

The refinement of NMR structures by molecular dynamics simulation

Andrew E. Torda¹ and Wilfred F. van Gunsteren¹

Laboratory of Physical Chemistry, University of Groningen, Nijenborgh 16, 9747 AG Groningen, The Netherlands

Received 26 January 1990

We discuss the use of molecular dynamics simulations as a tool for the refinement of structures based on NMR data. The procedure always involves the construction of a pseudo-energy term to model the experimental data and we consider the various approaches to this problem. We detail recent work where we account for the time averaging implicit in NMR measurements and attempt to model the experimental data more realistically. Finally, we discuss the problems and approximations involved in this work, the lack of consensus as to refinement methods and the scope for future developments.

1. Introduction

In its purest form, a molecular dynamics (MD) simulation involves some representation of a physical system. Force fields should be as accurate as possible, velocities should be carefully integrated and, hopefully, physical processes or properties will be reproduced. More recently, however, it has become popular to use MD and related methods as a tool for the refinement of molecular structures with respect to experimental data, especially from nuclear magnetic resonance (NMR) [1–4] or X-ray crystallographic measurements [5,6]. In this article, we shall concentrate on the specific case of data from nuclear magnetic resonance (NMR) measurements and discuss some of the problems with current procedures, recent improvements and future directions.

NMR measurements provide two main kinds of structural information. Firstly, *J*-coupling constants between protons separated by three bonds reflect the size of the included dihedral angle. The bulk of NMR information, however, consists of

nuclear Overhauser enhancement (NOE) measurements, corresponding to interproton distances, usually between sites less than about 5 Å apart [7]. Unfortunately, the distances are not precise and the set of data is usually not complete. This means that there is no analytical method which can generate structures consistent with the experimental data. Furthermore, the lack of data means that there is not even a single solution to the structural problem. Instead, there are one or more regions of conformational space containing structures consistent with the data. At the moment, the best way to estimate the location and size of this space is simply to generate families of plausible structures.

Most generally, MD simulations can be used for refinement whenever one can construct a potential-energy term whose value rises as the system deviates more from some experimentally measured parameter. When the system is simulated, it will tend to run downhill with respect to both real and artificial energetic terms. As long as the system's kinetic energy is regulated in some way, it should reach a state better in agreement with the experimental data and, ideally, also with respect to the physical terms in the force field. Compared to other optimisation methods, MD has the advantage that it is capable of temporarily

¹ Present address: Laboratory of Physical Chemistry, ETH Zentrum, CH-8092 Zürich, Switzerland. Email: torda @ igc.ethz.ch, wfvgn@igc.ethz.ch.

moving up potential-energy barriers and finding better regions of space with respect to the energetic terms.

The use of MD as a refinement tool introduces some considerations not normally present in a MD simulation. Firstly, neither force fields nor experimental data are perfect, so minima with respect to the two terms may not coincide exactly. One must then decide on the relative weight to give to the artificial terms. In the case of refinement using X-ray crystallographic data [5,6], it may be justified to have the artificial terms stronger than those of the normal force field. In the case of NMR refinement, where data is usually not as accurate, one may prefer for the physical terms in the force field to dominate. At the same time, the starting state of the system may be far away from the desired final state. Under these circumstances, it may be desirable to change the force field and use a simulation protocol so as to make energetic barriers more readily surmountable [8,9]. Continuing in this vein, it may be useful to consider dynamics schemes which are no longer Newtonian and merely serve as some means to drive the system into the desired state. We will also describe a recent innovation where we account for the fact that measured NMR properties reflect an average through time and not instantaneous values. This has led us to the introduction of a potential energy term which does not even conserve energy.

2. Modelling of the NOE with potential energy terms

Because the NOE is due to dipolar interactions between nuclei, the measured intensity is proportional to r^{-6} where r is the distance between a specified pair of protons. In practice, unknown distances are usually estimated by comparison of NOE intensity or buildup rates with those from protons at covalently fixed distances. This assumes that the reference and estimated interproton distances are subject to the same motions [10] and that only pairwise interactions contribute to the measured intensities or buildup rates [11]. A further complication arises, since molecules un-

dergo thermal motions and interproton distances will fluctuate on a time scale shorter than that for cross-relaxation processes. Then, one must remember that what is actually measured is a function of $\langle r^{-6} \rangle$ where the angle brackets denote an average over time. In the next section, we will address the problem of trying to model this non-linear time average in a MD force field.

Bearing in mind the nature of the NOE, the penalty function or potential-energy term to enforce the experimental data should be chosen. This term should be simple and computationally cheap while still driving the system to agree with the experimental data. The simplest choice for this term is quadratic with respect to the size of the violation of the distance constraint [1,12], so

$$V_{\text{dc}}(r) = \begin{cases} 0, & \text{if } r \leq r^0, \\ \frac{1}{2}K_{\text{dc}}(r - r^0)^2, & \text{if } r^0 < r < r^0 + \Delta r, \\ K_{\text{dc}}(r - r^0 - \frac{1}{2}\Delta r) \Delta r, & \text{if } r^0 + \Delta r \leq r, \end{cases} \quad (1)$$

where $V_{\text{dc}}(r)$ is the potential due to the distance-restraint term for a given pair of atoms, r is the instantaneous distance between the cross-relaxing nuclei and r^0 is the distance constraint calculated from the measured NOE. The force constant, K_{dc} , controls the relative strength of this artificial term in the force field.

Equation (1) actually describes three regimes. Firstly, the potential energy is zero in the ideal state where the instantaneous distance r is less than the constraint distance r^0 . Next, the term is quadratic for small violations of the restraint. Finally, if the distance r is larger than the sum of r^0 and Δr , the potential energy increases only linearly. This serves to put an upper limit on the size of the artificial force.

Fry et al. [13] used a similar term to eq. (1), but also performed some final refinements using a fourth-power term so

$$V_{\text{dc}}(r) = K(r - r^0)^4, \quad (2)$$

where K included a small series of constants and weights. Most recently, a form similar to eq. (2)

was used, but with a sixth-power term [14,15]. Certainly this approach can be used to produce steeper potential wells in the artificial energy term, but this again raises the issue of the degree to which one wants to balance the real and artificial terms in the force field.

Scarsdale et al. [16] proposed a penalty function which more accurately reflects the physical nature of the NOE, so

$$V_{\text{dc}}(r) = K \left[\left(r^{-3} - (r^0)^{-3} \right)^2 - (r^0)^{-6} \right]. \quad (3)$$

The purpose of this form can be seen by noting that for each distance constraint, $(r^0)^{-6}$ is a constant, so it does not contribute to the derivative. Equation (3) is thus appealing since it is quadratic with respect to a function that more closely reflects the measured NOE instead of the distance, a derived quantity.

3. Modelling the NOE as a time average

The pseudo-energy term described by eq. (3) is a better approximation to the NOE than, for example, eq. (1), but there is a more fundamental problem in this approach to enforcing experimental constraints. As described in section 2, the measured NOE is a weighted average of all conformations visited by a molecule on the NMR time scale. In the worst case, the molecule might be jumping between discrete states, none of which individually contain the distances used for the r^0 terms in eq. (1)–(3). Attempting to force structures to agree with the average NOE data may push the system into regions of conformational space which may hardly be populated in solution. In the case of small peptides where structural analysis is somewhat simpler, such discrete conformations have been identified and experimental data could only be explained by combinations of states [17,18]. For the refinement of a small oligosaccharide, Scarsdale et al. [16] actually performed simulations using a model based on two identifiable conformations. In the case of larger molecules like proteins, it is more difficult to identify individual conformations, since different parts of the molecule may be jumping between their own states.

Superimposed on these motions, there will be the normal thermal motions common to the whole molecule. In this situation, one cannot identify contributing conformers by inspection. The conformational space is best regarded as a continuum.

Recently [19], we proposed an alternative to the use of (1)–(3). Instead of forcing the individual distances r to agree with the experimental data, we used a potential energy based on the time-averaged distance $\bar{r}(t)$, so now

$$V_{\text{dc}}(\bar{r}(t)) = \begin{cases} 0, & \text{if } \bar{r}(t) \leq r^0, \\ \frac{1}{2} K_{\text{dc}} (\bar{r} - r^0)^2, & \\ 0, & \text{if } \bar{r}(t) > r^0. \end{cases} \quad (4)$$

This means that we only require a structure to satisfy the constraints as an average over time. Furthermore, it is possible to use the correct averaging of r to account for the power dependence of the NOE. So one could define

$$\bar{r}(t) = \langle r^{-6} \rangle^{-1/6}, \quad (5)$$

where, again, the angle brackets denote an average over time. The time scale of a simulation is usually much shorter than the correlation time for overall molecular tumbling, so one can neglect the influence of angular fluctuations [17]. Tropp [20] showed that under these circumstances, the NOE is effectively a function of r^{-3} , so we define

$$\bar{r}(t) = \langle r^{-3} \rangle^{-1/3}. \quad (6)$$

This can then be written in a form suitable for summation over the course of a MD trajectory,

$$\bar{r}(t) = \left(\frac{1}{t} \int_0^t r(t')^{-3} dt' \right)^{-1/3}. \quad (7)$$

Equation (7) is an average over the course of a whole trajectory, and it is this which must agree with the experimental data. It is thus the correct way to analyse a MD trajectory. It would, however, not be suitable as the basis for a pseudo-energy term in a simulation. If one were to use eq. (7), the averaging would be over an ever-growing time period and would become progressively less sensitive to changes in the system. In order to keep the system responsive to changes in the in-

stantaneous r , it is preferable to use some kind of running average. This is best done by using a memory function with a built-in exponential decay, so we define

$$\bar{r}(t) = \left(\frac{1}{\tau} \int_0^t e^{-t'/\tau} [r(t-t')]^{-3} dt' \right)^{-1/3}, \quad (8)$$

where τ is the decay constant for the exponential decay. This form of averaging results in an $\bar{r}(t)$ which does not feel the influence of short fluctuations in the system, but is still responsive to trends in behaviour. The degree of responsiveness is controlled by the parameter τ . In the limiting case of $\tau = 0$, there is no averaging. As τ becomes longer, the pseudo-energy term becomes less sensitive to fluctuations.

Originally, $\bar{r}(t)$ was defined by eq. (8), substituted directly into eq. (4) and the appropriate force constructed by taking the derivative with respect to r . This method worked well for a very simple model system [19] and for some small peptides, up to 12 residues (unpublished results). In the case of larger proteins and larger experimental data sets, however, large forces were occasionally generated. This came about since the force derived from eqs. (4) and (8) contained a fourth-power term with respect to $\bar{r}(t)/r(t)$.

This problem was avoided by adopting an unusual approach to enforcing experimental data in an MD simulation. No pseudo-energy term was defined at all. Instead, only a force was defined [21], so

$$F_i(t) = \begin{cases} 0, & \text{if } \bar{r}(t) \leq r^0, \\ -K_{dc} [\bar{r}_{ij}(t) - r^0] \frac{r_{ij}(t)}{r_{ij}(t)}, & \\ \text{if } \bar{r}(t) > r^0, \end{cases} \quad (9)$$

where $F_i(t)$ is the force on atom i due to atom j and $r_{ij} = r_i - r_j$. Equations (4–9) introduce some unusual properties into a MD simulation. The force is no longer simply a function of coordinates, but also, because of the use of $\bar{r}(t)$, a function of all previous configurations. This in turn means that the force field is no longer con-

servative. The consequences of this are discussed in section 5.

Furthermore, the use of eq. (9) means that it is no longer meaningful to refer to the pseudo-energy associated with a particular structure since, strictly, no such measure exists. Instead, it only makes sense to look at the properties of a system over some period of time. We then use eq. (9) to enforce experimental constraints, but we judge the success of a structural refinement by calculating the trajectory average given by eq. (7).

4. Application of time-averaged NOE constraints

The first tests of time-averaged distance constraints were performed on a very artificial system consisting of only three Lennard–Jones particles [19]. Two of these were fixed in space while the third was free to move without even periodic boundary conditions to restrain the accessible space. Two distance constraints were then imposed so as to require the mobile particle to be close to both of the fixed particles simultaneously. No single conformation existed which could satisfy both constraints instantaneously, so the constraints could only be satisfied, on average, if the mobile particle moved between the two fixed particles. This was analogous to a molecule having to jump between distinct conformations which individually could not explain experimental data.

This small system served to demonstrate the feasibility of the method and provided some indication of the effect of the adjustable parameters. Firstly, it was necessary to establish reasonable values for τ , the decay constant for the memory function in eq. (8). When $\tau = 0$, the system had no memory and no averaging of the calculated distances. As expected, the imposed constraints were not satisfied either instantaneously or as an average over the trajectory. As τ was increased and the averaging was over longer periods, the mobility of the system increased and the free particle was able to spend time close to both of the fixed particles in turn. Although at any instant at least one of the distance constraints was violated, both could be satisfied as an average over the trajectories. Most importantly, it was clear that the size of

was not critical, as long as it was longer than the time required for the system to visit all the conformations necessary to explain the experimental data. Furthermore, it could be seen that in order to achieve reasonable averaging, the length of a trajectory should be approximately an order of magnitude larger than τ .

The model system also served to highlight the increased reliance on the physical terms in the force field. If τ was too long, the distance constraints did not restrain the mobile particle. Because there were no covalent terms in the force field, the system could move far away from either of the two desired conformations and, on average, violate the distance restraints. If anything, this was encouraging for application to real molecules. It showed that as τ was increased, the influence of the artificial constraints decreased. It suggested that if one had a good conformation with respect to both physical and pseudo-energy terms, then the refinement procedure would become closer to a realistic molecular simulation.

Time-averaged distance constraints were subsequently applied to the refinement of a large molecule, the protein tendamistat [21]. This was an important test for several reasons. Firstly, the original structure from the Zürich group had served as a demonstration of the ability of NMR data to determine a solution structure [22]. With additional experimental information, tendamistat became one of the most precisely determined solution structures up to that time [23,24]. The structures, however, were the result of extensive distance-geometry calculations using the variable target function method [25], so they were static solutions to the structural problem. Because there was no evidence of conformational heterogeneity, it was of particular interest to see the additional conformational space that would be allowed using time averaging on the distance constraints.

Two of the published distance-geometry structures were selected for MD refinement on the basis of their agreement with the distance constraints. For each structure parallel simulations were run with normal MD refinement ($\tau = 0$) and with τ increased to a final value of 1.5 ps. All results were judged in terms of averages over 20 ps simulation trajectories.

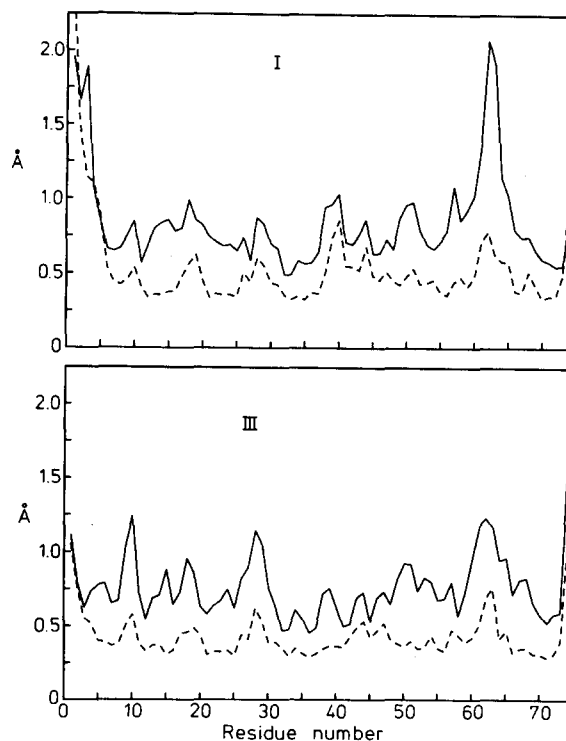


Fig. 1. Root mean square positional fluctuations of α -carbons in tendamistat over 20 ps trajectories. Solid lines are from runs using time averaged NOE's, dashed lines from runs using conventional MD refinement. I and III refer to the structure names used in ref. [23]. Taken from ref. [21].

MD refinement of the structures, with or without time-averaged distance constraints, resulted in a large improvement in potential energies of the structures. There were, however, significant differences in both residual violations of the experimental constraints and in the mobility of the structures during the simulations. Trajectories using time-averaged constraints had average violations typically 70–80% of those generated by normal refinement. The most remarkable difference, however, was in the mobility of the molecules. Considering backbone α -carbons, root-mean-square positional fluctuations were consistently larger and, at some sites, double those using conventional refinement (fig. 1).

Aside from a general increase in mobility, it was clear that certain parts of the molecule behaved quite differently under the influence of time-averaged distance constraints. The most

striking example was the extra motion of the sidechain of Tyr 15. Figure 2 shows the neighbouring peptide segment from two simulation snapshots superimposed on the starting distance-geometry structure. The experimental data contained 26 NOE restraints for this residue and the distance geometry calculations suggested that it had a rather rigid and well-defined location. The MD simulations, however, suggested an alternative explanation for the experimental data. The large number of NOE constraints may actually have resulted from the rapid motions of the sidechain bringing it near a series of other sites in turn. These rapid motions would result in only average NMR resonances and averaged NOE's would be detected in the experiments. This explanation is also consistent with the fact that no electron density was observed for the sidechain in the X-ray crystallographic structure determination [26].

The use of time-averaged distance constraints has advantages beyond simply better reproducing the experimental data. Firstly, because a structure is not required to satisfy all the restraints simultaneously, it is not necessary to use large force constants for the pseudo-energy terms. This should result in less distorted structures. The greater mo-

bility of the structures means that they spread through a greater region of conformational space. This is desirable when a refinement is viewed as a search of conformational space. Finally, trajectories produced by this method include a wide range of possible conformations rather than simply a cluster of structures centred about some artificial average. This extra realism should become important as NMR based structures become more frequently used for purposes such as drug design.

5. Future improvements in structural refinement

Although the introduction of time-averaged distance constraints appears to be an improvement over the use of static constraints, there are still some unresolved problems. The force on a particle can change over time, even if coordinates do not change. This results in heating of the structure. Unfortunately, this heating is a result of the distance constraints, so it need not be evenly distributed over the molecule. This means that individual atom-temperature coupling may be more appropriate than the currently used overall coupling to a temperature bath [27]. This will lead to more realistic distribution of the system over its

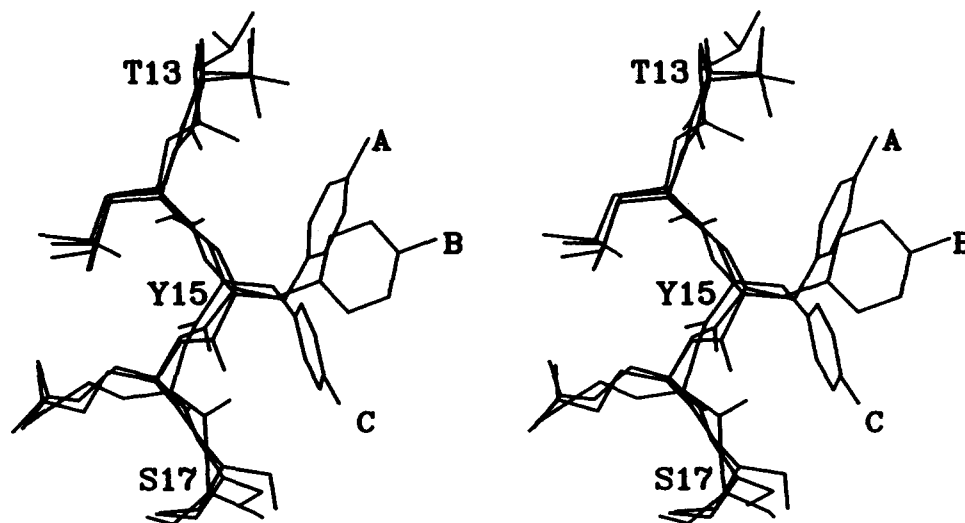


Fig. 2. Mobility of Tyr 15 in a 20 ps simulation of tendamistat. A stereo view of the peptide segment from residues 13 to 17 is shown from (A) 9 ps into a simulation trajectory, (B) distance-geometry starting structure (C), 16.2 ps into trajectory. Both MD structures were least-squares fitted to (B) on the basis of all backbone heavy atoms. Taken from ref. [21].

energy surface, but to a less Newtonian dynamics scheme.

The method is also limited by the finite time of a simulation. Ideally, one should simulate long enough to average over the conformational space covered by a molecule on the NMR time scale. This would typically be closer to milliseconds rather than the picoseconds currently usually simulated.

More generally, other fundamental changes should be considered for the modelling of NOE data. As a matter of principle, one should not construct the pseudo-energy term as a function of distances which are a derived quantity. Instead, one should write

$$V_{\text{nmr}} = K(\text{NOE}_{\text{obs}} - \text{NOE}_{\text{calc}})^2, \quad (10)$$

where NOE_{obs} and NOE_{calc} are the observed and calculated NOE's, respectively. Unfortunately, eq. (10) does not specify what approximations should be used for calculating the NOE. Certainly, the method should consider averaging through time, but it should also account for multiple spin cross-relaxation pathways [28,29]. This could be done by an iterative method [30,31], but Yip and Case did actually incorporate an expression including multi-spin effects into a MD pseudo-energy term [32]. As described, the procedure would need to be simplified in order to be computationally practical, but it is not clear what compromises would be necessary.

Aside from better modelling of the NOE, there are even more fundamental questions associated with the generation of structures from NMR data. There is no consensus as to the relative importance one should attach to the physical force field and the experimental data. At one extreme, one may use distance-geometry structures without any MD refinement or with refinement in a force field based on very simplified non-bonded interactions [9]. Alternatively, approaches such as time-dependent distance constraints attempt to minimise the influence of pseudo-energy terms.

Furthermore, there is still no agreed method for quantifying the degree to which a structure is defined. Clearly the root-mean-square positional differences within a family of structures is not a

sufficient criterion. Several authors have noted that different procedures introduce different biases and will produce families of structures centred about different averages with different spreads [33–35]. Moreover, if one enforces distance constraints as time-averaged quantities, additional restraints may actually increase the space covered by a structure, rather than making it appear better defined [21].

Finally, the issue of how best to treat J -coupling constants must also be addressed. These measurements do provide information on dihedral angles, but are also subject to averaging and considerations such as relative weighting in the force field.

In conclusion, the use of NMR for determining solution structures has become a popular field, but it is still possible for different groups to generate different solutions given the same experimental data. The questions discussed in this article thus seem relevant when considering the reliability, reporting and use of NMR structures.

Acknowledgements

These investigations were supported in part by the Netherlands' Foundation for Chemical Research (SON) with financial aid from the Netherlands's Technology Foundation (STW). We also wish to thank ICI for their financial support and interest in the project. We are especially grateful to Drs. I.D. Kuntz, R.M. Scheek and H.J.C. Berendsen for their suggestions and stimulating discussions.

References

- [1] W.F. van Gunsteren, R. Kaptein and E.R.P. Zuiderweg, Proc. NATO/CECAM Workshop on Nucleic Acid Conformation and Dynamics, Orsay (1984), W.K. Olsen, ed., pp. 79–82.
- [2] R. Kaptein, E.R.P. Zuiderweg, R.M. Scheek, R. Boelens and W.F. van Gunsteren, *J. Mol. Biol.* 182 (1985) 179.
- [3] G.M. Clore, A.M. Gronenborn, A.T. Brünger and M. Karplus, *J. Mol. Biol.* 191 (1985) 523.
- [4] A.T. Brünger, G.M. Clore, A.M. Gronenborn and M. Karplus Proc. Nat. Acad. Sci. USA. 83 (1986) 3801.

- [5] A.T. Brünger, J. Kuriyan and M. Karplus, *Science* 235 (1987) 458.
- [6] M. Fujinaga, P. Gros and W.F. van Gunsteren, *J. Appl. Crystallogr.* 22 (1989) 1.
- [7] K. Wüthrich, *NMR of Proteins and Nucleic Acids* (Wiley, New York, 1986).
- [8] M. Nilges, A.M. Gronenborn, A.T. Brünger and G.M. Clore, *Protein Eng.* 2 (1988) 27.
- [9] M. Nilges, G.M. Clore and A.M. Gronenborn, *FEBS Lett.* 229 (1988) 317.
- [10] J.H. Noggle and R.E. Schirmer, *The Nuclear Overhauser Effect* (Academic, New York, 1971).
- [11] A. Kalk and H.J.C. Berendsen, *J. Magn. Reson.* 24 (1976) 343.
- [12] W.F. van Gunsteren and H.J.C. Berendsen, *Groningen Molecular Simulation (GROMOS) library manual* (Biosmos, Groningen, 1987).
- [13] D.C. Fry, V.S. Madison, D.R. Bolin, D.N. Greeley, V. Toome and B.B. Wegrzynski, *Biochemistry* 28 (1989) 2399.
- [14] H. Widmer, M. Billeter and K. Wüthrich, *Proteins* 6 (1989) 357.
- [15] M. Billeter, Th. Schaumann, W. Braun and K. Wüthrich, *Biopolymers* 29 (1989) 695.
- [16] J.N. Scarsdale, P. Ram, J.H. Prestegard and R.K. Yu, *J. Comput. Chem.* 9 (1988) 133.
- [17] H. Kessler, C. Griesinger, J. Lutz, A. Müller and W.F. van Gunsteren, *J. Am. Chem. Soc.* 110 (1988) 3393.
- [18] H. Pepermans, D. Tourwé, G. van Binst, R. Boelens, R.M. Scheek, W.F. van Gunsteren and R. Kaptein, *Biopolymers* 27 (1988) 323.
- [19] A.E. Torda, R.M. Scheek and W.F. van Gunsteren, *Chem. Phys. Lett.* 157 (1989) 289.
- [20] J. Tropp, *J. Chem. Phys.* 72 (1980) 6035.
- [21] A.E. Torda, R.M. Scheek and W.F. van Gunsteren, *J. Mol. Biol.* 214 (1990) 223.
- [22] A.D. Kline, W. Braun and K. Wüthrich, *J. Mol. Biol.* 189 (1986) 377.
- [23] A.D. Kline, W. Braun and K. Wüthrich, *J. Mol. Biol.* 204 (1988) 675.
- [24] K. Wüthrich, *Science* 243 (1989) 45.
- [25] W. Braun and N. Gö, *J. Mol. Biol.* 186 (1985) 611.
- [26] M. Billeter, A.D. Kline, W. Braun, R. Huber and K. Wüthrich, *J. Mol. Biol.* 206 (1989) 677.
- [27] H.J.C. Berendsen, J.P.M. Postma, W.F. van Gunsteren, A. DiNola and J.R. Haak, *J. Chem. Phys.* 81 (1984) 3684.
- [28] J.W. Keepers and T.L. James, *J. Magn. Reson.* 57 (1984) 404.
- [29] B.A. Borgias and T.L. James, *J. Magn. Reson.* 79 (1988) 493.
- [30] R. Boelens, T.M.G. Koning, G.A. van der Marel, J.H. van Boom and R. Kaptein, *J. Magn. Reson.* 82 (1989) 290.
- [31] K.M. Banks, D.R. Hare and B.R. Reid, *Biochemistry* 28 (1989) 6996.
- [32] P. Yip and D.A. Case, *J. Magn. Reson.* 83 (1989) 643.
- [33] W.J. Metzler, D.R. Hare and A. Pardi, *Biochemistry* 28 (1989) 7045.
- [34] T.F. Havel, *Biopolymers* (1990) in print.
- [35] R.M. Levy, P.A. Bassolino, D.B. Kitchen and A. Pardi, *Biochemistry* 28 (1989) 9361.