2 order parameter is
263 equal to the long-time plateau value of an auto-correlation function involving
264 the vector connecting two nuclei. Thus S^2 values cannot be related to a single
265 configuration \vec{r}^N , but only to an average, see Table 1. When using such derived
266 quantities in a restraining function, one must make allowance for the uncertainty
267 introduced into the Q^{der} values due to the procedure f .

268 2.2 Insufficient experimental data

269 If the number of items of experimental data on a given biomolecular system is
270 lower than the number of its degrees of freedom, or if these items are correlated,
271 it may not be possible to uniquely determine the conformation of a molecule
272 that dominates the configurational ensemble. An example is the β -octa-peptide
273 discussed in the previous section.^[37] The 40 NOE distance bounds and 12 3J -
274 couplings are insufficient to uniquely determine the dominant conformer. The
275 set of 20 NMR model structures is distinct from the MD trajectory. While
276 the latter is dominated by a 2.5_{12} - P -helix, the former corresponds to a 2_8 - P -
277 helix.^[38] Other examples can be found in Refs. 39 and 40. The less independent
278 experimental values of Q^{obs} that are available, the more the set or ensemble of
279 structures derived from them will depend on the quality of the molecular model
280 or force field $V^{\text{phys}}(\vec{r}^N)$ used.

281 **2.3 Accuracy of the (physical) force field**

282 Of all the experimental techniques mentioned in the Introduction, high-resolution
283 X-ray diffraction of crystals yields the highest information density. A protein
284 may easily give rise to thousands or ten thousands of individual reflections. Due
285 to this abundance of restraints in $V^{Q,\text{restr}}$, the force field V^{phys} used in crystallo-
286 graphic structure refinement can be rather simple, e.g. containing only covalent
287 bond-length, bond-angle and torsional-angle terms in conjunction with short-
288 range (repulsive or van der Waals) interaction terms. The lack of long-range
289 electrostatic and hydrogen-bonding interaction terms is assumed to be compen-
290 sated by the large number of structure-factor amplitude restraints in $V^{Q,\text{restr}}$.
291 However, other experimental techniques have a much lower information density.
292 This means inappropriate structures may be generated if an inadequate force
293 field is used. This can again be illustrated by the case of the β -octa-peptide
294 discussed in the previous sections. The use of a simple volume-based force field
295 lacking electrostatic or hydrogen-bonding interaction terms in vacuo and the
296 requirement that the experimental data be satisfied by a single-structure during
297 the refinement procedure led to a 2_8 - P -helix be proposed. This is despite the
298 fact that such a helix had not been observed previously. In contrast the ensemble
299 generated by MD simulation using a thermodynamically calibrated force field for
300 this peptide in solution indicated the well-known 2.5_{12} - P -helix was dominant.
301 Based on the performance of the latter force field regarding the reproduction
302 of folding equilibria of peptides,^[41,42] the 2.5_{12} - P -helix is much more likely to
303 represent reality than a 2_8 - P -helix. The other experimental techniques men-
304 tioned in the Introduction have an even lower information density than NMR
305 spectroscopy. This makes the quality of the molecular model or force field used

306 in the derivation of structure from such experimental data progressively more
307 important and often the dominant factor.

308 2.4 Accuracy and consistency of the experimental data

309 In many cases the process of deriving structural information from experimental
310 data will indicate that the experimental data contain inaccuracies or inconsis-
311 tencies.^[43] In the case of redundant experimental data, simulation techniques
312 may be used to trace the source of these inconsistencies. An example is the
313 investigation of 103 ${}^3J_{N-H\beta}$ and 94 ${}^3J_{H\alpha-H\beta}$ coupling data for the 107 residue
314 FK506-binding protein FKBP.^[43] The 3J -couplings related to the χ_1 torsional
315 angle of the residues. These were separated into χ_1 angles for which 0, 1, 2, 3 or
316 4 measured 3J -couplings were available. Using Karplus relations to connect the
317 χ_1 -angle to the different 3J -couplings, it was possible to examine whether there
318 was a single χ_1 -angle value that would reproduce, allowing for an uncertainty
319 of 1 Hz in the Karplus relations, the different experimental 3J -coupling values
320 belonging to that torsional angle. No single χ_1 -angle value could satisfy all the
321 available data in 23 of the 63 residues.^[43] Thus, a single-structure refinement
322 procedure based on the 3J -couplings for these 23 residues would result in erro-
323 neous χ_1 -angle values. Because experimental 3J -couplings are averages and the
324 Karplus relation is non-linear in χ_1 , averaging of the 3J -couplings may lead to
325 ranges of χ_1 -angle values that are compatible with the observed 3J -couplings:
326 for 18 of the 23 residues such a range of χ_1 -angle values could be found. How-
327 ever, for 5 of these χ_1 -angles the only conclusions that could be drawn were that
328 either one or more of the experimental 3J -couplings for such a χ_1 -angle must
329 be incorrect or the use of the Karplus relation was inappropriate.^[43]

330 So, if redundant experimental data are available for a particular degree of
331 freedom or for a small set of degrees of freedom, these may be used to detect
332 inconsistencies between different experimental data, which in turn hint at mis-
333 takes in the experiment, e.g. in the assignment of NOE signals to atom pairs or
334 in the structures. We note, however, that experimental data $\langle Q \rangle_{\text{exp}}$ from dif-
335 ferent measurements may seem conflicting when they result from different time
336 windows of sensitivity of the different measurements. In addition, compensa-
337 tion of errors may occur. Unfortunately, most of the experimental techniques
338 mentioned in the Introduction rarely yield redundant experimental data for a
339 particular degree of freedom or a set of degrees of freedom of a system. Gener-
340 ally, indications that some or all of the experimental data may be problematic
341 becomes evident when structures display a very high internal potential energy
342 as a consequence of forcing the molecule to satisfy a set of inappropriate or
343 conflicting restraints Q^{exp} , see Section 2.7.

344 2.5 Accuracy of the function $Q(\vec{r})$ relating observable Q to 345 structure

346 The function $Q(\vec{r})$ that relates an observable quantity Q to a structure \vec{r} even
347 if postulated based on quantum-chemical theory or a set of physically reason-
348 able assumptions will still involve approximations and rely on a given set of
349 parameters and parameter values. For example, a 3J -coupling can be related
350 to a particular torsional angle θ through the Karplus relation (Table 1). The
351 parameters a , b and c of the Karplus relation depend on the atoms involved
352 in the 3J -coupling and their covalently bound neighbours. These are normally
353 calibrated using a set of measured 3J -couplings thought to be associated with

354 particular θ -values in a given structure or molecular fragment. Figure 1 shows
355 various parametrisations available in the literature. As can be seen, they lead to
356 differences of up to 3 Hz for particular angles θ .^[23,43,44] Although 3J -couplings
357 generally show a low experimental uncertainty, lower than 0.1 Hz, use of the
358 Karplus relation may introduce an uncertainty of 1-2 Hz. Again allowance
359 should be made for such uncertainties when using such data in structure deter-
360 mination.

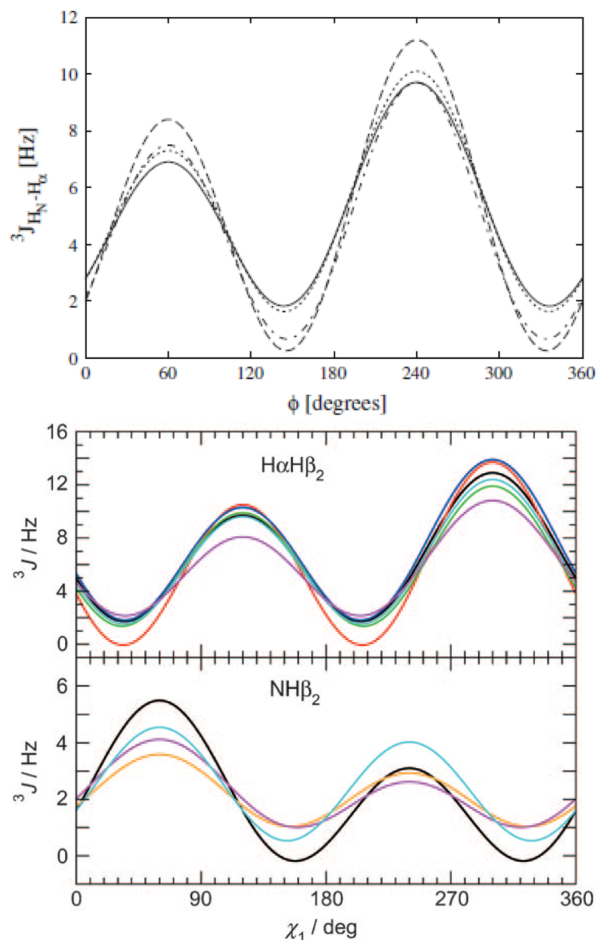


Figure 1: The variety of Karplus relations for the ${}^3J_{H_N H_\alpha}$ - (upper panel), ${}^3J_{H_\alpha H_\beta^-}$, and ${}^3J_{H_N H_\beta}$ -coupling values (lower panels) describing the 3J -coupling as function of the backbone ϕ -angle and side chain χ_1 -angle of an amino acid. Upper panel:^[44] The curves shown are from (solid line) Pardi et al.,^[46] (dashed line) Brüschweiler and Case,^[47] (dotted line) Wang and Bax,^[48] and (dash-dotted line) Schmidt et al.^[49] Lower panels:^[43] The curves shown are from (black) DeMarco et al.,^[50,51] (red) Abraham et al.,^[52] (green) Deber et al.,^[53] (blue) Cung et al.,^[54] (cyan) Fischman et al.,^[55] (magenta) Pérez et al.^[56] (NMR-based parameterization), and (orange) Pérez et al.^[56] (X-ray-based parameterization).

361 If different types of 3J -couplings related to a particular torsional angle can
 362 be measured with high precision,^[45] leading to a larger number of independent
 363 data points, it is possible to investigate the effect of varying (additional) model

364 parameters on the fit of 3J -couplings calculated using the Karplus relation to
 365 measured 3J -coupling values for a particular protein.^[45] However, increasing
 366 the number of changeable parameters of a model may also mask the approxima-
 367 tions and assumptions that affect the calculated 3J -coupling values. By allowing
 368 variations in parameters of the function $Q(\vec{r})$, the force field used or the struc-
 369 tures \vec{r}^N used to calculate Q -values, the inability of a particular function $Q(\vec{r})$
 370 to match all Q^{exp} values may be alleviated, but not eliminated.

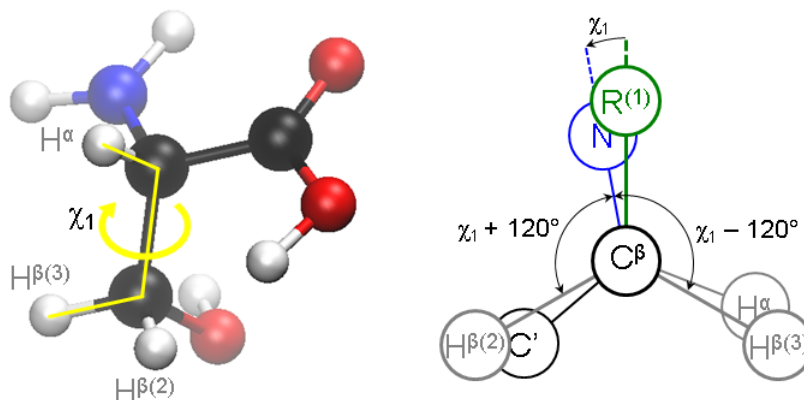


Figure 2: The relationship between the amino-acid side-chain torsional angle χ_1 , the hydrogen atom H_α bound to the C_α atom and the hydrogen atoms $H_{\beta 2}$ and $H_{\beta 3}$ bound to the C_β atom. Nitrogen: blue; oxygen: red; carbon: black; hydrogen: grey; remaining atoms in side chain (R): green.

371 Other functions $Q(\vec{r})$ will also suffer from inaccuracy or uncertainty. When
 372 calculating NMR chemical shifts for a variety of MD trajectory structures of a
 373 β -hepta-peptide in solution using semi-empirical quantum-chemical theory, the
 374 chemical shifts obtained displayed little sensitivity to the underlying conforma-
 375 tion compared to the variation due to different quantum-chemical approxima-
 376 tions.^[57] This can make the use of chemical shift values in structure refinement
 377 unreliable. Despite this chemical shift data in combination with empirical re-

378 lations between chemical shift and secondary structure is increasingly used in
379 protein structure refinement.^[58, 59]

380 Förster Resonance Energy Transfer (FRET) depends on the distance and
381 relative orientation of the transition dipole vectors of a donor and an acceptor
382 fluorophore. In most applications, isotropic orientational tumbling of the dipole
383 vectors is assumed. This leads the FRET signal to be interpreted in terms of
384 a single distance dependence. This greatly limits the accuracy of the distance
385 distributions that are derived from measured FRET intensities.^[15, 16]

386 2.6 Choice of the biasing function $V^{Q,\text{rest}}(\vec{r}^N)$ to guide $\langle Q \rangle_{\text{sim}}$ 387 towards Q^{exp}

388 The goal of adding a biasing or restraining function $V^{Q,\text{restr}}(Q(\vec{r}^N); Q^\circ)$ to the
389 physical interaction function $V^{\text{phys}}(\vec{r}^N)$ representing the molecular model is to
390 keep the value of the quantity $Q(\vec{r}^N)$ close to the target value Q° , which is
391 related to the experimentally observed value Q^{exp} ,

$$Q^\circ = Q^{\text{exp}} + \Delta Q^{\text{corr}}. \quad (6)$$

392 ΔQ^{corr} represents a correction to the Q^{exp} value that depends on the quantity
393 Q and its relation to molecular structure (see below).

394 The functional form of $V^{Q,\text{restr}}(Q; Q^\circ)$ should meet the following conditions
395 (Figure 3).^[60]

396 1. It should be a continuous function with a continuous derivative in order
397 to obtain continuous energy and force in an MD simulation.

398 2. It should have two or three ranges that display a different behaviour of

399 the restraining force as function of Q .

400 (a) a flat-bottom region of zero restraining energy and zero restraining
401 force around Q° , e.g. for $[Q^\circ - \Delta Q^{\text{fb}}, Q^\circ + \Delta Q^{\text{fb}}]$.

402 (b) a region in which the restraining energy and force increase with in-
403 creasing deviation of Q from Q° , e.g. for $[Q^\circ \pm \Delta Q^{\text{fb}}, Q^\circ \pm \Delta Q^{\text{fb}} \pm$
404 $\Delta Q^{\text{h}}]$.

405 (c) a region in which the restraining force reaches a constant maximum
406 in order to avoid too large restraining forces distorting the molecular
407 structure if the deviation of Q from Q° gets large, e.g. for $[Q^\circ \pm$
408 $\Delta Q^{\text{fb}} \pm \Delta Q^{\text{h}}, \pm\infty]$.

409 3. The restraining can be either purely attractive, that is apply when $Q > Q^\circ$,
410 e.g. as in the case of NOE atom-atom distance upper bound restrain-
411 ing,^[30,31] purely repulsive, that is apply when $Q < Q^\circ$, e.g. when repre-
412 senting the absence of NOE intensity in a spectrum^[61] or a combination
413 of attractive and repulsive restraining, e.g. as in the case of 3J -coupling
414 restraining.^[62]

The simplest representation of such a three-stage restraining function is a flat-bottom, then harmonic, then linear function

$$\begin{aligned}
V^{Q,\text{restr}}(Q; K^{Qr}, Q^\circ, \Delta Q^{\text{fb}}, \Delta Q^{\text{h}}) &= \frac{1}{2} K^{Qr} (Q - Q^\circ - \Delta Q^{\text{fb}})^2 \\
&\cdot H(Q; Q^\circ + \Delta Q^{\text{fb}}) [1 - H(Q; Q^\circ + \Delta Q^{\text{fb}} + \Delta Q^{\text{h}})] \\
&+ K^{Qr} \left(Q - Q^\circ - \Delta Q^{\text{fb}} - \frac{1}{2} \Delta Q^{\text{h}} \right) \Delta Q^{\text{h}} \\
&\cdot H(Q; Q^\circ + \Delta Q^{\text{fb}} + \Delta Q^{\text{h}}).
\end{aligned}$$

for $Q > Q^\circ$, (7)

and

$$\begin{aligned}
V^{Q, \text{restr}}(Q; K^{Qr}, Q^\circ, \Delta Q^{\text{fb}}, \Delta Q^{\text{h}}) &= \frac{1}{2} K^{Qr} (Q - Q^\circ + \Delta Q^{\text{fb}})^2 \\
&\cdot [1 - H(Q; Q^\circ - \Delta Q^{\text{fb}})] H(Q; Q^\circ - \Delta Q^{\text{fb}} - \Delta Q^{\text{h}}) \\
&- K^{Qr} \left(Q - Q^\circ + \Delta Q^{\text{fb}} + \frac{1}{2} \Delta Q^{\text{h}} \right) \Delta Q^{\text{h}} \\
&\cdot [1 - H(Q; Q^\circ - \Delta Q^{\text{fb}} - \Delta Q^{\text{h}})] \\
&\text{for } Q < Q^\circ, \tag{8}
\end{aligned}$$

where $\Delta Q^{\text{fb}} \geq 0$ and $\Delta Q^{\text{h}} \geq 0$, and the Heaviside step function $H(x; x_0)$ is defined by

$$\begin{aligned}
H(x; x_0) &= 0 \quad \text{for } x < x_0 \\
&= 1 \quad \text{for } x \geq x_0. \tag{9}
\end{aligned}$$

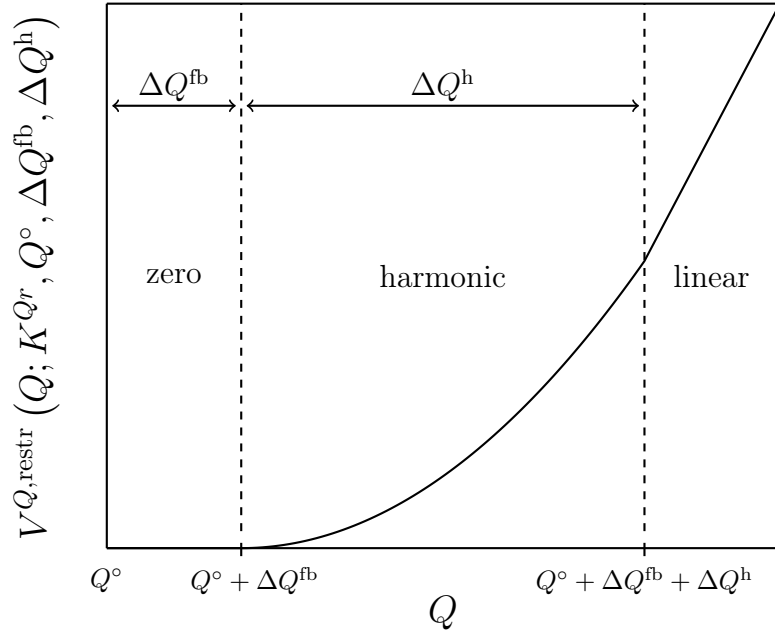


Figure 3: Potential energy term $V^{Q, \text{restr}}(Q; K^{Qr}, Q^\circ, \Delta Q^{\text{fb}}, \Delta Q^{\text{h}})$ for restraining the quantity Q .^[60]

415 The corresponding force along Q , i.e. the negative of the derivative of
 416 $V^{Q, \text{restr}}$ with respect to Q is then for $Q > Q^\circ$,

$$\begin{aligned}
 f^{Q, \text{restr}} &= 0 && \text{for } Q < Q^\circ + \Delta Q^{\text{fb}} \\
 &= -K^{Qr} (Q - Q^\circ - \Delta Q^{\text{fb}}) && \text{for } Q^\circ + \Delta Q^{\text{fb}} \leq Q \leq Q^\circ + \Delta Q^{\text{fb}} + \Delta Q^{\text{h}} \\
 &= -K^{Qr} \Delta Q^{\text{h}} && \text{for } Q > Q^\circ + \Delta Q^{\text{fb}} + \Delta Q^{\text{h}}, \quad (10)
 \end{aligned}$$

and for $Q < Q^\circ$,

$$\begin{aligned}
 f^{Q,\text{restr}} &= 0 && \text{for } Q > Q^\circ - \Delta Q^{\text{fb}} \\
 &= -K^{Qr} (Q - Q^\circ + \Delta Q^{\text{fb}}) && \text{for } Q^\circ - \Delta Q^{\text{fb}} - \Delta Q^{\text{h}} \leq Q \leq Q^\circ - \Delta Q^{\text{fb}} \\
 &= +K^{Qr} \Delta Q^{\text{h}} && \text{for } Q < Q^\circ - \Delta Q^{\text{fb}} - \Delta Q^{\text{h}}. \quad (11)
 \end{aligned}$$

417 To obtain the force on particle i , these expressions are multiplied by $\partial Q(\vec{r}^N)/\partial \vec{r}_i$,
 418 the derivative of Q with respect to \vec{r}_i . Of course, other functional forms of the
 419 restraining function that match the basic requirements, such as a Gaussian form,
 420 may be used.

421 In certain NMR experiments, the signals from different hydrogen atoms may
 422 not be distinguishable, e.g. for the three hydrogen atoms of a CH_3 group.
 423 In such a case the (distance) restraint is applied to a so-called pseudo-atom
 424 position.^[63] Such a position is defined in terms of the positions of the hydrogen
 425 atoms and the atoms these are bound to, see Figure 4. The use of pseudo
 426 atoms in the restraining requires corrections ΔQ^{corr} to the distance bounds Q^{exp}
 427 derived from experiment.^[60, 63–65] Figure 4 shows three sets of ΔQ^{corr} values for
 428 NOE distance bounds that are used in structure determination based on NOE
 429 data. An alternative is to use so-called ambiguous restraints, in which the sum
 430 of the $r^{-1/6}$ weighted distances to the different indistinguishable hydrogens is
 431 used in the restraining function.^[66, 67]

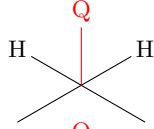
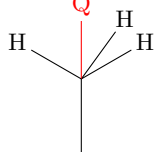
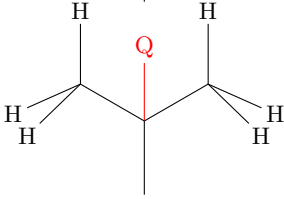
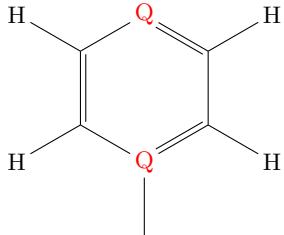
Geometry	Wüthrich (1983) correction (nm)	Fletcher (1996) correction (nm)	GROMOS (1996) correction (nm)
	+0.10	+0.07	+0.09
	+0.15	+0.04	+0.10
	+0.29	+0.15	+0.22
	+0.20	+0.20	+0.21

Figure 4: Geometries and NOE distance bound corrections for pseudo atoms Q , from Wüthrich et al.,^[63] Fletcher et al.,^[65] and van Gunsteren et al., GROMOS.^[64] Fletcher et al.^[65] also apply a multiplicity correction when an NOE signal is known to arise from an averaging over several hydrogens. For an NOE between groups of hydrogen atoms A and B, the upper bound is then scaled by a factor $(n_A n_B)^{\frac{1}{p}}$, where p corresponds to the power (3 or 6) of the distance averaging and n_A and n_B are the number of equivalent hydrogen atoms in each group.

432 **2.7 Weight of the restraining function $V^{Q,\text{restr}}$ relative to**
433 **that of the force field energy V^{phys}**

434 The relative weight of the biasing or restraining term $V^{Q,\text{restr}}$ compared to that
435 of the physical force field term V^{phys} in the total potential energy (Eq. 4) of
436 the system is determined by the size of the force constant K^{Qr} in $V^{Q,\text{restr}}$.

437 The larger the weight or K^{Qr} , the smaller will be the deviation of Q from Q° .
438 However, if the agreement with the experimental data is forced excessively this
439 will result in a higher V^{phys} of the system or a higher internal energy of the
440 molecule to which the restraints are applied. The molecule will simply adopt a
441 more strained conformation in order to match the values Q° .

442 An early example and analysis of the effect of restraining too strongly can
443 be found in Ref. 68. While deriving a spatial structure of the headpiece of the
444 protein lac repressor, two structures were obtained, structure I in which a loop
445 was incorrectly folded (total NOE bounds violation: 3.242 nm), and structure
446 II in which the loop was correct (total NOE bounds violation: 0.405 nm).^[61]
447 The structure refinement procedure was based on MD simulation followed by
448 energy minimization (EM) in the presence of 215 attractive NOE atom-atom
449 distance restraints with $K^{Qr} = 17 \cdot 10^3 \text{ kJ mol}^{-1} \text{ nm}^{-2}$, then followed by EM,
450 5 ps of MD and again MD simulation with $K^{Qr} = 0$.^[68] The results in Table 2
451 show that MD refinement with a large weight, $K^{Qr} = 17 \cdot 10^3 \text{ kJ mol}^{-1} \text{ nm}^{-2}$,
452 results in the incorrect structure having lower NOE bound violations, 0.180 nm,
453 than the correct structure, 0.258 nm. However, this results in a higher intra-
454 protein potential energy, $-2996 \text{ kJ mol}^{-1}$ compared to $-3053 \text{ kJ mol}^{-1}$. When
455 the restraining term $V^{Q,\text{restr}}$ is switched off, $K^{Qr} = 0$, the molecule relaxes
456 resulting in much larger NOE bound violations for the structure containing the
457 wrongly folded loop. This example illustrates that molecular structures derived
458 from experimental data should not only match that data, but must also have
459 a low energy. Otherwise one may satisfy the experimental data but obtain a
460 strained, unstable conformation.

Table 2: MD refinement of lac repressor structures.^[68] A set of 215 NOE distance bounds for pairs of atoms were used in the half-harmonic attractive restraining term of the potential energy function with a weight K^{Qr} . The GROMOS biomolecular force field^[76,77] was used as interaction function. The two initial structures I and II of the protein were taken from Ref. 61.

	K^{Qr} (kJ mol ⁻¹ nm ⁻²)	Structure I Loop wrongly folded		Structure II Loop correctly folded		energy (kJ mol ⁻¹)	
		constraint violation (nm)		constraint violation (nm)			
		sum	average	sum	average		
Initial structure		3.242	0.015	-2115	0.405	0.002	-3092
After 5 ps MD + EM	17000	0.180	0.001	-2996	0.258	0.001	-3053
After EM	0	0.523	0.002	-3083	0.461	0.002	-3100
After 5 ps MD + EM	0	3.297	0.015	-3032	1.823	0.008	-3102

461 This leads to the question how low should be the energy of a molecule. This
462 depends on the force field used as V^{phys} and on the size and composition of the
463 molecule. Table 3 presents the intra-protein potential energies of a number of
464 proteins for which a high-resolution X-ray crystal structure is available and for
465 some proteins the structures of which have been derived on the basis of NMR
466 NOE atom-atom distance bounds.^[68] The NOE bound based structures of the
467 lac repressor and of tendamistat both have energies comparable to those of X-
468 ray diffraction based structures of proteins of similar size. In contrast, the last
469 five structures in Table 3 display a relatively high energy, at least about 1000
470 kJ mol^{-1} higher than an X-ray structure of comparable size. This is most likely
471 due to the use of a large value for the restraining force constant K^{Qr} , which
472 lowers the distance bound violations at the expense of a high internal energy.
473 The discrepancy in energy cannot be due to the application of different force
474 fields. The biomolecular force fields used, CHARMM^[78] and GROMOS,^[76,77]
475 consist of comparable interaction terms: while the CHARMM force field yields
476 an energy of $-2247 \text{ kJ mol}^{-1}$ for the crambin X-ray structure, the GROMOS
477 force field yields a comparable value of $-2161 \text{ kJ mol}^{-1}$. The structures of the
478 last five proteins in Table 3 are likely to be highly strained.

479 This example illustrates that reporting and depositing structures of doubtful
480 quality is not a new phenomenon. It is a consequence of the understandable ten-
481 dency of researchers to push experimental and theoretical methods to the limit
482 of what they can do. Currently many structures of large protein assemblies are
483 being determined using electron microscopy and then used as initial structures
484 for MD simulations. Here it is likely that similar problems as in the example
485 given will be found in some of these structures. This state of affairs calls for a

486 continuous effort of the scientific community towards validation of structures in
487 data bases^[79–81] and if needed their improvement.^[81]

488 In this regard it is of interest to track the fate of the highly strained structures
489 discussed above in Table 3. None of the structures reported in Refs. 71–75
490 are currently available in the PDB. So they were either never deposited or
491 removed. The structure of hirudin^[71] was deposited in 1990 and has been
492 updated multiple times, the latest in 2009.

493 If different sets of experimental data, such as NOEs, RDCs, ³*J*-couplings
494 or X-ray structure factor amplitudes are used simultaneously when restraining,
495 there is a need to provide the relative weights of the different sets of Q^{exp} values
496 originating from different experiments.^[82] However, it is not possible to do this
497 in an objective fashion.

Table 3: Energies of X-ray and NMR protein structures.^[68] N_{Qr} : number of distance restraints. K^{Qr} : weight of the restraining function. The GROMOS^[76,77] or the CHARMM^[78] biomolecular force field was used as interaction function. For tendamistat the highest and lowest energy for a set of structures is given.

	Number of residues	Source of data or structure	N_{Qr}	K^{Qr} (kJ mol ⁻¹ nm ⁻²)	Average violation (nm)	Protein energy (kJ mol ⁻¹)	Force field used
aPP	36	X-ray ^[35]				-2180	GROMOS
crambin	46	X-ray ^[35]				-2161	GROMOS
BPTI	58	X-ray ^[35]				-3529	GROMOS
Phospholipase A2	123	X-ray ^[35]				-7848	GROMOS
Lac repressor	51	NMR ^[61]	215	4000	0.003	-3091	GROMOS
tendamistat	74	NMR ^[69]	842			-3140	GROMOS
						-2834	GROMOS
crambin	46	X-ray/NMR ^[70]	240	5000	0.033	-2247	CHARMM
hirudin	56	NMR ^[71]	359	17000	0.016	-1138	CHARMM
Histone H5	79	NMR ^[72]	307	17000	0.015	-1527	CHARMM
CPI	39	NMR ^[73]	309	33000	0.007	-724	CHARMM
phoratoxin	46	NMR ^[74]	331	33000	0.010	-1029	CHARMM
A1-purothionin	46	NMR ^[75]	310	33000	0.023	-498	CHARMM

498 **2.8 Averaging over space and time**

499 Values of Q^{exp} of experimentally measured properties correspond in general to
 500 averages of a quantity Q over both space (or molecules) and time,

$$Q^{\text{exp}} = \left\langle \left\langle Q \right\rangle_{\text{space}} \right\rangle_{\text{time}} \quad , \quad (12)$$

501 where the angular brackets denote averaging (Eq. 2) over the distribution $P(\vec{r}^N)$
 502 of configurations of the system, e.g. the Boltzmann distribution (Eq. 1).

503 Time averaging when restraining was first introduced^[83] in relation to the
 504 use of NOE derived atom-atom (i, j) distance, r_{ij} , information with $Q = r_{ij}^{-3}$ or
 505 $Q = r_{ij}^{-6}$. This was later followed by application to crystallographic structure
 506 factor amplitude restraining,^[84-87] chemical shift restraining,^[58] 3J -coupling re-
 507 straining^[62,88] and S^2 NMR order parameter restraining.^[89] Time-averaging
 508 restraining methods are characterized by two parameters, the force constant
 509 or weight K^{Qr} of the restraining function $V^{Q,\text{restr}}$ and the memory relaxation
 510 time τ_Q representing the time span over which $Q(\vec{r}^N(t))$ is to be averaged. The
 511 time-averaged value of Q is commonly exponentially damped

$$\langle Q(\vec{r}^N) \rangle_t = [\tau_Q (1 - \exp(-t/\tau_Q))]^{-1} \int_0^t \exp(-(t-t')/\tau_Q) Q(\vec{r}^N(t')) dt', \quad (13)$$

512 in the restraining function $V^{Q,\text{restr}}$ (Eqs. 7 and 8) in order to avoid that the
 513 restraining force progressively approaches zero with time.^[83] Time-averaging
 514 restraining has been investigated as function of K^{Qr} and τ_Q for a variety of
 515 systems.^[44,83-96] The force constant K^{Qr} should be taken as small as possible in
 516 order to avoid a restraining bias that puts strain into the molecule, as discussed

517 in the previous section, while being large enough to force $\langle Q \rangle_{\text{time}}$ to be close
 518 to Q° . In addition, too strong restraints may destroy the proper Boltzmann
 519 weighting (Eq. 1) of configurations of the trajectory. The memory relaxation
 520 time should be of the order of the experimental averaging time that determines
 521 Q^{exp} , but at least an order of magnitude shorter than the length of the MD
 522 simulation in order to secure sufficient statistics when averaging $Q(\vec{r}^N(t))$ over
 523 t . In addition, the heating of the system due to the non-conservative force
 524 resulting from the time-averaging restraining term should be small.

525 The averaging over molecules can be accounted for by simulating N_m in-
 526 dependent systems in parallel for which the initial coordinates and momenta
 527 $(\vec{r}^N(t_0), \vec{p}^N(t_0))$ are Boltzmann-distributed in regard to V^{phys} . This is easily
 528 achieved in the case of the momenta \vec{p}^N , because the kinetic energy term in
 529 the Hamiltonian of the system is – in the absence of constraints – quadratic
 530 in the momenta, but rather difficult if not impossible for configurations \vec{r}^N of
 531 a biomolecular system, because the potential energy term of the Hamiltonian
 532 is generally a complex function of the coordinates \vec{r}^N .^[97] Relaxation times of
 533 biomolecular systems generally exceed the time scale accessible using MD sim-
 534 ulation which means that if the initial configurations $\vec{r}^N(t_0)$ of the N_m systems
 535 are not Boltzmann distributed, the averages $\langle Q \rangle$ over the N_m systems are very
 536 likely to reflect non-Boltzmann averaging. Molecule-averaging restraints were
 537 again first introduced in relation to NOE-derived atom-atom distance informa-
 538 tion.^[98,99] Here, the molecule-averaged value of Q ,

$$\langle Q(\vec{r}^N) \rangle_{\text{molecules}} = \sum_{n=1}^{N_m} p_n Q(\vec{r}_n^N), \quad (14)$$

539 with

$$p_n = N_m^{-1}, \quad (15)$$

540 or

$$p_n = \frac{\exp(-V^{\text{phys}}(\vec{r}_n^N)/(k_B T))}{\sum_{n'=1}^{N_m} \exp(-V^{\text{phys}}(\vec{r}_{n'}^N)/(k_B T))}, \quad (16)$$

541 is used, instead of the single molecule value $Q(\vec{r}^N)$, in the restraining function
542 $V^{Q,\text{restr}}$ (Eqs. 7 and 8).^[99] Molecule-averaging restraining methods are also
543 characterized by two parameters, the force constant or weight K^{Qr} of the re-
544 straining function $V^{Q,\text{restr}}$ and the number of molecules or systems N_m over
545 which the averaging is performed.^[100,101] However, problems can arise if the
546 initial configurations $\vec{r}^N(t_0)$ are not Boltzmann distributed over the N_m sys-
547 tems and are given equal weight (Eq. 15) in the average $\langle Q \rangle$ (Eq. 14) unless
548 the simulation is much longer than the system relaxation time. In such cases
549 the average $\langle Q \rangle_{\text{molecules}}$ will only be Boltzmann weighted in the limit of a large
550 number of molecules N_m being used in combination with large restraining force
551 constants K^{Qr} ,^[102-104] and in the absence of noise in the target values Q° .^[105]
552 Unfortunately, neither condition is normally met in practice: (i) The number
553 of systems N_m is generally kept low, i.e. 10-100, in order to minimize the com-
554 putational effort; (ii) the force constant of the restraining function should be
555 as small as possible, as discussed in the previous section; (iii) the experimental
556 data Q^{exp} and the function $Q(\vec{r})$ relating Q to structure are generally not error
557 free.

558 The problem of non-Boltzmann averaging over molecules may be alleviated
559 by applying Hamiltonian replica-exchange techniques^[106-108] that swap config-
560 urations between the N_m molecules or systems based on the Boltzmann proba-

561 bility of the respective configurations.^[109,110]

562 Since most structures of proteins deposited in the protein data bank were
563 determined by single-structure refinement, i.e. not accounting for averaging, one
564 may ask whether accounting for averaging in restraining is important. Unless
565 the configurational ensemble of a molecule is very narrow and centered on a
566 single dominant structure, the average over a Boltzmann distribution (Eq. 2) of
567 structures \vec{r}^N , $\langle Q \rangle_{\vec{r}^N}$, will generally not be equal to the value of Q calculated
568 for a single structure \vec{r}^N

$$\langle Q \rangle_{\vec{r}^N} \neq Q(\vec{r}^N), \quad (17)$$

569 and it will also not be equal to the value of Q calculated for the mean structure
570 $\langle \vec{r}^N \rangle$,

$$\langle Q \rangle_{\vec{r}^N} \neq Q(\langle \vec{r}^N \rangle). \quad (18)$$

571 This is due to the non-linearity of the function $Q(\vec{r}^N)$ with respect to config-
572 urations \vec{r}^N or the non-linearity of the Boltzmann weighting in Eq. 2. As a
573 consequence, single-structure refinement, i.e. not accounting for averaging, may
574 lead to highly distorted structures. An example can be found in Figure 5.^[111]
575 The NOE atom-atom distance bounds derived from NMR experiments on a
576 β -hexa-peptide in methanol turned out to belong to two rather different confor-
577 mations, a left-handed 3_{14} helix and right-handed $2.7_{10/12}$ helix which are both
578 populated in solution. While single-structure refinement led to a distorted and
579 strained structure (Figure 5, middle), MD simulation without any restraints
580 sampled both helices and as a consequence satisfied all NOE bounds.^[111] For
581 proteins, accounting for averaging is particularly important in regard to the
582 mobility of side chains. Single-structure refinement of the protein tendamistat

583 based on NOE data resulted in the side chain of 15 Tyr being placed in a single
 584 conformation that resulted in small NOE bound violations to spatially adjacent
 585 side chains. However, in the structure determined crystallographically of this
 586 protein no electron density was observed for the side chain of 15 Tyr, indicating
 587 it was highly mobile. This illustrates how single-structure refinement can lead
 588 to the underestimation of the structural variability within a molecule.^[112,113]
 589 This applies to X-ray crystallographic single-structure refinement as well as to
 590 that based on NMR.^[86,87]

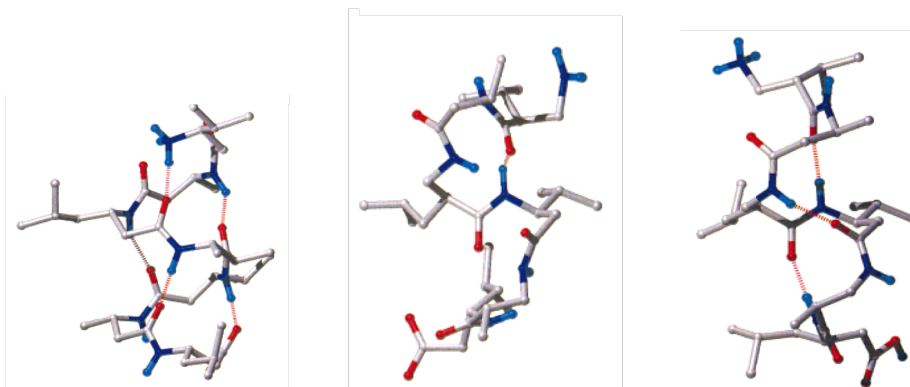


Figure 5: Left: Molecular model for a β -heptapeptide derived from NMR data obtained in methanol. Left-handed 3_{14} -helix with hydrogen bonds: $\text{NH}(i) - \text{O}(i+2)$.^[114] Middle: Molecular model for a β -hexapeptide derived using single-structure refinement using 34 distance restraints derived from NOE NMR data obtained in methanol. Average structure: Right-handed helix in MeOH with only one hydrogen bond: $\text{NH}(4) - \text{O}(1)$. Distorted, due to 3 NOEs characteristic for a left-handed 3_{14} -helix.^[115] Right: Molecular model for this β -hexapeptide derived from NMR data obtained in pyridine. Right-handed helix with hydrogen bonds: $\text{NH}(i) - \text{O}(i+1, i-3)$.^[116] Hydrogen bonds (with a maximum proton - acceptor distance of 0.25 nm and a minimum donor-proton-acceptor angle of 135°) are shown with red dashed lines.

591 When applied to 3J -coupling restraining, the use of the time-average $\langle Q \rangle_t$
 592 (Eq. 13) in the quadratic part of the restraining function $V^{Q,\text{restr}}$ (Eqs. 7, 8)
 593 led to large structural fluctuations.^[96] This is due to the fact that the average
 594 $\langle Q \rangle_t$ of Q can lag behind the instantaneous value of $Q(t)$ of Q . If the restraining

595 force only depends on $\langle Q \rangle_t$, it may drive $Q(t)$ away from Q° in case $\langle Q \rangle_t$ and
 596 $Q(t)$ are at different sides of Q° . This effect was not observed in time-averaging
 597 NOE distance bound restraining because (i) only attractive distance restraints
 598 are used, (ii) the r^{-3} or r^{-6} distance dependence of the NOE signal gives
 599 most weight to short distances when averaging thereby reducing the difference
 600 between $Q(t)$ and $\langle Q \rangle_t$ for short distances $Q(t)$, and (iii) the van der Waals
 601 repulsion between atoms in V^{phys} prevents very small atom-atom distances $Q(t)$
 602 being sampled.

603 The problem of $\langle Q \rangle_t$ lagging behind $Q(t)$ can be solved by using a biquadratic
 604 restraining function $V^{Q,\text{restr}}$ that depends on both $Q(t)$ and $\langle Q \rangle_t$,^[88]

$$\begin{aligned}
 V_{\text{att}}^{Q,\text{restr}}(Q(t), \langle Q \rangle_t; K^{Qr}, Q^\circ, \Delta Q^{\text{fb}}) &= \frac{1}{2} K^{Qr} (Q(t) - Q^\circ - \Delta Q^{\text{fb}})^2 \\
 &\cdot H(Q(t); Q^\circ + \Delta Q^{\text{fb}}) (\langle Q \rangle_t - Q^\circ - \Delta Q^{\text{fb}})^2 \\
 &\cdot H(\langle Q \rangle_t; Q^\circ + \Delta Q^{\text{fb}}), \quad (19)
 \end{aligned}$$

605 for $Q(t) > Q^\circ$ and $\langle Q \rangle_t > Q^\circ$

$$\begin{aligned}
 V_{\text{rep}}^{Q,\text{restr}}(Q(t), \langle Q \rangle_t; K^{Qr}, Q^\circ, \Delta Q^{\text{fb}}) &= \frac{1}{2} K^{Qr} (Q(t) - Q^\circ + \Delta Q^{\text{fb}})^2 \\
 &\cdot [1 - H(Q(t); Q^\circ - \Delta Q^{\text{fb}})] (\langle Q \rangle_t - Q^\circ + \Delta Q^{\text{fb}})^2 \\
 &\cdot [1 - H(\langle Q \rangle_t; Q^\circ - \Delta Q^{\text{fb}})], \quad (20)
 \end{aligned}$$

for $Q(t) < Q^\circ$ and $\langle Q \rangle_t < Q^\circ$. The corresponding force along Q , i.e. the
 negative of the derivative of $V^{Q,\text{restr}}$ with respect to Q , is then for the case

$Q(t) > Q^\circ$ and $\langle Q \rangle_t > Q^\circ$,

$$\begin{aligned}
f^{Q,\text{restr}} &= 0 && \text{for } Q(t) < Q^\circ + \Delta Q^{\text{fb}} \text{ or } \langle Q \rangle_t < Q^\circ + \Delta Q^{\text{fb}} \\
&= -K^{Qr} \left\{ (Q(t) - Q^\circ - \Delta Q^{\text{fb}}) (\langle Q \rangle_t - Q^\circ - \Delta Q^{\text{fb}})^2 \right. \\
&\quad \left. + (Q(t) - Q^\circ - \Delta Q^{\text{fb}})^2 (\langle Q \rangle_t - Q^\circ - \Delta Q^{\text{fb}}) \cdot (1 - \exp(-\Delta t/\tau_Q)) \right\} \\
&&& \text{for } Q(t) > Q^\circ + \Delta Q^{\text{fb}} \text{ and } \langle Q \rangle_t > Q^\circ + \Delta Q^{\text{fb}},
\end{aligned} \tag{21}$$

606 where the derivative of the average $\langle Q \rangle_t$ (Eq. 13) with respect to $Q(t)$ has been
607 expressed in terms of discrete MD time steps Δt ,

$$\frac{\partial \langle Q \rangle_t}{\partial Q(t)} = (1 - \exp(-\Delta t/\tau_Q)). \tag{22}$$

608 For the case $Q(t) < Q^\circ$ and $\langle Q \rangle_t < Q^\circ$ we have

$$\begin{aligned}
f^{Q,\text{restr}} &= 0 && \text{for } Q(t) > Q^\circ - \Delta Q^{\text{fb}} \text{ or } \langle Q \rangle_t > Q^\circ - \Delta Q^{\text{fb}} \\
&= -K^{Qr} \left\{ (Q(t) - Q^\circ + \Delta Q^{\text{fb}}) (\langle Q \rangle_t - Q^\circ + \Delta Q^{\text{fb}})^2 \right. \\
&\quad \left. + (Q(t) - Q^\circ + \Delta Q^{\text{fb}})^2 (\langle Q \rangle_t - Q^\circ + \Delta Q^{\text{fb}}) \cdot (1 - \exp(-\Delta t/\tau_Q)) \right\} \\
&&& \text{for } Q(t) < Q^\circ - \Delta Q^{\text{fb}} \text{ and } \langle Q \rangle_t < Q^\circ - \Delta Q^{\text{fb}}.
\end{aligned} \tag{23}$$

609 To obtain the force on particle i , these expressions must be multiplied by
610 $\partial Q(\vec{r}^N(t))/\partial \vec{r}_i(t)$. Using this biquadratic form, the restraining function only gener-
611 ates a force when both the instantaneous value $Q(t)$ of Q and the time-averaged
612 value $\langle Q \rangle_t$ of Q lie outside the flat bottom of the restraining function $V^{Q,\text{restr}}$.

613 If both, attractive ($Q > Q^\circ$) and repulsive ($Q < Q^\circ$) restraints are applied, the
 614 restraining is harmonic with a flat bottom. If $Q(t)$ and $\langle Q \rangle_t$ both lie outside
 615 the flat bottom and on the same side of Q° , the two terms in Eqs. 21 and 23
 616 originating from the derivatives of $Q(t)$ and $\langle Q \rangle_t$ yield contributions to the force
 617 $f^{Q,\text{restr}}$ that drive $Q(t)$ towards Q° . If $Q(t)$ and $\langle Q \rangle_t$ lie on different sides of
 618 Q° , i.e. $Q(t) < Q^\circ - \Delta Q^{\text{fb}}$ and $\langle Q \rangle_t > Q^\circ + \Delta Q^{\text{fb}}$ or $Q(t) > Q^\circ + \Delta Q^{\text{fb}}$ and
 619 $\langle Q \rangle_t < Q^\circ - \Delta Q^{\text{fb}}$, no restraining force is generated.

620 **2.9 How to deal with a multiple-valued relation between** 621 **structure and restrained quantity Q ?**

622 For some of the quantities Q listed in Table 1, there is more than one configura-
 623 tion \vec{r}^N that corresponds to the same Q value. An example of such a function
 624 $Q(\vec{r}^N)$ is the Karplus relation ${}^3J(\theta)$ that relates a 3J -coupling to a torsional
 625 angle θ , see Figure 1. There are up to four values of the angle θ that may corre-
 626 spond to a single 3J -coupling value. This implies that the inverse of ${}^3J(\theta)$, the
 627 function $\theta({}^3J)$, is multiple-valued. Depending on the initial value of the angle
 628 θ , the restraining function $V^{Q,\text{restr}}$ may drive the angle θ in different directions
 629 to one of the different angle values that correspond to ${}^3J^\circ$. This is illustrated
 630 in the upper panels of Figure 6. If the target ${}^3J^\circ$ -value lies beyond one maxi-
 631 mum or minimum of the function ${}^3J(\theta)$, but not beyond another maximum or
 632 minimum, the restraining function $V^{Q,\text{restr}}$ may drive the angle θ in the wrong
 633 direction. This is illustrated in the lower right panel of Figure 6. Thus if the
 634 function $\vec{r}^N(Q)$ is multiple-valued, the restraining techniques discussed so far
 635 may not lead to structures that match the target Q° values. This problem can
 636 be solved by driving the structure out of the range of θ -values for which $V^{Q,\text{restr}}$

637 has a low energy but where these energies cannot become zero.

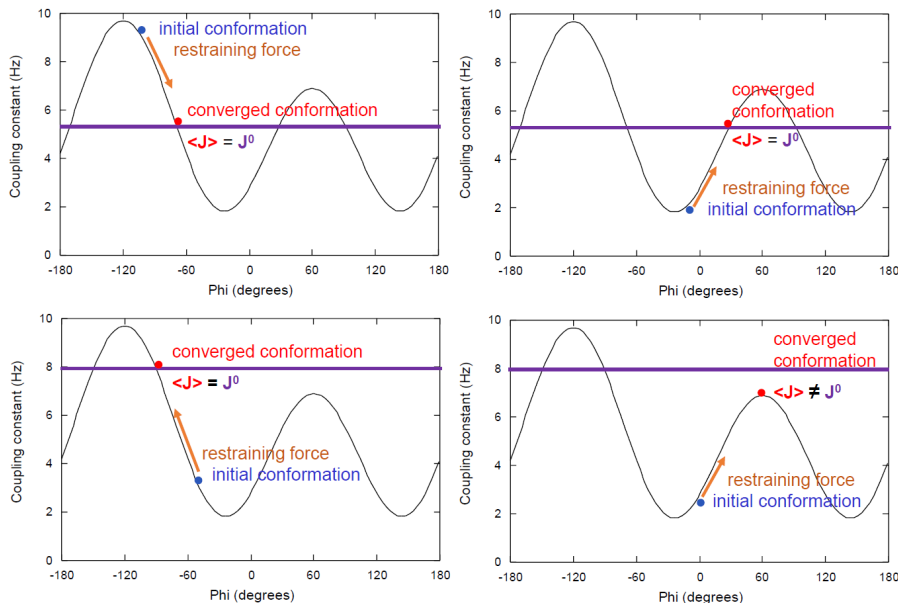


Figure 6: Consequences of a multiple-valued function $\phi(3J)$ for restraining an angle ϕ to a particular $3J$ -coupling value J^0 when starting from different initial values of ϕ in structure refinement.

638 A technique that achieves this, local-elevation sampling, has been proposed
639 more than two decades ago.^[117] It adds a memory function $V^{le}(\phi^M)$ depending
640 on a set $\phi^M \equiv (\phi_1, \phi_2, \dots, \phi_M)$ of torsional or dihedral angles ϕ_k defined by
641 four atoms or particles ($k_1 - k_2 - k_3 - k_4$) to the physical interaction function
642 V^{phys} . This memory function can have the form of a Gaussian of width $\Delta\phi^o$
643 centered at the value ϕ_i^o . For a given torsional angle ϕ_k , $V_k^{le}(\phi_k)$ is a sum of
644 N_{le} Gaussians centered at generally equidistant values ϕ_i^o , which are the same
645 for each k , with $i = 1, 2, \dots, N_{le}$ along the 360° range of ϕ_k values,

$$V_k^{le}(\phi_k(t)) = \sum_{i=1}^{N_{le}} V_k^{le}(\phi_k(t); \phi_i^o), \quad (24)$$

646 with

$$V_k^{\text{le}}(\phi_k(t); \phi_i^\circ) = K^{\text{le}} \omega(\phi_k(t); \phi_i^\circ) \exp\left(-\frac{(-\phi_k(t) - \phi_i^\circ)^2}{2(\Delta\phi^\circ)^2}\right). \quad (25)$$

647 K^{le} is the overall weight of the local-elevation contribution to the potential
 648 energy, while the weight factor $\omega(\phi_k(t); \phi_i^\circ)$ is enlarged by a given amount $\Delta\omega$
 649 at every time point t for which

$$\phi_i^\circ - \frac{1}{2}\Delta\phi^\circ \leq \phi_k(t) \leq \phi_i^\circ + \frac{1}{2}\Delta\phi^\circ. \quad (26)$$

650 In this way the total potential energy surface for the molecular configuration
 651 at the point defined by the local-elevation variables, i.e. torsional angles ϕ^M ,
 652 i.c. at position ϕ_i° , is lifted or locally elevated every time the molecule visits the
 653 region around ϕ_i° during the MD simulation. In this way the system is driven
 654 away from regions of configurational space that were visited previously. Multiple
 655 equivalent schemes have been proposed. For example, in 2002 the local-elevation
 656 sampling technique was republished under the name ‘‘metadynamics’’.^[118]

657 The local-elevation sampling technique can also be used when restraining
 658 structure \vec{r}^N or θ towards Q° or ${}^3J^\circ$ values when the function $\vec{r}^N(Q)$ or $\theta({}^3J)$
 659 is multiple-valued.^[119] In this case the local-elevation function V_k^{le} of Eqs. 24
 660 and 25 is used as biasing or restraining function $V^{Q,\text{restr}}$, but the way the weight
 661 $\omega(\phi_k(t); \phi_i^\circ)$ in Eq. 25 is assigned is varied. The weight is only increased when
 662 both $Q(t)$ and $\langle Q \rangle_t$ lie outside the flat bottom of $V^{Q,\text{restr}}$. With $Q = {}^3J_k$ we
 663 have, using a biquadratic attractive ($Q(t) > Q^\circ$ and $\langle Q \rangle_t > Q^\circ$) and repulsive
 664 ($Q(t) < Q^\circ$ and $\langle Q \rangle_t < Q^\circ$) flat-bottom function of $Q(t)$ and $\langle Q \rangle_t$,

$$\begin{aligned}
\omega(\phi_k(t); \phi_i^\circ) &= t^{-1} \int_0^t \delta_{\phi_k(t')\phi_i^\circ} \\
&\cdot \left\{ V_{\text{att}}^{Q, \text{restr}}(Q(t'), \langle Q \rangle_{t'}; K^{Qr}, Q^\circ, \Delta Q^{\text{fb}}) \right. \\
&\left. + V_{\text{rep}}^{Q, \text{restr}}(Q(t'), \langle Q \rangle_{t'}; K^{Qr}, Q^\circ, \Delta Q^{\text{fb}}) \right\} dt', \quad (27)
\end{aligned}$$

665 with

$$\delta_{\phi_k(t)\phi_i^\circ} = \begin{cases} 1 & \text{if } \phi_i^\circ - \frac{\Delta\phi^\circ}{2} \leq \phi_k(t) \leq \phi_i^\circ + \frac{\Delta\phi^\circ}{2} \\ 0 & \text{otherwise} \end{cases}. \quad (28)$$

666 This local-elevation biasing technique has been shown to generate structural
667 ensembles that match the Q^{exp} data used in the biasing function.^[44, 120, 121]

668 Local-elevation biasing of a MD simulation based on target Q° values of a
669 quantity Q can also be used to detect inconsistencies between different exper-
670 imental Q^{exp} values.^[43] This is illustrated in Figure 7 for the case of four 3J -
671 coupling values, ${}^3J_{H_\alpha H_{\beta 2}}$ (4.0 Hz), ${}^3J_{H_\alpha H_{\beta 3}}$ (2.0 Hz), ${}^3J_{NH_{\beta 2}}$ (1.8 Hz), ${}^3J_{NH_{\beta 3}}$
672 (1.8 Hz) belonging to one torsional angle $\phi = \chi_1$ of the side chain of residue 8
673 of the protein FKBP.^[43] Using local-elevation biasing for the four 3J -couplings
674 of the χ_1 -angle, this torsional angle sampled the complete range of $(0^\circ, 360^\circ)$,
675 upper left panel in Figure 7. For no combination of values of χ_1 could all four
676 3J -couplings be reproduced to within $\Delta Q^{\text{fb}} = 1$ Hz of the above mentioned
677 target values Q° . Instead there was a continual increase of the local-elevation
678 biasing or restraining energy V_k^{le} for the χ_1 -angle (upper middle panel in Fig-
679 ure 7). Thus one or more of the reported 3J -coupling values must be erroneous.
680 Analysis of the available data suggested that the value ${}^3J_{NH_{\beta 3}} = 1.8$ Hz was
681 probably incorrect. After removing this restraint from the restraining function,

682 the local-elevation potential energy did not increase with time and a range of
 683 χ_1 -angle values compatible with the three target 3J -coupling values could be
 684 identified (lower panels in Figure 7).

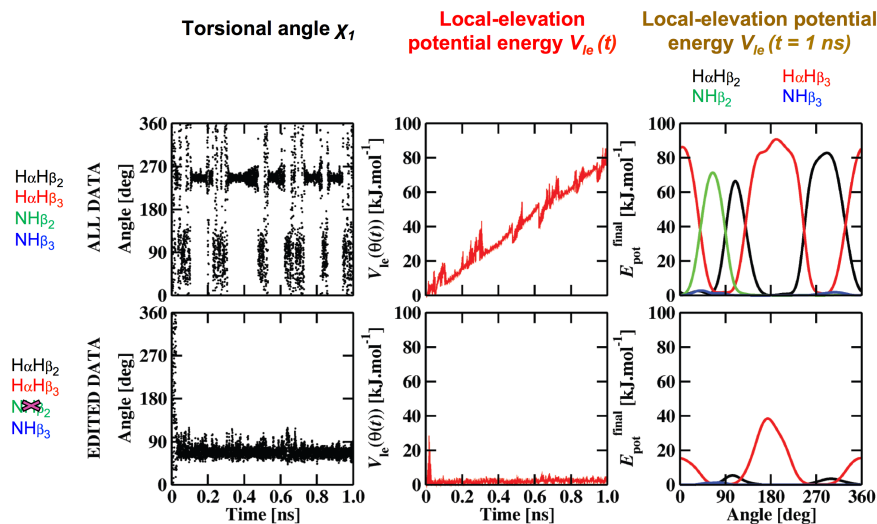


Figure 7: Detection of inconsistency of the (redundant) experimental data.^[43] The torsional angle χ_1 of residue 8 of FK506 binding protein (FKBP) determines four 3J -coupling constants. Through local-elevation search inconsistency between the four 3J -coupling values is detected. Upper panels: four 3J -coupling restraints. Lower panels: three 3J -coupling restraints. Left to right: torsional angle value $\chi_1(t)$, local-elevation potential energy $V_e(t)$, local-elevation potential energy $V_e(\chi_1)$ at $t = 1$ ns.

685 Multiple-valuedness of the inverse of the function $Q(\vec{r}^{2N})$ may also occur for
 686 other quantities Q than 3J -couplings or residual dipolar couplings (Table 1).
 687 The inverse of the crystallographic structure factors $F_{hkl}(\vec{r}^{2N})$ depends not only
 688 on the structure factor amplitudes, which are proportional to the square root of
 689 the scattering intensities I_{hkl} , but also on the phases ϕ_{hkl} . The phases cannot be
 690 experimentally measured and are instead inferred from a proposed model. The
 691 uncertainty with respect to the phases ϕ_{hkl} implies that different structures, i.e.
 692 electron densities, may map onto the same structure factor amplitudes, leading

693 to multiple-valuedness of the relation $\bar{r}^{2N}(I_{hkl})$.

694 The question why and in which cases local-elevation biased sampling is bet-
695 ter than instantaneous restraining as applied in single-structure refinement has
696 been discussed in Ref. 44. It is summarised in Figure 8 for a one-dimensional
697 system with $Q(x)$ and $V^{\text{phys}}(x)$. $V^{\text{phys}}(x)$, the molecular model or force field
698 used in the MD simulation, is represented by the dashed line, and the real,
699 correct potential energy surface of the system $V^{\text{real}}(x)$ by the solid line. Three
700 MD sampling techniques are compared, (i) free MD, that is using V^{phys} without
701 any restraining function $V^{Q,\text{restr}}$, (ii) instantaneous restraining (IR) MD using
702 $V^{Q,\text{restr}}$ but without accounting for averaging, so equivalent to single-structure
703 refinement, and (iii) local-elevation (LE) sampling MD using $V^{Q,\text{restr}}$ inclusive
704 averaging as described above. Depending on the levels of the different energy
705 wells of $V^{\text{phys}}(x)$ and $V^{\text{real}}(x)$ and the height of the barriers between them
706 in relation to the kinetic energy $k_{\text{B}}T$ available to the coordinate x in an MD
707 simulation, the different sampling techniques succeed or fail sampling the config-
708 urations x matching $\langle Q \rangle$. In free MD without restraints, the simulation may fail
709 to sample configurations with appropriate weight due to limitations in the force
710 field or molecular model (upper middle and lower left panels). Instantaneous
711 restraining or single-structure refinement may fail when different configurations
712 x contribute to the average $\langle Q \rangle$ such that $\langle Q \rangle = Q(x)$ corresponds to a config-
713 uration x with a high energy $V^{\text{phys}}(x)$ or $V^{\text{real}}(x)$ (right-hand panels in Figure 8).
714 Only local-elevation MD sampling yields ranges of configurations that reproduce
715 $\langle Q \rangle$ and are of low energy V^{phys} .

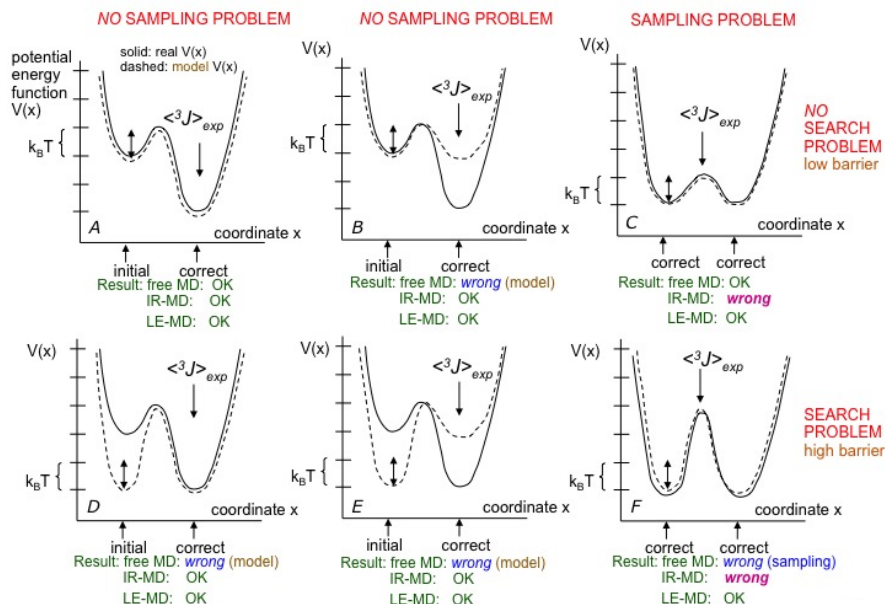


Figure 8: Schematic representation of six different real (solid line) and model (dashed line) potential energy functions illustrating the force-field problem of unrestrained MD simulations (examples *B*, *D*, *E*), the sampling problem of MD simulations which occurs when instantaneous restraints are applied (examples *C* and *F*), and the search problem of MD simulations due to high-energy barriers between different conformations (examples *D*, *E*, *F*). The double arrow indicates the thermal energy ($\frac{1}{2}k_B T$) associated with the degree of freedom x . If the thermal energy is comparable to the barrier height, transitions are easy, whereas a higher barrier leads to rare transitions. If the measured 3J -value, $\langle ^3J \rangle_{\text{exp}}$, corresponds according to the non-linear Karplus relation to a torsional-angle coordinate x for which the potential energy is larger (examples *C* and *F*), instantaneous restraining will lead to an unrealistic configuration x .

716 2.10 Boltzmann sampling or weighting of configurations

717 It follows from statistical mechanics that a set of molecular configurations should
718 approximate a statistical-mechanical ensemble, e.g. a Boltzmann-weighted set
719 of structures approximating at constant temperature a canonical ensemble. The
720 weights of the individual structures or configurations \vec{r}^N should be proportional
721 to $\exp(-V^{\text{real}}(\vec{r}^N)/(k_B T))$. Unfortunately, the real potential energy function

722 of a biomolecular system is not known. At the atomic level of resolution it
723 can be approximated by a model function or force field $V^{\text{phys}}(\vec{r}^N)$, and the
724 weights of a set of structures or conformations should then be proportional
725 to $\exp(-V^{\text{phys}}(\vec{r}^N)/(k_{\text{B}}T))$. Using Monte Carlo simulation or MD simulation
726 at constant temperature this relative weighting of trajectory structures is well
727 approximated. Yet, for most biomolecular structures reported in the litera-
728 ture and deposited in structural data banks, this Boltzmann weighting is not
729 achieved. Normally such structures result from energy minimisation or temper-
730 ature annealing single-structure refinement procedure, which by definition are
731 not Boltzmann weighted. Such sets of structures do not constitute an experi-
732 mental ensemble. In fact, structures based directly on experimental data should
733 not be called an ensemble, because their relative weights do not correspond to
734 or even approximate a well defined statistical-mechanical ensemble. For this
735 reason it is also inappropriate to use protein structures deposited in the Pro-
736 tein Data Bank (PDB) as a means to directly calibrate parameters for use in a
737 biomolecular force field. The variation in bond angles for example, represents
738 the variation in the minimum energy position between models, not the variation
739 which would be expected due to thermal motion.

740 **3 What should be avoided when deriving structural infor-** 741 **mation from measured data ?**

742 The considerations in the previous sections lead to a number of choices and
743 procedures that are to be avoided when deriving structural information from
744 measured data for biomolecular systems.

- 745 1. Use of derived “experimental” data in the biasing or restraining function

746 $V^{Q,\text{restr}}$. Unless the procedure f (Eq. 5) that converts measured values
747 $\langle Q^{\text{obs}} \rangle_{\text{exp}}$ of an observable quantity Q^{obs} into derived values of a non-
748 observable quantity Q^{der} involves well founded assumptions and does not
749 increase the uncertainty of the values Q^{der} , derived data should be avoided.

750 2. Use of a poor molecular model or force field, e.g. without electrostatic in-
751 teraction terms or without explicit solute-solvent interactions. The quality
752 of the force field is particularly important in situations where the number
753 of independent experimental values Q^{exp} is small compared to the number
754 of degrees of freedom of the system.

755 3. Use of poorly defined experimental data or Q^{exp} values beset with a large
756 uncertainty in a restraining function $V^{Q,\text{restr}}$. In some cases the effects
757 of uncertainty in the experimental data can be accounted for by using
758 a flat-bottom potential energy term for restraining the range of which
759 corresponds to the uncertainty.

760 4. The use of any relation $Q(\vec{r}^N)$ between an observable Q and structure
761 \vec{r}^N for which the inverse function $\vec{r}^N(Q)$ is multiple-valued without an
762 appropriate sampling technique such as local-elevation. This applies not
763 only to 3J -coupling restraints, but also to restraints based on e.g. Residual
764 Dipolar Couplings or crystallographic structure factor amplitudes.

765 5. Use of instantaneous restraining or single-structure refinement, unless it
766 is *a priori* clear that the system of interest adopts a single conformation.

767 6. Use of a too large force constant K^{Qr} or weight of the restraining potential
768 energy term $V^{Q,\text{restr}}$ compared to the physical potential energy V^{phys} .

- 769 7. Use of molecule averaging in restraining without using a statistical-mechanically
770 based weighting criterion, e.g. canonical or microcanonical.^[122] For exam-
771 ple, the use of replica-exchange protocols between concurrently simulated
772 copies of the molecular system may ensure correct weighting.
- 773 8. Use of a poor or a non-Boltzmann sampling algorithm to generate a set
774 of structures.

775 We note that these issues play a role when deriving structure from measured
776 values of all observables mentioned in the Introduction and that these issues still
777 play a role in contemporary structure determination.^[123–125] It depends on the
778 information density of the particular measurement technique, that is the ratio
779 of the number of independent measured values Q^{exp} and the number of degrees
780 of freedom of the molecule whose structure is to be determined, how much the
781 choices made in the structure refinement procedure will influence the resulting
782 molecular structure. The majority of the measurement techniques mentioned in
783 the Introduction yield a rather low information density, which means that the
784 structures modelled based on the experimental data will be dominated by the
785 choice of molecular model and refinement techniques and parameters.

786 **4 Flawed structural interpretation of measured data due** 787 **to inappropriate choices made when deriving structure** 788 **from experimental data**

789 Since many structures reported in the literature have been derived based on
790 procedures and assumptions that ideally should be avoided, one might wonder to
791 what extent biomolecular structures that result from adhering to the standards

792 listed in the previous section lead to a different structural interpretation of the
793 experimental data. This, of course, will depend on a variety of factors involved
794 in the structure derivation, the type of molecules considered and the type of
795 experimental data available. Below, a few examples from the literature are
796 presented to illustrate potential consequences.

797 Structures of the 16-residue C-terminal half GCN4p16-31 of the leucine zip-
798 per GCN4 were derived on the basis of 172 NOE atom-atom distance bounds
799 and 15 $^3J_{H_N H_\alpha}$ -couplings measured by NMR.^[44,121] Single-structure refinement
800 using the program X-PLOR based on (i) 172 NOE distance restraints, (ii) 14
801 α -helical hydrogen bond restraints (2 per hydrogen bond) between the peptide
802 H and N atoms of residue $i + 4$ and the carbonyl oxygen of residue i with
803 $i = 18 - 24$, and (iii) 8 backbone ϕ -angle restraints for residues 17-24, led to a
804 set of 20 NMR model structures through simulated annealing.^[126] However, this
805 set of NMR model structures showed five NOE bound violations larger than 0.05
806 nm and two 3J -coupling deviations larger than 3 Hz. It was decided, therefore to
807 perform local-elevation MD time-averaging structure refinement using the GRO-
808 MOS program^[109] and the GROMOS force field 53A6^[127] based on (i) the 172
809 NOE distance restraints using time-averaging and (ii) the 15 $^3J_{H_N H_\alpha}$ -couplings
810 using local-elevation sampling.^[44] This led to a greater variability between the
811 structures but overall the deviations from the experimental data were smaller:
812 no NOE bound violations beyond 0.03 nm and no 3J -coupling deviations larger
813 than 1 Hz. The network of polar side chain contacts in the LE-MD structural
814 ensemble was more extensive than in the set of NMR model structures obtained
815 through single-structure refinement using X-PLOR. This is likely to be due to
816 the force field used lacking electrostatic and hydrogen bonding terms.^[121] In

817 other words, this set of 20 NMR model structures, deposited in the Protein Data
818 Bank, does not capture the expected contacts between the polar side chains of
819 the helix. In this case an inappropriate choice of hydrogen-bond and torsional-
820 angle restraints based on values of derived, non-observable quantities Q^{der} , the
821 neglect of averaging by applying instantaneous restraining or single-structure
822 refinement, and the use of a force field without electrostatic interaction terms
823 led to a set of structures with a lower structural variability than indicated by
824 the measured data, and thus not only to a poor set of structural models but to
825 the misinterpretation of the data on which these models were based.

826 An example of the combined effect of inaccuracies of the force field used in
827 single-structure refinement and a limited set of atom-atom distance bounds de-
828 rived from NMR data can be found in Ref. 128. Although the proposed structure
829 of the complex of the barley lipid transfer protein 1 (LTP) in complex with the
830 fatty acid ligand palmitate was based on the use of 25 protein-ligand NOE re-
831 straints,^[128] these were not sufficient to identify the details of the ligand binding
832 site using the calculation protocol applied. Indeed, MD simulation starting from
833 the NMR model structure indicated this structure was not stable. The palmi-
834 tate ligand moved out of the internal cavity and became exposed to solvent.^[129]
835 Using different sets of restraints (omitting hydrogen-bond restraints), different
836 protonation states and force-field parameters two alternative palmitate binding
837 modes characterised by different hydrogen-bonding patterns, were identified for
838 palmitate.^[129] The two modes have similar protein-ligand interaction energies
839 suggesting that both are significantly populated. In these binding modes the
840 palmitate head group is not involved in a salt bridge or hydrogen bond to the
841 side chain of Cys 9 as had been previously suggested.^[128] Instead the ligand

842 formed hydrogen bonds with main-chain amide groups of residues in helix A.
843 The ligand binding exploits irregularities in the helical hydrogen-bonding pat-
844 tern in this helix due to the presence of Pro 12. Simulations of single mutant
845 P12V and double mutant P12V,P20V variants of barley LTP also suggested
846 that the presence of Proline at position 12 in helix A is needed to stabilize the
847 binding of palmitate to barley LTP.

848 Another example of the incomplete picture of ligand binding based on X-ray
849 and NMR data as obtained by single-structure refinement based on a limited
850 set of experimental data can be found in Ref. 130. The data suggested the oc-
851 currence of two ligand binding modes of the ligand caprate to maize LTP. Using
852 different protonation states and configurational sampling strategies MD simu-
853 lation^[131] indicated that one particular binding mode (M) is preferred in maize
854 LTP rather than a mixture of two different binding modes (M and B) suggested
855 previously. The mobility of the caprate ligand explained the absence of electron
856 density for the ligand head group in the X-ray structure determination.^[132]

857 Currently the study of so-called intrinsically disordered proteins is rather
858 popular. For such proteins the ratio of the amount of independent data that
859 can be accurately measured compared to the number of protein configurational
860 degrees of freedom that allow for disorder is very, very low. This implies that
861 the reliability of the generated structures is almost entirely determined by the
862 combination of molecular model and force field used in the structure generation
863 rather than by the measured data. This applies to most of the observable
864 quantities and measurement techniques mentioned in the Introduction.

865 Even when using a high information density technique such as X-ray diffrac-
866 tion crystallography the interpretation of the electron density of a molecule may

867 be flawed by ignoring the motion of atoms in the structure refinement process.
868 This is illustrated by a time-averaging structure factor amplitude restraining
869 MD simulation of the protein basic pancreatic trypsin inhibitor (BPTI).^[87] The
870 X-ray model structure derived using single-structure refinement based on X-
871 ray diffraction data contained four internal water molecules without contact to
872 bulk solvent due to the presence of the side chain of Glu 7 blocking their ac-
873 cess to the protein surface. Allowing for averaging through the time-averaging
874 structure refinement procedure it was found that the electron density at the pro-
875 posed position of the Glu 7 side chain could be due to a water molecule. While
876 the side chain actually does not have its own well-defined electron density, it
877 exists in many different conformations during the time-averaging refinement-
878 simulation.^[87] These conformations are consistent with the crystallographic
879 data, because no significant density difference exists. In contrast to the single-
880 structure refinement result, the new water site bridges the four internal water
881 sites to the bulk water outside. Other examples of misinterpretation of experi-
882 mental data can be found in Ref. 79–81, 133, 134. The risk of misinterpretation
883 of experimental data is even higher for the measurement techniques of low in-
884 formation density mentioned in the Introduction.

885 **5 Pollution of structural databases by flawed structures**

886 During the past years, a number of protein structures deposited in the Protein
887 Data Bank (PDB) were corrected or retracted, because the interpretation of the
888 corresponding X-ray diffraction data was flawed.^[79, 80] Considering the number
889 of things that can go wrong in the process of structure refinement,^[81, 135, 136]
890 the number of partially flawed or non-optimal structures is likely to be large.^[81]

891 This highlights the importance of being able to correct structures deposited in
892 the PDB as more robust procedures to derive structure from measured values
893 of observable quantities^[134] or to check the derived structures with respect to
894 geometric and energetic criteria^[81] become available and especially when these
895 result in a better match to the measured values of observables^[137] than existing
896 structures. It also highlights the need for all structures to be deposited together
897 with values of the observables such as NOE bounds, ³*J*-couplings and struc-
898 ture factor amplitudes, force fields and computational protocols used for their
899 generation. The primary criteria for whether a structure is included in these
900 critical public databases must be the extent to which the models can account
901 for the available experimental data. Limiting the ability of researchers other
902 than the original depositors to correct structures will lead to the progressive
903 pollution of the PDB with partially flawed structures. Structures should be
904 stored independently of the observed data from which they were derived, and
905 the key choices and procedures of the structural derivation should be deposited
906 with the structures. Recent moves by the PDB and other repositories to enforce
907 this should be encouraged.^[138]

908 In fact, the quality of structures is an on-going concern to the PDB, and the
909 PDB has held a series of meetings to address concerns regarding the deposition
910 of flawed structures.^[139] While the ease with which data can be acquired has
911 increased and data processing has improved, the extent to which partially auto-
912 mated tools are used by less experienced persons to solve larger structures has
913 led to flawed structures.^[140] Solving structures for more complex systems even
914 with modern tools can be just as challenging (i.e. error prone) as for simpler
915 systems with older tools. For MD simulation, the quality of force fields and

916 sampling has improved, but so has the complexity of the systems studied. It
917 would be very wrong to assume all recent simulation studies are of high quality.

918 **6 Conclusion**

919 The process of deriving structural information from measured data for biomolec-
920 ular systems is a complex one involving many assumptions and approximations.
921 These may introduce uncertainty or error in the derived structures, which in turn
922 may then misrepresent reality. This may be avoided by carefully considering the
923 various choices to be made in the refinement process. These considerations re-
924 gard many types of observables and measurement techniques, not only NMR
925 spectroscopy, but also X-ray and neutron diffraction, small-angle and wide-
926 angle X-ray scattering, free-electron laser imaging, cryo-electron microscopy,
927 EPR, Infrared and Raman spectroscopy, CD, FRET, ATM and ion-mobility
928 mass spectrometry. The lower the information density of a measurement tech-
929 nique, the larger will be the chance that problematic choices of approximations
930 and parameters of the structure refinement process will go undetected.

931 Due to the complexity of the process of deriving structures from experimen-
932 tally measured data and the different characteristics of the variety of experimen-
933 tal data that can be used, it is not possible to assert what is “best practice” and
934 expect this to be applicable to all cases. What procedure to choose will depend
935 on *(i)* the type, quality and number of experimental data, *(ii)* the relationship
936 between these data and structure, *(iii)* the choice of the degrees of freedom
937 of the molecular model used, *(iv)* the quality of the molecular model and force
938 field used, *(v)* the range of the model parameter values that is sampled, and *(vi)*
939 the extent of configurational sampling and biasing executed. While it would of

940 course be nice to have a straightforward checklist for determining structure qual-
941 ity, the factors at play and their relative importance are likely to vary on a
942 case-by-case basis, thus such a list will always be a simplification and likely to
943 lead to even further problems. Unfortunately, a high-dimensional problem often
944 only has a high-dimensional answer.

945 **7 Acknowledgements**

946 This work was financially supported by grant No 20020-137827 of the Swiss
947 National Science Foundation and by grant No 228076 of the European Research
948 Council (ERC), which is gratefully acknowledged.

949 **References**

- 950 [1] J. C. Kendrew, G. Bodo, H. M. Dintzis, R. G. Parrish, H. Wyckoff, D. C.
951 Phillips, *Nature* **1958**, 181, 662-666
- 952 [2] C. C. F. Blake, D. F. Koenig, G. A. Mair, A. C. T. North, D. C. Phillips,
953 V. R. Sarma, *Nature* **1965**, 206, 757-761
- 954 [3] J. Drenth, C. M. Enzing, K. H. Kalk, J. C. A. Vessies, *Nature* **1976**, 264,
955 373-377
- 956 [4] A. Wlodawer, J. Walter, R. Huber, L. Sjölin, *J. Mol. Biol.* **1984**, 180,
957 301-329
- 958 [5] D. I. Svergun, M. H. J. Koch, *Rep. Prog. Phys.* **2003**, 66, 1735-1782
- 959 [6] K. J. Gaffney, H. N. Chapman, *Science* **2007**, 316, 1444-1448
- 960 [7] M. Walczak, H. Grubmüller, *Phys. Rev. E* **2014**, 90, 022714
- 961 [8] X.-C. Bai, G. McMullan, S. H. W. Scheres, *Trends Biochem. Sci.* **2015**,
962 40, 49-57
- 963 [9] R. R. Ernst, G. Bodenhausen, A. Wokaun, *Principles of Nuclear Mag-*
964 *netic Resonance in One and Two Dimensions*, Clarendon Press, Oxford,
965 England, **1990**
- 966 [10] K. D. Berndt, P. Güntert, L. P. M. Orbons, K. Wüthrich, *J. Mol. Biol.*
967 **1992**, 227, 757-775
- 968 [11] C. C. Jao, B. G. Hegde, J. Chen, I. S. Haworth, R. Langen, *Proc. Natl.*
969 *Acad. Sci. USA* **2008**, 105, 19666-19671
- 970 [12] S. Bagchi, S. G. Boxer, M. D. Fayer, *Phys. Chem. B* **2012**, 116, 4034-4042
- 971 [13] Y. Hong, L. Meng, S. Chen, C. W. T. Leung, L.-T. Da, M. Faisal, D.-A.
972 Silva, J. Liu, J. W. Y. Lam, X. Huang, B. Z. Tang, *J. Am. Chem. Soc.*
973 **2012**, 134, 1680-1689
- 974 [14] C. A. E. Hauser, R. Deng, A. Mishra, Y. Loo, U. Khoe, F. Zhuang, D. W.
975 Cheong, A. Accardo, M. B. Sullivan, C. Riek, J. Y. Ying, U. A. Hauser,
976 *Proc. Natl. Acad. Sci. USA* **2011**, 108, 1361-1366
- 977 [15] M. Hoefling, N. Lima, D. Haenni, C. A. M. Seidel, B. Schuler, H. Grub-
978 müller, *PLoS ONE* **2011**, 6, e19791
- 979 [16] M. M. Reif, C. Oostenbrink, *J. Comput. Chem.* **2014**, 35, 2319-2332
- 980 [17] A. T. Brunger, P. Strop, M. Vrljic, S. Chu, K. R. Weninger, *J. Struct.*
981 *Biol.* **2011**, 173, 497-505
- 982 [18] T. Ando, *Nanotechnology* **23**, 062001
- 983 [19] Z. Hall, A. Politis, M. F. Bush, L. J. Smith, C. V. Robinson, *J. Am.*
984 *Chem. Soc.* **2012**, 134, 3429-3438

- 985 [20] J. Drenth, *Principles of Protein X-ray Crystallography*, Springer, New
986 York, USA, **1994**
- 987 [21] R. Boelens, T. M. G. Koning, R. Kaptein, *J. Mol. Struct.* **1988**, 173,
988 299-311
- 989 [22] M. Karplus, *J. Chem. Phys.* **1959**, 30, 11-15
- 990 [23] D. Steiner, J. R. Allison, W. F. van Gunsteren, *J. Biomol. NMR* **2012**,
991 53, 223-246
- 992 [24] B. Han, Y. Liu, S. W. Ginzinger, D. S. Wishart, *J. Biomol. NMR*, **2011**,
993 50, 43-57
- 994 [25] W. F. van Gunsteren, A. M. J. J. Bonvin, X. Daura, L. J. Smith, in:
995 *Structure Computation and Dynamics in Protein NMR, Biol. Magnetic*
996 *Resonance Vol. 17* (Eds.: R. N. Krishna and L. J. Berliner), Plenum
997 Publishers, New York, USA, **1999**, pp. 3-35.
- 998 [26] J. H. Lee, F. Li, A. Grishaev, A. Bax *J. Am. Chem. Soc.*, **2015**, 137,
999 1432-1435
- 1000 [27] U. Stocker, D. Juchli, W. F. van Gunsteren, *Mol. Sim.*, **2003**, 29, 123-138
- 1001 [28] A. T. Brünger, *X-PLOR. A System for X-ray Crystallography and NMR*,
1002 Yale University Press, New Haven, CT, USA, **1992**.
- 1003 [29] W. A. Hendrickson, J. H. Konner, in: *Biomolecular Structure, Conforma-*
1004 *tion, Function and Evolution, Volume 1: Diffraction and Related Studies*
1005 (Ed.: R. Srinivasan), Pergamon, Oxford, U.K., **1981**, pp. 43-57.
- 1006 [30] W. F. van Gunsteren, R. Kaptein, E. R. P. Zuiderweg, in: *Proceedings*
1007 *NATO/CECAM workshop on nucleic acid conformation and dynamics*
1008 (Ed.: W.K. Olson), CECAM, Orsay, France, **1984**, pp. 79-92.
- 1009 [31] R. Kaptein, E. R. P. Zuiderweg, R. M. Scheek, R. Boelens, W. F. van
1010 Gunsteren, *J. Mol. Biol.* **1985**, 182, 179-182
- 1011 [32] A. T. Brünger, J. Kuriyan, M. Karplus, *Science* **1987**, 235, 458-460
- 1012 [33] M. Fujinaga, P. Gros, W. F. van Gunsteren, *J. Appl. Cryst.* **1989**, 22, 1-8
- 1013 [34] G. F. Schröder *Curr. Opin. Struct. Biol.* **2015**, 31, 20-27
- 1014 [35] H. M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T. N. Bhat, H. Weiss-
1015 sig, I. N. Shindyalov, P. E. Bourne, *Nucleic Acids Res* **2000**, 28, 235-242
- 1016 [36] W. F. van Gunsteren, J. Dolenc, A. E. Mark, *Curr. Opin. Struct. Biology*
1017 **2008**, 18, 149-153
- 1018 [37] A. Glättli, W. F. van Gunsteren, *Angew. Chem.* **2004**, 116, 6472-6476;
1019 *Angew. Chem. Int. Ed.* **2004**, 43, 6312-6316
- 1020 [38] K. Gademann, A. Häne, M. Rueping, B. Jaun, D. Seebach, *Angew. Chem.*
1021 **2003**, 115, 1573-1575; *Angew. Chem. Int. Ed.* **2003**, 42, 1534-1537

- 1022 [39] B. Zagrovic, G. Jayachandran, I. S. Millett, S. Doniach, V. S. Pande, *J.*
1023 *Mol. Biol.* **2005**, 353, 232-241
- 1024 [40] B. Zagrovic, W. F. van Gunsteren, *Proteins Struct. Funct. Bioinf.* **2006**,
1025 63, 210-218
- 1026 [41] W. F. van Gunsteren, R. Bürgi, C. Peter, X. Daura, *Angew. Chem.* **2001**,
1027 113, 363-367; *Angew. Chem. Int. Ed.* **2001**, 40, 351-355
- 1028 [42] W. F. van Gunsteren, D. Bakowies, R. Baron, I. Chandrasekhar, M. Christen,
1029 X. Daura, P. Gee, D. P. Geerke, A. Glättli, P. H. Hünenberger, M. A.
1030 Kastenholz, C. Oostenbrink, M. Schenk, D. Trzesniak, N. F. A. van der
1031 Vegt, H.B. Yu, *Angew. Chem.* **2006**, 118, 4168-4198; *Angew. Chem. Int.*
1032 *Ed.* **2006**, 45, 4064-4092
- 1033 [43] J. R. Allison, W. F. van Gunsteren, *ChemPhysChem* **2009**, 10, 3213-3228
- 1034 [44] J. Dolenc, J. H. Missimer, M. O. Steinmetz, W. F. van Gunsteren, *J.*
1035 *Biomol. NMR* **2010**, 47, 221-235
- 1036 [45] B. Vögeli, J. Ying, A. Grishaev, A. Bax, *J. Am. Chem. Soc.* **2007**, 129,
1037 9377-9385
- 1038 [46] A. Pardi, M. Billeter, K. Wüthrich, *J. Mol. Biol.* **1984**, 180, 741-751
- 1039 [47] R. Brüschweiler, D. A. Case, *J. Am. Chem. Soc.* **1994**, 116, 11199-11200
- 1040 [48] A. C. Wang, A. Bax, *J. Am. Chem. Soc.* **1996**, 118, 2483-2494
- 1041 [49] J. M. Schmidt, M. Blümel, F. Löhr, H. Rüterjans, *J. Biomol. NMR* **1999**,
1042 14, 1-12
- 1043 [50] A. DeMarco, M. Linàs, K. Wüthrich, *Biopolymers* **1978**, 17, 617-636
- 1044 [51] A. DeMarco, M. Linàs, K. Wüthrich, *Biopolymers* **1978**, 17, 2727-2742
- 1045 [52] R. J. Abraham, K. A. McLauchlan, *Mol. Phys.* **1962**, 5, 513-523
- 1046 [53] C. M. Deber, E. R. Blout, D. A. Torchia, *J. Am. Chem. Soc.* **1971**, 93,
1047 4893-4897
- 1048 [54] M. T. Cung, M. Marraud, J. Nèel, in: *10th Prague IUPAC Micro-*
1049 *Symposium on Macromolecules*, Prague, **1972**, p. C-3.
- 1050 [55] A. J. Fischman, D. H. Live, H. R. Wyssbrod, W. C. Agosta, D. Cowburn,
1051 *J. Am. Chem. Soc.* **1980**, 102, 2533-2539
- 1052 [56] C. Pérez, F. Löhr, H. Rüterjans, J. M. Schmidt, *J. Am. Chem. Soc.* **2001**,
1053 123, 7081-7093
- 1054 [57] X. Daura, I. Antes, W. F. van Gunsteren, W. Thiel, A. E. Mark, *Proteins*
1055 *Struct. Funct. Genet.* **1999**, 36, 542-555
- 1056 [58] T. S. Harvey, W. F. van Gunsteren, in: *Techniques in Protein Chemistry*
1057 *IV*, Academic Press, **1993**, pp. 615-622.
- 1058 [59] D. A. Case, *Curr. Opin. Struct. Biol.* **2013**, 23, 172-176

- 1059 [60] W. F. van Gunsteren, R. Boelens, R. Kaptein, R. M. Scheek, E. R. P.
1060 Zuiderweg, in: *Molecular Dynamics and Protein Structure* (Ed.: J. Her-
1061 mans), Polycrystal Book Service, P.O. Box 27, Western Springs, Ill. 60558,
1062 USA, **1985**, pp. 92-99.
- 1063 [61] J. de Vlieg, R. Boelens, R. M. Scheek, R. Kaptein, W. F. van Gunsteren,
1064 *Isr. J. Chem.* **1986**, 27, 181-188
- 1065 [62] A. E. Torda, R. M. Brunne, T. Huber, H. Kessler, W. F. van Gunsteren,
1066 *J. Biomol. NMR* **1993**, 3, 55-66
- 1067 [63] K. Wüthrich, M. Billeter, W. Braun, *J. Mol. Biol.* **1983**, 169, 949-961
- 1068 [64] W. F. van Gunsteren, S. R. Billeter, A. A. Eising, P. H. Hünenberger, P.
1069 Krüger, A. E. Mark, W. R. P. Scott, I. G. Tironi, Biomolecular Simulation:
1070 The GROMOS96 Manual and User Guide, Vdf Hochschulverlag AG an der
1071 ETH Zürich, Zürich, Switzerland, **1996**.
- 1072 [65] C. M. Fletcher, D. N. M. Jones, R. Diamond, *J. Biomol. NMR* **1996**, 8,
1073 292-310
- 1074 [66] M. Nilges, *Proteins* **1993**, 17, 297-309
- 1075 [67] M. Nilges, M. J. Macias, S. I. O'Donoghue, H. Oschkinat, *J. Mol. Biol.*
1076 **1997**, 269, 408-422
- 1077 [68] W. F. van Gunsteren, in: *Studies in Physical and Theoretical Chemistry,*
1078 *Vol 71, Modelling of Molecular Structures and Properties* (Ed.: J.-L. Ri-
1079 vail), Elsevier, Amsterdam, **1990**, pp. 463-478.
- 1080 [69] A. D. Kline, W. Braun, K. Wüthrich, *J. Mol. Biol.* **1988**, 204, 675-724
- 1081 [70] G. M. Clore, A. T. Brünger, M. Karplus, A. M. Gronenborn, *J. Mol. Biol.*
1082 **1986**, 191, 523-551
- 1083 [71] G. M. Clore, D. K. Sukumaran, M. Nilges, J. Zarbock, A. M. Gronenborn,
1084 *EMBO J.* **1987**, 6, 529-537
- 1085 [72] G. M. Clore, A. M. Gronenborn, M. Nilges, D. K. Sukumaran, J. Zarbock,
1086 *EMBO J.* **1987**, 6, 1833-1842
- 1087 [73] G. M. Clore, A. M. Gronenborn, M. Nilges, C. A. Ryan, *EMBO J.* **1987**,
1088 26, 8012-8023
- 1089 [74] G. M. Clore, D. K. Sukumaran, M. Nilges, A. M. Gronenborn, *Biochem-*
1090 *istry* **1987**, 26, 1732-1745
- 1091 [75] G. M. Clore, M. Nilges, D. K. Sukumaran, A. T. Brünger, M. Karplus, A.
1092 M. Gronenborn, *EMBO J.* **1986**, 5, 2729-2735
- 1093 [76] W. F. van Gunsteren, M. Karplus, *Macromolecules* **1982**, 15, 1528-1544
- 1094 [77] J. Hermans, H. J. C. Berendsen, W. F. van Gunsteren, J. P. M. Postma,
1095 *Biopolymers* **1984**, 23, 1513-1518

- 1096 [78] B. R. Brooks, R. E. Bruccoleri, B. D. Olafson, D. J. States, S. Swami-
1097 nathan, M. Karplus, *J. Comput. Chem.* **1983**, 4, 187-217
- 1098 [79] A. M. Davis, S. A. St-Gallay, G. J. Kleywegt, *Drug Disc. Today* **2008**,
1099 13, 831-841
- 1100 [80] J. Agirre, G. Davis, K. Wilson, K. Cowtan, *Nature Chem. Biol.* **2015**, 11,
1101 303
- 1102 [81] W. G. Touw, R. P. Joosten, G. Vriend, *J. Mol. Biol.* **2016**, 428, 1375-1393
- 1103 [82] C. A. Schiffer, R. Huber, K. Wüthrich, W. F. van Gunsteren, *J. Mol. Biol.*
1104 **1994**, 241, 588-599
- 1105 [83] A. E. Torda, R. M. Scheek, W. F. van Gunsteren, *Chem. Phys. Lett.* **1989**,
1106 157, 289-294
- 1107 [84] P. Gros, W. F. van Gunsteren, W. G. J. Hol, *Science* **1990**, 249, 1149-1152
- 1108 [85] P. Gros, W. F. van Gunsteren, *Mol. Sim.* **1993**, 10, 377-395
- 1109 [86] C. A. Schiffer, P. Gros, W. F. van Gunsteren, *Acta Cryst. D* **1995**, 51,
1110 85-92
- 1111 [87] C. A. Schiffer, W. F. van Gunsteren, *Proteins Struct. Funct. Bioinf.* **1999**,
1112 36, 501-511
- 1113 [88] W. R. P. Scott, A. E. Mark, W. F. van Gunsteren, *J. Biomol. NMR* **1998**,
1114 12, 501-508
- 1115 [89] N. Hansen, F. Heller, N. Schmid, W. F. van Gunsteren, *J. Biomol. NMR*
1116 **2014**, 60, 169-187
- 1117 [90] D. A. Pearlman, P. A. Kollman, *J. Mol. Biol.* **1991**, 220, 457-479
- 1118 [91] U. Schmitz, A. Kumar, T. L. James, *J. Am. Chem. Soc.* **1992**, 114, 10654-
1119 10656
- 1120 [92] U. Schmitz, B. Ulyanov, A. Kumar, T. L. James, *J. Mol. Biol.* **1993**, 234,
1121 373-389
- 1122 [93] D. A. Pearlman, *J. Biomol. NMR* **1994**, 4, 1-16
- 1123 [94] D. A. Pearlman, *J. Biomol. NMR* **1994**, 4, 279-299
- 1124 [95] A. P. Nanzer, W. F. van Gunsteren, A. E. Torda, *J. Biomol. NMR* **1995**,
1125 6, 313-320
- 1126 [96] A. P. Nanzer, A. E. Torda, C. Bisang, C. Weber, J. A. Robinson, W. F.
1127 van Gunsteren, *J. Mol. Biol.* **1997**, 267, 1012-1025
- 1128 [97] A. P. E. Kunz, J. R. Allison, D. P. Geerke, B. A. C. Horta, P. H. Hünen-
1129 berger, S. Riniker, N. Schmid, W. F. van Gunsteren, *J. Comput. Chem.*
1130 **2012**, 33, 340-353

- 1131 [98] R. M. Scheek, A. E. Torda, J. Kemmink, W. F. van Gunsteren, in: *Com-*
1132 *putational Aspects of the Study of Biological Macromolecules by Nuclear*
1133 *Magnetic Resonance Spectroscopy* (Eds.: J. C. Hoch, F. M. Poulsen, C.
1134 Redfield), NATO ASI Series A225, Plenum Press, New York, USA, **1991**,
1135 pp. 209-217.
- 1136 [99] J. Fennen, A. E. Torda, W. F. van Gunsteren, *J. Biomol. NMR* **1995**, 6,
1137 163-170
- 1138 [100] T. Huber, W. F. van Gunsteren, *J. Phys. Chem. A* **1998**, 102, 5937-5943
- 1139 [101] B. Hess, R. M. Scheek, *J. Magn. Reson.* **2003**, 164, 19-27
- 1140 [102] J. W. Pitera, J. D. Chodera, *J. Chem. Theory Comput.* **2012**, 8, 4335-3451
- 1141 [103] B. Roux, J. Weare, *J. Chem. Phys.* **2013**, 138, 084107
- 1142 [104] A. Cavalli, C. Camilloni, M. Vendruscolo, *J. Chem. Phys.* **2013**, 138,
1143 094112
- 1144 [105] S. Olsson, J. Frellsen, W. Boomsma, K. V. Mardia, T. Hamelryck, *PLoS*
1145 *One* **2013**, 8, e79439
- 1146 [106] Y. Sugita, Y. Okamoto, *Chem. Phys. Lett.* **1999**, 314, 141-151
- 1147 [107] Y. Sugita, A. Kitao, Y. Okamoto, *J. Chem. Phys.* **2000**, 113, 6042-6051
- 1148 [108] H. Fukunishi, O. Watanabe, S. Takada, *J. Chem. Phys.* **2002**, 116, 9058-
1149 9067
- 1150 [109] N. Schmid, J. R. Allison, J. Dolenc, A. P. Eichenberger, A.-P. E. Kunz,
1151 W. F. van Gunsteren, *J. Biomol. NMR* **2011**, 51, 265-281
- 1152 [110] W. Rieping, M. Habeck, M. Nilges, *Science* **2005**, 309, 303-306
- 1153 [111] X. Daura, K. Gademann, B. Jaun, D. Seebach, W. F. van Gunsteren,
1154 A. E. Mark, *Angew. Chem.* **1999**, 111, 249-253; *Angew. Chem. Int. Ed.*
1155 **1999**, 38, 236-240
- 1156 [112] A. E. Torda, R. M. Scheek, W. F. van Gunsteren, *J. Mol. Biol.* **1990**,
1157 214, 223-235
- 1158 [113] C. A. Schiffer, R. Huber, K. Wüthrich, W. F. van Gunsteren, *J. Mol. Biol.*
1159 **1994**, 241, 588-599
- 1160 [114] D. Seebach, P. E. Ciceri, M. Overhand, B. Jaun, D. Rigo, L. Oberer, U.
1161 Hommel, R. Amstutz, H. Widmer, *Helv. Chim. Acta* **1996**, 79, 2043-2066
- 1162 [115] D. Seebach, S. Abele, K. Gademann, G. Guichard, T. Hintermann, B.
1163 Jaun, J. L. Matthews, J. V. Schreiber, L. Oberer, U. Hommel, H. Widmer,
1164 *Helv. Chim. Acta* **1998**, 81, 932-982
- 1165 [116] D. Seebach, K. Gademann, J. V. Schreiber, J. L. Matthews, T. Hinter-
1166 mann, B. Jaun, L. Oberer, U. Hommel, H. Widmer, *Helv. Chim. Acta*
1167 **1997**, 80, 2033-2038

- 1168 [117] T. Huber, A. E. Torda, W. F. van Gunsteren, *J. Comput.-Aided Mol.*
1169 *Design* **1994**, 8, 695-708
- 1170 [118] A. Laio, M. Parrinello, *Proc. Natl. Acad. Sci. USA* **2002**, 99, 12562-12566
- 1171 [119] M. Christen, B. Keller, W. F. van Gunsteren, *J. Biomol. NMR* **2007**, 39,
1172 265-273
- 1173 [120] Z. Gattin, J. Zaugg, W. F. van Gunsteren, *ChemPhysChem* **2010**, 11,
1174 830-835
- 1175 [121] J. H. Missimer, J. Dolenc, M. O. Steinmetz, W. F. van Gunsteren, *Prot.*
1176 *Sci.* **2010**, 19, 2462-2474
- 1177 [122] Z. Lin, W. F. van Gunsteren, *J. Chem. Phys.* **2015**, 143, 034110
- 1178 [123] A. Kuzmanic, N. S. Pannu, B. Zagrovic, *Nat. Commun.* **2014**, 5, 3220
- 1179 [124] Z. Dauter, A. Wlodawer, W. Minor, M. Jaskolski, B. Rupp, *IUCrJ* **2014**,
1180 1, 179-193
- 1181 [125] H. van den Bedem, J. S. Fraser, *Nat. Meth.* **2015**, 12, 307-318
- 1182 [126] M. O. Steinmetz, I. Jelesarov, W. M. Matousek, S. Honnappa, W. Jahnke,
1183 J. H. Missimer, S. Frank, A. T. Alexandrescu, R. A. Kammerer *Proc. Natl.*
1184 *Acad. Sci. USA* **2007**, 104, 7062-7067
- 1185 [127] C. Oostenbrink, A. Villa, A.E. Mark, W.F. van Gunsteren, *J. Comp.*
1186 *Chem.* **2004**, 25, 1656-1676
- 1187 [128] M. H. Lerche, F. M. Poulsen, *Protein Sci.* **1998**, 7, 2490-2498
- 1188 [129] L. J. Smith, W. F. van Gunsteren, J. R. Allison, *Protein Sci.* **2013**, 22,
1189 56-64
- 1190 [130] L. F. Pacios, C. Gómez-Casado, L. Tordesillas, A. Palacin, R. Sánchez-
1191 Monge, A. Díaz-Perales, *J. Comput. Chem.* **2012**, 33, 1831-1844
- 1192 [131] L. J. Smith, Y. Roby, J. R. Allison, W. F. van Gunsteren, *Biochemistry*
1193 **2013**, 52, 5024-5038
- 1194 [132] G. W. Han, J. Y. Lee, H. K. Song, C. S. Chang, K. Min, J. Moon, D. H.
1195 Shin, M. L. Kopka, M. R. Sawaya, H. S. Yuan, T. D. Kim, J. Choe, D.
1196 Lim, H. J. Moon, S. W. Suh, *J. Mol. Biol.* **2001**, 308, 263-278
- 1197 [133] D. Trzesniak, W. F. van Gunsteren, *Prot. Sci.* **2006**, 15, 2544-2551
- 1198 [134] B. T. Burnley, P. V. Afonine, P. D. Adams, P. Gros, *eLife* **2012**, 1, 1,
1199 e00311
- 1200 [135] A. Wlodawer, W. Minor, Z. Dauter, M. Jaskolski, *FEBS Journal* **2008**,
1201 275, 1-21
- 1202 [136] A. Wlodawer, W. Minor, Z. Dauter, M. Jaskolski, *FEBS Journal* **2013**,
1203 280, 5705-5736

- 1204 [137] R. P. Joosten, K. Joosten, G. N. Murshudov, A. Perrakis, *Acta Cryst.*
1205 **2012**, D68, 484-496
- 1206 [138] G. T. Montelione, M. Nilges, A. Bax, P. Güntert, T. Herrmann, J. S.
1207 Richardson, C. D. Schwieters, W. F. Vranken, G. W. Vuister, D. S.
1208 Wishart, H. M. Berman, G. J. Kleywegt, J. L. Markley, *Structure* **2013**,
1209 21, 1563-1570
- 1210 [139] P. D. Adams, K. Aertgeerts, C. Bauer, J. A. Bell, H. M. Berman, T.
1211 N. Bhat, J. M. Blaney, E. Bolton, G. Bricogne, D. Brown, S. K. Burley,
1212 D. A. Case, K. L. Clark, T. Darden, P. Emsley, V. A. Feher, Z. Feng,
1213 C. R. Groom, S. F. Harris, J. Hendle, T. Holder, A. Joachimiak, G. J.
1214 Kleywegt, T. Krojer, J. Marcotrigiano, A. E. Mark, J. L. Markley, M.
1215 Miller, W. Minor, G. T. Montelione, G. Murshudov, A. Nakagawa, H.
1216 Nakamura, A. Nicholls, M. Nicklaus, R. T. Nolte, A. K. Padyana, C. E.
1217 Peishoff, S. Pieniazek, R. J. Read, C. Shao, S. Sheriff, O. Smart, S. Soisson,
1218 J. Spurlino, T. Stouch, R. Svobodova, W. Tempel, T. C. Terwilliger, D.
1219 Tronrud, S. Velankar, S. C. Ward, G. L. Warren, J. D. Westbrook, P.
1220 Williams, H. Yang, J. Young, *Structure* **2016**, 24, 502-508
- 1221 [140] G. Miller, *Science* **2006**, 314, 1856-1857