

Bayesian multi-state multi-condition modeling of a protein structure from X-ray crystallography

Matthew Hancock^{a,c,1}, James S Fraser^a, Paul D Adams^a, and Andrej Sali^{b,1,2}

This manuscript was compiled on June 25, 2024

A model of a protein structure at atomic resolution is key to rationalizing and predicting its biological function. Many such models are computed from a diffraction pattern from X-ray crystallography. Despite the protein crystal containing billions of protein molecules that independently sample the energy landscape during data collection, most models computed from X-ray data depict a single set of atomic coordinates. A model with multiple sets of atomic coordinates (multi-state) may improve the satisfaction of the X-ray data and is a more accurate, precise, and informative depiction of the protein. However, computing a multi-state model is challenging on account of a low data-to-parameter ratio. X-ray datasets collected for the same system under distinct experimental conditions (eg, temperature) may provide additional observations, thereby improving the data-to-parameter ratio. Here, we develop, benchmark, and illustrate MultiXray: Bayesian multi-state multi-condition modeling for X-ray crystallography. The input information is J X-ray datasets collected under distinct conditions and a molecular mechanics force field. The model representation is N independent atomic models of the protein structure (states) and the weight of each state under each condition. A Bayesian posterior model density quantifies the match of the model with all X-ray datasets and molecular mechanics. A sample of models is drawn from the posterior model density using molecular dynamics simulations. We benchmark MultiXray on simulated X-ray data and analyze the impact of additional states and conditions on the scoring function and model search. We illustrate MultiXray on temperature-dependent X-ray datasets collected for SARS-CoV-2 M^{pro} and compute multi-state multi-condition models that improve the R^{free} relative to the PDB model by up to 0.05. MultiXray is implemented in our open-source *Integrative Modeling Platform*, relying on integration with *Phenix*, thus making it easily applicable to other systems.

Keyword 1 | Keyword 2 | Keyword 3 | ...

A. A protein crystal is a heterogeneous mix of structural states. X-ray crystallography is an important experimental technique for obtaining structural models of proteins at atomic resolution. In X-ray crystallography, the protein crystal contains between $10^6 - 10^{15}$ copies of the protein (1–3). Due to the high solvent content of the crystal, protein molecules within the crystal can fluctuate nearly independently throughout data collection and adopt distinct structural states based on the energy landscape (4–6). Despite the resulting structural heterogeneity, a majority of models computed from X-ray datasets describe a single structural state with Gaussian isotropic B-factors (7). Fitting a B-factor to an X-ray dataset convolutes all sources of experimental uncertainty, including the heterogeneous mix of structural states, random thermal motion, and long-range crystallographic disorder (8, 9). The inability of a B-factor to describe the heterogeneous mix of structural states found within the protein crystal contributes to the inability of protein models to satisfy an X-ray dataset within its theoretical uncertainty (10, 11). A model that depicts anharmonic conformational substates found in a protein crystal will improve the satisfaction of X-ray data and reveal the structural basis of important biological properties at atomic resolution, such as allosteric networks or hidden cavities for small molecule binding (cryptic pockets) (12, 13).

B. Approaches to modeling a heterogeneous mix of structural states. There are 2 approaches for computing models that depict multiple structural states from X-ray data. First, conformational substates can be captured by computing single-state models that independently satisfy the X-ray dataset, as is seen in the modeling of Nuclear Magnetic Resonance spectra (14). For example, phenix.ensemble_refinement computes an ensemble from snapshots of a molecular simulation restrained by a time-averaged X-ray target function (15). Such approaches may not find weakly occupied states in a reasonable amount of computation time, as they are dependent on

Significance Statement

Authors must submit a 120-word maximum statement about the significance of their research paper written at a level understandable to an undergraduate educated scientist outside their field of specialty. The primary goal of the significance statement is to explain the relevance of the work in broad context to a broad readership. The significance statement appears in the paper itself and is required for all research papers.

Author affiliations: ^aAffiliation One; ^bAffiliation Two; ^cAffiliation Three

Please provide details of author contributions here.

Please declare any competing interests here.

²To whom correspondence should be addressed. E-mail: salis@ilab.org

overcoming potentially large barriers in the scoring function. Alternatively, a model can depict conformational substates by introducing additional structural variables (3). All variables are then collectively fit against the X-ray data. The extent and detail of the additional degrees of freedom depend on the available computational power and data quality (16). For example, qFit-3 avoids introducing excessive structural parameters by representing side chains as ensembles of one or more rotameric states (17). Methods that refine multiple fully parameterized atomic models have been limited to systems that diffract to ultra-high resolution (6, 18–21).

C. Multi-condition crystallography. Often, multiple X-ray datasets are collected for the same system under distinct experimental conditions. One example is multi-temperature crystallography, where data collection is performed at temperatures from cryogenic to near-physiological (22). Data collection at higher temperatures has been shown to dramatically modify protein dynamics within the protein crystal (23–25). More recently, models computed from higher temperature datasets have revealed a fuller set of structural states within the protein’s energy landscape at atomic resolution (26–28). Comparisons of models of the same system computed at distinct temperatures show similar structural states but shifted thermodynamic equilibrium (26, 29, 30). Therefore, multiple X-ray datasets under distinct conditions containing mutual structural information could increase the data-to-parameter ratio and inform a more accurate, precise, and complete multi-state model.

D. Computing a multi-state multi-condition model. Here, we seek to compute a multi-state model from multiple X-ray diffraction patterns collected under distinct conditions (1a). We can fit a multi-state model using multiple X-ray datasets without needing ultra-high-resolution X-ray data. We accommodate for the thermodynamic shift of distinct experimental conditions by computing each state’s weight under each condition. In other words, we assume the sets of structural states are consistent across all experimental conditions with varied weights, allowing a single multi-state model to be informed by all X-ray datasets and massively boosting the data-to-parameter ratio. The assumption of a consistent set of structural states under all conditions is not limiting because the structural states may include sparsely populated or completely unpopulated states. We formulate a Bayesian posterior model density for the multi-state multi-condition model. We draw a sample from the Bayesian posterior model density using a molecular dynamics algorithm where all structural states are jointly restrained by the satisfaction of all X-ray datasets in addition to molecular mechanics. We benchmarked the method using synthetic X-ray datasets simulated from multi-state SARS-CoV-2 M^{Pro} structural models and showed that multiple structural states are needed to satisfy the data and that by including multiple X-ray datasets, all datasets are better satisfied individually. We illustrate our method to compute a multi-state model of the SARS-CoV-2 viral target M^{Pro} from multi-temperature X-ray under 6 conditions and show that all datasets are satisfied better by increasing the number of states and X-ray datasets.

1. Methods

A. Overview of modeling method. A model is a depiction of our knowledge about a system or process. We wish to inform the model based on the input information, which generally includes experimental data and prior models (eg, physical theories and statistical preferences). A model can then rationalize past observations and predict future ones. Modeling is the search for a set of models consistent with the input information. We aim to find all models that satisfy the input information, reflecting the uncertainty of the input information and the modeling process. It is convenient to divide modeling into three steps: (i) specifying all model variables (representation), (ii) ranking alternative models by their agreement with the input information (scoring), and (iii) generating a sample of good-scoring models (sampling). A model should be validated before being interpreted, often by examining the satisfaction of input information withheld from the modeling process. Multiple iterations of gathering input information, modeling, and validation are often necessary to compute a sufficiently precise model (31, 32).

B. Input information. The input information may be used to inform any step of modeling: representation, scoring, and sampling. Here, we use the following input information: 1) X-ray datasets D_1, \dots, D_J collected under distinct experimental conditions (eg, at distinct temperatures) and corresponding PDB models. Each X-ray dataset is a set of observed structure factor amplitudes $\{|F_{\vec{s}}^O|\}_{\vec{s} \in S}$ indexed by a scattering vector \vec{s} . The X-ray datasets are used in scoring. The PDB model is used in scoring and sampling. 2) The CHARMM19 force field parameters (33). The force field parameters are used in scoring to evaluate the prior.

C. Representation. The representation defines the degrees of freedom whose values are to be determined by modeling. The multi-state multi-condition model M includes the following variables: (i) N states M_i (N is selected before modeling). Each state is an independently parameterized model of a protein structure. The number and type of atoms (composition) of each state are based on the PDB models (Input Information). Only non-hydrogen atoms from the protein are included, however, the representation may be generalized to include additional atoms/molecules (eg, hydrogen, solvent, ion, and ligand). As we are computing a model containing a set of discrete structural states, all atomic B-factors are set to 15 and atomic occupancies to 1. (ii) Weight matrix $W_{N \times J}$ containing the weight of M_i under condition j , w_{ij} . The weight matrix is subject to the constraint that the sum of the weights of all states under each condition (ie, columns) is 1. (iii) Nuisance variables, taken from established modeling methods, to improve the fit of the X-ray datasets by the model (34, 35). For clarity, we do not include nuisance parameters in the notation.

D. Scoring.

D.1. Bayesian posterior model density. A scoring function assesses the match of the model to the input information. Bayesian inference is one approach to assess this fit. In Bayesian inference, the posterior model density $p(M|D, I)$ describes the relationship between model M , prior information I , and data D . The posterior model density may be factored

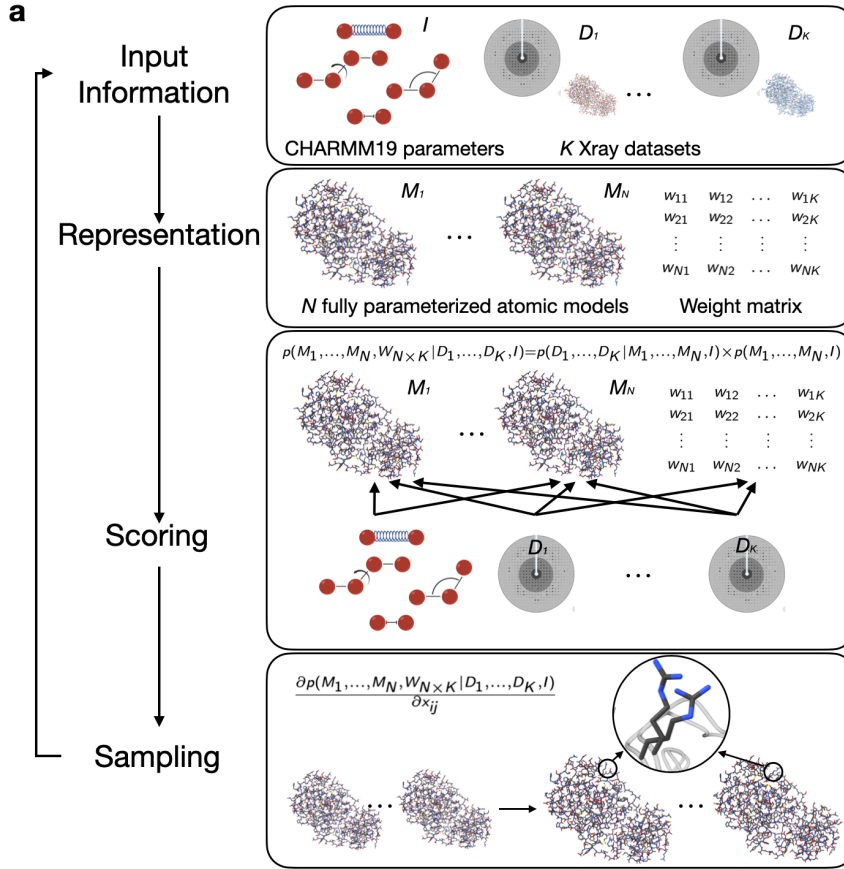


Fig. 1. modeling can be framed as a model search given some input information. Here, the input information is the CHARMM19 force field parameters and J X-ray datasets collected for the same system under distinct experimental conditions (eg, temperature). The representation is the N atomic states containing all the heavy atoms of the system along with the weight matrix that parameterizes the weight of each state under each condition. All states and the weight matrix are scored collectively against each X-ray dataset. Each state is individually scored against the molecular mechanics force field. A sample is drawn from the posterior model density using molecular dynamics. All states are initialized by a starting structure and the force on the atoms is computed from the satisfaction of all X-ray datasets along with the molecular mechanics, and the weights are stochastically sampled.

into a prior $p(M|I)$ that quantifies the state of knowledge prior to the observation of data and a likelihood $p(D|M, I)$ that quantifies what is learned from the observation of data (36):

$$p(M|D, I) = p(D|M, I) \times p(M|I) \quad [1]$$

It is often helpful to decompose the likelihood into a forward model $f(M)$ that simulates a noiseless data observation and a noise model $N(f(M); D, \sigma)$ that quantifies the difference between the observed and simulated data based on a model of the experiment's noise parameterized by σ (31).

D.2. Joint-likelihood. The likelihood of observing all X-ray diffraction patterns from a multi-state multi-condition model is:

$$p(D_1, \dots, D_J | M_1, \dots, M_N, W_{N \times J}) \quad [2]$$

The likelihood may be factored into the likelihood of observing each condition's dataset:

$$\begin{aligned} p(D_1, \dots, D_J | M_1, \dots, M_N, W_{N \times J}) \\ = \prod_{j=1}^J p(D_j | M_1, \dots, M_N, W_{N \times j}) \end{aligned} \quad [3]$$

The likelihood for each condition's dataset depends only on the set of states and the weights representing the states under the corresponding condition (matrix column):

$$\begin{aligned} p(D_j | M_1, \dots, M_N, W_{N \times j}) \\ = p(D_j | M_1, \dots, M_N, w_{1j}, \dots, w_{Nj}) \end{aligned} \quad [4]$$

The likelihood of an X-ray dataset is factored into the likelihood of observing each structure factor amplitude F_s^O (35):

$$\begin{aligned} p(D_j | M_1, \dots, M_N, w_{1j}, \dots, w_{Nj}) = \\ \prod_{\vec{s} \in S} p(F_s^O | M_1, \dots, M_N, w_{1j}, \dots, w_{Nj}) \end{aligned} \quad [5]$$

The likelihood of a structure factor amplitude given a multi-state multi-conditional model is a noise model that quantifies the difference between the observed structure factor amplitude F_s^O and the model structure factor amplitude F_s^M simulated by the forward model from the multi-condition multi-state model (35):

$$p(F_s^O | M_1, \dots, M_N, w_{1j}, \dots, w_{Nj}, I) = p(F_s^O | F_s^M) \quad [6]$$

The noise model is based on the assumption that the real and imaginary parts of a complex structure factor \mathbf{F} are sampled from a two-dimensional Gaussian with variance $\epsilon\beta$ and scale parameter $\alpha_{\vec{s}}$ (35):

$$\begin{aligned} p(\text{Re}\mathbf{F}, \text{Im}\mathbf{F}) = \frac{1}{\pi\epsilon\beta_{\vec{s}}} \exp(-(\text{Re}\mathbf{F} - \alpha_{\vec{s}}F^M \cos\phi^M)^2 \\ + (\text{Im}\mathbf{F} - \alpha_{\vec{s}}F^M \sin\phi^M)^2)/\epsilon\beta_{\vec{s}} \end{aligned} \quad [7]$$

The likelihood for a single structure factor amplitude is (35):

$$p(F_{\vec{s}}^O | F_{\vec{s}}^M, \alpha_{\vec{s}}, \beta_{\vec{s}}) = \begin{cases} \frac{2F_{\vec{s}}^O}{\epsilon\beta_{\vec{s}}} \exp\left(-\frac{(F_{\vec{s}}^O)^2 + \alpha_{\vec{s}}^2 (F_{\vec{s}}^M)^2}{\epsilon\beta_{\vec{s}}}\right) \times \\ I_0\left(\frac{2\alpha_{\vec{s}} F_{\vec{s}}^O F_{\vec{s}}^M}{2\epsilon\beta_{\vec{s}}}\right) & \dots \text{if } \vec{s} \text{ acentric} \\ \left(\frac{2}{\epsilon\pi\beta_{\vec{s}}}\right)^{1/2} \exp\left(-\frac{(F_{\vec{s}}^O)^2 + \alpha_{\vec{s}}^2 (F_{\vec{s}}^M)^2}{2\epsilon\beta_{\vec{s}}}\right) \times \\ \cosh\left(\frac{2\alpha_{\vec{s}} F_{\vec{s}}^O F_{\vec{s}}^M}{2\epsilon\beta_{\vec{s}}}\right) & \dots \text{if } \vec{s} \text{ centric} \end{cases} \quad [8]$$

where I_0 is the hyperbolic Bessel function of the first kind ($\alpha = 0$) and \cosh is the hyperbolic cosine function.

D.3. Forward model. The forward model $F_{\vec{s}}^M$ is the magnitude of the complex model structure factor $\mathbf{F}_{\vec{s}}^M$:

$$F_{\vec{s}}^M = |\mathbf{F}_{\vec{s}}^M| \quad [9]$$

The model structure factor is the sum of the protein $\mathbf{F}_{\vec{s}}^C$ and bulk solvent structure factor $\mathbf{F}_{\vec{s}}^B$ with an overall scale factor k_{total} and bulk solvent scaling factor k_{mask} (37).

$$\mathbf{F}_{\vec{s}}^M = k_{\text{total}}(\mathbf{F}_{\vec{s}}^C + k_{\text{mask}}\mathbf{F}_{\vec{s}}^B) \quad [10]$$

It is well-established how to compute the scattering from a model of a protein structure assuming a perfect crystal (38). Here, we compute the scattering of a model of a protein structure with multiple weighted structural states. It is important to keep in mind that different structures and positions of molecules in a crystal can result in the same X-ray diffraction data, even when the data are noiseless, complete, and collected instantaneously. We show that the scattering in the 5 cases is equivalent (2ab).

The degeneracy between (i) and (ii) relates the scattering of a multi-state crystal to a set of single-state crystals. The structure factor for a multi-state model may be computed from the structure factor $\mathbf{F}_{i\vec{s}}$:

$$\mathbf{F}_{\vec{s}} = \sum_{i=1}^{\text{states}} \mathbf{F}_{i\vec{s}} \quad [11]$$

The degeneracy of (i), (iv), and (v) complicates the interpretation of a match between a multi-state model and an X-ray dataset. For example, imagine we compute an N -state model with N rotamers for residues A and B each. Without exhaustive sampling of all possible multi-state models, we do not know if the rotamer pair always occurs in the same molecule or is one of many possible combinations of the rotameric states in multiple molecules. In other words, we cannot rule out any combination of the rotameric states of A and B in one molecule. As an aside, the degeneracy of the X-ray data might be broken by using additional information for modeling; for example, a potential energy function may rule out combinations of rotameric states that result in overlapping atoms.

D.4. Prior. The prior for the multi-state multi-condition model is the product of a prior for each state:

$$p(M|I) = \prod_i^N p(M_i|I) \quad [12]$$

The prior for a state is the Boltzmann distribution corresponding to the potential energy of the state, including terms for bond lengths b , bond angles θ , dihedral angles ϕ , improper dihedral angles w , and non-bonded interactions r_{ij} :

$$p(M_i|I) = -\frac{1}{Q} \left(\exp \sum_{\text{bond}} k_b (b - b_0)^2 + \sum_{\text{ang}} k_{\theta} (\theta - \theta_0)^2 \right. \quad [13] \\ \left. + \sum_{\text{dih}} k_{\phi} [1 + \cos(n\phi - \delta)] + \sum_{\text{imp}} k_w (\omega - \omega_0)^2 \right. \\ \left. + \sum_{\text{nb}} \left(\epsilon \left[\left(\frac{R_{\text{min}_{ij}}}{r_{ij}} \right)^{12} - \left(\frac{R_{\text{min}_{ij}}}{r_{ij}} \right)^6 \right] \right) \right)$$

where the parameters k_b , b_0 , k_{θ} , θ_0 , k_{ϕ} , δ , k_w , w_0 , ϵ , and $R_{\text{min}_{ij}}$ are obtained from the CHARMM19 molecular mechanics force field (33) and Q is the partition function, which is ignored in practice.

D.5. Scoring function. A Bayesian scoring function for ranking alternative solutions based on the available information is:

$$S = -\log p(D|M, I) - \log p(M|I) \quad [14]$$

E. Sampling. The purpose of sampling is to find all models consistent with the input information. In Bayesian modeling, this is achieved by computing a sufficiently converged estimate of $p(D|M, I)$.

E.1. Atomic coordinates. The atomic positions for each state are sampled via molecular dynamics (MD). The atomic positions of a state are initialized to one of the input PDB models. The velocity of each atomic coordinate is sampled from a Boltzmann distribution with $T = 300\text{K}$. Because of the likelihood, the force on the atoms is non-conservative, and a thermostat is used to maintain a simulation temperature of 300K (15). We do not include explicit solvent molecules in the simulations, but the likelihood restrains each state to well-folded conformations. To account for discrepancies in the X-ray datasets, we reset the center of the mass to the center of mass of the corresponding PDB model before scoring against an X-ray dataset. The force computed on each atom is the partial derivative of the scoring function with respect to the atomic position of atom \vec{x} in state i , \vec{x}_i :

$$\frac{\partial S}{\partial \vec{x}_i} = - \left(\frac{\partial \log p(D_1|M_1, \dots, M_N, w_{11}, \dots, w_{N1}, I)}{\partial \vec{x}_i} + \dots \right. \quad [15] \\ \left. + \frac{\partial \log p(D_J|M_1, \dots, M_N, w_{1J}, \dots, w_{NJ}, I)}{\partial \vec{x}_i} \right. \\ \left. - \frac{\partial \log p(M_i|I)}{\partial \vec{x}_i} \right)$$

To calculate $\frac{\partial \log p(D_j|M_1, \dots, M_N, w_{1j}, \dots, w_{Nj}, I)}{\partial \vec{x}_i}$ more efficiently, a quadratic approximation of $p(F_{\vec{s}}^O | F_{\vec{s}}^M)$ is used (35). By computing the partial derivative of all atomic positions for a multi-state model, the atoms are restrained to satisfy a heterogeneous X-ray dataset collectively. Furthermore, all atoms of all states are restrained to satisfy all J X-ray datasets collectively.

Because the magnitudes of the derivative of the likelihood and prior can vary significantly, we introduce 2 scaling parameters w_{xray} and w_{auto} :

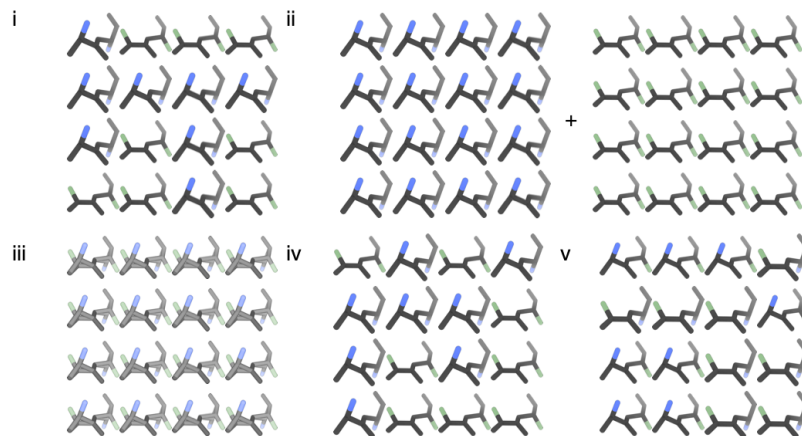
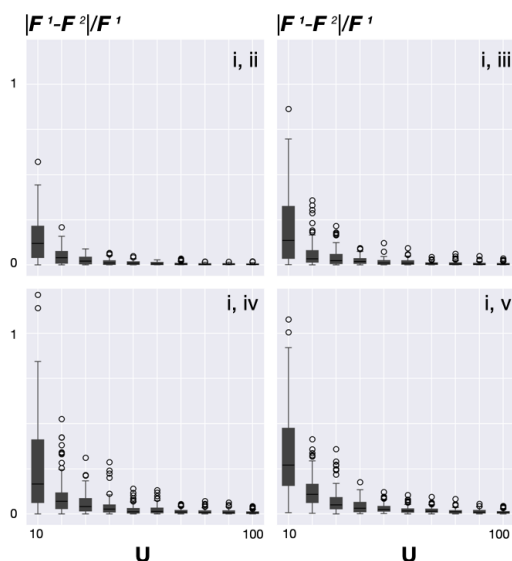
a**b**

Fig. 2. Five cases where the scattering is degenerate based on a multi-state model containing 2 states M_1 and M_2 with w_1 and w_2 . (i) Scattering from a heterogeneous crystal containing unit cells occupied by either M_1 or M_2 with probability w_1 and w_2 , respectively. (ii) A weighted combination of the scattering from homogeneous crystals containing identical unit cells with a single copy of M_1 and M_2 respectively. (iii) Scattering from a homogeneous crystal that contains the occupancy-weighted average of M_1 and M_2 (ie, the atomic occupancy of an atom from M_1 and M_2 is w_1 and w_2 , respectively). (iv) Scattering from a heterogeneous crystal the same as i where the states are re-indexed (ie, M_1 becomes M_2 and w_1 becomes w_2 and vice versa). (v) If $w_1 = w_2$, the scattering from a heterogeneous crystal except the states contain an intermixing of the atomic positions of M_1 and M_2 . **b**, as the number of unit cells goes to infinity, the reciprocal lattice points simulated from a crystal in ii, iii, iv, and v converge to those for the crystal in i. We construct 100 random reference crystals with U^3 unit cells where each unit cell is randomly populated by M_1 or M_2 with probability w_1 , w_2 respectively. M_1 , M_2 each contain 2 scatterers with random atomic position and scattering type. We compute the reciprocal lattice points up to a maximum miller index of $hkl=(3,3,3)$ via a discrete Fourier Transform (DFT). For (ii), (iii), (iv), and (v), we construct 100 random crystals as described and compute the same reciprocal lattice points. For each reciprocal lattice point, we compute the Euclidean distance for each structure factor to the reference structure factor U from 10 to 100. As U grows larger, the average Euclidean distance between the structure factors and reference structure factors converges to 0.

$$\begin{aligned} \frac{\partial S}{\partial \vec{x}_i} = & -w_{\text{xray}}(w_{\text{auto}}^0 \frac{\partial \log p(D_1|M_1, \dots, M_N, w_{11}, \dots, w_{N1}, I)}{\partial \vec{x}_i} + \dots \\ & [16] \\ & + w_{\text{auto}}^J \frac{\partial \log p(D_J|M_1, \dots, M_N, w_{1J}, \dots, w_{NJ}, I)}{\partial \vec{x}_i}) \\ & - \frac{\partial \log p(M_i|I)}{\partial \vec{x}_i} \end{aligned}$$

w_{auto}^j is computed automatically to ensure the average force on the atoms from the likelihood of D_j is approximately equal to that from the prior (21). To account for the presence of multiple X-ray datasets, w_{xray} is selected before modeling. We select $w_{\text{xray}} \in \{2^{-i}; -5 \leq i \leq 5\}$ by computing a trial sample using each w_{xray} and identifying the w_{xray} that results in a sample that best satisfies the data (R^{free}) after filtering out w_{xray} that produce non-physiological conformations (fig:fig3c).

We perform a short refinement minimization of $p(M_i|I)$ for each state of each model.

E.2. Weight matrix. For every 100 steps of sampling of the atomic positions, we propose 10 sets of weights per condition. A proposal weight is accepted for a condition if it improves the likelihood for the corresponding X-ray dataset. Each proposal is sampled from a Normal distribution with a mean equal to the current weight and a standard deviation of 0.05 and normalized.

E.3. Nuisance parameters. The nuisance variables are not stochastically sampled to enhance computational efficiency but optimized by minimizing the least squares residual between the model and the observed X-ray diffraction data (34).

F. Software availability. The software, input files, and output files are freely available at as part of our open-source *Integrative Modeling Platform* (<https://integrativemodeling.org/2.20.0/doc/manual/>) (39). The software relies on integration with *Phenix* (<https://phenix-online.org/>) (37), as described previously (40).

Digital Figures. EPS, high-resolution PDF, and PowerPoint are preferred formats for figures that will be used in the main manuscript. Authors may submit PRC or U3D files for 3D images; these must be accompanied by 2D representations in TIFF, EPS, or high-resolution PDF format. Color images must be in RGB (red, green, blue) mode. Include the font files for any text.

Images must be provided at final size, preferably 1 column width (8.7cm). Figures wider than 1 column should be sized to 11.4cm or 17.8cm wide. Numbers, letters, and symbols should be no smaller than 6 points (2mm) and no larger than 12 points (6mm) after reduction and must be consistent.

Figures and tables should be labelled and referenced in the standard way using the `\label{}` and `\ref{}` commands.

Figure 3 shows an example of how to insert a column-wide figure. To insert a figure wider than one column, please use the `\begin{figure*}...\end{figure*}` environment. Figures wider than one column should be sized to 11.4 cm or 17.8 cm wide. Use `\begin{SCfigure*}...\end{SCfigure*}` for a wide figure with side legends.

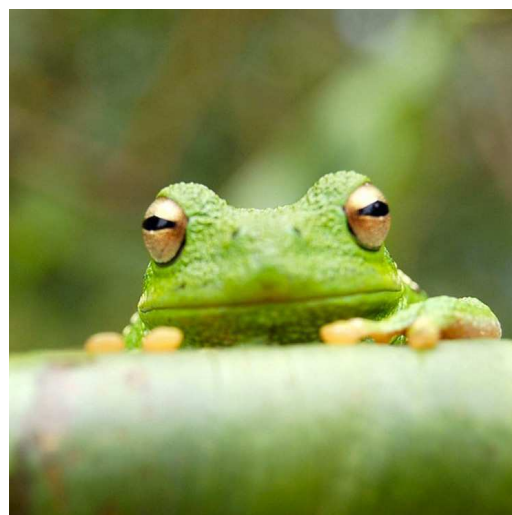


Fig. 3. Placeholder image of a frog with a long example legend to show justification setting.

Tables. Tables should be included in the main manuscript file and should not be uploaded separately.

Single column equations. Authors may use 1- or 2-column equations in their article, according to their preference.

To allow an equation to span both columns, use the `\begin{figure*}...\end{figure*}` environment mentioned above for figures.

Note that the use of the `\widetext` environment for equations is not recommended, and should not be used.

Supporting Information Appendix (SI). Authors should submit SI as a single separate SI Appendix PDF file, combining all text, figures, tables, movie legends, and SI references. SI will be published as provided by the authors; it will not be edited or composed. Additional details can be found in the [PNAS Author Center](#). The PNAS Overleaf SI template can be found [here](#). Refer to the SI Appendix in the manuscript at an appropriate point in the text. Number supporting figures and tables starting with S1, S2, etc.

Authors who place detailed materials and methods in an SI Appendix must provide sufficient detail in the main text methods to enable a reader to follow the logic of the procedures and results and also must reference the SI methods. If a paper is fundamentally a study of a new method or technique, then the methods must be described completely in the main text.

SI Datasets. Supply .xlsx, .csv, .txt, .rtf, or .pdf files. This file type will be published in raw format and will not be edited or composed.

SI Movies. Supply Audio Video Interleave (avi), Quicktime (mov), Windows Media (wmv), animated GIF (gif), or MPEG files. Movie legends should be included in the SI Appendix file. All movies should be submitted at the desired reproduction size and length. Movies should be no more than 10MB in size.

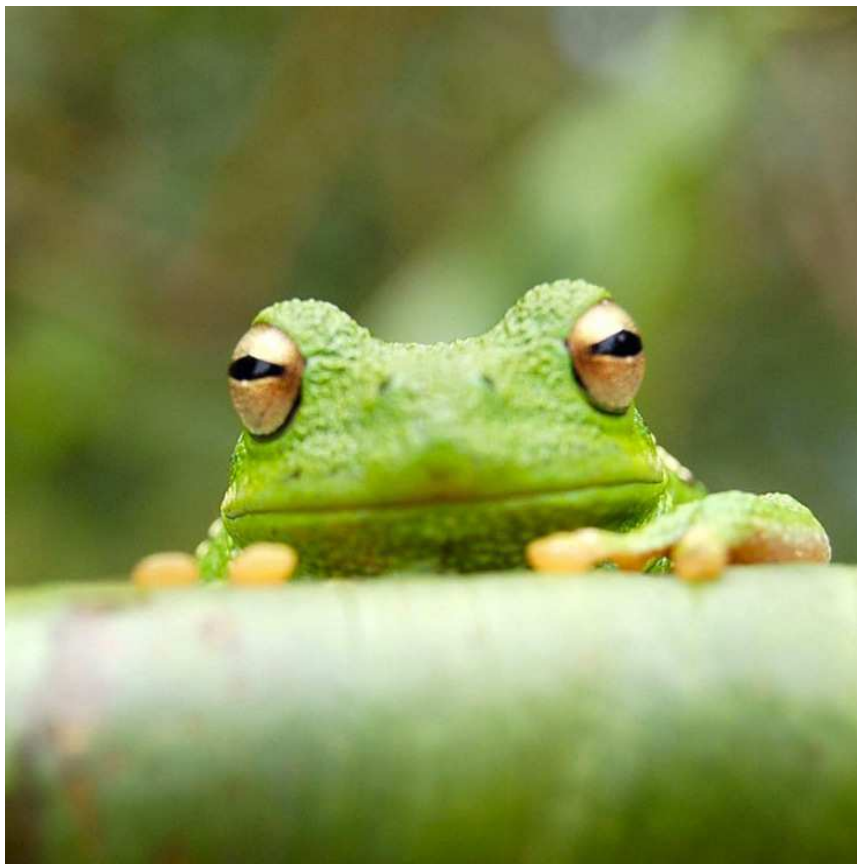


Fig. 4. This legend would be placed at the side of the figure, rather than below it.

$$\begin{aligned}
 (x + y)^3 &= (x + y)(x + y)^2 \\
 &= (x + y)(x^2 + 2xy + y^2) \\
 &= x^3 + 3x^2y + 3xy^2 + x^3.
 \end{aligned}
 \tag{17}$$

Materials and Methods

Please describe your materials and methods here. This can be more than one paragraph, and may contain subsections and equations as required.

Subsection for Method. Example text for subsection.

1. PA Rejto, ST Freer, Protein conformational substates from x-ray crystallography. *Prog. Biophys. Mol. Biol.* **66**, 167–196 (1996).
2. CA Smith, et al., Population shuffling of protein conformations. *Angew. Chem. Int. Ed Engl.* **54**, 207–210 (2015).

ACKNOWLEDGMENTS. Please include your acknowledgments here, set in a single paragraph. Please do not include any acknowledgments in the Supporting Information, or anywhere else in the manuscript.

3. RA Woldeyes, DA Sivak, JS Fraser, E pluribus unum, no more: from one crystal, many conformations. *Curr. Opin. Struct. Biol.* **28**, 56–62 (2014).
4. MA DePristo, PIW de Bakker, TL Blundell, Heterogeneity and inaccuracy in protein structures solved by x-ray crystallography. *Structure* **12**, 831–838 (2004).
5. LH Jensen, Refinement and reliability of macromolecular models based on x-ray diffraction data. *Methods Enzym.* **277**, 353–366 (1997).
6. D Ringe, GA Petsko, Study of protein dynamics by x-ray diffraction. *Methods Enzym.* **131**, 389–433 (1986).
7. Z Sun, Q Liu, G Qu, Y Feng, MT Reetz, Utility of B-Factors in protein science: Interpreting rigidity, flexibility, and internal motion and engineering thermostability. *Chem. Rev.* **119**, 1626–1665 (2019).
8. A Kuzmanic, NS Pannu, B Zagrovic, X-ray refinement significantly underestimates the level of microscopic heterogeneity in biomolecular crystals. *Nat. Commun.* **5**, 3220 (2014).
9. J Kuriyan, GA Petsko, RM Levy, M Karplus, Effect of anisotropy and anharmonicity on protein crystallographic refinement. an evaluation by molecular dynamics. *J. Mol. Biol.* **190**, 227–254 (1986).
10. JM Holton, S Classen, KA Frankel, JA Tainer, The r-factor gap in macromolecular crystallography: an untapped potential for insights on accurate structures. *FEBS J.* **281**, 4046–4060 (2014).
11. D Vitkup, D Ringe, M Karplus, GA Petsko, Why protein r-factors are so large: a self-consistent analysis. *Proteins* **46**, 345–354 (2002).

869	12. H van den Bedem, JS Fraser, Integrative, dynamic structural biology at atomic resolution—it's about time. <i>Nat. Methods</i> 12 , 307–318 (2015).	931
870	13. K Henzler-Wildman, D Kern, Dynamic personalities of proteins. <i>Nature</i> 450 , 964–972 (2007).	932
871	14. CA Schiffer, Time-averaging crystallographic refinement in <i>Computer Simulation of Biomolecular Systems: Theoretical and Experimental Applications</i> , eds. WF van Gunsteren, PK Weiner, AJ Wilkinson. (Springer Netherlands, Dordrecht), pp. 265–269 (1997).	933
872	15. BT Burnley, PV Afonine, PD Adams, P Gros, Modelling dynamics in protein crystal structures by ensemble refinement. <i>Elife</i> 1 , e00311 (2012).	934
873	16. SA Wankowicz, et al., Uncovering protein ensembles: Automated multiconformer model building for x-ray crystallography and Cryo-EM. <i>bioRxiv</i> (2024).	935
874	17. BT Riley, et al., qfit 3: Protein and ligand multiconformer modeling for x-ray crystallographic and single-particle cryo-EM density maps. <i>Protein Sci.</i> 30 , 270–285 (2021).	936
875	18. MA Wilson, AT Brunger, The 1.0 Å crystal structure of Ca(2+)-bound calmodulin: an analysis of disorder and implications for functionally relevant plasticity. <i>J. Mol. Biol.</i> 301 , 1237–1256 (2000).	937
876	19. EJ Levin, DA Kondrashov, GE Wesenberg, GN Phillips, Jr, Ensemble refinement of protein crystal structures: validation and application. <i>Structure</i> 15 , 1040–1052 (2007).	938
877	20. J Kuriyan, et al., Exploration of disorder in protein structures by x-ray restrained molecular dynamics. <i>Proteins</i> 10 , 340–358 (1991).	939
878	21. FT Burling, AT Brünger, Thermal motion and conformational disorder in protein crystal structures: Comparison of multi-conformer and time-averaging models. <i>Isr. J. Chem.</i> 34 , 165–175 (1994).	940
879	22. MC Thompson, Combining temperature perturbations with x-ray crystallography to study dynamic macromolecules: A thorough discussion of experimental methods. <i>Methods Enzymol.</i> 688 , 255–305 (2023).	941
880	23. H Frauenfelder, GA Petsko, D Tsernoglou, Temperature-dependent x-ray diffraction as a probe of protein structural dynamics. <i>Nature</i> 280 , 558–563 (1979).	942
881	24. H Frauenfelder, SG Sligar, PG Wolynes, The energy landscapes and motions of proteins. <i>Science</i> 254 , 1598–1603 (1991).	943
882	25. RF Tilton, Jr, JC Dewan, GA Petsko, Effects of temperature on protein structure and dynamics: X-ray crystallographic studies of the protein ribonuclease-a at nine different temperatures from 98 to 320 K. <i>Biochemistry</i> 31 , 2469–2481 (1992).	944
883	26. JS Fraser, et al., Hidden alternative structures of proline isomerase essential for catalysis. <i>Nature</i> 462 , 669–673 (2009).	945
884	27. B Halle, Biomolecular cryocrystallography: structural changes during flash-cooling. <i>Proc. Natl. Acad. Sci. U. S. A.</i> 101 , 4793–4798 (2004).	946
885	28. DA Keedy, et al., Crystal cryocooling distorts conformational heterogeneity in a model Michaelis complex of DHFR. <i>Structure</i> 22 , 899–910 (2014).	947
886	29. A Ebrahim, et al., The temperature-dependent conformational ensemble of SARS-CoV-2 main protease (Mpro). <i>bioRxiv</i> (2021).	948
887	30. S Du, et al., Refinement of multiconformer ensemble models from multi-temperature x-ray diffraction data. <i>bioRxiv</i> (2023).	949
888	31. MP Rout, A Sali, Principles for integrative structural biology studies. <i>Cell</i> 177 , 1384–1403 (2019).	950
889	32. A Sali, From integrative structural biology to cell biology. <i>J. Biol. Chem.</i> 296 , 100743 (2021).	951
890	33. BR Brooks, et al., CHARMM: the biomolecular simulation program. <i>J. Comput. Chem.</i> 30 , 1545–1614 (2009).	952
891	34. PV Afonine, RW Grosse-Kunstleve, PD Adams, A Urzhumtsev, Bulk-solvent and overall scaling revisited: faster calculations, improved results. <i>Acta Crystallogr. D Biol. Crystallogr.</i> 69 , 625–634 (2013).	953
892	35. VY Lunin, PV Afonine, AG Urzhumtsev, Likelihood-based refinement. i. irremovable model errors. <i>Acta Crystallogr. A</i> 58 , 270–282 (2002).	954
893	36. R McElreath, Statistical rethinking: A bayesian course with examples in R and stan. (2015).	955
894	37. D Liebschner, et al., Macromolecular structure determination using x-rays, neutrons and electrons: recent developments in phenix. <i>Acta Crystallogr D Struct Biol</i> 75 , 861–877 (2019).	956
895	38. A Guinier, X-ray diffraction in crystals. <i>Imperfect Crystals, Amorph. Bodies. Dorer</i> (1963).	957
896	39. D Russel, et al., Putting the pieces together: integrative modeling platform software for structure determination of macromolecular assemblies. <i>PLoS Biol.</i> 10 , e1001244 (2012).	958
897	40. M Hancock, et al., Integration of software tools for integrative modeling of biomolecular systems. <i>J. Struct. Biol.</i> 214 , 107841 (2022).	959
898		960
899		961
900		962
901		963
902		964
903		965
904		966
905		967
906		968
907		969
908		970
909		971
910		972
911		973
912		974
913		975
914		976
915		977
916		978
917		979
918		980
919		981
920		982
921		983
922		984
923		985
924		986
925		987
926		988
927		989
928		990
929		991
930		992