

Refinement of Multidomain Protein Structures by Combination of Solution Small-Angle X-ray Scattering and NMR Data

Alexander Grishaev,^{*,†} Justin Wu,[‡] Jill Trewhella,[§] and Ad Bax^{*,†}

Contribution from the Laboratory of Chemical Physics, NIDDK, National Institutes of Health, Bethesda, Maryland 20892-0520, Department of Biochemistry, The Ohio State University, Columbus, Ohio 43210, and Department of Chemistry, University of Utah, Salt Lake City, Utah 84112-0850

Received June 30, 2005; E-mail: grishaev@speck.niddk.nih.gov; bax@nih.gov

Abstract: Determination of the 3D structures of multidomain proteins by solution NMR methods presents a number of unique challenges related to their larger molecular size and the usual scarcity of constraints at the interdomain interface, often resulting in a decrease in structural accuracy. In this respect, experimental information from small-angle scattering of X-ray radiation in solution (SAXS) presents a suitable complement to the NMR data, as it provides an independent constraint on the overall molecular shape. A computational procedure is described that allows incorporation of such SAXS data into the mainstream high-resolution macromolecular structure refinement. The method is illustrated for a two-domain 177-amino-acid protein, γ S crystallin, using an experimental SAXS data set fitted at resolutions from ~ 200 Å to ~ 30 Å. Inclusion of these data during structure refinement decreases the backbone coordinate root-mean-square difference between the derived model and the high-resolution crystal structure of a 54% homologous γ B crystallin from 1.96 ± 0.07 Å to 1.31 ± 0.04 Å. Combining SAXS data with NMR restraints can be accomplished at a moderate computational expense and is expected to become useful for multidomain proteins, multimeric assemblies, and tight macromolecular complexes.

Introduction

Determination of the three-dimensional structures of large proteins by solution NMR techniques presents a number of unique challenges. Increased line width resulting from slower rotational diffusion leads to a decrease in signal-to-noise ratio, increased resonance overlap, and larger uncertainty of the resonance positions. These effects decrease the number of observable NMR signals and complicate the process of their assignment. One way to address this problem is by combining ^{13}C and ^{15}N enrichment with perdeuteration, where the majority of ^1H nuclei are replaced by the effectively NMR-invisible ^2H .^{1,2} When complemented by transverse relaxation-optimized spectroscopy (TROSY)-based pulse sequence techniques, such labeling leads to a dramatic simplification of the NMR spectra, narrower resonance signals, and increased signal-to-noise ratios.³ Perdeuteration, however, also has a downside: since it effectively makes sparse the set of NMR observables, it decreases the intrinsic information content of the NMR data. Additional difficulties arise due to the nonglobular nature of many multidomain proteins. Even though the conformations and relative orientations of the individual domains can be determined

accurately by using backbone–backbone nuclear Overhauser effects (NOEs) and extensive sets of residual dipolar couplings (RDCs), relative positioning of the individual domains can remain challenging as protein perdeuteration eliminates the majority of the resonances necessary for defining the requisite side-chain-mediated interdomain NOE contacts.

Any source of experimental data that can compensate for the decrease in NOE restraint information associated with the application of NMR to large, multidomain proteins is therefore expected to be invaluable. In particular, information is needed that complements restraints derived from the common types of NMR data, including short-range interproton distances derived from NOEs,^{4–7} dihedral angles derived from J couplings,^{8,9} and orientations derived from residual dipolar couplings.^{10,11} It is well recognized that such complementary information is contained in the profiles of small-angle scattering of X-ray radiation by macromolecules in solution (SAXS).¹² Previously, SAXS data have been used in ad hoc calculations to complement NMR

[†] NIDDK, National Institutes of Health.

[‡] The Ohio State University.

[§] University of Utah.

- (1) Grzesiek, S.; Anglister, J.; Ren, H.; Bax, A. *J. Am. Chem. Soc.* **1993**, *115*, 4369–4370.
- (2) Tugarinov, V.; Hwang, P. M.; Kay, L. E. *Annu. Rev. Biochem.* **2004**, *73*, 107–146.
- (3) Salzmann, M.; Wider, G.; Pervushin, K.; Senn, H.; Wuthrich, K. *J. Am. Chem. Soc.* **1999**, *121*, 844–848.

- (4) Wuthrich, K. *NMR of Proteins and Nucleic Acids*; John Wiley & Sons: New York, 1986.
- (5) Kaptein, R.; Boelens, R.; Scheek, R. M.; van Gunsteren, W. F. *Biochemistry* **1988**, *27*, 5389–5395.
- (6) Clore, G. M.; Gronenborn, A. M. *Crit. Rev. Biochem. Mol. Biol.* **1989**, *24*, 479–564.
- (7) Wagner, G. *J. Biomol. NMR* **1993**, *3*, 375–385.
- (8) Biamonti, C.; Rios, C. B.; Lyons, B. A.; Montelione, G. T. *Adv. Biophys. Chem.* **1994**, *4*, 51–120.
- (9) Bax, A. *Methods Enzymol.* **1994**, *239*, 79–125.
- (10) Tolman, J. R.; Flanagan, J. M.; Kennedy, M. A.; Prestegard, J. H. *Proc. Natl. Acad. Sci. U.S.A.* **1995**, *92*, 9279–9283.
- (11) Tjandra, N.; Bax, A. *Science* **1997**, *278*, 1111–1114.
- (12) Svergun, D. I.; Koch, M. H. *J. Rep. Prog. Phys.* **2003**, *66*, 1735–1782.

data in solving the solution structures of modular proteins (e.g., the Gla-EGF domain of the blood coagulation factor Xa protein¹³ and a calmodulin/trifluoperazine complex¹⁴) essentially by evaluating which NMR-derived relative domain positions are in best agreement with the SAXS data or by a grid search for a 3D translation vector between the rigidly held domains. However, the potential for combining the two types of data has never been fully exploited directly in NMR structure calculation.

The SAXS intensity curve, recorded as a function of the scattering angle, is essentially a Fourier transform of the distribution of the interatomic distances within the macromolecule. Since the latter is known to encode both the overall molecular shape and the nonuniform distribution of the protein's atomic density,¹⁵ incorporation of this information into macromolecular structure refinement can compensate for the deficiency of the translational information derived from interdomain NOEs. Other advantages of using SAXS in the context of NMR-based structure determination are its independence of isotopic labeling, the high speed of data acquisition at the conditions that can be matched to those used for the solution NMR experiments, and smaller sample volumes ($\sim 15 \mu\text{L}$ per sample) compared to those required for NMR measurements. The main experimental challenges in applying SAXS methodology are the following: (i) sample conditions have to be carefully optimized to prevent aggregation, (ii) subtraction of the solvent contribution to the scattering must be done with high precision, and (iii) the sample can suffer radiation damage.

Here we demonstrate that direct incorporation of SAXS data in NMR structure calculation is readily feasible, and at moderate computational expense. The combination of NMR data, recently used for determining the solution structure of the eye lens protein γS crystallin, with SAXS data results in considerably closer agreement with the X-ray structures of homologous members of the γ -crystallin family than the original NMR structure.

Materials and Methods

Protein Sample Preparation. A uniformly ^{15}N -enriched sample of γS crystallin was used for collecting the SAXS data. Enrichment of the protein in ^{15}N was used only because the sample initially was intended for NMR studies, and does not affect the protein stability or its scattering profile. Protein preparation details have been described elsewhere.¹⁶ To minimize oxidation-induced dimerization through the Cys residues on the surface of the protein, the sample was dialyzed against 100 mL of buffer containing fresh reducing agent (dithiothreitol, DTT) for 6 h under the flow of N_2 on-site, immediately prior to data acquisition. The sample composition was 9 mg/mL protein, 0.04% NaN_3 , 5 mM DTT, 25 mM imidazole, pH 6.0. An aliquot of the dialysate was used to measure the solvent blank, which must be subtracted from the sample measurement in order to determine the scattering from the protein molecules alone. This same dialysate was also used for diluting the sample, to evaluate the concentration dependence of the SAXS profile.

SAXS Data Acquisition and Processing. Each 12 μL sample was centrifuged at ~ 1000 rpm into a glass capillary mounted on a brass holder, which was used to position the capillary precisely and reproducibly in the focused X-ray beam. Scattering data were acquired

with the sample cooled to 291.4 K using the X-ray instrument at the University of Utah, described in a previous publication.¹⁷ The instrument uses a sealed tube source (Cu $K\alpha$ -edge giving 1.542 Å wavelength) and a slit geometry with a one-dimensional position-sensitive detector. The sample-to-detector distance was 0.64 m, corresponding to an accessible q range of 0.0054–0.3192 Å⁻¹. Individual detector channels were mapped onto the momentum transfer axis using the 50.1 ± 0.1 Å d spacing of the (100) reflection of the polycrystalline cholesterol myristate sample. To prevent oxidation of the sample by air during the measurement, N_2 was flowing around the capillary throughout the experiment. Scattering data were acquired for 12 h per sample at two protein concentrations: 9.0 and 4.5 mg/mL. Data normalization, correction for the detector sensitivity, and subtraction of the solvent scattering were done as described previously.¹⁷ Preliminary data analysis was done using Guinier formalism and $P(r)$ analysis based on an indirect Fourier transform; it uses a $\sin(x)/x$ series expansion and is implemented in the program P_of_R that includes beam geometry corrections.¹⁸ The $P(r)$ analysis was also carried out using the program GNOM^{19,20} which, along with the beam geometry corrections, utilizes a regularized indirect transform and thus avoids the potential for systematic oscillations in the calculated $P(r)$. For the acquired γS crystallin data, both programs gave essentially the same result, indicating that the scattering data are of good quality in that they have a robust $P(r)$ solution, independent of the details of the Fourier transform. The contribution to the scattering arising from the hydration layer at the surface of the protein was calculated for a given structure by fitting the desmeared scattering data to the structure in question using the program CRYSOLO.²¹ The globboc correction was calculated from the structural coordinates using scattering profile simulation software written in-house, and available upon request from the authors.

Structure Calculation Protocol. γS crystallin structure models were generated by a restrained molecular dynamics simulated annealing protocol using the CNS package.²² The force field included the usual empirical energy terms: bonds, angles, improper angles, and a repulsive-only quartic nonbonded term with all van der Waals radii scaled down by a factor of 0.8, as well as a backbone–backbone hydrogen-bonding potential of mean force.²³ Additional terms included those for the NOEs, experimental dihedral angles, and RDCs, and were identical to those used previously for calculating the γS crystallin structure in the absence of SAXS data (Protein Data Bank (PDB) entries 1ZWM and 1ZWO). The temperature was linearly decreased from 2000 K to 1 K in 200 stages of 200 steps each, with the $\text{H}^{\text{N}}\text{--N}$ RDC force constant ramped up from 0.01 to 0.40 kcal/Hz². NOE and backbone dihedral angle force constants were fixed throughout the calculations at 50 kcal/Å² and 10 kcal/rad², respectively. All statistics were extracted from the ensembles of 20 calculated structures, starting from the structures previously calculated and deposited in the absence of SAXS data. In all cases, reference calculations were run in exactly the same way, but with the SAXS data fit term inactivated. The original NMR structure of γS crystallin was based primarily on backbone one-bond dipolar couplings, supplemented by a moderate number of easily accessible $\text{H}^{\text{N}}\text{--H}^{\text{N}}$ and $\text{CH}_3\text{--CH}_3$ NOE data. A total of 179 $\text{H}^{\text{N}}\text{--H}^{\text{N}}$ NOEs and 70 $\text{CH}_3\text{--CH}_3$ NOEs were available, 15 of them between the N- and C-terminal domains. The dipolar restraints include an extensive set of couplings recorded in two media, and comprise 291 N--H^{N} , 303 C--C^{α} , 273 $\text{N--C}'$, and 246 $\text{C}^{\alpha}\text{--C}^{\beta}$ RDCs. Backbone

(17) Heidorn, D. B.; Trewthella, J. *Biochemistry* **1988**, *27*, 909–915.

(18) Moore, P. B. *J. Appl. Crystallogr.* **1980**, *13*, 168–175.

(19) Svergun, D. I.; Petoukhov, M. V.; Koch, M. H. J. *Biophys. J.* **2001**, *80*, 2946–2953.

(20) Svergun, D. I. *J. Appl. Crystallogr.* **1992**, *25*, 495–503.

(21) Svergun, D.; Barberato, C.; Koch, M. H. J. *J. Appl. Crystallogr.* **1995**, *28*, 768–773.

(22) Brunger, A. T.; Adams, P. D.; Clore, G. M.; DeLano, W. L.; Gros, P.; Grosse-Kunstleve, R. W.; Jiang, J. S.; Kuszewski, J.; Nilges, M.; Pannu, N. S.; Read, R. J.; Rice, L. M.; Simonson, T.; Warren, G. L. *Acta Crystallogr. D, Biol. Crystallogr.* **1998**, *54*, 905–921.

(23) Grishaev, A.; Bax, A. *J. Am. Chem. Soc.* **2004**, *126*, 7281–7292.

(13) Sunnerhagen, M.; Olah, G. A.; Stenflo, J.; Forsen, S.; Drakenberg, T.; Trewthella, J. *Biochemistry* **1996**, *35*, 11547–11559.

(14) Mattinen, M. L.; Paakkonen, K.; Ikonen, T.; Craven, J.; Drakenberg, T.; Serimaa, R.; Waltho, J.; Annala, A. *Biophys. J.* **2002**, *83*, 1177–1183.

(15) Koch, M. H. J.; Vachette, P.; Svergun, D. I. *Q. Rev. Biophys.* **2003**, *36*, 147–227.

(16) Wu, Z.; Delaglio, F.; Wyatt, K.; Wistow, G.; Bax, A. *Protein Sci.* **2005**, *14*, 3101–3114.

dihedral angles (ϕ, ψ) are restrained by values derived from the previously described molecular fragment replacement (MFR) database search procedure,^{16,24} which is based on the observed dipolar couplings and yields a total of 318 torsion restraints. Restraints for 71 χ^1 and 11 χ^2 side-chain angles, extracted from ${}^3J_{C\gamma C'}$ and ${}^3J_{C\gamma N}$ couplings, were also used.

Results and Discussion

SAXS Data Analysis in the Context of High-Resolution Structure Refinement. X-rays are scattered by electrons, and the intensity of the radiation scattered by the macromolecules in solution depends on the electron scattering density difference, or “contrast”, between the macromolecule and the bulk solvent. An additional contribution to the scattering arises from a thin layer of solvent at the macromolecular surface which can have an electron density different from that of the bulk solvent. The existence of the latter hydration layer effect has been demonstrated in a number of experimental and computational studies.^{21,25,26}

In isotropic conditions, the scattering intensity is averaged over all orientations of the macromolecule with respect to the incident radiation beam. The scattering vector $q = 4\pi(\sin \theta)/\lambda$ denotes the momentum transfer between the incident beam of wavelength λ and the radiation scattered at the angle 2θ . In the absence of macromolecular aggregation, the intensity of the scattered beam can be represented as²¹

$$I(q) = \langle |A_m(\mathbf{q}) - \rho_s A_s(\mathbf{q}) + \delta\rho A_l(\mathbf{q})|^2 \rangle_\Omega \quad (1)$$

Here $\langle \rangle_\Omega$ denotes the solid angle average over all orientations of the momentum transfer vector \mathbf{q} for the fixed norm q , $A_m(\mathbf{q})$, $A_s(\mathbf{q})$, and $A_l(\mathbf{q})$ are the scattering amplitudes of the macromolecule, solvent displaced by the macromolecular volume, and its hydration layer, respectively, and ρ_s and $\delta\rho$ are the bulk solvent electron density ($0.334 \text{ e}/\text{\AA}^3$) and the density of the hydration layer ($0.00\text{--}0.07 \text{ e}/\text{\AA}^3$).²¹ At a given orientation of the momentum transfer vector \mathbf{q} with respect to the molecular frame, the scattering amplitude of the macromolecule is a Fourier transform of the atomic coordinates \mathbf{r}_j over its N atoms, weighted by the atomic X-ray form factors f_j :

$$A_m(\mathbf{q}) = \sum_{j=1}^N f_j(q) \exp(i\mathbf{q}\mathbf{r}_j) \quad (2)$$

The scattering of the solvent displaced by the macromolecule can be approximated by placing dummy solvent atoms at all atomic positions within the macromolecule with the form factors given by²⁷

$$g_j(q) = G(q)V_j \exp\left(-\frac{q^2 V_j^{2/3}}{4\pi}\right) \quad (3)$$

Here, V_j are the volumes of the solvent displaced by each atom represented by the Gaussian spheres of previously tabulated²⁷ radii r_j . The expansion factor $G(q)$ is given by^{21,25,26}

$$G(q) = \left(\frac{r_0}{r_m}\right)^3 \exp\left(-\frac{q^2(r_0^2 - r_m^2)}{(36\pi)^{1/3}}\right) \quad (4)$$

Here, r_0 is the average atomic radius in the macromolecule and r_m is the adjustable parameter that allows one to vary the average displaced solvent volume per atomic group. Here, we set $r_m = r_0$, which makes the expansion factor equal to one. The total scattering amplitude of the contrast between the macromolecule and the displaced solvent can then be conveniently expressed as the Fourier transform of the macromolecular coordinates weighted by the solvent-corrected form factors f_j^s :

$$\begin{aligned} A_m(\mathbf{q}) - \rho_s A_s(\mathbf{q}) &= \sum_{j=1}^N [f_j(q) - \rho_s g_j(q)] \exp(i\mathbf{q}\mathbf{r}_j) \\ &= \sum_{j=1}^N f_j^s(q) \exp(i\mathbf{q}\mathbf{r}_j) \end{aligned} \quad (5)$$

We will restrict our treatment to the range of $q < 1 \text{ \AA}^{-1}$, where this approximate procedure can be expected to work reasonably well.

There are two common approaches to solid angle averaging over the $\exp(i\mathbf{q}\mathbf{r}_j)$ terms, one exploiting the favorable properties of their spherical harmonics expansion^{21,25,26} and the other relying on application of the Debye formula.^{28,29} Both involve a comparable computational overhead for proteins of up to ~ 300 residues. We chose the Debye formula for its mathematical simplicity, representing the spherical average in eq 1 as

$$I(q) = \sum_{i=1}^N \sum_{j=1}^N f_i^s(q) f_j^s(q) \frac{\sin(qr_{ij})}{qr_{ij}} \quad (6)$$

The quality of the fit between the experimental scattering data and those predicted from the model is described by the χ^2 statistics over the set of N_q experimental values:

$$\chi^2 = \frac{1}{N_q - 1} \sum_{k=1}^{N_q} \left[\frac{I_{\text{expt}}(q_k) - c_k I_{\text{calc}}(q_k)}{\sigma(q_k)} \right]^2 \quad (7)$$

Here, c_k are scattering vector-dependent correction factors described in more detail below and $\sigma(q_k)$ are the uncertainties of each experimental data point q_k . Fitting SAXS data would thus involve simulation of the model-based scattering intensity $I_{\text{calc}}(q_k)$ for all q_k , correction of the latter by the c_k factors, calculation of the χ^2 statistics, and finally, differentiation of χ^2 with respect to the current atomic coordinates to yield a set of atomic forces that aim to minimize χ^2 . When added to an empirical force field used in the molecular dynamics (MD)-based structure refinement, these forces should allow a refinement against SAXS data in combination with other data sources (in this case, a set of NMR-generated restraints). The gradient of the χ^2 with respect to the atomic coordinates r_j can be expressed as

(24) Kontaxis, G.; Delaglio, F.; Bax, A. *Methods Enzymol.* **2005**, *394*, 42–78.

(25) Merzel, F.; Smith, J. C. *Acta Crystallogr. D, Biol. Crystallogr.* **2002**, *58*, 242–249.

(26) Svergun, D. I. *Biophys. J.* **1999**, *76*, 2879–2886.

(27) Fraser, R. D. B.; Macrae, T. P.; Suzuki, E. *J. Appl. Crystallogr.* **1978**, *11*, 693–694.

(28) Chacon, P.; Moran, F.; Diaz, J. F.; Pantos, E.; Andreu, J. M. *Biophys. J.* **1998**, *74*, 2760–2775.

(29) Walther, D.; Cohen, F. E.; Doniach, S. *J. Appl. Crystallogr.* **2000**, *33*, 350–363.

$$\nabla_{r_j}[\chi^2] \approx \sum_{k=1}^{N_q} \frac{I_{\text{expt}}(q_k) - c_k I_{\text{calc}}(q_k)}{\sigma_k^2} \sum_{i=1}^N \sum_{j \neq i}^N f_i^s(q_k) f_j^s(q_k) \times \left[\cos(q_k r_{ij}) - \frac{\sin(q_k r_{ij})}{q_k r_{ij}} \right] \frac{\mathbf{r}_{ij}}{r_{ij}^2} \quad (8)$$

Hence, fitting SAXS data involves evaluation of eqs 6–8 at each step of molecular dynamics/energy minimization. Because the number of operations necessary for these calculations scales as $N_q N^2$, it is clear that one problem that has been preventing incorporation of SAXS data into structure refinement is its enormous computational overhead. For example, calculation of the χ^2 and its gradients takes tens of seconds of CPU time on a modern Pentium-class processor per step, for proteins between 100 and 200 residues in length. Since MD trajectories commonly used in high-resolution structure refinement may involve 10^4 – 10^5 such steps, the challenges are quite apparent.

The solution to this problem is hinted at by the form of the $N_q N^2$ expression: a suitable approximation to eqs 6–8 with smaller values of N_q and N will alleviate the computational burden. Starting with N^2 -dependent terms, it is known that the shapes of the spherically averaged scattering form factors of small, closely proximal sets of atoms do not show a pronounced dependence on the exact atomic geometries below ~ 3 Å resolution.³⁰ The resulting “globbic approximation”, in which an all-atom representation of the macromolecular structure is coarse-grained into a smaller number of spatially proximal “globs”, has been widely used in the interpretation of the low-resolution X-ray crystallographic³¹ and SAXS^{19,28} data. Following this strategy, we have split protein structures into sets of small fragments, each involving 3–9 heavy atoms, along with their associated H’s (see the Supporting Information for the definition of the “globs”). We have then recalculated the spherically averaged scattering form factors for each glob as

$$f_k^{\text{glob}}(q) = \left[\sum_{i=1}^N \sum_{j=1}^N f_i^s(q) f_j^s(q) \frac{\sin(q r_{ij})}{q r_{ij}} \right]^{1/2} \quad (9)$$

One can then approximate the scattering intensity curve with the sum in eq 6 running over the set of globs, positioned at the coordinates weighted by the atomic electron number counts within each glob, and using the globbic form factors instead of the atomic ones. Since our specification reduces N input heavy atoms into approximately $N/3$ globs, the required CPU time is reduced by about an order of magnitude. The procedure, however, has a drawback: the approximated scattering intensity curves show small but systematic differences with respect to the “exact” ones, obtained from all-atom calculations. We address this problem via an approach used by others:³¹ derivation of “globbic” correction factors $c_k = c(q_k)$ as ratios between the “exact” scattering curves and the globbically approximated ones. Figure 1 shows the average and standard deviation of this correction, calculated over a large set of protein structures in the 100–200 residue size range. Application of such a correction will decrease the systematic errors of our approximation to

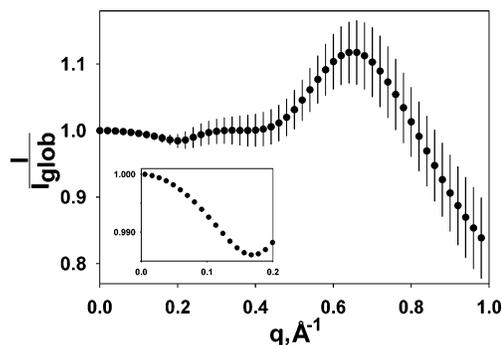


Figure 1. Globbic correction factor calculated as the ratio between the atomic and globbic scattering curves, $I(q)$ and $I_{\text{glob}}(q)$, respectively. The mean and standard deviation of the curve points are calculated on the basis of 538 single-chain protein X-ray structures of 100–200 residues length, solved at resolutions of 1.8 Å or better. The calculations were carried out according to eqs 6 and 9. The inset shows the average correction factor calculated from the γ S crystallin models used in the final round of structure refinement.

values comparable to the error bars indicated within the figure. Notice that since our globs are smaller than the “dummy residues” usually employed in SAXS data analysis, the average correction factors and their variances are smaller than the ones obtained in those approaches (compare to Figure 2 of ref 19). In fact, we have adjusted the size and composition of the globs to provide a conservative compromise between the computational speed-up and the magnitudes of the systematic errors resulting from the approximate nature of the calculation. The shape and overall features of the globbic correction curve are largely independent of the size and secondary structure content of the protein, while showing a pronounced dependence on the glob size, especially in the higher resolution range (see Supporting Information for details).

In practice, these correction factors are calculated from the current structural model, and re-estimated after each successive cycle of structure refinement until convergence is reached. Such a procedure will, in general, ensure that the approximated globbic correction curve approaches the exact one as the refined structure approaches the correct model. The calculated scattering intensity curves are also corrected for the effect of the bound solvent layer using CRY SOL,²¹ taking as input the entire family of structures prior to every cycle of structure refinement and fitting the bound solvent density as the only adjustable parameter.

The second part of our strategy involves reducing N_q , the number of experimental points to be fitted. For proteins of up to a few hundred residues, the maximum curvature of the simulated scattering curves, ca. 10^{-2} \AA^{-1} , is much smaller than the scattering vector step of the oversampled experimental data (typically ca. 10^{-3} \AA^{-1}). Reduction of the fitted data set to fewer points within the same q interval is thus expected to speed-up the calculation by an amount proportional to the ratio of the number of points in the original data to that in the “sparsened” data set. If the separation in q between the sparsened data points is substantially smaller than the distance between the features of the scattering curve, sparsening is not expected to have any detrimental effects on the accuracy of the data representation. We have performed a regularized fit of the oversampled, desmeared experimental data set using the package GNOM^{20,32}

(30) Guo, D. Y.; Blessing, R. H.; Langs, D. A.; Smith, G. D. *Acta Crystallogr. D, Biol. Crystallogr.* **1999**, *55*, 230–237.

(31) Guo, D. Y.; Blessing, R. H.; Langs, D. A. *Acta Crystallogr. D, Biol. Crystallogr.* **2000**, *56*, 1148–1155.

(32) Svergun, D. I. *Biophys. J.* **1991**, *24*, 485–592.

and sparsened the smoothed data fit by a factor of 8. The combination of these two procedures results in an overall speed-up factor of ~ 80 , placing the time for a single-point SAXS pseudo-energy/forces calculation to less than $\sim 1/3$ of a second for a protein of up to ~ 180 residues, when fitting up to 30 SAXS data points on a 2.8 GHz Pentium 4 processor. This gain makes it possible to conduct regular-length MD structure refinement in a reasonable amount of time (ca. 6 h per structure for 40 000 MD steps). The SAXS data fitting module was coded into the CNS structure refinement package²² with the corresponding energy term introduced by the "SAXS" keyword.

Application to γ S Crystallin. We demonstrate the utility of the solution scattering data in NMR structure refinement of murine γ S crystallin, a two-domain eye lens protein of 177 residues. The N- and C-terminal domains are topologically similar, each consisting of two four-strand β -sheets arranged in Greek key motifs, linked by a Tyr corner. The entire protein shares 54% sequence identity with bovine γ B crystallin, for which a 1.1 Å resolution X-ray structure is available (PDB code 1AMM³³), and 50% sequence identity with human γ D crystallin (PDB code 1HK0³⁴). In addition, a crystal structure is available for a dimer formed by the C-terminal domains of bovine γ S crystallin (PDB code 1A7H³⁵). The primary sequence of γ S crystallin can be aligned to these entries without any gaps or insertions within each individual domain.

The NMR structure for γ S crystallin was recently determined by molecular fragment replacement (MFR) methodology,²⁴ using primarily dipolar couplings as input restraints, supplemented by small numbers of H^N-H^N and CH_3-CH_3 NOE restraints.¹⁶ The two globular domains of the recent NMR structure of γ S crystallin are very similar to those seen in the homologous γ B crystallin (backbone root-mean-square deviation (rmsd) 0.63 and 1.09 Å for the N- and C-terminal domains, respectively). The relative orientation of the two domains in γ S crystallin is also very similar to that seen in other crystallin structures, but the two domains are farther apart in the NMR structure, presumably as a result of the scarcity of interdomain restraints. This situation is encountered more frequently, in particular in protein-protein complexes, and in larger proteins where interdomain NOEs tend to be relatively sparse, but relative orientations of domains are tightly defined by RDCs.^{36,37} Therefore, the SAXS data present an ideal complement for determining an accurate solution structure of such systems.

The SAXS data for γ S crystallin at 4.5 mg/mL protein concentration were minimally affected by aggregation, as determined by the linearity of the Guinier plot (see Supporting Information) and $P(r)$ analysis. The latter yields a gyration radius (R_g) value of 18.3 ± 0.2 Å, a maximum linear dimension (D_{max}) of 54–57 Å, and an estimated molecular volume of $(25.2 \pm 0.7) \times 10^3$ Å³, approximated from the total intensity under the measured scattering profile and using the Porod invariant.³⁸ The same parameters determined using the 1AMM, 1A7H, and

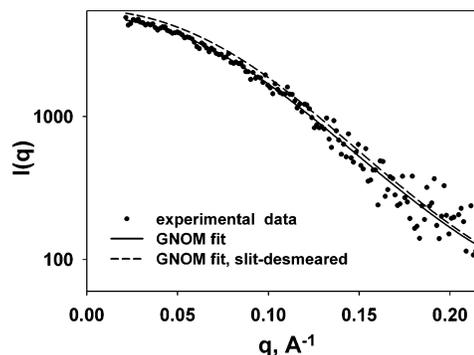


Figure 2. Experimental scattering data recorded for the 4.5 mg/mL γ S crystallin sample. The solid line shows regularized data fit from the GNOM program. The dashed line corresponds to the slit-desmeared data fit. A total of 16 points of this curve, equally spaced between ~ 0.02 and ~ 0.22 Å⁻¹, are subsequently used for structure calculation.

1HK0 crystal structures and the program CRY SOL²¹ are $R_g = 16.6$ – 16.8 Å, $D_{max} = 55.2$ – 56.5 Å, and a molecular volume of $(25.5$ – $25.8) \times 10^3$ Å³. The observed difference in R_g is likely to be a consequence of a thin surface layer of solvent with a density higher than that of the bulk solvent, a phenomenon often leading to an increase of the apparent SAXS-extracted R_g values by 1–2 Å with respect to the numbers calculated from the atomic coordinates. A weak tail is seen in the $P(r)$ distribution that appears to have a D_{max} of ca. 80 Å, which likely reflects a small amount of dimerized protein in the sample volume. The presence of seven reduced Cys residues in γ S crystallin, of which surface-exposed Cys²⁴ and Cys²⁶ are particularly reactive, promotes dimerization and formation of higher-order multimers under oxidizing conditions. $I(0)$ analysis of the data, using lysozyme as a standard, indicates that the dimers account for less than 8.5% of the total protein. The raw data as well as the regularized GNOM fits are shown in Figure 2. Even though the recorded scattering intensity extends up to 0.32 Å⁻¹, the uncertainty in our data precludes interpretation beyond about 0.22 Å⁻¹. The increased uncertainty is due in part to the fact that the SAXS instrument used has a one-dimensional detector and hence captures an increasingly smaller percentage of the solid angle of the circularly averaged scattering pattern at larger angles; a much higher signal-to-noise ratio can be attained using a synchrotron source coupled with an area detector, providing the sample can withstand the high radiation levels.

A total of five cycles of structure refinement were necessary to make globular and surface solvent layer corrections consistent with the ensemble of refined structures. The density of the bound solvent layer, assumed to be 3.5 Å thick, was determined from CRY SOL fits to be 0.025 e/Å³ higher than the bulk solvent density, which is within the expected range for a typical protein in solution.

The accuracy of the atomic coordinates of the refined models was evaluated with respect to the high-resolution X-ray structures of γ B, γ D, and C-terminal γ S crystallins (PDB entries 1AMM, 1A7H, and 1HK0). The γ B and γ D crystallins share ca. 50% sequence identity with γ S crystallin, and 74% with one another. With a two-domain backbone rmsd of 0.69 Å, the crystal structures of γ B and γ D crystallins exhibit very close similarity, despite crystallization in two different space groups. When comparing relative domain positions in γ B and γ D (keeping their N-terminal domains superimposed), the orientations of their C-terminal domains differ primarily by a 5.5°

(33) Kumaraswamy, V. S.; Lindley, P. F.; Slingsby, C.; Glover, I. D. *Acta Crystallogr. D, Biol. Crystallogr.* **1996**, *52*, 611–622.

(34) Basak, A. K.; Bateman, O.; Slingsby, C.; Pande, A.; Asherie, N.; Ogun, O.; Benedek, G. B.; Pande, J. *J. Mol. Biol.* **2003**, *328*, 1137–1147.

(35) Basak, A. K.; Kroone, R. C.; Lubsen, N. H.; Naylor, C. E.; Jaenicke, R.; Slingsby, C. *Protein Eng.* **1998**, *11*, 337–344.

(36) Clore, G. M. *Proc. Natl. Acad. Sci. U.S.A.* **2000**, *97*, 9021–9025.

(37) Tugarinov, V.; Choy, W. Y.; Orekhov, V. Y.; Kay, L. E. *Proc. Natl. Acad. Sci. U.S.A.* **2005**, *102*, 622–627.

(38) Glatter, O.; Kratky, O. *Small-Angle X-ray Scattering*; Academic Press: New York, 1982.

Table 1. Impact of Inclusion of SAXS Data as Restraints during Structure Calculation

	no SAXS data	with SAXS data
backbone rmsd to 1AMM, Å		
N-terminal domain (6–85)	0.63 ± 0.05	0.56 ± 0.05
C-terminal domain (94–175)	1.09 ± 0.09	0.90 ± 0.04
both domains (6–85, 94–175)	1.96 ± 0.07	1.31 ± 0.04
backbone rmsd to 1HK0, Å		
N-terminal domain (6–85)	0.70 ± 0.05	0.63 ± 0.05
C-terminal domain (94–175)	1.13 ± 0.08	0.95 ± 0.04
both domains (6–85, 94–175)	1.89 ± 0.08	1.18 ± 0.05
backbone rmsd to 1A7H, Å (94–175)	1.07 ± 0.08	0.87 ± 0.05
rmsd to mean, Å		
backbone atoms (6–85, 94–175)	0.26 ± 0.07	0.25 ± 0.04
all heavy atoms (6–85, 94–175)	0.83 ± 0.06	0.82 ± 0.05
Procheck Ramachandran statistics, %		
most favored	89.7	89.0
allowed	9.7	10.8
generous	0.6	0.2
steric clashes/100 residues	2.0 ± 1.1	4.4 ± 1.3
χ of SAXS data fit	1.1 ± 0.1	0.25 ± 0.02

rotation and exhibit no detectable translation. The packing at the hydrophobic interface in the homodimer of the C-terminal γ S crystallin domain is similarly tight, but shows a 23° rotation relative to γ B. In contrast, in our previously determined solution structure of γ S crystallin, the backbone rmsd relative to γ B and γ D is dominated by translation, not by relative domain orientation, and presumably results from insufficient interdomain NOE restraints.¹⁶ Therefore, this backbone rmsd presents a reasonable measure for the error in the relative position of the two domains of γ S crystallin.

Table 1 lists the values of the backbone rmsd for the ordered regions of the protein, comprising residues 6–85 and 94–175. It is clear from the data presented in Table 1 that inclusion of the SAXS data in the refinement brings on a considerably better agreement between the NMR structure of γ S and X-ray structures of γ B, γ D, and γ S crystallins. Inspection of the structure shows that the SAXS data fit results in the predicted tighter packing of the two domains with respect to each other. Remarkably, however, inclusion of the SAXS data also results in a small lowering of the individual domain rmsd values relative to the homologous X-ray structures of γ crystallins, even though these were already quite small without SAXS refinement. This result suggests that even for smaller, globular systems, SAXS data can improve the quality of NMR structures due to the constraint it provides on the overall molecular shape. Relative orientations of the two domains in the family of the calculated γ S structures are also rather similar to those in the X-ray structures. With the N-terminal domain of γ S NMR structure again best-fitted to the N-terminal domain of γ B crystallin, the orientation of the C-terminal domain differs by an \sim 11.5° “twisting” rotation about an axis that deviates by 20° from the long axis of the molecule. Relative to the γ D crystallin structure, the corresponding rotation is only 7°, resulting in the slightly lower two-domain backbone rmsd values relative to 1HK0 (Table 1). The small differences in orientation result primarily from the RDC restraints, and are minimally affected by the SAXS data: inclusion of the SAXS data in the structure refinement affects the relative domain orientation in γ S crystallin by less than 1°. This result highlights the complementary nature of SAXS and RDC restraints. The calculated family of structures

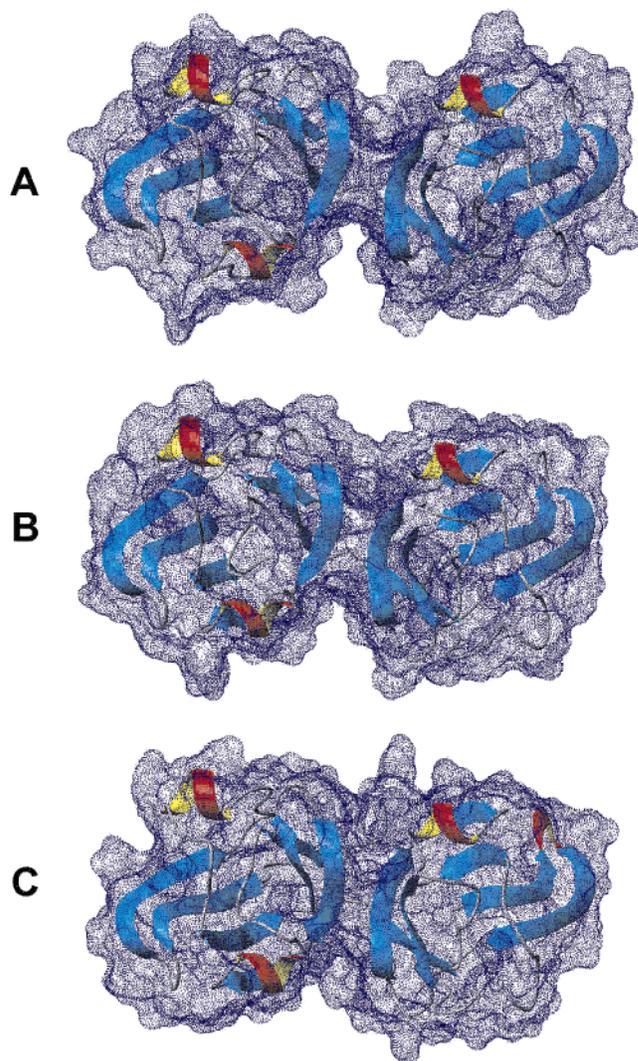


Figure 3. Impact of the SAXS data fit on the overall geometry. Protein backbones are shown in ribbon representation, and the molecular surfaces are calculated by sliding a 1.4-Å radius sphere over the molecule. (A) A representative model before SAXS data fit, (B) a representative model after SAXS data fit, and (C) the X-ray structure of γ B crystallin (1AMM) used to evaluate the accuracy of the atomic coordinates. The figure was generated using the program MOLMOL.⁴⁴

is deposited into the RSCB Protein Data Bank with the accession number 2A5M.

The impact of including the SAXS data in the structure calculation is illustrated in Figure 3. Incorporation of the scattering data clearly has the effect of bringing the two domains together, closer to their relative position in the 1AMM and other X-ray models. The same effect can be seen by comparing the interatomic distance distribution curves before and after SAXS refinement shown in Figure 4. When the SAXS data are not included in the structure calculation, the $P(r)$ distribution of the NMR structure is typical for a well-separated two-domain system, and shows a shoulder around 33 Å, roughly corresponding to the separation between the centers of the two domains. Inclusion of the experimental scattering data results in a more globular distribution, removing the “neck” separating the two domains. The scattering profile remains noticeably asymmetric, characteristic of a prolate ellipsoid shape. A substantial difference between the two calculated curves underscores the information content in the experimental SAXS

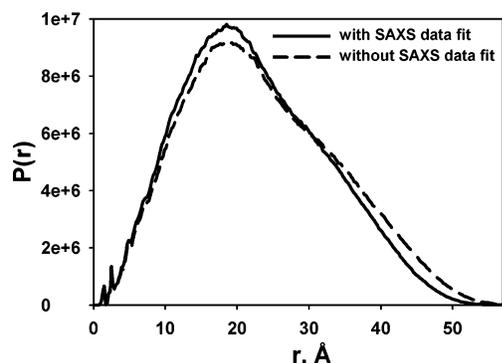


Figure 4. Impact of inclusion of SAXS data in the structure calculation on the distribution of the interatomic distances in the obtained structures. Each distribution is weighted by the products of the atomic numbers for each atom pair.

data, notwithstanding their modest signal-to-noise ratio and resolution range. Interestingly, inclusion of the SAXS data in the structure calculation does not simply push the N- and C-terminal domains as close as possible. Beta strands $\beta 6$ of the N-terminal and $\beta 14$ of the C-terminal domain remain separated by a distance that is not short enough to form backbone–backbone hydrogen bonds. Instead, the relative position of residues Met⁵⁸ and Gln¹⁴⁸, located near the closest point of interdomain backbone–backbone approach, suggests the presence of a water-separated pair of backbone–backbone hydrogen bonds, and residues Gln¹⁴⁸ and Ile⁶⁰ are likely to be linked by a backbone–side chain hydrogen bond Gln¹⁴⁸:O^{ε2}–Ile⁶⁰:H^N, as seen in other crystallin structures. A comparison between the $P(r)$ distribution of the NMR family of structures and those from the available X-ray models indicates that crystal structures are slightly more compact and globular in shape, consistent with their smaller gyration radii (see the Supporting Information).

To test the dependence of the calculated structures on the number of SAXS points being fitted, N_q , we have doubled it within the same resolution interval, i.e., reduced the sparsening of the original data. The resulting structures are very similar to the 2A5M bundle, with the backbone rmsd to the mean of 2A5M in the 0.2–0.3 Å range. This result indicates that detrimental effects on the structural quality produced by our sparsening of SAXS data are negligible.

In addition to using the NMR restraints listed above, we have also investigated several other model scenarios. First, we have completely removed all NMR-derived distance restraints between the C- and N-terminal domains. Perhaps surprisingly, when SAXS data are used in the refinement, removal of the interdomain NOE restraints has a negligible effect on the difference between the obtained structure and the structure of γB crystallin. This result underscores the utility of the SAXS data for relative positioning of the two domains. On the other hand, the rmsd values relative to the X-ray structures for the C-terminal domain alone increase slightly as a result of the loss of the specific interdomain connectivity restraints. This latter result reflects small structural changes in the vicinity of residues 151–155 and 131–135, for which no backbone amides could be observed due to intermediate time scale conformational exchange, and for which therefore no dipolar or backbone torsion restraints were available. In contrast, when no SAXS data are included, removal of the interdomain distance restraints results in substantial translations of the (oriented) domains

relative to one another, and considerable deviations from the homologous γB and γD crystallin structures (see the Supporting Information).

To further test the limits of how uniquely the SAXS data define the relative domain positions, we have, in addition to removing any interdomain distance constraints, also severed the C–N bond between residues 89 and 90, along with all its associated bonds, angles, and stereochemical restraints in the empirical force field used for the structure calculation. As a result, the relative translational position of the two domains is completely unspecified by either the NMR data or the chain connectivity, while their relative orientation remains tightly defined by the dipolar couplings from the two alignment media. Such a scenario simulates the case of docking of a tight complex between two independent macromolecular entities, based on a combination of RDCs and SAXS data only. The starting geometries for the MD runs contained the C-terminal domain randomly translated on a sphere of 50 Å radius around the N-terminal one. Our results, outlined in detail in the Supporting Information, illustrate the limitations inherent in our data: the quality of the final fit is only weakly correlated with the backbone rmsd to the γB reference structure. In our case, the shape of the individual domains is too globular, i.e., has insufficient unique features to unambiguously establish the correct solution from the scattering data at hand. Higher signal-to-noise ratios and/or higher resolution may aid such discrimination. For cases where the individual domains are less symmetric than those of γS crystallin, one also may expect the scattering data to be more successful in at least limiting the potential solution set.

To gauge the dependence of our results on the amount of the intradomain information input, we have repeated all structure calculations while also including additional distance restraints for the 89 hydrogen bonds that could be determined in a consensus manner by sequence alignment to the homologous 1AMM, 1HKO, and 1A7H structures. As expected, addition of the corresponding H-bond restraints results in a further decrease of the backbone rmsd of both individual domains and the two-domain construct with respect to all X-ray models (see Supporting Information for details).

We have also investigated the effect of the decrease in the amount of the orientational NMR restraints on the quality of structures resulting from SAXS data fit. In one such test, we have deactivated all RDC restraints and compared the accuracy of the resulting coordinates with and without SAXS data fitted. In a second test, we deactivated all RDC restraints except for the N–HN RDCs from one alignment medium (gelled Pf1). The results, outlined in detail in the Supporting Information, show the importance of having at least one set of dipolar couplings in addition to the SAXS data as these restraints are crucial for the correct positioning of the two domains.

It is perhaps interesting to consider the information content of the SAXS data alone, in the absence of any NMR data. SAXS data alone clearly provide insufficient restraints for independent structure determination. However, we have attempted fold recognition instead, by submitting our experimental SAXS data to the server DARA,³⁹ which takes the SAXS scattering profile and the molecular mass of the protein as input. The server

(39) Sokolova, A. V.; Volkov, V. V.; Svergun, D. I. *J. Appl. Crystallogr.* **2003**, *36*, 865–868.

returned the 10 highest-scoring hits out of 223 searched in the range of molecular masses from 19.3 to 22.3 kDa (PDB codes 1AMM, 1AWD, 1DJ7, 1E7N, 1G8Q, 1MJS, 1N0Q, 1SPH, 2GCR, and 262L). Remarkably, 3 among those 10 were members of the crystallin family (PDB codes 1AMM, 1E7N, and 2GCR).

The dependence of small-angle scattering intensity on the square of the molecular weight of the scattering particle results in a scattering profile that is quite sensitive to small amounts of aggregation. In contrast, NMR is relatively insensitive to minor degrees of aggregation in the sample. Thus, combining NMR and scattering data could be problematic if the procedure were intolerant to even weak degrees of self-association. Considering that γ S crystallin has a tendency to self-associate, as judged by the steeper than expected increase in rotational correlation time with volume fraction, and to form covalent homodimers through oxidation of the solvent-exposed Cys²⁴ and Cys²⁶ residues, it presents a challenging case for SAXS refinement. Therefore, the fact that we obtained a considerable improvement in structural accuracy for this rather challenging system bodes well for the future utility of this technique. It is also encouraging that significant gains in structural accuracy can be made even with the relatively modest statistical quality of our SAXS data, which were obtained using a simple laboratory-based instrument that uses a sealed tube X-ray source. Scattering profiles extending to much higher angles and at much higher signal-to-noise ratios can be recorded at synchrotrons for favorable systems, such as larger proteins and nucleic acids.⁴⁰

Our structure refinement procedure is based on the assumption of a single, well-defined conformation. However, it is important to bear in mind that SAXS data represent an average over all conformations sampled by the molecule in solution. In the application to γ S crystallin, the assumption of a single well-defined conformation is supported by a variety of NMR data, including ¹⁵N backbone dynamics measurements and the indistinguishable values of the alignment tensors of the two domains. However, there is no a priori reason that prevents application of the SAXS refinement procedure to a multiconformer refinement of a more dynamic complex.

Another issue of potential interest is whether including SAXS data in the refinement, as done in the current study, has any advantages over calculating a family of structures and then selecting from these the subset with the lowest χ^2 of the SAXS data fit, a task that can easily be performed with existing software.^{13,14} We have generated a family of 166 structures without inclusion of SAXS data, and evaluated SAXS χ^2 on those models (see Supporting Information). Our results indicate that, while selection by the lowest SAXS χ^2 will lower the rmsd to 1AMM, the decrease is considerably smaller than when the SAXS data are fitted directly. This outcome results in part from the commonly used “repulsive-only” nonbonded interactions, and underscores the limitations in providing sufficient sampling of conformational space during the structural refinement, which can be overcome by including the SAXS data as restraints in the structure calculation.

Concluding Remarks

In this study we have demonstrated the utility of solution X-ray scattering data as a component of high-resolution NMR

structure refinement. The obtained improvements in accuracy are very encouraging, particularly given the limited effective resolution range of only up to ~ 30 Å spanned by our acquired scattering data. SAXS data present an ideal complement to NMR data sets rich in orientational restraints, such as those contained in residual dipolar couplings, but lacking a large number of accurate translational restraints, such as NOEs. Use of the SAXS data clearly will be most advantageous for defining the solution structure of larger macromolecules, where the number of restraints per residue tends to be sparse, but where dipolar couplings are still readily accessible. Higher informational content within the same resolution range and higher signal-to-noise ratios for SAXS data when applied to these systems is well suited to offset the decrease of the density of the NMR-based structural constraints.^{41–43}

To date, the usage of SAXS data in structural biology has mainly been limited to (i) de novo low-resolution shape reconstruction, (ii) testing previously derived high-resolution structural models, and (iii) rigid-body refinement of multiunit macromolecular assemblies. With the substantial improvements in the formalism connecting the observed data to the underlying structural model that has occurred in the past few years, this situation is likely to change. The direct fitting approach described in the current study is intended to facilitate a more routine usage of this key data source during macromolecular structure refinement.

Acknowledgment. We thank Peter Flynn (University of Utah) for the use of his biochemistry laboratory facilities for on-site sample preparation for the SAXS experiments. This research was supported by the Intramural Research Program of the NIDDK, NIH, and by the Intramural Antiviral Target Program of the Office of the Director, NIH.

Supporting Information Available: A table defining the “globs” used to represent the molecular structure; a Guinier plot; a figure with the results of the refinement with no interdomain chain connectivity; tables with structural statistics for additional test cases where the number and nature of NMR restraints are varied; a figure showing the difference between the results of the joint NMR–SAXS refinement and the procedure in which the SAXS data would be used to filter the family of structures generated from NMR data alone; figures showing the dependence of the globbic scattering curve on the protein parameters and the size of the globs; a figure comparing scattering curves from all-atom and globbic calculations; a figure showing an all-atom fit of the representative final structure to the experimental SAXS data; and a figure showing $P(r)$ distributions calculated from the NMR family of structures compared with those of the X-ray structures of γ -crystallins (PDF). This material is available free of charge via the Internet at <http://pubs.acs.org>. Programs for calculation of the globbic form factors and globbic correction factors as well as source code for SAXS data fitting routines in CNS and Xplor-NIH are available from the authors.

JA054342M

- (41) Tugarinov, V.; Kay, L. E. *J. Mol. Biol.* **2003**, *327*, 1121–1133.
- (42) Lukin, J. A.; Kontaxis, G.; Simplaceanu, V.; Yuan, Y.; Bax, A.; Ho, C. *Proc. Natl. Acad. Sci. U.S.A.* **2003**, *100*, 517–520.
- (43) Williams, D. C.; Cai, M. L.; Clore, G. M. *J. Biol. Chem.* **2004**, *279*, 1449–1457.
- (44) Koradi, R.; Billeter, M.; Wuthrich, K. *J. Mol. Graph.* **1996**, *14*, 51–55.

(40) Zuo, X. B.; Tiede, D. M. *J. Am. Chem. Soc.* **2005**, *127*, 16–17.