

# PrISM: Precision for Integrative Structural Models

Varun Ullanat<sup>1,+</sup>, Nikhil Kasukurthi<sup>1,+</sup>, and Shruthi Viswanath<sup>1,\*</sup>

<sup>1</sup>National Center for Biological Sciences, Tata Institute of Fundamental Research,  
Bangalore, India.

+ Equal contribution.

\* Corresponding author. [shruthiv@ncbs.res.in](mailto:shruthiv@ncbs.res.in)

Short title: Precision for integrative models

*Keywords:* model precision, integrative structure modeling, validation of integrative models, density map, 3D grid, spatial clustering, thematic mapping.

## Abstract

**Summary** A single precision value is currently reported for an integrative model. However, precision may vary for different regions of an integrative model owing to varying amounts of input information. We develop PrISM (Precision for Integrative Structural Models), to efficiently identify high and low-precision regions for integrative models.

**Availability and Implementation** PrISM is available at <https://github.com/isblab/prism>; benchmark data used in this paper is available at doi:10.5281/zenodo.6241200

**Contact** shruthiv@ncbs.res.in

**Supplementary information** This article contains Supplementary data online.

# Introduction

Integrative modeling has emerged as the method of choice for determining the structures of large macromolecular assemblies which are challenging to characterize using a single experimental method [1]–[6]. Several assemblies have been determined by this approach, yielding insights on transcription [7], gene regulation and DNA repair [8], [9], intra-cellular transport [10], [11], cell cycle progression [12], [13], immune response and metabolism [14], [15]. Integrative modeling often relies on sparse, noisy, and ambiguous data from heterogenous samples [16]. Consequently, usually, more than one model (structure) that satisfies the data. Therefore, an important attribute of an integrative model is its precision, defined as the variability among all the models that satisfy the input data. The precision defines the uncertainty of the structure and is a lower bound on its accuracy. Importantly, downstream applications of the structure are limited by its precision. For example, a protein model of 20 Å precision cannot be used to accurately identify binding sites for drug molecules. Precision aids in making informed choices for future modeling, including altering the representation, degrees of freedom, or the amount of sampling [13], [17].

Currently, a single precision value is reported for the integrative model. However, there can be varying amounts of input information for different regions in the model, resulting in different precisions for different regions [18]. It would be useful to identify regions of high and low-precision in the model. For instance, low-precision regions can be used to suggest where the next set of experimental data would be most impactful. High-precision regions can be used for further analysis such as identifying binding interfaces, rationalizing known mutations, and suggesting new mutations.

Several methods have been proposed for detecting substructure similarities and determining flexible/rigid regions in molecular dynamics simulations [19]–[24]. However, they are not directly applicable to integrative models of macromolecular assemblies. First, these methods rely on the input being a set of atomic structures with known secondary structure. In contrast, integrative models are encoded by a more complex representation of an ensemble of multi-scale, multi-state, time-ordered models, and can comprise of regions with unknown atomic structure [18], [25], [26]. Second, these methods identify rigid substructures without quantifying precision for all parts of the structure. Finally, these methods analyze structures with a small number of proteins and have not been demonstrated to be scalable for tens of thousands of models of large assemblies.

Validation of integrative models, including assessment of model precision, is an open research challenge and timely due to the new PDB archive for integrative structures [25]–[29] (<http://pdb-dev.wwpdb.org>). Here, we demonstrate PrISM, an efficient method to visualize high and low-precision regions of an integrative model. Methods like PrISM are expected to improve the utility of deposited integrative structures.

## PrISM inputs and outputs

The input for PrISM is a set of structurally superposed integrative models (Fig. 1). Commonly, these models are encoded by a multi-scale representation. [18], [25], [26]. Each protein is represented by a sequence of spherical beads; each bead corresponds to a number of contiguous residues along the protein sequence. Coarse-grained bead representations are necessary since large assemblies cannot be efficiently and exhaustively sampled in atomic detail [5], [6], [18]. Regions with atomic structure are represented at higher-resolution (*e.g.*, one residue per bead), while other

regions are usually further coarse-grained (*e.g.*, thirty residues per bead). The most common input would be the models from the most populated cluster from integrative modeling analysis [6], [17], [30]. Additional optional user inputs include the voxel size for bead grids and the number of high and low-precision classes.

The outputs from PrISM are regions (‘patches’) of high and low-precision. They are visualized as a bi-polar color map overlaid on a representative model, with high-precision patches in shades of green and low-precision patches in shades of red.

## PrISM algorithm

The method implemented in PrISM is described here. Alternate design choices are also discussed (Supplementary Methods).

### Obtaining density maps for beads

A coarse-grained bead is the smallest primitive, *i.e.*, unit of representation, of an integrative model. We first compute a density map for each bead. A density map is a projection or rasterization of the beads onto a 3D grid, storing a density value for each grid element (voxel). We use a spherical kernel projection since it explicitly considers the bead mass and radius. The contribution to density  $d_k^i$  to voxel  $k$ , centered at  $v_k$ , in a grid with voxel spacing  $V$ , from bead  $i$  of model  $M_j$ , with centre coordinates  $b_i^j$ , mass  $m_i$ , and radius  $r_i$  is given by:

$$d_k^i = \begin{cases} \frac{m_i V^3}{\frac{4}{3} \pi r_i^3}, & \|b_i^j - v_k\|_2^2 \leq r_i^2 \\ 0, & \text{Otherwise} \end{cases}$$

The densities at each voxel are subsequently normalized by the number of input models to obtain average density at a voxel.

Since the density map for each bead can be independently computed, this step is trivially parallelized. The density map allows for the comparison of beads of different sizes using a uniform representation.

### Computing bead spread

We define the spread of a bead, a measure of its precision, as its density-weighted RMSF from its center of density. That is, the center of density  $c_i$  for bead  $i$  is:

$$c_i = \sum_{k=0}^{N_V} v_i d_k^i$$

where  $N_V$  is the total number of voxels in the grid. We then compute the bead spread  $s$  using the following:

$$s_i = \sqrt{\frac{\sum_{k=0}^{N_V} (\|v_i - c_i\|_2)^2 d_k^i}{\sum_{k=0}^V d_k^i}}$$

This step is also parallelized.

### Classifying beads by spread

Next, we use the Jenks Natural Breaks algorithm to classify beads into high and low-precision classes, given the required number of high- and low-precision classes [31]. This algorithm produces the classification that optimizes a goodness-of-variance measure, similar to  $k$ -means clustering, and is used in thematic mapping for clustering one-dimensional data [31].

### Obtaining patches for each class of beads

Next, we detect beads with concerted localization by identifying correlations in the positions and precisions of beads. We define a *patch* as a set of beads in the same precision class, computed by the Natural Breaks algorithm above, which are also proximal to each other in the set of models. A pair of beads is proximal if the average distance between their surfaces, across the input set of models, is less than 10 Å. Patches represent a further grouping of beads in each class. We obtain patches efficiently by deriving connected components of the graph representing the average contact map of the input set of models.

A naïve implementation of the above algorithm is prohibitively expensive in terms of runtime and memory even for small complexes. Several enhancements were therefore implemented (Supplementary Methods).

### PrISM evaluation and usage

PrISM is benchmarked on twelve systems and shown to be fast (Supplementary Results, Table S1) [6], [8], [9], [17], [18], [30], [32]. The annotated precision is shown to be consistent with measures such as root mean-square fluctuation (RMSF) and localization density maps, providing more fine-grained information than the latter in some cases (Supplementary Results, Table S2-S3, Fig. S1-S5). We also recommended

Detailed usage information is at <https://github.com/isblab/prism>.

### Conclusion

PrISM is an efficient method for annotating precision for integrative models of large assemblies. A limitation is that it is applicable for integrative models generated by the Integrative Modeling Platform (IMP, <https://integrativemodeling.org>). In contrast to atomic structural models, these models are multi-scale, *i.e.* coarse-grained at multiple levels by primitives such as spherical beads. In future, the approach could be extended to other model ensembles of coarse-grained models. Methods such as PrISM are expected to improve the utility of deposited integrative structures in the new Protein Data Bank archive for integrative structures [25], [26], [28], [29] (<http://pdb-dev.wwpdb.org>).

### Availability

The source code is available at <https://github.com/isblab/prism>

**Data availability:** The benchmarks underlying this article are available at doi:10.5281/zenodo.6241200

## Acknowledgements

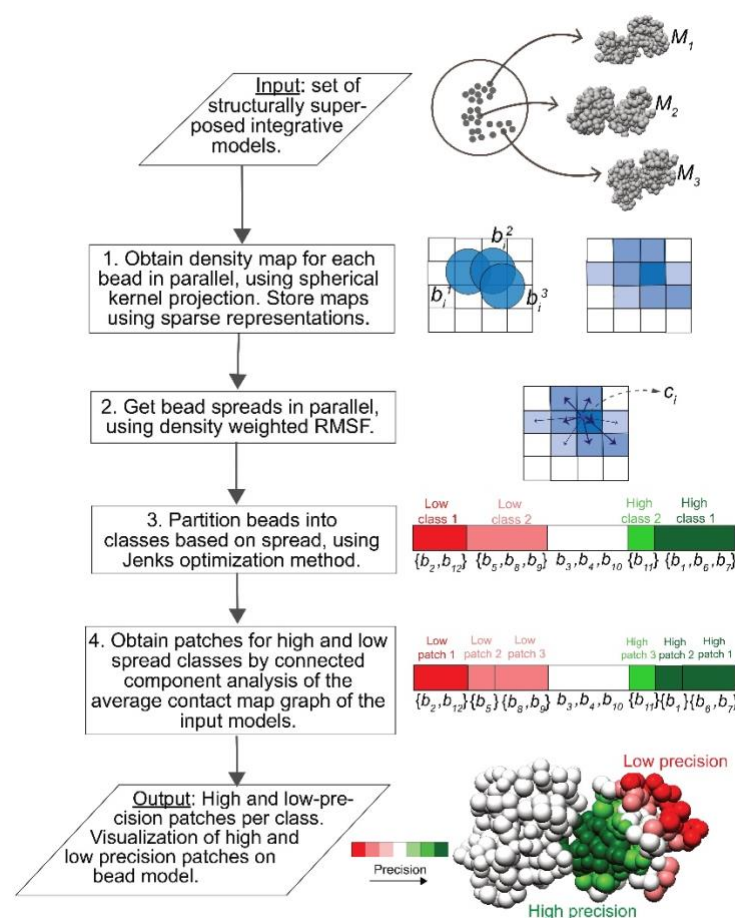
We thank lab members Shreyas Arvindekar, Satwik Pasani, Aditi Pathak, Kartik Majila, and Praveen Roy DS for feedback on the method and manuscript and testing out early versions of the method.

## Funding

This work has been supported by the Department of Science and Technology SERB grant SPG/2020/000475 and Department of Atomic Energy (DAE) TIFR grant RTI 4006 from the Government of India.

*Conflict of Interest:* none declared.

## Figures



**Fig. 1 Schematic of PrISM.** The input is a set of structurally superposed integrative models. Models of three protein-protein complexes are depicted here as  $M_1$ ,  $M_2$  and  $M_3$ . First, a density map is obtained for each bead. Three beads corresponding to bead  $i$  from the three models,  $M_1$ ,  $M_2$  and  $M_3$ , are projected onto the grid. The obtained density map has blue-colored squares; the color intensity corresponds to the density of the square. Next, the bead spread is computed from the density map as the deviation of densities around the center of density  $c_i$ . Subsequently, the Jenks method is used to classify beads into high and low-precision classes. In the example shown, there are two low and two high-precision classes. These classes are further partitioned into patches. The output is a set of high and low-precision patches per class. It is visualized on a representative model as a bi-polar colormap, with shades of green (red) corresponding to high-precision (low-precision) patches.

## References

- [1] F. Alber *et al.*, “The molecular architecture of the nuclear pore complex,” vol. 450, p. 7, 2007.
- [2] D. Russel *et al.*, “Putting the pieces together: integrative structure determination of macromolecular assemblies,” *PLoS Biol*, vol. 10, no. 1, p. e1001244, 2012.
- [3] A. B. Ward, A. Sali, and I. A. Wilson, “Integrative Structural Biology,” *Science*, vol. 339, no. 6122, pp. 913–915, Feb. 2013, doi: 10.1126/science.1228565.
- [4] B. Webb *et al.*, “Integrative structure modeling with the Integrative Modeling Platform,” *Prot Sci*, vol. 27, pp. 245–258, 2018.

- [5] M. P. Rout and A. Sali, “Principles for Integrative Structural Biology Studies,” *Cell*, vol. 177, pp. 1384–1403, 2019.
- [6] D. J. Saltzberg, S. Viswanath, I. Echeverria, I. E. Chemmama, B. Webb, and A. Sali, “Using Integrative Modeling Platform to compute, validate, and archive a model of a protein complex structure,” *Protein Sci. Publ. Protein Soc.*, vol. 30, no. 1, pp. 250–261, Jan. 2021, doi: 10.1002/pro.3995.
- [7] P. Robinson *et al.*, “Molecular architecture of the yeast Mediator complex,” *eLife*, p. 10.7554/eLife.08719, 2015.
- [8] J. Luo *et al.*, “Architecture of the human and yeast general transcription and DNA repair factor TFIIF,” *Mol Cell*, vol. 59, no. 5, pp. 794–806, 2015.
- [9] S. Arvindekar, M. J. Jackman, J. K. K. Low, M. J. Landsberg, J. P. Mackay, and S. Viswanath, “Molecular architecture of nucleosome remodeling and deacetylase sub-complexes by integrative structure determination.” bioRxiv, p. 2021.11.25.469965, Dec. 11, 2021. doi: 10.1101/2021.11.25.469965.
- [10] S. J. Kim *et al.*, “Integrative Structure and Functional Anatomy of a Nuclear Pore Complex,” *Nature*, vol. 555, no. 7697, pp. 475–482, 2018.
- [11] S. J. Ganesan *et al.*, “Integrative structure and function of the yeast exocyst complex,” *Protein Sci. Publ. Protein Soc.*, vol. 29, no. 6, pp. 1486–1501, Jun. 2020, doi: 10.1002/pro.3863.
- [12] S. Viswanath *et al.*, “The molecular architecture of the yeast spindle pole body core determined by Bayesian integrative modeling,” *Mol Biol Cell*, vol. 28, no. 23, pp. 3298–3314, 2017.
- [13] S. Pasani and S. Viswanath, “A Framework for Stochastic Optimization of Parameters for Integrative Modeling of Macromolecular Assemblies,” *Life*, vol. 11, no. 11, Art. no. 11, Nov. 2021, doi: 10.3390/life11111183.
- [14] K. Lasker *et al.*, “Molecular architecture of the 26S proteasome holocomplex determined by an integrative approach,” *Proc Natl Acad Sci USA*, vol. 109, pp. 1380–1387, 2012.
- [15] C. Gutierrez *et al.*, “Structural dynamics of the human COP9 signalosome revealed by cross-linking mass spectrometry and integrative modeling,” *Proc. Natl. Acad. Sci. U. S. A.*, vol. 117, no. 8, pp. 4088–4098, Feb. 2020, doi: 10.1073/pnas.1915542117.
- [16] D. Schneidman-Duhovny, R. Pellarin, and A. Sali, “Uncertainty in Integrative Structural Modeling,” *Curr Opin Struct Biol*, vol. 28, pp. 96–104, 2014.
- [17] S. Viswanath, I. Chemmama, P. Cimermancic, and A. Sali, “Assessing Exhaustiveness of Stochastic Sampling for Integrative Modeling of Macromolecular Structures,” *Biophys J*, vol. 113, no. 11, Art. no. 11, 2017.
- [18] S. Viswanath and A. Sali, “Optimizing model representation for integrative structure determination of macromolecular assemblies,” *Proc. Natl. Acad. Sci.*, vol. 116, no. 2, pp. 540–545, Jan. 2019, doi: 10.1073/pnas.1814649116.
- [19] W. Wriggers and K. Schulten, “Protein domain movements: detection of rigid domains and visualization of hinges in comparisons of atomic coordinates,” *Proteins Struct. Funct. Bioinforma.*, vol. 29, no. 1, pp. 1–14, 1997, doi: [https://doi.org/10.1002/\(SICI\)1097-0134\(199709\)29:1<1::AID-PROT1>3.0.CO;2-J](https://doi.org/10.1002/(SICI)1097-0134(199709)29:1<1::AID-PROT1>3.0.CO;2-J).
- [20] K. Kedem, L. P. Chew, and R. Elber, “Unit-vector RMS (URMS) as a tool to analyze molecular dynamics trajectories,” *Proteins Struct. Funct. Bioinforma.*, vol. 37, no. 4, pp. 554–564, 1999, doi: [https://doi.org/10.1002/\(SICI\)1097-0134\(19991201\)37:4<554::AID-PROT6>3.0.CO;2-1](https://doi.org/10.1002/(SICI)1097-0134(19991201)37:4<554::AID-PROT6>3.0.CO;2-1).



- [21] D. J. Jacobs, A. J. Rader, L. A. Kuhn, and M. F. Thorpe, “Protein flexibility predictions using graph theory,” *Proteins Struct. Funct. Bioinforma.*, vol. 44, no. 2, pp. 150–165, 2001, doi: <https://doi.org/10.1002/prot.1081>.
- [22] C. Pfleger, S. Radestock, E. Schmidt, and H. Gohlke, “Global and local indices for characterizing biomolecular flexibility and rigidity,” *J. Comput. Chem.*, vol. 34, no. 3, pp. 220–233, 2013, doi: <https://doi.org/10.1002/jcc.23122>.
- [23] L. Martínez, “Automatic identification of mobile and rigid substructures in molecular dynamics simulations and fractional structural fluctuation analysis,” *PloS One*, vol. 10, no. 3, p. e0119264, 2015, doi: [10.1371/journal.pone.0119264](https://doi.org/10.1371/journal.pone.0119264).
- [24] F. Cazals and R. Tetley, “Characterizing molecular flexibility by combining least root mean square deviation measures,” *Proteins Struct. Funct. Bioinforma.*, vol. 87, no. 5, pp. 380–389, 2019, doi: <https://doi.org/10.1002/prot.25658>.
- [25] A. Sali *et al.*, “Outcome of the First wwPDB Hybrid/Integrative Methods Task Force Workshop,” *Structure*, vol. 23, no. 7, pp. 1156–67, 2015.
- [26] B. Vallat, B. Webb, J. Westbrook, A. Sali, and H. M. Berman, “Development of a prototype system for archiving integrative/hybrid structure models of biological macromolecules,” *Structure*, vol. 26, no. 6, pp. 894–904.e2, 2018.
- [27] H. M. Berman *et al.*, “Federating Structural Models and Data: Outcomes from A Workshop on Archiving Integrative Structures,” *Structure*, vol. 27, no. 12, pp. 1745–1759, Dec. 2019, doi: [10.1016/j.str.2019.11.002](https://doi.org/10.1016/j.str.2019.11.002).
- [28] B. Vallat, B. Webb, J. Westbrook, A. Sali, and H. M. Berman, “Archiving and Disseminating Integrative Structure Models,” *J Biomol NMR*, vol. 73, pp. 385–98, 2019.
- [29] B. Vallat *et al.*, “New system for archiving integrative structures,” *Acta Crystallogr. Sect. Struct. Biol.*, vol. 77, no. 12, Art. no. 12, Dec. 2021, doi: [10.1107/S2059798321010871](https://doi.org/10.1107/S2059798321010871).
- [30] D. Saltzberg *et al.*, “Modeling biological complexes using Integrative Modeling Platform,” *Methods Mol Biol*, vol. 2022, pp. 353–77, 2019.
- [31] G. F. Jenks, “The Data Model Concept in Statistical Mapping,” in *International Yearbook of Cartography*, 7th ed., 1967, pp. 186–190.
- [32] A. F. Brilot *et al.*, “CM1-driven assembly and activation of yeast  $\gamma$ -tubulin small complex underlies microtubule nucleation,” *eLife*, vol. 10, p. e65168, doi: [10.7554/eLife.65168](https://doi.org/10.7554/eLife.65168).