# A large data set comparison of protein structures determined by crystallography and NMR: Statistical test for structural differences and the effect of crystal packing

Michael Andrec,[1] David A. Snyder,[2] Zhiyong Zhou,[1] Jasmine Young,[3] Gaetano T. Montelione,[2] and Ronald M. Levy[1]*

[1] BioMaPS Institute for Quantitative Biology, Northeast Structural Genomics Consortium and Department of Chemistry and Chemical Biology, The State University of New Jersey, Piscataway, New Jersey 08854

[2] Center for Advanced Biotechnology and Medicine, Northeast Structural Genomics Consortium and Department of Molecular Biology and Biochemistry, The State University of New Jersey, Piscataway, New Jersey 08854

[3] RCSB Protein Data Bank and Department of Chemistry and Chemical Biology Rutgers, The State University of New Jersey, Piscataway, New Jersey 08854

## ABSTRACT

*The existence of a large number of proteins for which both nuclear magnetic resonance (NMR) and X-ray crystallographic coordinates have been deposited into the Protein Data Bank (PDB) makes the statistical comparison of the corresponding crystal and NMR structural models over a large data set possible, and facilitates the study of the effect of the crystal environment and other factors on structure. We present an approach for detecting statistically significant structural differences between crystal and NMR structural models which is based on structural superposition and the analysis of the distributions of atomic positions relative to a mean structure. We apply this to a set of 148 protein structure pairs (crystal vs NMR), and analyze the results in terms of methodological and physical sources of structural difference. For every one of the 148 structure pairs, the backbone root-mean-square distance (RMSD) over core atoms of the crystal structure to the mean NMR structure is larger than the average RMSD of the members of the NMR ensemble to the mean, with 76% of the structure pairs having an RMSD of the crystal structure to the mean more than a factor of two larger than the average RMSD of the NMR ensemble. On average, the backbone RMSD over core atoms of crystal structure to the mean NMR is approximately 1 Å. If non-core atoms are included, this increases to 1.4 Å due to the presence of variability in loops and similar regions of the protein. The observed structural differences are only weakly correlated with the age and quality of the structural model and differences in conditions under which the models were determined. We examine steric clashes when a putative crystalline lattice is constructed using a representative NMR structure, and find that repulsive crystal packing plays a minor role in the observed differences between crystal and NMR structures. The observed structural differences likely have a combination of physical and methodological causes. Stabilizing attractive interactions arising from intermolecular crystal contacts which shift the equilibrium of the crystal structure relative to the NMR structure is a likely physical source which can account for some of the observed differences. Methodological sources of apparent structural difference include insufficient sampling or other issues which could give rise to errors in the estimates of the precision and/or accuracy.*

## INTRODUCTION

The growth in the size of the Protein Data Bank[1] (PDB) in recent years has substantially increased the number of proteins for which both nuclear magnetic

resonance (NMR) and crystal structures have been deposited. In an October 2005 survey of the database available on the PDB web site (www.pdb.org/pdb/clusterExp-Methods.do), a total of 288 sets of PDB-deposited structural models were found in which each member of the set has high sequence similarity with the other members, and where both crystal and NMR structures are represented. The size of this set is increasing due to progress in traditional structural biology as well as from the substantial investment in protein structure determination by the Protein Structure Initiative in structural genomics. Given this large set of structures determined using both experimental methods, it becomes possible to compare NMR and X-ray crystal structures on a much larger scale than was possible previously and to examine statistically the possible effects of crystal vs solution environments on protein structure.

In this article, we explore the structural differences between NMR and crystal structural models of the same protein using a data set of 148 such structure pairs. Conformational differences are detected using a method based on the FindCore method of Snyder & Montelione[2] for defining core atom sets for purposes of structural superposition. We make use of the ensemble nature of the NMR structural model to determine the statistical significance of deviations in atomic positions between the crystal and mean NMR model by asking about the likelihood of the position of a given atom in the crystal structure relative to the positions of the corresponding atom in the NMR ensemble. This is similar in spirit to methods for the statistical description of structural ensembles developed by Gerstein and co-workers[3,4]; however, it differs in that we do not attempt to fit full three-dimensional normal (Gaussian) distributions for each atom. Instead, we use a statistic based on the distance of each atom from the mean. In this way, it is similar to the distance-derived pseudo $B$-factor that has been proposed for crystallographic phase determination via molecular replacement using NMR structures.[5] Our method provides both local information about which particular atoms have changed position, as well as "aggregated" statistics that provide a global summary of the overall agreement of the two structures.

We apply this statistical methodology to a large subset of the available protein structure pairs in the PDB to study the nature of the differences between NMR and crystal structures of the same protein. Possible methodological and physical origins for the structural differences are examined, including the use of more modern data and structure determination protocols, structure quality, and similarity of structure determination conditions. We explore the role of crystal contacts using a systematic procedure that replaces the protein structure in the crystal lattice with the NMR solution structure and analyzes the resulting system for steric clashes.

The question of the consistency of two structural models must be based on a proper estimation of precision and accuracy, which allows the user of the model to assess the reproducibility and correctness, respectively, of a given result. When the fundamental data is interpreted in terms of a model, uncertainties, ambiguities, and errors in the raw data must be "propagated" through the modeling process in order to obtain estimates of precision and accuracy in the model parameters (in this case, the atomic coordinates). The proper assessment of precision and accuracy is necessary in structural biology to evaluate whether differences between different structural models are or are not significant, and, by implication, whether such differences contain functional information.

In small molecule crystallography, the quantitation of precision is well-established, leading to anisotropic $B$-factors that represent thermal motion via the Debye-Waller relationship.[6] In macromolecular crystallography, the situation is more complex, as there is typically insufficient data to estimate anisotropic $B$-factors. The reported $B$-factors represent the scaling of an atom's contribution to the diffraction pattern required to optimally match the back-calculated and observed diffraction patterns. While the Debye-Waller equation enables the interpretation of such indirectly calculated $B$-factors as coordinate uncertainties, they do not conceptually represent the same quantity as the variability in atomic position across an ensemble of NMR-derived structures. Recent results indicate that $B$-factors and the variances in NMR-derived atomic position may be compatible measures of precision.[7]

In NMR structure determination, there exist only nascent methods for the statistically rigorous estimation of precision, and active research is ongoing in this area.[8] Traditionally, NMR structures in the PDB have been deposited as ensembles which are the result of replicating the structure determination procedure several times. The conformational variability across this ensemble can then be used as a measure of reproducibility, and therefore the precision. Measurement of the divergence of an ensemble of structures typically involves calculating a superposition, minimizing the root-mean-square deviation (RMSD) from an average or representative set of atomic coordinates. This approach suffers from pitfalls, however, in that the set of atoms being superimposed must be well-chosen.[2,9] Furthermore, this strategy also implicitly assumes that the ensemble deposited in the PDB well-characterizes the uncertainty inherent in the data.[10–14]

The estimation of accuracy also has complicating factors, not the least of which is the fact that the "correct" structure is typically unknown. In the absence of knowledge of the true structure, one can attempt to estimate accuracy using internal measures. This is now routinely done in macromolecular X-ray crystallography using the free $R$-factor,[15,16] which provides a cross-validated estimate of the goodness-of-fit that is correlated with the

phase error.[16] The direct application of such a strategy in NMR structure determination is difficult because of the combined effect of low data density and high information content of individual data points[17] (e.g. critical long-range nuclear Overhauser effect (NOE) restraints), but extensions based on a jackknife approach have been suggested.[18] More sophisticated approaches based on the assessment of the degree to which the model satisfies the raw data have also been described.[19,20]

While crystal and NMR structures occasionally differ from each other because of specific structural rearrangements arising from interactions in the crystalline environment (e.g. salt bridges, steric interactions),[21–24] one more generally sees more "diffuse" differences that are reflected in the statistical distributions of structural properties. For example, it has been observed that "knowledge-based" potentials can differ depending on whether they are derived from NMR or crystal structure databases.[25,26] Furthermore, packing and global conformational properties of crystal and NMR structures have been shown to exhibit significant statistical differences. For example, Garbuzynskiy *et al.*[27] have examined the statistical distribution of the distance of closest approach for all pairs of residues in a set of 60 proteins the structures of which have been determined by both crystallography and NMR. They found that this distribution is skewed toward shorter distances in NMR structures relative to crystal structures, that is short residue–residue distances are overrepresented while longer distances are underrepresented in NMR structures relative to crystal structures. In addition, they observed differences in hydrogen bonding patterns. The authors of that study concluded that the differences are likely due to the methodology used to determine the structural models, as the observed effects are correlated with the structure determination software used to do the NMR structural modeling, and become less prominent after re-refinement using explicit solvent. The view that deviations in the global properties of NMR structures from those of crystal structures represents a lack of accuracy in the NMR structures has been used as the basis for the introduction of additional constraints into the NMR structure determination process to make the resulting structures have more of the properties of high-resolution crystal structures.[28,29]

The investigation of specific structural differences between crystal and NMR structure pairs has been common since the earliest days of NMR-based structure determination. The most recent systematic reviews of such comparisons are now a decade old,[21–24] with later comparisons scattered throughout the structural biology literature in the form of reports on individual proteins. This paper will not attempt to review this body of work. Rather, we will provide a systematic study of overall structural differences in a large sample of structure pairs using the statistical methodologies for superposition and structural comparison described below.

## MATERIALS AND METHODS

### Superposition of structures using a core atom set

To evaluate the structural variability of the NMR-derived ensemble and to determine whether the crystal structure is significantly different from the NMR-derived ensemble of structures, it is necessary to superimpose the ensemble and the crystal structure in a common reference frame. To perform such a superposition one must choose which atoms will be used in optimizing the translational and rotational degrees of freedom that relate the structural models to each other. Although such superpositions are routinely performed by the depositors of structures to the PDB based on subsets of residues which they consider to be ordered, the criteria used to which determine this subset are not uniform between labs or across time.[2] Using depositor-defined "core" atoms or residues would not allow us to make consistent comparisons across the range of 148 NMR-crystal structure pairs. Instead, we use the FindCore method[2] to identify a core set of atoms (which need not include all of the atoms in any given residue) for which superpositions of the NMR-derived ensemble and the crystal structure are well-defined.

Briefly, the FindCore method takes an ensemble of structures and for selected pairs of atoms in the protein calculates the variance of the distance between that pair of atoms across the ensemble. Typically, only a subset of atoms are considered, for example the backbone heavy atoms (N, C, and $C^\alpha$) or all heavy atoms (backbone and side-chain). The row vectors of variances associated with each atom is transformed into "order parameters" by counting the number of elements less than a critical variance, and atoms are categorized into core and noncore classes using least squares clustering (see Ref. 2 for further details). It should be noted that the "core" determined by FindCore is a reflection only of the degree of order within the NMR ensemble, and need not correspond to any physical property of the protein (such as a hydrophobic core). In addition, FindCore can detect the presence of multiple domains, which are defined as a set of substructures which when subjected to rigid body motion reduce the apparent disorder (as opposed to biological structural domains). In all of the proteins studied here, FindCore determines that only one such domain is present.

To define a mean NMR structure onto which to superimpose both the NMR-derived ensemble and the crystal structure, all pairwise superpositions of the models in the NMR ensemble are first calculated. The "central model" is defined to be the model $i$ which minimizes $S(i) = \sum_j R_{ij}$, where the sum is over all models in the NMR-derived ensemble of structures and where $R_{ij}$ represents the RMSD (calculated over core atoms) between models $i$ and $j$. All models other than the central model are then

superimposed onto the central model, and the mean coordinates for all atoms under consideration, both core and noncore, are calculated based on this superimposition. The structures of the NMR-derived ensemble and the crystal structure are then superimposed to these mean coordinates.

### A statistical test for difference in atomic position between crystal and NMR structures

Given the superpositions of both the crystal structure and the NMR-derived ensemble to the mean NMR-derived coordinates calculated using the methods described earlier, we define the statistic

$$\chi_\psi^2(i) = \frac{\|\mathbf{a}_{\mathrm{crystal}}(i) - \mathbf{m}(i)\|^2}{\frac{1}{N}\sum_j \|\mathbf{a}_{\mathrm{NMR}}(i,j) - \mathbf{m}(i)\|^2} \quad (1)$$

which measures the deviation of coordinates of atom $i$ in the crystal structure $\mathbf{a}_{\mathrm{crystal}}(i)$ from the mean NMR-derived coordinates of that atom $\mathbf{m}(i)$ normalized by the variation in that atom's coordinates across the NMR-derived ensemble (where $\mathbf{a}_{\mathrm{NMR}}(i,j)$ represents the coordinates of the $i$-th atom in the $j$-th model and $N$ is the number of models in the NMR-derived ensemble). The $\chi_\psi^2$ statistic defined in Eq. (1) increases as the position of atom $i$ in the crystal moves farther away from the mean NMR position (for a given NMR ensemble), and also increases as the spread of atom positions in the NMR ensemble decreases. Therefore, it agrees with our intuition concerning the significance of a difference in atomic position: the position in the crystal structure is deemed to be different if its distance from the NMR mean is large compared to the width of the NMR bundle itself. We can quantify the degree of significance using the standard method of the $P$-value,[30] which we define to be the probability that the $\chi_\psi^2$ statistic under the null hypothesis is as big or bigger than that observed.
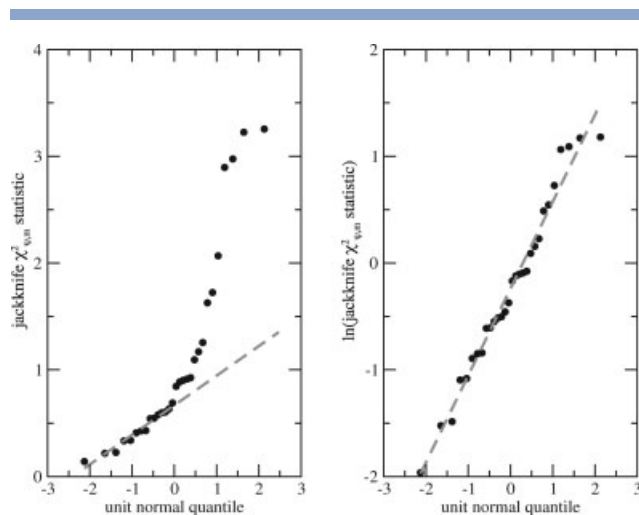
Although the $\chi_\psi^2$ statistic is formally similar to the $\chi^2$ and $F$ statistics of elementary statistical theory,[30] its distribution under the null hypothesis may be quite different from that of its more standard relatives since the quantities in the numerator and denominator arise from variables that are likely not normally distributed. It is possible to show that the distribution of $\chi_\psi^2$ under the null hypothesis if we assume that the atomic positions are distributed according to an isotropic three-dimensional normal distribution with variance $\sigma^2$ is a gamma distribution with mean 1 and variance 2/3.[31] However, the distribution of atomic positions is unlikely to be isotropic[3,4] and may not even be normally distributed. Therefore, we have chosen to estimate the $P$-values using an empirical method.

We evaluate whether a given value of $\chi_\psi^2(i)$ is significant by explicitly constructing a sample of possible $\chi_\psi^2(i)$

values under the null hypothesis using a jackknife procedure as follows. We repeat the core finding step using the NMR-derived ensemble only and then treat each member $m$ of the NMR-derived ensemble as the "crystal structure", recalculating the mean coordinates $\mathbf{m}_{(m)}$ of the NMR-derived ensemble having left out model $m$, superimposing the NMR-derived ensemble onto that mean and then calculating the statistic

$$\chi_{\psi,m}^2(i) = \frac{\|\mathbf{a}_{\mathrm{NMR}}(i,m) - \mathbf{m}_{(m)}(i)\|^2}{\frac{1}{N-1}\sum_{j\neq m}\|\mathbf{a}_{\mathrm{NMR}}(i,j) - \mathbf{m}_{(m)}(i)\|^2} \quad (2)$$

for each atom $i$ and model $m$. The collection of statistics $\{\chi_{\psi,m}{}^2(i)\}_{m=1\ldots N}$ constitutes a sample taken under the null hypothesis. Since there are typically only 10–20 members in NMR ensembles deposited in the PDB, this sample is not sufficiently large to accurately estimate a $P$-value directly. However, we have empirically found that the distribution of $\chi_{\psi,m}^2(i)$ for a given core atom $i$ is approximately log-normal (Fig. 1). The choice of a log-normal is conservative in the sense that it leads to a higher likelihood (smaller statistical significance) of more distant outliers. Therefore, we use the sample mean and variance of the log-transformed values $\mu_\psi(i) = \frac{1}{N}\sum_m \log(\chi_{\psi,m}^2(i))$ and $\sigma_\psi^2(i) = \frac{1}{N-1}\sum_m[\log(\chi_{\psi,m}^2(i))$



**Figure 1**

*Normal Q-Q plot of the jackknife statistics $\chi_{\psi,m}^2$ [Eq. (2)] for the $C_\alpha$ atom of residue 41 over the 30-member NMR structural ensemble 1IKM. The normal Q-Q plot is constructed from a set of values $S = \{x_i\}_{i=1\ldots N}$ by plotting each value $x_j$ vs the unit normal quantile $q_j$, which is defined as the value for which the standard normal cumulative probability $\int_{-\infty}^{q_j} (2\pi)^{-1/2}exp(-x^2/2)\mathrm{d}x$ equals the empirical fraction of values in $S$ smaller than $x_j$.[32] (A) Normal Q-Q plot of the untransformed jackknife statistics. The dashed line is a least-squares fit to the smallest 15 points, and the deviation from linearity indicates a non-normal distribution. (B) Similar plot of the log-transformed jackknife statistics, showing the agreement with a log-normal distribution. The dashed line is the least-squares fit to all 30 points. Similar behavior is seen for other residues and for other proteins in the test set.*

$-\mu_\psi(i)]^2$ to calculate the $z$-score

$$z(i) = \frac{\log(\chi_\psi^2(i)) - \mu_\psi(i)}{\sigma_\psi(i)}. \qquad (3)$$

If the underlying distribution is log-normal, then the $z$-score will follow a standard normal distribution and the $P$-value can be calculated using the error function.[33] In summary, our measure of the statistical significance of structural differences is based on the assumption that the structural ensemble provided in the PDB entry represents its precision, that the distribution of the $\chi_\psi^2$ statistic is log-normal under the null hypothesis for core atoms, and that the uncertainty in the atomic positions of the crystal structure (as estimated from its $B$ factors) is small compared to that of the NMR structure. In the following section, we comment on the effects of uncertainties in atomic positions in the the crystal structure on this analysis. It is possible that the use of alternative sampling strategies for generating the NMR ensemble[11,13,14,34] may not result in log-normal distributions under the null hypothesis, and the form of the distribution would have to be re-evaluated for these sampling methods.

### Analysis of steric clashes arising from crystal packing

In addition to assessing the location and extent of the structural differences between crystal and NMR structures of the same protein, we explore one possible origin of those differences which is related to crystal contacts by determining how consistent the NMR structure is with the crystalline environment. A reliable automatic superposition method is essential for such a study. This is because the orientation of the NMR structure relative to the $(x,y,z)$ coordinate axes of the PDB file are arbitrary, while those of the crystal structure are intimately connected to the axes of the unit cell and symmetry axes of the crystal. Therefore, to replicate a crystalline environment of a given crystal structure using a different conformation of the asymmetric unit derived from an NMR structure, it is first necessary to translate and rotate the NMR structure to coincide as best as possible with the crystal structure. FindCore provides a method for performing this transformation. Once the transformation has been found, the crystal environment can be simulated by application of the appropriate symmetry operations and distances can be examined to find close contacts.

Specifically, we use the core atoms found by FindCore to bring the NMR ensemble into the same coordinate frame as the crystal structure, extract the central model as defined in the description of the FindCore algorithm above, build copies of that structure by applying the translation and rotation symmetry operators of the space group of the corresponding crystal structure, and search-

ing for "clashes" between heavy atoms on symmetry related copies of the central model of the NMR ensemble. This analysis was performed using the structure validation program used by PDB's ADIT validation server. Two heavy atoms are considered to clash if they are within 2.2 Å of each other.

### Application to a large test set of crystal/NMR structure pairs

As of October 2005, the PDB contained a total of 427 "clusters" containing structures determined by multiple methods (crystallography, NMR, and cryoelectron microscopy). These clusters are defined as groups of PDB entries which have 95% or greater sequence similarity with the other members of the cluster and where more than one experimental method is represented. Of these, 288 clusters correspond to proteins whose structure has been determined using both crystallography and NMR. We have partially hand-curated these clusters to select pairs of PDB entries (one X-ray, one NMR) that could be used for our structural comparison. Specifically, we eliminated structures that had only $C_\alpha$ traces, no NMR ensembles, or where the models were the result of pushing the limits of technology (e.g. NMR models of very large proteins). Furthermore, we eliminated pairs with large differences in structure determination conditions (e.g. differences in ligands, extra/missing domains, or participation in protein complexes), as well as pairs in which fewer than half of the residues include two or more backbone heavy atoms classified as core by FindCore. This filtering resulted in a total of 136 protein pairs remaining. In addition, we also collected and analyzed a more limited set of 35 proteins based on that used previously by Garbuzynskiy et al.[27] (Tables I–III). Of these, 23 proteins overlap with the larger set, resulting in a total of 148 unique protein structure pairs which constitute our large data set.

Analysis of RMSD and statistically significant atomic position differences was performed on the full set of 148 protein structure pairs (referred to as the "complete test set"). In addition, analysis of contributions to structural differences, including the analysis of crystal contacts, was performed on the 35-protein set shown in Tables I–III (the "reduced test set"). FindCore was run twice for each protein structure pair, once using only backbone heavy atoms and once using all heavy atoms. The former were used to obtain the RMSD statistics in Table I and the residue-aggregated statistics in Table III, while the latter were used to obtain the atomic-level statistics in Table II. The $P$-values assessing the significance of the difference in position between the crystal and NMR structure was evaluated only for core atoms. The reported $P$-values have been adjusted for multiple hypothesis testing[36] to control the false discovery rate at 0.05 (i.e. on average 5% of the positives in a given protein pair will be false positives).

**Table I**
*RMSD, Quality, Age and Condition Similarity Statistics for the Reduced Test Set Based on Superpositions Using Backbone Heavy Atoms*

| PDB ID | | Core backbone (Å)[a] | | All backbone RMS$_{cryst}$ (Å)[a] | Quality[b] | | NMR year[c] | Similar conditions[d] |
|---|---|---|---|---|---|---|---|---|
| X-ray | NMR | RMS$_{ens}$ | RMS$_{cryst}$ | | X-ray | NMR | | |
| 1AIE | 1PET | 0.20 | 0.61 | 0.59 | −5.30 | −13.5 | 1994 | * |
| 1AIL | 1NS1 | 0.39 | 0.45 | 0.93 | −3.37 | −4.85 | 1997 | * |
| 1BSY | 1DV9 | 0.45 | 0.70 | 1.72 | −3.08 | −6.33 | 2000 | |
| 1BV1 | 1BTV | 0.47 | 0.72 | 1.50 | −0.83 | −17.56 | 1997 | * |
| 1CHN | 1DJM | 0.38 | 1.16 | 1.71 | 0.59 | −2.6 | 1999 | |
| 1DUK | 1MYF | 0.42 | 0.60 | 1.20 | −3.75 | 0.41 | 1994 | |
| 1EKG | 1LY7 | 0.27 | 1.11 | 1.03 | −0.80 | −5.85 | 2002 | * |
| 1EY4 | 1JOR | 0.38 | 0.47 | 3.74 | −2.11 | −5.68 | 2001 | |
| 1FIL | 1PFL | 0.25 | 1.08 | 1.64 | −1.12 | −19.78 | 1994 | |
| 1FXD | 1F2G | 0.30 | 0.72 | 0.89 | −1.28 | −12.96 | 1998 | * |
| 1GPR | 1AX3 | 0.39 | 0.89 | 1.30 | −1.54 | −4.91 | 1997 | * |
| 1GSV | 3PHY | 0.27 | 0.71 | 1.28 | −0.80 | −20.08 | 1998 | |
| 1HB6 | 2ABD | 0.32 | 0.97 | 1.70 | −1.77 | −12.04 | 1993 | |
| 1HOE | 2AIT | 0.33 | 0.56 | 1.05 | −1.61 | −18.51 | 1989 | * |
| 1I27 | 1NHA | 0.46 | 0.58 | 1.09 | −1.44 | −2.90 | 2002 | |
| 1JF4 | 1VRE | 0.30 | 1.25 | 1.28 | 0.79 | −15.42 | 1999 | |
| 1KF5 | 2AAS | 0.23 | 0.65 | 1.23 | −0.32 | −2.13 | 1992 | |
| 1LDS | 1JNJ | 0.31 | 0.94 | 3.71 | −1.61 | −5.91 | 2001 | * |
| 1R69 | 1R63 | 0.41 | 0.59 | 0.87 | −0.20 | −1.83 | 1996 | * |
| 1RBV | 1RCH | 0.45 | 0.93 | 1.53 | −1.01 | −3.31 | 1995 | |
| 1RDG | 1E8J | 0.33 | 0.67 | 1.02 | −0.64 | −5.92 | 2000 | * |
| 1RNB | 1BNR | 0.29 | 0.68 | 0.95 | −4.62 | −16.29 | 1995 | |
| 1SHG | 1AEY | 0.25 | 0.31 | 1.08 | −2.89 | −3.02 | 1997 | * |
| 1UBI | 1D3Z | 0.07 | 0.20 | 0.39 | −1.68 | 0.16 | 1999 | * |
| 1WHO | 1BMW | 0.45 | 0.85 | 1.89 | −0.59 | −8.46 | 1998 | * |
| 1ZON | 1DGQ | 0.15 | 0.56 | 1.04 | −1.24 | −2.07 | 1999 | |
| 2CDV | 1IT1 | 0.12 | 0.77 | 1.27 | −1.93 | −22.37 | 2001 | |
| 2CI2 | 3CI2 | 0.36 | 0.71 | 1.50 | −4.49 | −12.56 | 1991 | * |
| 2CPL | 1OCA | 0.30 | 0.46 | 0.91 | −0.83 | −3.25 | 1997 | * |
| 2ERL | 1ERC | 0.26 | 0.92 | 1.24 | −2.25 | −3.43 | 1994 | |
| 2FFM | 1PQX | 0.31 | 1.20 | 2.26 | −4.05 | −4.14 | 2003 | * |
| 2OVO | 1TUR | 0.27 | 0.57 | 1.11 | −3.37 | −4.55 | 1994 | |
| 3ICB | 1CDN | 0.27 | 0.63 | 1.37 | −4.61 | 0.47 | 1995 | |
| 3IL8 | 1IKM | 0.15 | 1.01 | 1.37 | 0.16 | −6.15 | 1995 | * |
| 4RNT | 1YGW | 0.31 | 1.09 | 1.68 | −1.89 | −44.95 | 1996 | |

[a]"Core backbone" RMSDs are over all core backbone atoms. "All backbone" RMSDs use the core atom superposition and all backbone atoms, but exclude residues at the termini which do not contain core atoms.
[b]Worst PSVS Z-score[35] over all validation methods (Procheck, Verify3D, ProsaII, and MolProbity).
[c]Year in which the NMR structure was deposited in the PDB.
[d]Asterisk indicates that there is no indication of a difference in solution/crystalization conditions, ligands, or amino acid sequence that might cause a structural difference.

# RESULTS

## Comparison using RMSDs

Traditionally, the global comparison of a crystal and an NMR structure has been done by calculating the RMSD between the crystal and the mean (or representative) NMR structure (which we will denote as RMS$_{cryst}$), and the significance of the differences has been assessed by comparing it to the RMSD characterizing the spread of the NMR ensemble (e.g. the average RMSD to the mean NMR structure RMS$_{ens}$).[21] Table I shows the values of RMS$_{cryst}$ and RMS$_{ens}$ calculated over all backbone core atoms using the FindCore procedure as described above for the reduced test set, while Figures 2 and 3 and Supplementary Material Table 1 show the results for the complete test set. In Figure 2, we show the distributions of RMS$_{cryst}$ and RMS$_{ens}$ for the complete set of 148 protein structure pairs. The RMS$_{ens}$ has a rather narrow distribution centered around 0.4 Å, while RMS$_{cryst}$ shows considerably more variability, with a broad distribution peaked near 1 Å but extending beyond 3 Å. Of the 148 structure pairs, 61 (41%) have RMS$_{cryst}$ > 1 Å, and 13 (9%) have RMS$_{cryst}$ > 2 Å. It should be emphasized that these are backbone RMSDs based on the core atom set. The inclusion of atoms outside of that set results in an increase in the RMSDs, as can be seen in the "all backbone RMS$_{cryst}$" column of Table I, the average of which is 1.4 Å compared to 1 Å for the core RMS$_{cryst}$. A typical example of the structural differences that we observe is shown in Figure 4, which shows the core-atom superpo-

**Table II**
*Results of FindCore Superpositions, P-value Based Structural Difference, and Crystal Clash Analysis Using all Heavy Atoms for the Reduced Test Set*

| PDB ID | | Heavy atoms | Core heavy atoms | | Significant P-value[a] | | Clashes[b] | |
|---|---|---|---|---|---|---|---|---|
| X-ray | NMR | | Number | % | Number | $f_\Delta$ (%) | Total | Core |
| 1AIE | 1PET | 295 | 179 | 61 | 70 | 39 | 29 | 5 |
| 1AIL | 1NS1 | 558 | 359 | 64 | 20 | 6 | 23 | 0 |
| 1BSY | 1DV9 | 1286 | 647 | 50 | 11 | 2 | 63 | 0 |
| 1BV1 | 1BTV | 1278 | 713 | 56 | 14 | 2 | 24 | 2 |
| 1CHN | 1DJM | 966 | 578 | 60 | 353 | 61 | 83 | 0 |
| 1DUK | 1MYF | 1212 | 697 | 58 | 26 | 4 | 4 | 0 |
| 1EKG | 1LY7 | 940 | 545 | 58 | 222 | 41 | 26 | 0 |
| 1EY4 | 1JOR | 1110 | 651 | 59 | 11 | 2 | 190 | 0 |
| 1FIL | 1PFL | 1046 | 664 | 64 | 517 | 78 | 13 | 2 |
| 1FXD | 1F2G | 430 | 254 | 59 | 48 | 19 | 9 | 1 |
| 1GPR | 1AX3 | 1191 | 642 | 54 | 108 | 17 | 115 | 1 |
| 1GSV | 3PHY | 951 | 539 | 57 | 324 | 60 | 13 | 0 |
| 1HB6 | 2ABD | 698 | 413 | 59 | 286 | 69 | 13 | 0 |
| 1HOE | 2AIT | 558 | 328 | 59 | 20 | 6 | 59 | 0 |
| 1I27 | 1NHA | 616 | 133 | 22 | 1 | 0.8 | 93 | 0 |
| 1JF4 | 1VRE | 1056 | 639 | 61 | 438 | 69 | 5 | 0 |
| 1KF5 | 2AAS | 1055 | 205 | 19 | 117 | 57 | 22 | 0 |
| 1LDS | 1JNJ | 809 | 447 | 55 | 259 | 58 | 112 | 0 |
| 1R69 | 1R63 | 484 | 273 | 56 | 0 | 0 | 28 | 1 |
| 1RBV | 1RCH | 1234 | 669 | 54 | 21 | 3 | 60 | 1 |
| 1RDG | 1E8J | 398 | 231 | 58 | 1 | 0.4 | 44 | 0 |
| 1RNB | 1BNR | 907 | 199 | 22 | 79 | 40 | 20 | 0 |
| 1SHG | 1AEY | 427 | 265 | 56 | 8 | 3 | 11 | 0 |
| 1UBI | 1D3Z | 602 | 407 | 68 | 298 | 73 | 6 | 0 |
| 1WHO | 1BMW | 745 | 400 | 54 | 109 | 27 | 54 | 1 |
| 1ZON | 1DGQ | 1779 | 846 | 48 | 657 | 78 | 20 | 3 |
| 2CDV | 1IT1 | 801 | 471 | 59 | 422 | 90 | 3 | 0 |
| 2CI2 | 3CI2 | 521 | 274 | 53 | 29 | 11 | 16 | 0 |
| 2CPL | 1OCA | 1258 | 668 | 53 | 3 | 0.5 | 10 | 0 |
| 2ERL | 1ERC | 598 | 186 | 31 | 135 | 73 | 41 | 2 |
| 2FFM | 1PQX | 626 | 359 | 57 | 193 | 54 | 84 | 0 |
| 2OVO | 1TUR | 418 | 258 | 62 | 66 | 26 | 55 | 2 |
| 3ICB | 1CDN | 600 | 334 | 56 | 126 | 38 | 14 | 0 |
| 3IL8 | 1IKM | 558 | 315 | 56 | 299 | 95 | 116 | 0 |
| 4RNT | 1YGW | 776 | 449 | 58 | 390 | 87 | 33 | 0 |

[a]Number (or fraction) of core heavy atoms that are significantly different at the 0.05 significance level.
[b]Number of atoms (or core atoms) that have a contact of $\leq 2.2$ Å with an atom from a symmetry-related molecule after insertion of the NMR structure into the crystal frame.

sition of the crystal and NMR structures of Photoactive Yellow Protein ($RMS_{cryst} = 0.71$). Most of the structural variability is located in the noncore loop regions and the N-terminus. However, significant differences are also evident at the edges of secondary structural units, such as the start of the C-terminal β-strand.

There is little overlap between the $RMS_{cryst}$ and $RMS_{ens}$ distributions. This is seen clearly in Figure 3, which shows that $RMS_{cryst}$ is larger than $RMS_{ens}$ for every structure pair, as no points lie below the solid line of slope one passing through the origin. Furthermore, there is a great deal of variability in the distance of the points from that line, ranging from proteins for which $RMS_{cryst} \approx RMS_{ens}$ to ones where $RMS_{cryst}$ is a factor of 2–10

larger than $RMS_{ens}$. Of the 148 structure pairs, 112 (76%) have RMSDs that differ by more than a factor of 2, 68 (46%) by a factor of 3, and 42 (28%) by a factor of 4. Unlike some much earlier studies,[24,37] we do not see a linear correlation between $RMS_{cryst}$ and $RMS_{ens}$.

The differences between $RMS_{cryst}$ relative to $RMS_{ens}$ plotted in Figures 2 and 3 can be caused by methodological or physical effects (or a combination of the two). Since it is not easy to classify the structure pairs by the specifics of the methodologies used to generate the models, we have chosen to use two surrogate measures: the year in which the NMR structure was determined, and the overall quality of the models as measured by packing, rotomeric states, and other biophysical characteristics. The methodology of crystal structure determination is considered to be mature and has been stable in the time frame of the structures that we are considering here. NMR spectroscopic structure determination, on the other hand, has undergone substantial improvements in methodology, ranging from additional data sources to advances in data analysis methodology (not to mention orders of magnitude increases in computational power). Therefore, older NMR structures may be less accurate than more recent ones. However, there is no discernible correlation between the size of $RMS_{cryst}$ relative to $RMS_{ens}$ and the age of the NMR structure [Table I and Fig. 5(A)] for the reduced test set. An alternative measure of methodological effects is the overall "quality" of the structure based on biophysical prior knowledge. It is possible that structure pairs for which $RMS_{cryst}$ is much larger than $RMS_{ens}$ may simply arise from one (or both) of the structures being of poor quality. We have made use of the recently developed Protein Structure Validation Suite (PSVS), an automated metaserver that calculates a variety of structure quality measures (including Procheck,[38] ProSaII,[39] Verify3D,[40] and MolProbity[41]) and summarizes the results in terms of z-scores relative to a set of high-resolution crystal structures.[35] These z-scores represent how far away (in units of standard deviation) a given quality measure is from the average of that quality measure for the high-resolution crystal structures. In Table I, we report the most unfavorable (most negative) of the z-scores from among all of the quality measures, and the structure pairs for which this z-score is less than −6 (corresponding to poor quality structures) are indicated in Figure 5(B). We do see a larger fraction of higher-quality structures in the region where $RMS_{cryst} \approx RMS_{ens}$, and a slight overabundance of low-quality structures far away from the $RMS_{cryst} = RMS_{ens}$ line, however the correlation is weak. Thus, while at least some of the variability in "accuracy" has a methodological origin, its effect is relatively weak based on the measures which we examined (Fig. 5). Another possible source of methodological difference could result from the use of residual dipolar couplings (RDCs) in the NMR refinement.[42–44] In the data set studied here, only 8 of the

**Table III**
*Results of FindCore Superpositions, P-value Based Structural Difference, and Crystal Clash Analysis Using Backbone Atoms and Aggregated at the Residue Level for the Reduced Test Set*

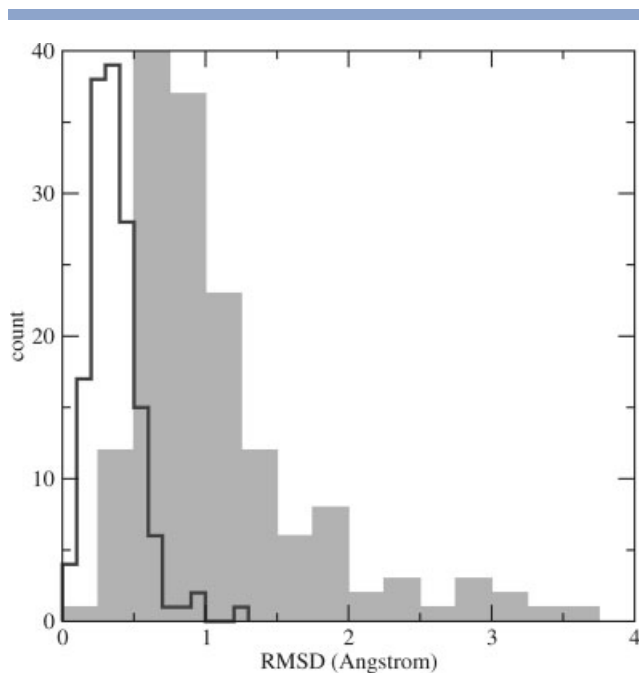| PDB ID | | | Core residues | | Significant P-value[a] | | Clashes[b] | | |
|---|---|---|---|---|---|---|---|---|---|
| X-ray | NMR | Residues | Number | % | Number | % | Total | Core | Core+ sig.diff. |
| 1AIE | 1PET | 31 | 17 | 55 | 8 | 47 | 8 | 1 | 1 |
| 1AIL | 1NS1 | 73 | 40 | 55 | 1 | 2 | 8 | 2 | 1 |
| 1BSY | 1DV9 | 162 | 97 | 60 | 0 | 0 | 7 | 0 | 0 |
| 1BV1 | 1BTV | 159 | 86 | 54 | 0 | 0 | 9 | 2 | 0 |
| 1CHN | 1DJM | 126 | 70 | 56 | 47 | 67 | 11 | 1 | 1 |
| 1DUK | 1MYF | 153 | 92 | 60 | 0 | 0 | 4 | 0 | 0 |
| 1EKG | 1LY7 | 121 | 82 | 68 | 33 | 40 | 8 | 3 | 1 |
| 1EY4 | 1JOR | 144 | 89 | 62 | 1 | 1 | 21 | 3 | 0 |
| 1FIL | 1PFL | 139 | 91 | 66 | 67 | 74 | 7 | 3 | 2 |
| 1FXD | 1F2G | 58 | 40 | 69 | 5 | 13 | 4 | 3 | 0 |
| 1GPR | 1AX3 | 159 | 91 | 57 | 12 | 13 | 12 | 2 | 1 |
| 1GSV | 3PHY | 122 | 69 | 57 | 41 | 59 | 16 | 3 | 1 |
| 1HB6 | 2ABD | 86 | 48 | 56 | 34 | 71 | 3 | 0 | 0 |
| 1HOE | 2AIT | 74 | 45 | 61 | 2 | 4 | 8 | 1 | 0 |
| 1I27 | 1NHA | 73 | 42 | 56 | 0 | 0 | 12 | 1 | 0 |
| 1JF4 | 1VRE | 147 | 109 | 74 | 63 | 58 | 2 | 2 | 2 |
| 1KF5 | 2AAS | 124 | 67 | 54 | 36 | 54 | 2 | 0 | 0 |
| 1LDS | 1JNJ | 100 | 57 | 57 | 32 | 56 | 11 | 0 | 0 |
| 1R69 | 1R63 | 63 | 37 | 59 | 0 | 0 | 4 | 1 | 0 |
| 1RBV | 1RCH | 155 | 89 | 57 | 0 | 0 | 9 | 4 | 0 |
| 1RDG | 1E8J | 53 | 29 | 55 | 3 | 10 | 9 | 0 | 0 |
| 1RNB | 1BNR | 111 | 66 | 60 | 19 | 29 | 1 | 0 | 0 |
| 1SHG | 1AEY | 57 | 36 | 63 | 1 | 3 | 3 | 0 | 0 |
| 1UBI | 1D3Z | 76 | 45 | 59 | 28 | 62 | 5 | 0 | 0 |
| 1WHO | 1BMW | 94 | 56 | 60 | 13 | 23 | 8 | 3 | 1 |
| 1ZON | 1DGQ | 184 | 112 | 61 | 90 | 80 | 7 | 3 | 3 |
| 2CDV | 1IT1 | 107 | 65 | 61 | 60 | 92 | 1 | 0 | 0 |
| 2CI2 | 3CI2 | 65 | 35 | 54 | 0 | 0 | 6 | 0 | 0 |
| 2CPL | 1OCA | 164 | 94 | 57 | 0 | 0 | 2 | 0 | 0 |
| 2ERL | 1ERC | 40 | 25 | 63 | 19 | 76 | 8 | 3 | 2 |
| 2FFM | 1PQX | 90 | 49 | 54 | 26 | 53 | 12 | 1 | 1 |
| 2OVO | 1TUR | 56 | 31 | 55 | 6 | 19 | 8 | 0 | 0 |
| 3ICB | 1CDN | 75 | 42 | 56 | 5 | 12 | 5 | 0 | 0 |
| 3IL8 | 1IKM | 68 | 40 | 59 | 38 | 95 | 14 | 1 | 1 |
| 4RNT | 1YGW | 104 | 61 | 59 | 50 | 82 | 5 | 1 | 1 |
| totals: | | 3613 | 2144 | 59 | 740 | 35 | 260 | 44 | 19 |

[a]Number (or fraction) of core residues that have at least two backbone heavy atoms that are significantly different at the 0.05 significance level.
[b]Number of residues (or core residues) that have a contact of ≤2.2 Å with an atom from a symmetry-related molecule after insertion of the NMR structure into the crystal frame. "Core+sig.diff." represents the number of core residues that also have significantly different backbone conformations.

NMR structures included RDC refinement (1D3Z, 1E8L, 1L3G, 1R36, 1U81, 1YCM, 1YJI, and 2EZN), with one additional structure (1E8E) that made use of paramagnetic pseudocontact shifts which also provide orientational information relative to a global reference frame.[43] Although this is too small of a sample to study the effect of RDCs in a comprehensive way, it is interesting to note that these 9 examples have relatively large $RMS_{cryst}$:$RMS_{ens}$ ratios, averaging 4.3 and ranging up to 6.9, implying substantial structural differences relative to the crystal structure. We will discuss the possible impact of RDC data on precision and accuracy in more detail below.

One possible source of variability that is of a physical origin is obvious: one structure in each pair was determined in solution and the other in a crystalline environment, and the impact of this difference may vary from protein to protein. This will be considered in detail below. However, physical differences can also arise from differences in ligands, amino acid sequence, oxidation state, and other changes in the environment between the NMR and crystal structure determination. Unfortunately, this information is not always apparent from the information given in the PDB headers. We have tried to identify protein pairs for which there is no information that

**Figure 2**

*Distributions of the RMSDs of the crystal structure to the mean NMR structure (filled gray) and the average RMSD of the structures in the NMR ensemble to the mean NMR structure (dark line) for the complete 148-protein test set (data in Supplementary Material Table 1).*
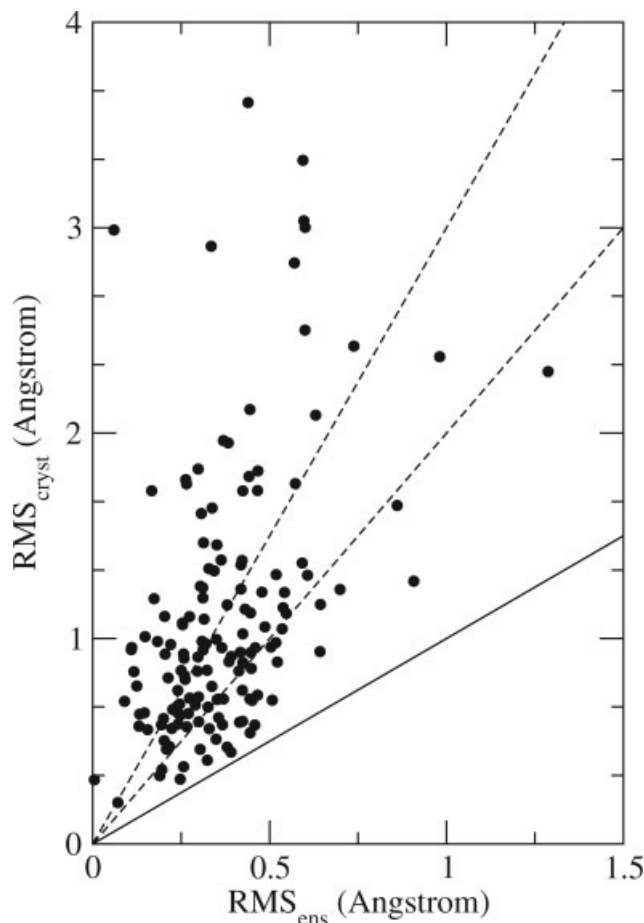
would indicate a physical difference in the protein targets other than the presence of a crystal lattice, and these have been indicated in Table I by an asterisk in the rightmost column. If we examine where these structure pairs lie in the ($RMS_{cryst}$, $RMS_{ens}$) plane, again there appears to be some correlation: structure pairs for which $RMS_{cryst} \approx RMS_{ens}$ tend to have been determined under similar conditions [Fig. 5(C)]. However, the correlation is not strong. There are many structure pairs where the NMR structure was determined recently, do not have poor quality scores, and were determined under similar conditions, but still have large $RMS_{cryst}$:$RMS_{ens}$ ratios (such as 2FFM/1PQX, 1EKG/1LY7, and 1LDS/1JNJ).

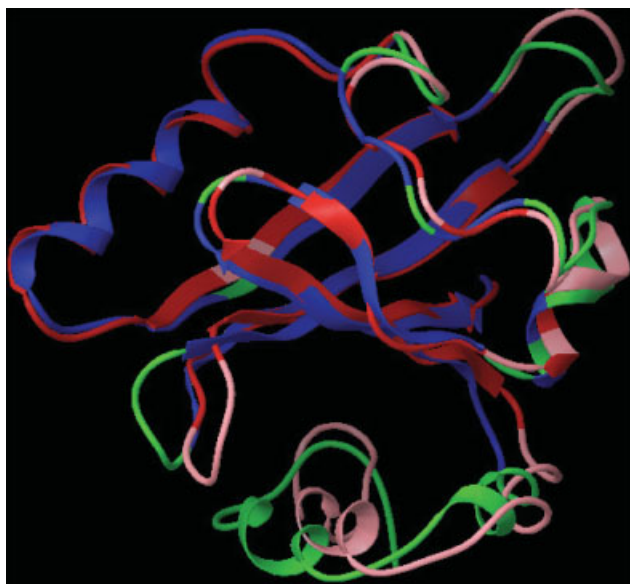## Comparison using a P-values

### Local measure of structural difference

The *P*-value based measure of structural similarity based on the $\chi^2_\psi$ statistic of Eq. (1) differs from a simple comparison of $RMS_{cryst}$ and $RMS_{ens}$ in that it is based on a statistical measure that leads to well-defined false positive and false negative rates relative to a null hypothesis that the crystal structure coordinates do not differ from those in the NMR ensemble. As we have constructed it, this measure is local, and therefore provides

information on an atom-by-atom basis about where the structures differ from each other in a statistically significant way. In Figure 6, we show the multiple testing corrected *P*-values for heavy atoms in the crystal and NMR structures of influenza A virus non-structural protein 1 N-terminal domain (NS1) (1AIL/1NS1),[45,46] ubiquitin (1UBI/1D3Z),[47,48] and the hypothetical protein SAV1430 from *S. aureus* (2FFM/1PQX). The *P*-value is a measure of how unlikely the observed value of $\chi^2_\psi$ would be if the position of a given atom in the crystal structure comes from the same distribution as the atomic positions in the NMR ensemble. Small values of *P* indicate that this hypothesis is likely false. In the case of NS1 [Fig. 6(A)], it is clear that almost all of the core heavy atoms have *P*-values greater than 0.1, and the
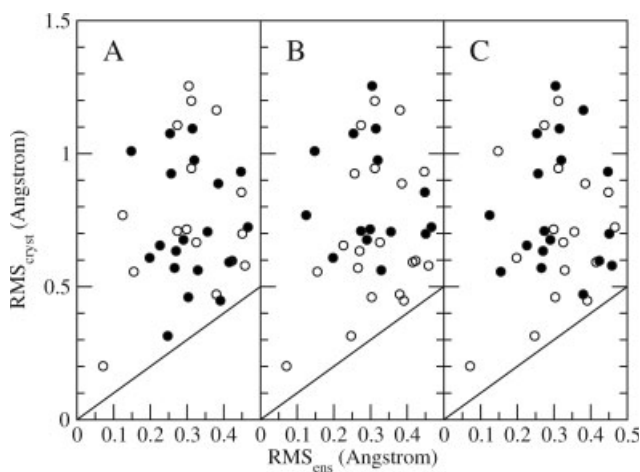


**Figure 3**

*Correlation between the RMSD of the crystal structure to the mean NMR structure and the average RMSD of the structures in the NMR ensemble to the mean NMR structure for the complete 148-protein test set (data in Supplementary Material Table 1). The solid line indicates the points for which $RMS_{cryst}$ equals $RMS_{ens}$, while the two dashed lines correspond to $RMS_{cryst}$:$RMS_{ens}$ ratios of 2 and 3, respectively.*

atoms for which there are significant *P*-values are highly localized in the vicinity of residues 8, 15, 29, and 51. This indicates that except for those regions, the crystal and NMR structures are identical to within the precision of the NMR structure.

By contrast, both the SAV1430 [Fig. 6(B)] and ubiquitin [Fig. 6(C)] structure pairs show large numbers of *P*-values less than 0.01 throughout the protein. The origin of these significant differences, however, are quite different in the two proteins. In SAV1430, RMS$_{cryst}$ is large (1.2 Å), therefore we might expect a large number of differences in atomic positions between the crystal and NMR structure that are significant. Ubiquitin, by contrast, has an RMS$_{cryst}$ of 0.2 Å, indicating that the structural differences between the crystal and mean NMR structures are quite small. Nonetheless, these very small differences appear to be statistically significant because the "width" of the NMR ensemble is even smaller (as suggested by the RMS$_{ens}$ of 0.07 Å). It should be noted that the 1D3Z ubiquitin structure makes extensive use of residual dipolar coupling data, which may contribute to its high apparent precision.



**Figure 4**

*Core atom superposition of crystal structure (1GSV) and central NMR structure (3PHY model 20) of Photoactive Yellow Protein. The crystal structural model is shown in blue (core residues) and green (non-core residues), while the NMR structural model is in red (core) and pink (non-core). A "core residue" is defined to be one in which at least 2 backbone atoms belong to the core. The overall backbone RMSD after core atom superposition omitting the disordered N-terminus (at the bottom of the figure) is 1.28 Å.*
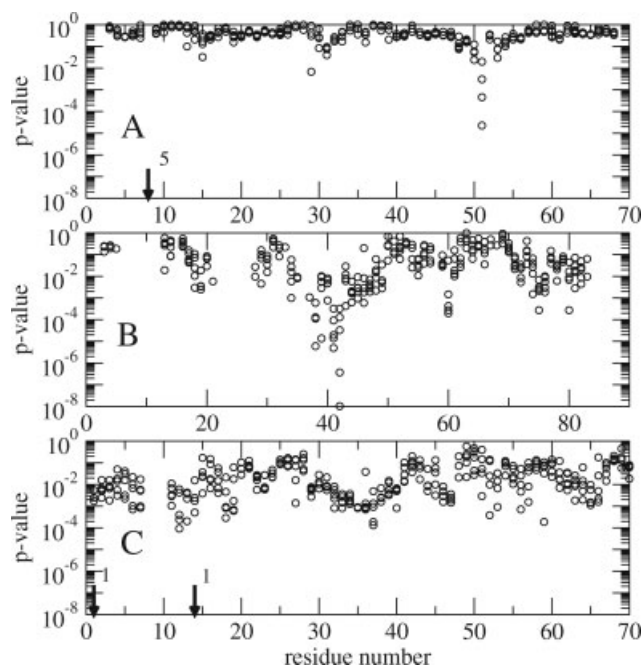


**Figure 5**

*Correlation between RMS$_{cryst}$ and RMS$_{ens}$ for the reduced 35-protein test set which have been color-coded by properties related to possible sources of methodological or physical difference (Table I). (A) Filled circles represent structure pairs in which the NMR structure was determined prior to 1998. (B) Red dots represent structure pairs in which either the crystal or NMR structure had a minimal PSVS z-score of less than −6. (C) Filled circles indicate structure pairs that upon manual curation were deemed to have significant differences in their experimental conditions. The solid line correspond to RMS$_{cryst}$ = RMS$_{ens}$.*



**Figure 6**

*Profile of structural difference p-values along the protein chain for three selected protein structure pairs 1AIL/1NS1 (A), 2FFM/1PQX (B), and 1UBI/1D3Z (C). Small P-values (e.g. <0.05) indicate that it is unlikely that the observed atomic position came from the same distribution as the atomic positions in the NMR ensemble. Values for all core atoms for which the positions are determined by backbone degrees of freedom (N, C, O, C$^\alpha$, and C$^\beta$) are shown. The arrows and numbers located at residue 8 in panel A and residues 1 and 14 in panel C indicate the number of atoms for which the P-values are less than $10^{-8}$.*

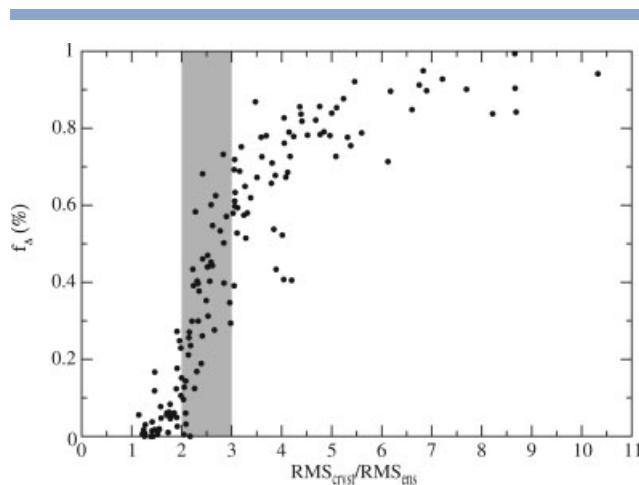### Global measure of structural difference using aggregate statistics

The differences seen in the three panels of Figure 6 suggest that we can use the fraction of core heavy atoms showing significant deviation between the crystal and NMR structures as an aggregated measure of the overall similarity of two structures. If we examine the fraction of core heavy atoms showing significant deviation $f_\Delta$ at a significance level of 0.05 (Table II and Supplementary Material Table 1), we see that there is considerable variability, with some structure pairs (such as NS1) having very few deviant core heavy atoms, while others (such as 2CDV/1T1 and 3IL8/1IKM) having essentially all of their core heavy atoms differing between the crystal and NMR structures. This is despite the fact that almost all of the structure pairs considered have a similar fraction of their heavy atoms in the core (50–65%).

There are a large number of structure pairs for which the differences between the crystal and NMR structural models appears to be significant: 73 of the complete set of 148-protein test structure pairs (49%) have more than 50% of their core heavy atoms in significantly different positions in the NMR and crystal structures at a confidence level of 0.05. This suggests that these are global differences in the results of the two structure determination methods for approximately half of the proteins in the large data set whose structures were determined by both NMR and X-ray crystallography. Furthermore, of the nine NMR structural models that make use of RDCs or pseudocontact shifts, all but one have $f_\Delta$ values greater than 50%. In addition to the global aggregation used to calculate $f_\Delta$, one could also aggregate on intermediate levels, such as residues or secondary structure units. We make use of aggregation at the residue level in the study of crystal contacts below. Aggregation at the secondary structure level could detect correlated motions of such units as rigid bodies.

### Comparison with RMSD-based structural difference measures

It is interesting to ask if $f_\Delta$ contains similar information to that contained in the more traditional comparison of $RMS_{cryst}$ and $RMS_{ens}$. In Figure 7, we plot $f_\Delta$ versus the ratio of $RMS_{cryst}$ to $RMS_{ens}$: it can be seen that structure pairs with very small $RMS_{cryst}:RMS_{ens}$ ratios (smaller than $\approx 2$) consistently have very small values of $f_\Delta$, while those with large $RMS_{cryst}:RMS_{ens}$ ratios (greater than $\approx 3$) have large values of $f_\Delta$. For RMSD ratios in the 2–3 range, however, we see that there is a "twilight zone": protein pairs with very similar RMSD ratios can have extremely different $f_\Delta$ values.

We can further explore the relationship between the RMSD ratio and atomic-level $\chi_\psi^2$-based $P$-values by interpreting the RMSD ratio in a statistical manner. $RMS_{ens}$ is defined as an average over an ensemble of



**Figure 7**

*Correlation between RMSD ratio and the fraction of core heavy atoms whose crystal vs NMR coordinate differences are significant ($f_\Delta$) for the complete 148-protein test set. The outlier pairs 1IW6/1R2N and 1BR9/2TMP have been deleted for clarity. The shaded region represents the "twilight zone" for the RMSD ratio measure as discussed in the text.*

RMSDs to the mean NMR structure. If we assume that the distribution of those individual RMSDs that contribute to $RMS_{ens}$ is a half-normal

$$P(\mathrm{RMSD}) = \begin{cases} \sqrt{\dfrac{2}{\pi\sigma^2}}\exp\left(\dfrac{-\mathrm{RMSD}^2}{2\sigma^2}\right) & \mathrm{RMSD} \geq 0 \\ 0 & \text{otherwise} \end{cases},$$

then we can calculate a $P$-value corresponding to the probability that a value $\mathrm{RMSD} \geq RMS_{cryst}$ would be observed. We choose $\sigma$ to be such that the mean of the half-normal is equal to $RMS_{ens}$, or $\sigma = \sqrt{\pi/2}\,RMS_{ens}$. To be significant at a confidence level of 0.05, $RMS_{cryst}$ must exceed $RMS_{ens}$ by a factor of 2.44. Of the 148 proteins studied here, 87 (59%) have RMSD ratios that are significant at a 0.05 level according to this measure. This is larger fraction than the 49% found using the globally aggregated atomic level $\chi_\psi^2$-based $P$-value analysis above, and is consistent with the log-normal distribution being more tolerant of outliers and the presence of a "twilight zone" in the relationship between $f_\Delta$ and the RMSD ratio.

These results together with the $f_\Delta$ analysis provide a statistical basis for the interpretation of RMSD ratios: if the ratio is less than 2 or greater than 3, then we can say with confidence that the structures are globally similar or different, respectively, and that these differences are statistically significant. On the other hand, if the RMSD ratio is in the 2–3 range, then the globally aggregated $P$-value statistic $f_\Delta$ provides additional information about the significance of the differences between the crystal and

NMR models that complements that provided by the RMSD ratio.

### Effect of crystal structure uncertainty

The analysis performed above depends on the assumption that all of the structural uncertainty is associated with the NMR model and none of it with the crystallographic model. This is clearly not the case, as crystal structures have atomic $B$-factors that can be understood as a measure of precision. We have found for our test set that the mean-square displacement derived from a simple Debye-Waller analysis of the $B$-factors for core atoms as well as the corresponding variability in position of core atoms in the NMR ensemble are comparable in magnitude and are on the order of 0.1–0.3 $\text{Å}^2$. It might seem that since the atom-level precisions are similar, the uncertainty in the crystal structure cannot be ignored. However, what is important is whether inclusion of the crystallographic uncertainty will convert a statistically significant position difference into an insignificant one. This will happen when the distance between the mean NMR coordinates and the crystal coordinates for a given atom is on the order of tenths of an Ångström, since if it is shorter than that, the difference will already be insignificant, while for longer distances the inclusion of the crystallographic uncertainty will cause a small $P$-value to increase but still remain small.

A rough calculation using NMR variability and $B/8\pi^2$ both 0.2 $\text{Å}^2$ and a typical log-normal $\chi^2_{\psi,m}$ distribution indicates that inclusion of the crystallographic uncertainty causes an uncorrected $P$-value of 0.02 to increase to 0.05. Therefore, we expect that inclusion of $B$-factor uncertainties will only have a consequential effect on $P$-values in the 0.01–0.05 range. On average, proteins in our test set have less than 20% of their heavy atoms with $P$-values of this magnitude. While this is a non-negligible fraction, it does not have a large impact on our conclusions. Inclusion of $B$-factor based uncertainties will cause the points in Figure 7 to shift downwards on average by 0.2, but is unlikely to change the shape or correlation with the RMSD ratio. Similarly, our statement that 73 of the 148 proteins have more than half of their core atoms in significantly different positions remains qualitatively correct. For example, only 14 of the 73 proteins have $f_\Delta$ in the range of 0.5–0.6. If crystallographic uncertainty is taken into account, then the fraction of proteins having more than half of their core atoms in significantly different positions decreases from 49% to approximately 40%.

### Analysis of crystal-induced clashes

We have deferred the discussion of crystal contacts as sources of conformational difference until now. As described in the Methods section earlier, it is possible to use the FindCore superposition methodology to simulate the presence of a crystalline environment by translating and rotating the NMR ensemble such that its mean coincides as best as possible with the crystal structure and then constructing the corresponding crystal environment by application of the appropriate symmetry operations to the central model. The close contact statistics obtained using this analysis are shown in Table II. It should be noted that the "clashes" discussed below are not actual physical clashes that are seen in any structural model in the PDB. Instead, they are "virtual collisions" that would arise if a crystal lattice was constructed using representative NMR structures.

For the reduced test set, the number of close contacts ranges widely from 3 to 190, or 1 to 21% of the heavy atoms in the protein. While these represent a small fraction of the total atoms in each protein, in absolute terms they are substantial. However, there is no discernible correlation between the fraction of close contacts and the fraction of core heavy atoms occupying significantly different positions $f_\Delta$. Furthermore, the number of clashes that occur between pairs of core atoms is remarkably small both in fractional and absolute terms, with 24 of the 35 proteins having no core-core clashes and another 9 having only one or two core clashes, as seen in Table II. From this point of view, it would seem that crystal contacts play only a small role in the differences between crystal and NMR structural models.

The above analysis was performed on an atomic level. It is also possible to examine the clash data from a residue-based viewpoint. We consider an entire residue to be "core" if the backbone-based FindCore analysis determined that at least two of its backbone heavy atoms belong to the core. Similarly, a residue is considered to be in a significantly different position in the crystal and NMR structures if at least two of its backbone heavy atoms have significant $P$-values. A residue is counted as being involved in a crystal clash if any of its atoms (core and non-core, backbone and sidechain) are involved in a close contact.

The results of this residue-aggregated analysis shown in Table III are similar to what was found in the atomic-level case. In particular, only 2% (44/2144) of the core residues (those with well-ordered backbone atoms) are involved in a clash, and only 17% (44/260) of the residues that are involved in clashes also have well-ordered backbones. As before, this suggests that incompatibility with the crystal environment plays only a minor role in giving rise to the observed differences between core residues in the crystal and NMR structures. This result is supported by a comparison of the rate at which statistically significant position differences occur depending on whether the residues is or is not involved in a clash: of the 44 core residues involved in clashes, 19 (43%) have significantly changed backbone atomic positions. On the other hand, of the $2144 - 44 = 2100$ core residues *not* involved in a clash, $740 - 19$

= 721 (34%) have significantly changed backbone atomic positions. The increase in structural difference rate from 34 to 43% is not large or even statistically significant, and this result again points out the relatively small role played by steric clashes.

## DISCUSSION

On the basis of biophysical principles, it is expected that a crystal and NMR structure of the same protein should agree with each other to within appropriate "error bounds" that reflect both the dynamic nature of the protein and the lack of knowledge because of imprecise and sparse experimental data. If they do not agree, then the differences must be attributable to either physical or methodological factors, or a combination of the two. Using both the $\chi_\psi^2$ and RMSD-based methodologies described earlier, we have found that there are widespread statistically significant structural differences within pairs of crystal and NMR structures of the same protein. We would like to understand the origins of these differences.

The observation that the apparent precision of NMR structures is greater than the observed structural differences (i.e. $RMS_{ens} < RMS_{cryst}$) is a phenomenon that has been seen previously in the results of experimental structure determinations.[18,24] It is possible that this indicates real differences between crystal and NMR structures arising from differences between the crystal and solution environments[21] or, alternatively, these structural differences may reflect systematic bias introduced by differences in the structural constraints and modeling methods used by the two approaches.[27,29] However, this phenomenon was also observed in early studies involving error analyses performed using synthetic data where the true accuracy can be determined directly (since the structure which generated the data is known).[37,49,50] The apparent precision was greater than the difference between the calculated and true structure even in studies where the synthetic data were subjected to relatively mild "noise",[37] and only became worse when more realistic models of the noise were used.[50]

It should be noted that several of these studies are at least a decade old, and it is possible that changes in NMR refinement protocols have ameliorated the problem. In fact, more recent studies with synthetic data give a less consistent picture: Chalaoux et al. report a "bias" similar to that observed in the older studies,[51] while Bassolino-Klimas, et al. do not.[52] This may be due to the differences in the methods used to add "noise" to the NOE data or differences in the form of the NOE restraint pseudopotential used.

It has been shown that the crystal structures of the same protein in different crystal forms or as multiple copies in an asymmetric unit can have significant differences (e.g. see Refs. 53 and 54), indicating that crystal contacts do play a role in determining the observed conformations. Although we see weak evidence for structural differences between crystal and NMR structures because of the NMR structure's incompatibility with the crystalline environment, most of the significant differences cannot be explained by clashes in the immediate vicinity of the difference. It is possible that clashes could give rise to nonlocal changes in structure, although we see no direct evidence for this.

Forces other than repulsive interactions, such as hydrogen bonds, ion pairs, and hydrophobic surface contacts all play a role in stabilizing conformations. These attractive interactions could, for example, selectively stabilize conformations in the crystalline state that are rare in the solution ensemble, leading to observable differences between crystal and NMR structures. This possibility is supported by recent work using new experimental tools for the study of protein dynamics by NMR.[55] In particular, elegant studies from the Kern laboratory have reported large global conformational fluctuations on the microsecond timescale in enzymes, and they observe that this motion is coupled to catalysis.[56,57] Furthermore, they have shown that this dynamics is also present in the free enzyme (without bound substrate).[58]

It is not unreasonable to suppose that the crystalline environment could have the effect of "freezing out" one global conformation from the more diverse ensemble present in solution, leading to considerably different average structures in the solution state and the crystal. It is also possible that the structure which is favored by the crystalline environment also exists in solution, but as a minor conformer that could be missed in the NMR structure determination process. This may be exacerbated by the use of cryo-cooling techniques in crystallography that have become common in recent years.[59,60] The detection of such effects will require a more detailed and extensive analysis than we can provide here, and will be left for future research.

Alternatively, it may be that the apparent differences are in fact only apparent: that is, they appear to be significant because of methodological issues. For example, small differences in atomic position (on an absolute scale) can appear to be significant because the precision estimate is even smaller than the position difference [e.g. Fig. 6(C)]. One might argue that since an "expert" would likely not consider the differences between the two structures to be "real", the statistical measures of similarity are generating spurious results. However, measures of statistical significance are only as good as the estimates of precision (e.g. the diversity of the NMR ensemble) on which they are based. The methodology by which NMR ensembles are typically generated is not geared toward reliably sampling the full conformational distribution, but is instead based on repeated minimizations and selection of the

"best" structures that satisfy the constraints as much as possible. This strategy may lead to a systematic overestimation of the precision of NMR structures,[11] thereby rendering small apparent deviations significant. While the case illustrated in Figure 6(C) may not specifically be an example of such precision overestimation, it has the characteristics that would be symptomatic of such a situation.

Previous studies[11,12] have suggested that it is possible to find more diverse sets of structures that fit a given set of NMR data than those generated by currently used protocols. In particular, a large-scale re-refinement of NMR structures based on publicly-available data using a modern protocol found that the RMSD values of the ensembles increased significantly over those reported by the original authors. On average, a 0.4 Å increase was seen, while for 44 proteins (out of a total of 545) the ensemble RMSD increased by more than 1 Å.[12] Furthermore, recent work on the interpretation of crystallographic data has suggested that a single structure may not be the best representation of the solution to the crystallographic structure determination problem,[61,62] and that this may cause overestimation of precision in the crystal structure as well. Underestimation of precision can also be due to "conformational pinning" in which the lack of consideration of ensemble-averaging leads to the "trapping" of the conformation between conflicting minima in the constraint pseudoenergies.[10] It is possible that Bayesian sampling-based approaches to NMR structure determination[13,14] or other alternative sampling procedures[11,34] including explicit ensemble-based refinement[63–67] will lead to more realistic estimates of the precision and accuracy, and therefore more reliable identification of conformational differences between crystal and NMR structures.

Some of the differences observed here are too large to be explained by erroneous precision estimates (such as those in the long tail of the $RMS_{cryst}$ distribution in Fig. 2). These could be due to inaccuracies in the NMR or crystal structure. It has been observed, for example, that re-refinement of NMR structures results in deviation of the mean structures from those originally reported, and that the resulting average structures have improved structural quality[12] and are more consistent with corresponding crystal structures.[27] However, the changes observed on re-refinement are relatively small ($\approx 0.5$ Å)[12] and are unlikely to render many of our observed differences insignificant, particularly when coupled with the increase in coordinate uncertainty that they observe. The large differences observed here may be examples of true structural differences that have a physical origin, or may be due to a combination of factors.

NOE distance constraints remain the dominant form of structural data used in the determination of NMR structural models. Within the past decade, however, the use of RDCs has become much more common as an additional source of data on macromolecular structure and dynamics.[42–44,68,69] RDCs arise from the incomplete averaging of the dipolar Hamiltonian due to the partial ordering of the molecule relative to the static magnetic field and provide information about the orientation of internuclear vectors relative to an alignment tensor describing the ordering. This dependence on the angle relative to a global reference frame makes RDCs highly complementary to the local distance information contained in NOEs. However, RDCs differ from NOEs in that a single RDC value is consistent with many bond vector orientations that can be visualized as a continuous curve on the surface of a unit sphere. On the other hand, even very small deviations off of that curve can result in significant disagreements with a measured RDC value.

This highly-constraining nature of RDCs will have an impact on both precision and accuracy: they could improve the accuracy of the resulting structural models, but could also exacerbate the overestimation of precision, in that they are particularly susceptible to conformational pinning, especially if steep quadratic constraint pseudo-potentials are used. It has been observed that the refinement of protein structures using RDCs can significantly increase the number of residues falling in the most favored regions of the Ramachandran plot[42]; however, it is not always clear if such refinement actually improves the accuracy of a structure.[44] In addition, such refined structures have higher apparent precision (based on the RMSD within the ensemble) by as much as a factor of 2–3.[42] It is not yet clear whether such increases in the ensemble tightness represent a true improvement in the precision or are the result of effects such as conformational pinning.

RDCs hold the promise of clarifying some of the possible sources of the structural differences that we observe here by improving the accuracy of the structural model. For example, suppose a protein undergoes a large-scale motion involving the closing of a flap (as in Ref. 56, for example), and that the open form has the flap too far away from the body of the protein to give rise to observable NOEs, while the closed form does lead to NOEs. Even if the closed form is a minor conformer, NOEs would still be observed from the flap to the body of the protein, since their $r^{-6}$-dependence causes them to be dominated by the distance of closest approach. The observation of such NOEs could lead to the incorrect conclusion that the closed form predominates in solution. On the other hand, the global orientation dependence of RDCs could be used to produce a more accurate picture of the structure and dynamics.

The structural interpretation of RDCs in the presence of such conformational dynamics is not straightforward.[68,69] Because RDCs are the result of direct spin-spin interactions and not relaxation, the effect of motion on RDCs is not "washed out" by the overall tumbling of the molecule in solution. Therefore, RDCs are subject to

conformational averaging from motions ranging from librational motions on the sub-picosecond timescale to the millisecond timescale of the NMR experiment itself. The dynamical interpretation of RDCs has been the subject of considerable controversy,[70–73] in that the same data can produce very different pictures of the magnitude of the conformational fluctuations depending on the type of analysis that is performed.[68] Many analyses of dynamics are based on discrepancies of the data relative to predictions based on a model that was itself derived on the basis of the assumption of a static structure,[69] and it is possible that the extensive refinement of (implicitly static) structures with respect to RDCs could partially absorb the effects of internal motion.[68]

The varying interpretations of RDC data—some suggesting more dynamics,[70,72] others less[73]—lead to different pictures of the native state of proteins, and have consequences for the interpretation of structural differences studied here. The resolution of these discrepancies will require further developments in refinement methodology and data interpretation, and will likely involve the generation of ensemble representations of protein conformation that simultaneously satisfy all available structural and dynamical data (including NOEs, scalar couplings, RDCs, and relaxation measurements), thereby avoiding the current artificial dichotomy between structure and dynamics.[65]

## CONCLUSIONS

We have performed a statistical comparison of the structural differences between crystal and NMR structural models of the same protein on a large data set of 148 structure pairs for which structural models derived from both crystallographic and NMR data have been deposited in the PDB. To the best of our knowledge, this is the largest set of proteins for which such an analysis has been performed. The detection of significant structural differences was performed in an automated fashion based on the FindCore method for structural superposition, and made use of *P*-value estimates of statistical significance based on the distribution of distances of atomic coordinates from their mean positions to obtain an estimate of the probability of observing the crystal structure given the NMR ensemble.

Using this statistical methodology, we have found pervasive statistically significant structural differences within pairs of crystal and NMR structures of the same protein. For all of the 148 structure pairs examined, the difference in RMSD between the crystal and the mean NMR structure was found to exceed the RMSD within the NMR ensemble, and in 73 of the 148 structure pairs more than half of the core heavy atoms are in significantly different positions.

While we would like to be able to attribute the large number of apparent structural differences observed to a particular physical or methodological source, this is not possible. We can, however, rule out the role of repulsive crystal packing effects as a dominant source of structural difference. Further clarification of the origins of the structural differences between crystal and NMR structures must await further methodology development in structural biology, including the development and study of advanced sampling methods both in NMR and crystallography,[11,13,14,34,61,74] coupled with ensemble methods for refinement.[63–66] The precision and accuracy of crystallographic and NMR structural models are more than academic issues, with ramifications for ligand docking and virtual screening,[62] and it is hoped that with careful experimental work and statistical analysis these questions can be resolved.

## REFERENCES

1. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE. The protein data bank. Nucleic Acids Res 2000;28:235–242.
2. Snyder DA, Montelione G. Clustering algorithms for identifying core atom sets and for assessing the precision of protein structure ensembles. Proteins 2005;59:673–686.
3. Altman RB, Hughes C, Gerstein MB. Methods for displaying macromolecular structural uncertainty: applications to the globins. J Mol Graphics 1995;13:142–152.
4. Gerstein M, Altman RB. Average core structures and variability measures for protein families: application to the immunoglobulins. J Mol Biol 1995;251:161–175.
5. Wilmanns M, Nilges M. Molecular replacement with NMR models using distance-derived pseudo *B* factors. Acta Cryst D 1996;52:973–982.
6. Willis BTM, Pryor AW. Thermal vibrations in crystallography. Cambridge: Cambridge University Press; 1975.
7. Clore GM, Schwieters CD. Concordance of residual dipolar couplings, backbone order parameters and crystallographic *B*-factors for a small α/β protein: a unified picture of high probability, fast atomic motions in proteins. J Mol Biol 2006;355:879–886.
8. Snyder DA, Bhattacharya A, Huang YJ, Montelione G. Assessing precision and accuracy of protein structures derived from NMR data. Proteins 2005;59:655–661.
9. Nilges M, Clore GM, Gronenborn AM. A simple method for delineating well-defined and variable regions in protein structures determined from interproton distance data. FEBS Letters 1987;219:11–16.
10. Tejero R, Bassolino-Klimas D, Bruccoleri RE, Montelione GT. Simulated annealing with restrained molecular dynamics using CONGEN: energy refinement of the NMR solution structure of epidermal and type-α transforming growth factors. Protein Sci 1996;5:578–592.
11. Spronk CAEM, Nabuurs SB, Bonvin AMJJ, Krieger E, Vuister GW, Vriend G. The precision of NMR structure ensembles revisited. J Biomol NMR 2003;25:225–234.
12. Nederveen AJ, Doreleijers JF, Vranken W, Miller Z, Spronk CAEM, Nabuurs SB, Güntert P, Livny M, Markley JL, Nilges M, Ulrich EL, Kaptein R, Bonvin AMJJ. RECOORD: a recalculated coordinate database of 500+ proteins from the PDB using restraints from the BioMagResBank. Proteins 2005;59:662–672.
13. Rieping W, Habeck M, Nilges M. Inferential structure determination. Science 2005;309:303–306.

14. Habeck M, Rieping W, Nilges M. Weighting of experimental evidence in macromolecular structure determination. Proc Natl Acad Sci USA 2006;103:1756–1761.

15. Brünger AT. Free *R* value: a novel statistical quantity for assessing the accuracy of crystal structures. Nature 1992;355:472–475.

16. Brünger AT. Free *R* value: cross-validation in crystallography. Methods Enzymol 1997;277:366–396.

17. Nabuurs SB, Spronk CAEM, Krieger E, Maassen H, Vriend G, Vuister GW. Quantitative evaluation of experimental NMR constraints. J Am Chem Soc 2003;125:12026–12034.

18. Brünger AT, Clore GM, Gronenborn AM, Saffrich R, Nilges M. Assessing the quality of solution nuclear magnetic resonance structures by complete cross-validation. Science 1993;261:328–331.

19. Gronwald W, Kirchhöfer R, Görler A, Kremer W, Ganslmeier B, Neidig K-P, and Kalbitzer HR. RFAC, a program for automated NMR R-factor estimation. J Biomol NMR 2000;17:137–151.

20. Huang YJ, Powers R, Montelione GT. Protein NMR recall, precision and F-measure scores (RPF scores): structure quality assessment measures based on information retrieval statistics. J Am Chem Soc 2005;127:1665–1674.

21. Billeter M. Comparison of protein structures determined by NMR in solution and by X-ray diffraction in single crystals. Q Rev Biophys 1992;25:325–377.

22. Wagner G, Hyberts SG, Havel TF. NMR structure determination in solution: a critique and comparison with X-ray crystallography. Annu Rev Biophys Biomol Struct 1992;21:167–198.

23. MacArthur MW, Laskowski RA, Thornton JM. Knowledge-based validation of protein structure coordinates derived by X-ray crystallography and NMR spectroscopy. Curr Opinion Struct Biol 1994;4:731–737.

24. Gronenborn AM. Clore GM. Structures of protein complexes by multidimensional heteronuclear magnetic resonance spectroscopy. Crit Rev Biochem Mol Biol 1995;30:351–385.

25. Godzik A, Koliński A, Skolnick J. Are proteins ideal mixtures of amino acids? Analysis of energy parameter sets. Protein Sci 1995;4:2107–2117.

26. Bastolla U, Farwer J, Knapp EW, Vendruscolo M. How to guarantee optimal stability for most representative structures in the Protein Data Bank. Proteins 2001;44:79–96.

27. Garbuzynskiy SO, Melnik BS, Lobanov MY, Finkelstein AV, Galzitskaya OV. Comparison of X-ray and NMR structures: is there a systematic difference in residue contacts between X-ray- and NMR-resolved protein structures? Proteins 2005;60:139–147.

28. Kuszewski J, Gronenborn AM, Clore GM. Improving the quality of NMR and crystallographic protein structures by means of a conformational database potential derived from structure databases. Protein Sci 1996;5:1067–1080.

29. Kuszewski J, Gronenborn AM, Clore GM. Improving the packing and accuracy of NMR structures with a pseudopotential for the radius of gyration. J Am Chem Soc 1999;121:2337–2338.

30. Bhattacharyya GK, Johnson RA. Statistical concepts and methods. New York: Wiley; 1977.

31. Johnson NL, Kotz S. Distributions in statistics: continuous univariate distributions—1. Boston: Houghton Mifflin; 1970.

32. Cleveland WS. Visualizing data. Summit, New Jersey: Hobart Press; 1993.

33. Abramowitz M, Stegun IA. Handbook of mathematical functions. New York: Dover Publications; 1972.

34. Chen J, Won H-S, Im W, Dyson HJ, Brooks CL III. Generation of native-like protein structures from limited NMR data, modern force fields and advanced conformational sampling. J Biomol NMR 2005; 31:59–64.

35. Bhattacharya A, Tejero R, Montelione GT. Evaluating protein structures determined by structural genomics consortia. Proteins 2007; 66:778–795.

36. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. J R Statist Soc B 1995;57:289–300.

37. Clore GM, Robien MA, Gronenborn AM. Exploring the limits of precision and accuracy of protein structures determined by nuclear magnetic resonance spectroscopy. J Mol Biol 1993;231:82–102.

38. Laskowski RA, MacArthur MW, Moss DS, Thornton JM. PRO-CHECK: a program to check the stereochemical quality of protein structures. J Appl Cryst 1993;26:283–291.

39. Sippl MJ. Recognition of errors in three-dimensional structures of proteins. Proteins 1993;17:355–362.

40. Lüthy R, Bowie JU, Eisenberg D. Assessment of protein models with three-dimensional profiles. Nature 1992;356:83–85.

41. Lovell SC, Davis IW, Arendall WB III, de Bakker PIW, Word JM, Prisant MG, Richardson JS, Richardson DC. Structure validation by Cα geometry: φ, ψ, and Cβ deviation. Proteins 2003;50:437–450.

42. Lipsitz RS. Tjandra N. Residual dipolar couplings in NMR structural analysis. Annu Rev Biophys Biomol Struct 2004;33:387–413.

43. Prestegard JH, Bougault CM, Kishore AI. Residual dipolar couplings in structure determination of biomolecules. Chem Rev 2004;104:3519–3540.

44. Bax A, Grishaev A. Weak alignment NMR: a hawk-eyed view of biomolecular structure. Curr Opin Struct Biol 2005;15:563–570.

45. Chien C-Y, Tejero R, Huang Y, Zimmerman DE, Ríos CB, Krug RM, Montelione GT. A novel RNA-binding motif in influenza A virus non-structural protein 1. Nat Struct Biol 1997;4:891–895.

46. Liu J, Lynch PA, Chien C-Y, Montelione GT, Krug RM, Berman HM. Crystal structure of the unique RNA-binding domain of the influenza virus NS1 protein. Nat Struct Biol 1997;4:896–899.

47. Alexeev D, Bury SM, Turner MA, Ogunjobi OM, Muir TW, Ramage R, Sawyer L. Synthetic, structural and biological studies of the ubiquitin system: chemically synthesized and native ubiquitin fold into identical three-dimensional structures. Biochem J 1994;299:159–163.

48. Cornilescu G, Marquardt JL, Ottiger M, Bax A. Validation of protein structure from anisotropic carbonyl chemical shifts in a dilute liquid crystalline phase. J Am Chem Soc 1998;120:6836–6837.

49. Havel TF, Wüthrich K. An evaluation of the combined use of nuclear magnetic resonance and distance geometry for the determination of protein conformations in solution. J Mol Biol 1985;182:281–294.

50. Zhao D, Jardetzky O. An assessment of the precision and accuracy of protein structures determined by NMR: dependence on distance errors. J Mol Biol 1994;239:601–607.

51. Chalaoux F-R, O'Donoghue SI, Nilges M. Molecular dynamics and accuracy of NMR structures: effects of error bounds and data removal. Proteins 1999;34:453–463.

52. Bassolino-Klimas D, Tejero R, Krystek SR, Metzler WJ, Montelione GT, Bruccoleri RE. Simulated annealing with restrained molecular dynamics using a flexible restraint potential: theory and evaluation with simulated NMR constraints. Protein Sci 1996;5:593–603.

53. Kossiakoff AA, Randal M, Guenot J, Eigenbrot C. Variability of conformations at crystal contacts in BPTI represent true low-energy structures: correspondence among lattice packing and molecular dynamics structures. Proteins 1992;14:65–74.

54. Eigenbrot C, Randal M, Kossiakoff AA. Structural effects induced by mutagenesis affected by crystal packing factors: the structure of a 30–51 disulfide mutant of basic pancreatic trypsin inhibitor. Proteins 1992;14:75–87.

55. Mittermaier A, Kay LE. New tools provide new insights in NMR studies of protein dynamics. Science 2006;312:224–228.

56. Eisenmesser EZ, Bosco DA, Akke M, Kern D. Enzyme dynamics during catalysis. Science 2002;295:1520–1523.

57. Wolf-Watz M, Thai V, Henzler-Wildman K, Hadjipavlou G, Eisenmesser EZ, Kern D. Linkage between dynamics and catalysis in a thermophilic-mesophilic enzyme pair. Nat Struct Mol Biol 2004;11:945–949.

58. Eisenmesser EZ, Millet O, Labeikovsky W, Korzhnev DM, Wolf-Watz M, Bosco DA, Skalicky JJ, Kay LE, Kern D. Intrinsic dynamics of an enzyme underlies catalysis. Nature 2005;438:117–121.

59. Juers DH. Matthews BW. Reversible lattice repacking illustrates the temperature dependence of macromolecular interactions. J Mol Biol 2001;311:851–862.

60. Juers DH, Matthews BW. Cryo-cooling in macromolecular crystallography: advantages, disadvantages and optimization. Q Rev Biophys 2004;37:105–119.

61. DePristo MA, de Bakker PIW, Blundell TL. Heterogeneity and inaccuracy in protein structures solved by X-ray crystallography. Structure 2004;12:831–838.

62. Furnham N, Blundell TL, DePristo MA, Terwilliger TC. Is one solution good enough? Nat Struct Mol Biol 2006;13:184–185.

63. Bonvin AMJJ, Brünger AT. Conformational variability of solution nuclear magnetic resonance structures. J Mol Biol 1995;250:80–93.

64. Dedmon MM, Lindorff-Larsen K, Christodoulou J, Vendruscolo M, Dobson CM. Mapping long-range interactions in α-synuclein using spin-label NMR and ensemble molecular dynamics simulations. J Am Chem Soc 2005;127:476–477.

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

66. Burling FT, Brünger AT. Thermal motion and conformational disorder in protein crystal structures: comparison of multi-conformer and time-averaging models. Israel J Chem 1994;34:165–175.

67. Knight JL, Zhou Z, Gallicchio E, Himmel DM, Friesner RA, Arnold E, Levy RM. Modeling maximal structural diversity in X-ray crystallographic refinement using protein local optimization by torsion angle sampling. 2007, submitted for publication.

68. Blackledge M. Recent progress in the study of biomolecular structure and dynamics in solution from residual dipolar couplings. Prog Nucl Magn Reson Spect 2005;46:23–61.

69. Tolman JR, Ruan K. NMR residual dipolar couplings as probes of biomolecular dynamics. Chem Rev 2006;106:1720–1736.

70. Tolman JR, Flanagan JM, Kennedy MA, Prestegard JH. NMR evidence for slow collective motions in cyanometmyoglobin. Nat Struct Biol 1997;4:292–297.

71. Bax A, Tjandra N. Are proteins even floppier than we thought? Nat Struct Biol 1997;4:254–256.

72. Peti W, Meiler J, Brüschweiler R, Griesinger C. Model-free analysis of protein backbone motion from residual dipolar couplings. J Am Chem Soc 2002;124:5822–5833.

73. Clore GM, Schwieters CD. How much backbone motion in ubiquitin is required to account for dipolar coupling data measured in multiple alignment media as assessed by independent cross-validation? J Am Chem Soc 2004;126:2923–2938.

74. DePristo MA, de Bakker PIW, Johnson RJK, Blundell TL. Crystallographic refinement by knowledge-based exploration of complex energy landscapes. Structure 2005;13:1311–1319.