

Persistent Sheaf Laplacian Analysis of Protein Flexibility

Nicole Hayes¹, Xiaoqi Wei^{1,*}, Hongsong Feng¹, Ekaterina Merkurjev^{1,2†} and Guo-Wei Wei^{1,3,4‡}

¹ Department of Mathematics,

Michigan State University, MI 48824, USA.

² Department of Computational Mathematics, Science and Engineering

Michigan State University, MI 48824, USA.

³ Department of Electrical and Computer Engineering,

Michigan State University, MI 48824, USA.

⁴ Department of Biochemistry and Molecular Biology,

Michigan State University, MI 48824, USA.

Abstract

Protein flexibility, measured by the B-factor or Debye-Waller factor, is essential for protein functions such as structural support, enzyme activity, cellular communication, and molecular transport. Theoretical analysis and prediction of protein flexibility are crucial for protein design, engineering, and drug discovery. In this work, we introduce the persistent sheaf Laplacian (PSL), an effective tool in topological data analysis, to model and analyze protein flexibility. By representing the local topology and geometry of protein atoms through the multiscale harmonic and non-harmonic spectra of PSLs, the proposed model effectively captures protein flexibility and provides accurate, robust predictions of protein B-factors. Our PSL model demonstrates an increase in accuracy of 32% compared to the classical Gaussian network model (GNM) in predicting B-factors for a dataset of 364 proteins. Additionally, we construct a blind machine learning prediction method utilizing global and local protein features. Extensive computations and comparisons validate the effectiveness of the proposed PSL model for B-factor predictions.

Keywords: Protein flexibility, persistent sheaf Laplacians, topological data analysis, machine learning, B-factor prediction.

*Current address: Department of Mathematics, North Carolina State University, Raleigh, NC.

†Corresponding author, Email: merkurje@msu.edu

‡Corresponding author, Email: weig@msu.edu

1 Introduction

Proteins are pivotal to life, playing an essential role in many biological processes, including signaling, gene regulation, transcription, translation, interaction with a protein or substrate molecule, etc.¹ They are composed of amino acids, which form polypeptide chains and fold into specific three-dimensional (3D) structures. There are four levels of protein structures: primary, secondary, tertiary, and quaternary. The primary structure is the linear sequence of amino acids, whereas the secondary structure refers to α -helices and β -sheets due to hydrogen bonds and electrostatic interactions. The tertiary structure corresponds to the 3D shape of a single polypeptide chain, while the quaternary structure describes the global arrangement of multiple polypeptide chains into a functional complex.²

Proteins have various functions; most notably, some of the functions of proteins include catalyzing metabolic reactions (enzymes), providing structural support (e.g., collagen in connective tissues), facilitating cellular communication (e.g., receptors and signaling molecules), and transporting molecules (e.g., hemoglobin for oxygen transport). These functions originate from their 3D structures. In particular, protein structure flexibility is a vital characteristic of protein structure that is essential to protein functions.³ Specifically, protein flexibility enables proteins to adapt to various shapes and conditions, which facilitate their interactions with other molecules, such as DNA, RNA, ions, co-factors, ligands, and other small molecules. Under physiological conditions, proteins undergo constant thermal fluctuation, which enables the proteins to bind substrates, catalyze reactions, and transmit signals. Enzymes, for example, exhibit an induced fit mechanism, where their active sites adapt complementary shapes to accommodate substrates, improving the catalytic efficiency. In a similar way, molecular motors, such as myosins and kinesins, utilize flexibility to enable directed movement during muscle contraction and intracellular transport.

Protein flexibility can be measured by the B-factor, also known as the Debye-Waller factor, which measures the attenuation of X-ray or neutron scattering due to thermal motion of atoms in protein crystallography. Specifically, the B-factor is defined according to the mean displacement of a scattering center in X-ray diffraction data.^{4,5} The B-factor is used to describe the flexibility of atoms and/or amino acids within a protein structure, and it further provides valuable information about the protein's thermal motion, structural stability, activity, and other protein functions.⁶

Protein flexibility has been intensively studied in computational biophysics in recent decades.⁷⁻¹⁰ In addition to the thoroughly investigated flexibility of proteins involved in folding, folded proteins (i.e., proteins in their native conformations) are also flexible and, in fact, exhibit internal motion in neighborhoods of their native conformations.^{11,12} In a seminal work, McCammon et al.¹¹ investigated such local motion in a small folded globular protein using a molecular dynamics (MD) approach, demonstrating the fluid-like characteristics of the internal motions. However, analyzing the dynamics of a large protein would require simulations at time scales that are intractable for the MD approach.¹³ Consequently, other methods have since emerged using a time-harmonic approximation¹⁴ to the protein's potential energy function used in MD, resulting in time-independent techniques. Such methods include normal mode analysis (NMA)¹⁴⁻¹⁸ and elastic network models (ENMs).¹⁹⁻²⁴

Some of the most popular methods^{13,25,26} for protein flexibility analysis include the Gaussian network model (GNM)^{21,27,28} and anisotropic network model (ANM),¹⁹ both of which are types of ENMs. The GNM approach treats the protein as a network, with the residues representing the junctions. B-factors are then approximated using the first few eigenvalues of the connectivity matrix, which correspond to the long-time dynamics of proteins that MD simulations are unable to capture.²⁹ Moreover, multiple methods have emerged as modifications of the original GNM and ANM models, including generalized GNM (gGNM), multiscale GNM (mGNM), and multiscale ANM (mANM).²⁶ Such methods attempt to improve the efficiency and accuracy of GNM and ANM. Due to their ability to capture multiscale information intrinsic to protein structures, mGNM and mANM models have been shown²⁶ to significantly improve B-factor predictions of proteins compared to the original GNM and ANM methods.

Other algorithms, such as the flexibility-rigidity index (FRI),¹³ which relies on the theory of continuum elasticity with atomic rigidity (CEWAR), have also improved results for B-factor prediction over the original GNM method. The FRI is based on the assumption that protein functions depend solely upon the protein's structure and environment, and therefore it assesses flexibility and rigidity by analyzing the topological connectivity and geometric compactness of protein structures. A benefit of the flexibility-rigidity index is that it bypasses the Hamiltonian interaction matrix and matrix diagonalization. Consequently, the FRI has significantly reduced computational complexity compared to other algorithms for protein flexibility analysis. Additional modifications, including fast FRI (fFRI),²⁵ anisotropic FRI (aFRI),²⁵ and multiscale FRI (mFRI),³⁰ have been developed to further improve the efficiency of FRI as well as its accuracy on structures that are difficult for the NMA, GNM, and FRI algorithms.³⁰

Recently, many machine learning approaches have been developed for protein flexibility analysis. For example, sequence-based predictions have been reported,³¹⁻³³ and other machine-learning-based predictions of protein flexibility have also been proposed.³³⁻³⁵ More recently, a method that utilizes both sequence information and structure information has been developed for protein B-factor prediction.³⁶

In 2019, persistent topological Laplacians (PTLs)^{37,38} were first introduced to overcome certain drawbacks of persistent homology, a key technique used in topological data analysis (TDA).^{39,40} Many PTLs have been proposed in the past few years, including the persistent combinatorial Laplacian, the persistent path Laplacian, the persistent

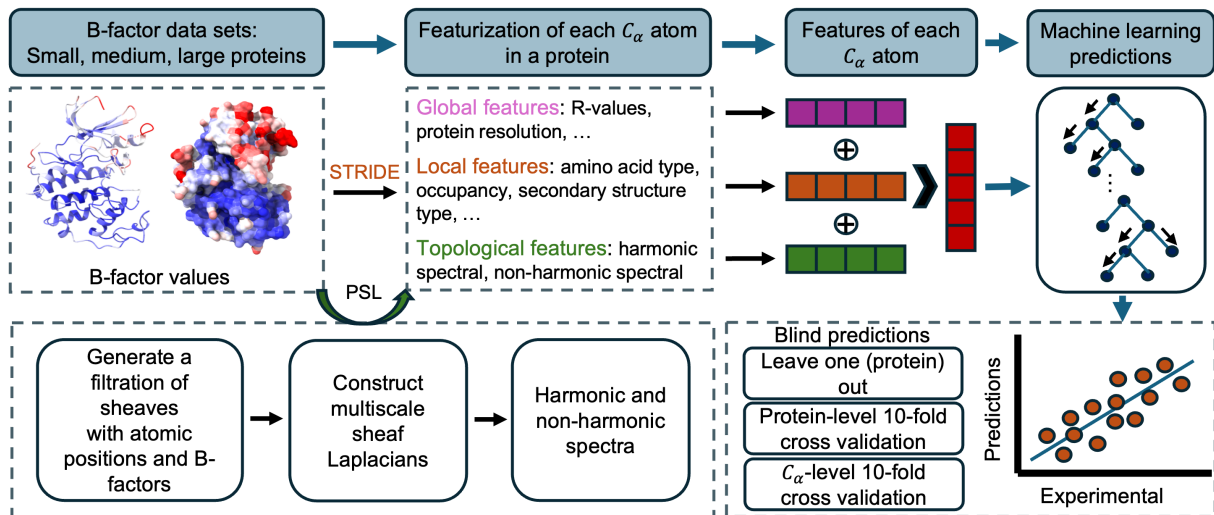


Figure 1: Outline of the methods used in our work. The blind B-factor prediction in Section 2.3 utilizes all pictured features, while the protein subset results from Section 2.1 include only the topological features generated using the persistent sheaf Laplacian (PSL) model.

sheaf Laplacian (PSL),⁴¹ the persistent directed graph Laplacian, and the persistent hyperdigraph Laplacian.⁴² Most of these algorithms are global, offering the topological and geometric descriptions of all objects in their topological space. In other words, they generate information about the protein as a whole. However, for protein flexibility analysis, one must have a method to describe the local properties of individual atoms. The PSL model serves such a function, as it allows the assignment of a specific weight at each node (or atom); thus, it provides local topological and geometric information in its spectra, making it suitable for protein flexibility analysis.

The aim of the present work is to demonstrate the utility of the PSL model for protein flexibility analysis via the prediction of protein B-factors. The remainder of this manuscript is organized in the following manner: all results of this work are given in Section 2. Section 2.1 summarizes our results on protein subsets from the literature, and Section 2.2 presents the performance of the PSL model on individual proteins that are challenging for the GNM. Section 2.3 details the results for blind machine learning prediction using the PSL model. In Section 3, we describe the algorithms used in this manuscript, including some background on persistent homology and cellular sheaves.

2 Results

In this section, we present our results for experiments applying the persistent sheaf Laplacian (PSL) model as outlined in the previous section. Figure 1 summarizes the methods used to generate the results throughout this section.

2.1 Results on protein subsets

2.1.1 Data sets

To demonstrate the persistent sheaf Laplacian model’s performance on proteins of various sizes, we conducted computational experiments on four data sets. Three of these data sets were constructed by Park et al.¹⁴ as sets of relatively small-, medium-, and large-sized protein structures. There are 33 proteins in the set of small-sized proteins, 36 in the set of medium-sized proteins, and 35 in the set of large-sized proteins. The fourth data set is a superset constructed by Opron et al.^{25,30} consisting of (1) the three aforementioned sets, (2) 40 proteins of varying sizes randomly selected from the Protein Data Bank (PDB),⁴³ and (3) 263 high-resolution protein structures used by Xia et al.¹³ in tests of their FRI algorithm, with the duplicates subsequently removed. (Note that in their earlier paper, Opron et al.²⁵ used a set of 365 proteins, but their later manuscript³⁰ excluded the protein with PDB ID 1AGN due to an unrealistic B-factor. The present paper utilizes the updated set consisting of 364 proteins.)

Additionally, all protein data sets used for B-factor prediction in the present study were preprocessed to contain only the C_α atoms from their respective proteins. As discussed by Xia et al.,¹³ the B-factor for an arbitrary atom in a protein is associated with that atom’s flexibility, but its B-factor may be affected by diffraction in data collection,

preventing a direct interpretation of flexibility. However, the B-factors of C_α atoms correlate directly with their atomic flexibility. Accordingly, our B-factor predictions in this work can be interpreted as atomic flexibility predictions.

Table 1 displays the results of the PSL model compared to other methods on the data sets of small, medium, and large proteins as well as the superset.

2.1.2 Parameters and results

For all PSL results in this section and Section 2.2, we utilized a filtration induced by three radii: 6Å, 9Å, and 12Å. For each radius, we generate a 0th persistent sheaf Laplacian matrix L_0 and compute its eigenvalues, then compute the maximum, minimum, mean, and median of the set of non-zero eigenvalues, as well as the number of zero eigenvalues. These quantities comprise five features for each radius, resulting in 15 features in total for each residue. To obtain the B-factor predictions in this section, we performed linear regression using the set of PSL features for the full set of 364 proteins as well as the subsets.

To better assess the performance of the PSL method relative to other approaches and to avoid overfitting, we did not perform an extensive search for the optimal filtration radii and eigenvalue statistic parameters for each task below. Rather, we conducted experiments on the set of 364 proteins with a few sets of parameters and chose those that yielded a good average Pearson correlation coefficient over the entire set. The above parameters may be tuned to further improve model performance for a given task—higher-order persistent sheaf Laplacian matrices and their respective eigenvalues may also be used to generate such features, and other statistics may be used as well, such as the standard deviation of the non-zero eigenvalues. Moreover, suitable filtration radii may be chosen to capture desired multiscale information for a given protein. Another example of PSL feature generation can be seen in Section 2.3.2.

Protein Set	PSL	ASPH (B) ⁵	ASPH (W) ⁵	opFRI ²⁵	pfFRI ²⁵	GNM ¹⁴	NMA ¹⁴
Small	0.927	0.85	0.86	0.667	0.594	0.541	0.480
Medium	0.728	0.69	0.69	0.664	0.605	0.550	0.482
Large	0.643	0.61	0.62	0.636	0.591	0.529	0.494
Superset	0.751	0.65	0.66	0.673	0.626	0.565	NA

Table 1: Average Pearson correlation coefficients for the PSL model compared to other methods. Experiments were conducted on the full set of 364 proteins as well as three subsets of small, medium, and large protein structures as described by Park et al.¹⁴ ASPH denotes the atom-specific persistent homology method developed by Bramer et al.,⁵ with results using Bottleneck (B) and Wasserstein (W) metrics displayed. Both sets of ASPH results used both an exponential and Lorenz kernel for least-squares fitting. opFRI and pfFRI results are from Opron et al.,²⁵ and GNM and NMA results are from Park et al.¹⁴

The PSL model achieves improved performance over all other compared methods on all data sets shown in Table 1. In particular, the PSL model improves the benchmark GNM by 32%.

2.2 Individual protein case studies

As Opron et al. discussed in their 2015 work,³⁰ the Gaussian network model (GNM) experiences difficulty in predicting B-factors for certain protein structures. In addition to the comparison shown in Table 1, in this section, we examine a few case studies of particular proteins to demonstrate the success of the PSL model on such structures. All protein structural visualizations were generated using the Visual Molecular Dynamics software (VMD),⁴⁴ and residues of each protein are assigned colors based on their experimental or predicted B-factors. Lower B-factors are shown as blue (corresponding to “colder” or more rigid residues), and higher B-factors are shown as red (corresponding to “warmer” or more flexible residues). All GNM results were obtained using the default GNM model with a cutoff of 7Å.

Calmodulin is a calcium detector within the cells and plays a significant role in numerous cellular pathways. Its flexibility allows it to interact with varied target proteins. Figure 2 displays the predicted and experimental B-factors for the calcium-binding protein calmodulin (PDB ID: 1CLL)⁴³ using our persistent sheaf Laplacian model as well as the Gaussian network model. We observe that the Gaussian network model produces a large error in B-factor prediction for residues from about 65-85. These residues correspond to a flexible hinge region of the protein.³⁰ The root mean square error (RMSE) for the PSL model is 9.14 for calmodulin, a 23% decrease from the GNM model’s RMSE of 11.9.

Next, we consider a monomeric cyan fluorescent protein (mTFP) that emits cyan light. It is used in biological experiments to visualize specific targets. Figure 3 shows experimental B-factors and predicted B-factors of the protein mTFP1 (PDB ID: 2HQK). Again, the predicted B-factors shown were computed using the Gaussian network model

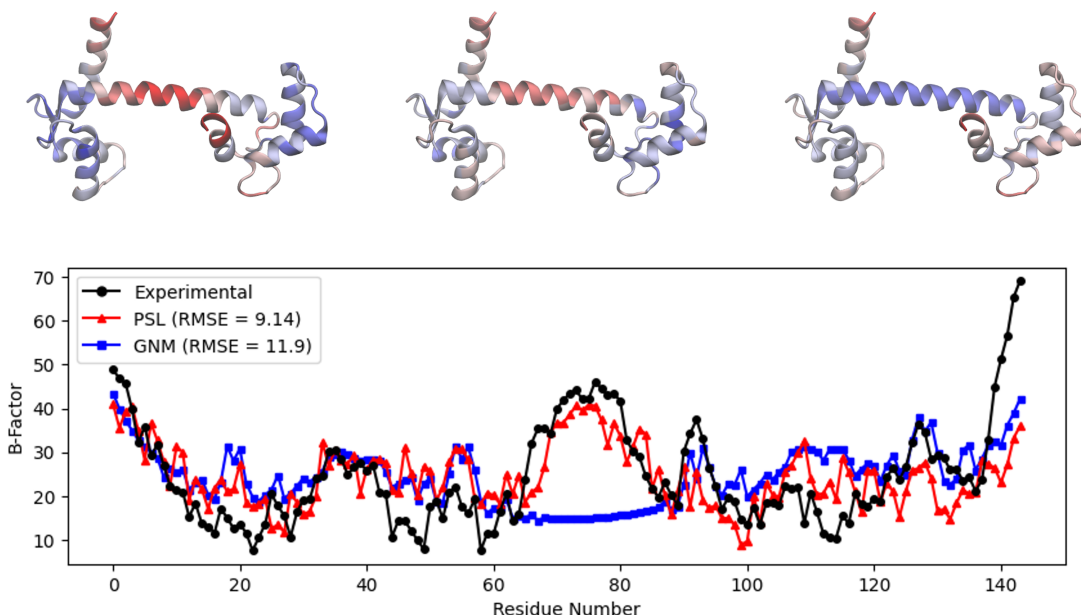


Figure 2: Top: visualization of the protein calmodulin (PDB ID: 1CLL) using Visual Molecular Dynamics (VMD),⁴⁴ with residues colored by experimental B-factors (left), B-factors predicted by PSL (center), and B-factors predicted by GNM (right). Bottom: experimental and predicted B-factors for each residue of the protein. The GNM result uses the default cutoff of 7Å. The GNM underestimates the B-factors for residues between about 65 and 85.

and our PSL model. As in the results for the protein calmodulin, the GNM is unable to correctly predict B-factors for one range of residues (around residues 50-60) in the protein mTFP1. Here, however, the Gaussian network model overestimates the B-factors in this region, visible in the GNM structural representation as the red α -helix in the center of the β -barrel.³⁰ Opron et al.³⁰ observed that using a cutoff of 8Å for GNM somewhat resolves this error, and they suggested that the GNM may experience difficulty in this region due to its use of hard thresholds based on connectivity parameters. The persistent sheaf Laplacian model is significantly more accurate in this region, likely due to the fact that it captures atom-specific information as well as molecular information at multiple scales. Overall, the PSL model improves the RMSE on mTFP1 to 3.43 from 8.74 for the GNM, a nearly 61% decrease.

We further consider a probable antibiotics synthesis protein from *Thermus thermophilus*. In Figure 4, we investigate the experimental and predicted B-factors of this protein (PDB ID: 1V70). On this protein, our persistent sheaf Laplacian model is able to predict the B-factors accurately across all residues of the protein, while the Gaussian network model experiences a high level of inaccuracy on residues from about 0-10. This vast overprediction contributes to a very high RMSE value for the GNM, at 17.9. Our PSL model achieves a significantly lower RMSE of 2.78 on the protein 1V70, 84% lower than that of the GNM.

Finally, we studied the ribosomal protein L14 (PDB ID: 1WHI),³⁰ one of the most conserved ribosomal proteins. It functions as an organizational component of the translational apparatus. In Figure 5, we show the experimental and predicted B-factors for the ribosomal protein L14. Again, we observe that the GNM overestimates the flexibility of some regions of this protein, most significantly for the residues around 60-80. The RMSE for the PSL model on this protein is nearly half that of the GNM model, whose RMSE is 6.59.

2.3 Blind machine learning prediction

2.3.1 Data sets

Two datasets, one from Opron et al.^{25,30} and the other from Park et al.¹⁴ are used in our work. The first dataset contains 364 proteins,^{25,30} and the second¹⁴ has three sets of proteins with small, medium, and large sizes, which are the subsets of the 364 protein set.

In our blind predictions, proteins 1OB4, 1OB7, 2OXL, and 3MD5 from the superset are excluded because the STRIDE software cannot generate features for these proteins. We exclude protein 1AGN due to the known problems

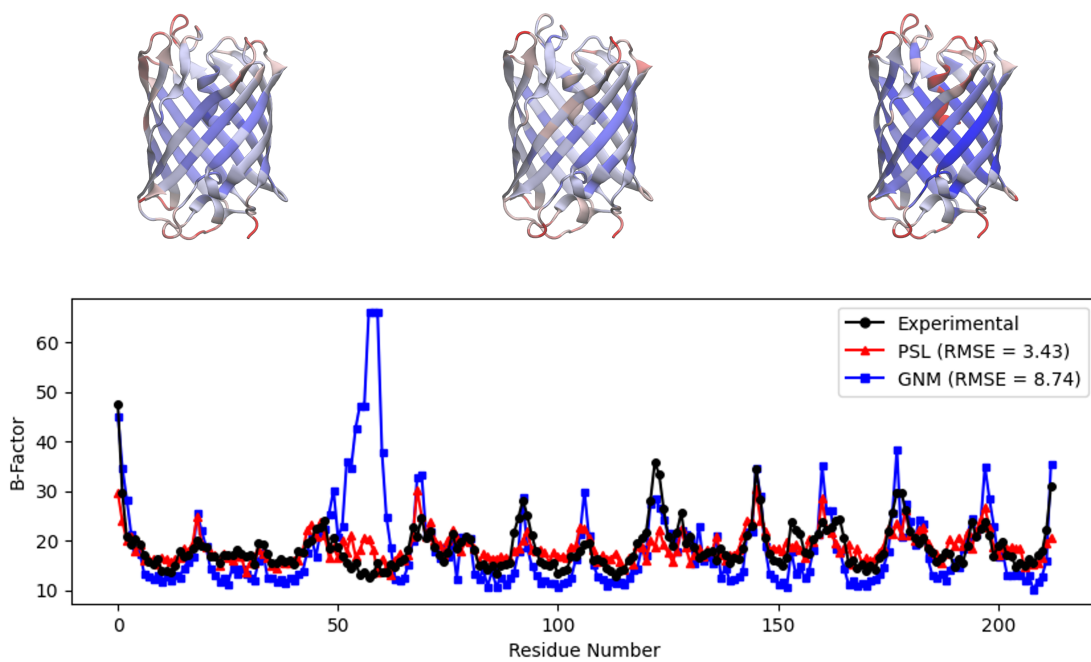


Figure 3: Top: visualization of the protein mTFP1 (PDB ID: 2HQB) using VMD,⁴⁴ with residues colored by experimental B-factors (left), B-factors predicted by PSL (center), and B-factors predicted by GNM (right). Bottom: experimental and predicted B-factors for each residue of the protein. The GNM result uses the default cutoff of 7Å. The GNM vastly overestimates the B-factors of residues around 50-60.

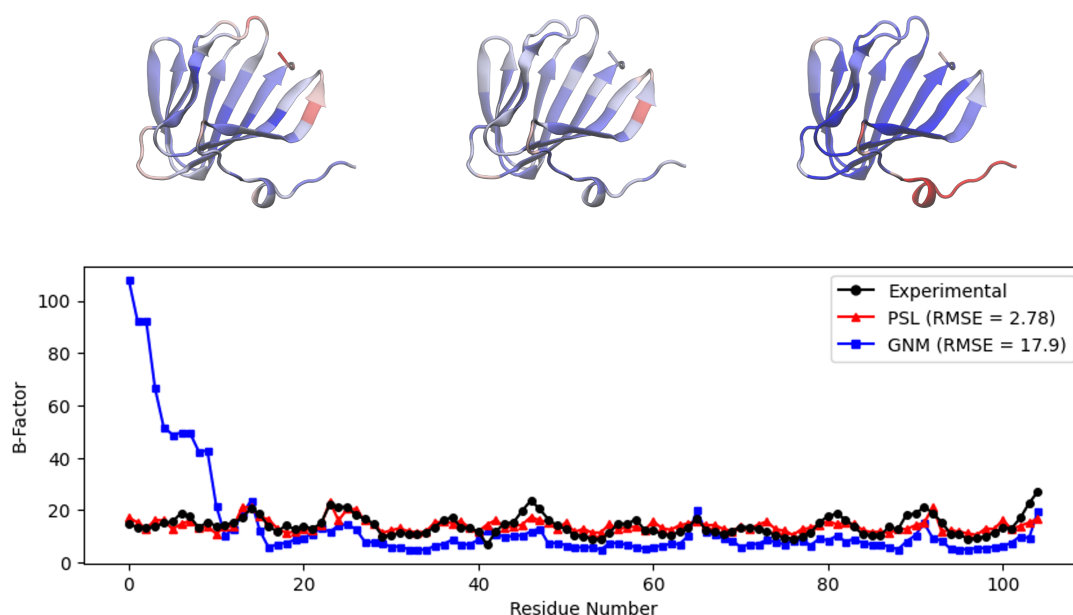


Figure 4: Top: visualization of the protein with PDB ID 1V70 using VMD,⁴⁴ with residues colored by experimental B-factors (left), B-factors predicted by PSL (center), and B-factors predicted by GNM (right). Bottom: experimental and predicted B-factors for each residue of the protein. The GNM result uses a cutoff of 7Å. The GNM vastly overestimates the B-factors for residues from about 0-10.

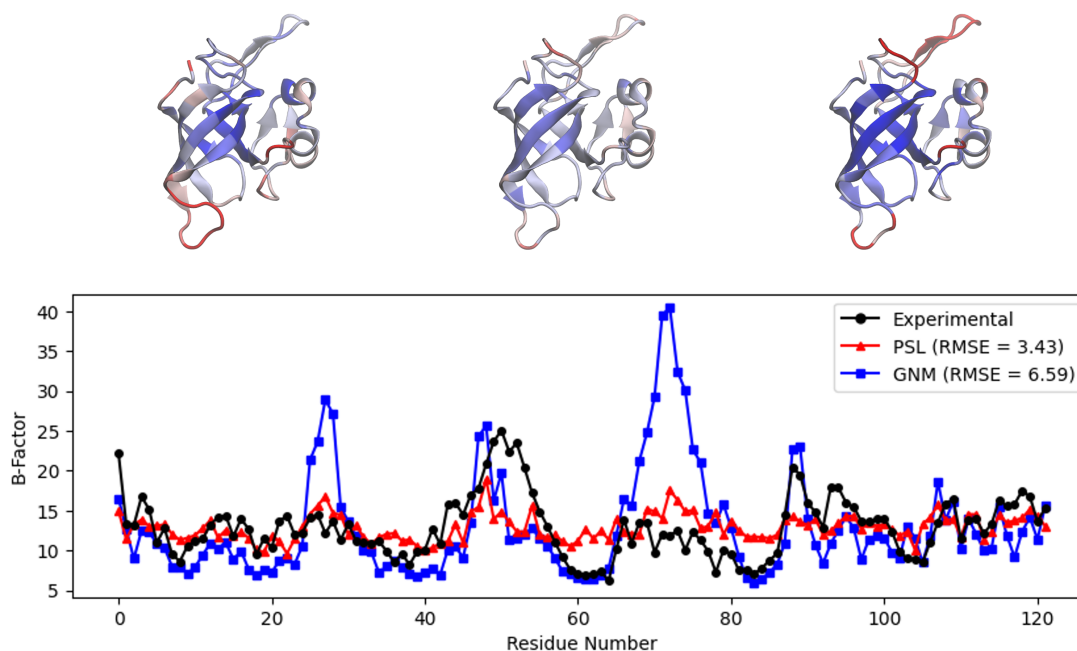


Figure 5: Top: visualization of the ribosomal protein L14 (PDB ID: 1WHI) using VMD,⁴⁴ with residues colored by experimental B-factors (left), B-factors predicted by PSL (center), and B-factors predicted by GNM (right). Bottom: experimental and predicted B-factors for each residue of the protein. The GNM result uses a cutoff of 7Å. The GNM overestimates the B-factors for residues between 60-80.

with this protein data.^{25,30} Additional proteins from the superset are also excluded. Proteins 1NKO, 2OCT, and 3FVA are excluded because these proteins have unphysical B-factors (i.e., zero values). We also excluded proteins 3DWV, 3MGN, 4DPZ, 2J32, 3MEA, 3AOM, 3IVV, 3W4Q, 3P6J, and 2DKO due to inconsistent protein data processed with STRIDE compared to original PDB data. A total of 346 proteins are used for blind predictions. Those data can be found in our provided GitHub repository.

2.3.2 PSL features

The second approach to B-factor prediction that we examined is a blind prediction for protein B-factors. We use PSL features as local descriptors of protein structures, applying three cutoff distances, i.e., 7, 10, and 13Å, to define the atom groups used to construct a sheaf Laplacian matrix. For each cutoff distance, we generate a sheaf Laplacian matrix, L_1 , with a filtration radius matching the cutoff distance. From each matrix, we extract five features: the count of zero eigenvalues, and the maximum, minimum, mean, and standard deviation of the non-zero eigenvalues. Together, these provide 15 PSL features for blind machine learning predictions.

2.3.3 Additional features

In addition to PSL features, we extract a range of global and local protein features for building machine learning models. Each PDB structure is associated with global features, such as the R-value, resolution, and the number of heavy atoms, which are extracted from the PDB files. These features enable the comparison of the B-factors in different proteins. The local characteristics of each protein consist of packing density, amino acid type, occupancy, and secondary structure information generated by STRIDE.⁴⁵ STRIDE provides comprehensive secondary structure details for a protein based on its atomic coordinates from a PDB file, classifying each atom into categories such as α -helix, 3-10-helix, π -helix, extended conformation, isolated bridge, turn, or coil. Furthermore, STRIDE provides ϕ and ψ angles and residue solvent-accessible area, contributing a total of 12 secondary features. In our implementation, we use one-hot encoding for both amino acid types and the 12 secondary features. The packing density of each C_α atom in a protein is calculated based on the density of surrounding atoms, with short, medium, and long-range packing density features defined for each C_α atom. The packing density of the i th C_α atom is defined as

$$p_i^d = \frac{N_d}{N}, \quad (1)$$

where d represents the specified cutoff distance in Å, N_d denotes the number of atoms within the Euclidean distance d from the i th atom, and N is the total number of heavy atoms in the protein. The packing density cutoff values used in this study are provided in Table 2. Our PSL features, combined with the global and local features provided

Short	Medium	Long
$d < 3$	$3 \geq d < 5$	$5 \leq d$

Table 2: Packing density parameter in distance d Å.

for each PDB file, offer a comprehensive feature set for each C_α atom in the protein. For blind predictions, we integrate these features with machine learning algorithms to build regression models. To evaluate the performance of our machine learning model on blind predictions, we conducted two validation tasks: 10-fold cross-validation and leave-one-(protein)-out validation. For 10-fold cross-validation, we designed two types of experiments—one based on splitting by PDB files and another on splitting by all C_α atoms collected from the PDB files. Our modeling and predictions are centered on the B-factors of C_α atoms.

2.3.4 Evaluation metrics

To assess our method for B-factor prediction, we use the Pearson correlation coefficient (PCC):

$$\text{PCC}(\mathbf{x}, \mathbf{y}) = \frac{\sum_{m=1}^M (B_m^e - \bar{B}^e)(B_m^t - \bar{B}^t)}{\sqrt{\sum_{m=1}^M (B_m^e - \bar{B}^e)^2 \sum_{m=1}^M (B_m^t - \bar{B}^t)^2}},$$

where $B_m^t, m = 1, 2, \dots, N$ are the predicted B-factors and $B_m^e, m = 1, 2, \dots, N$ are the experimental B-factors from the PDB file. Here \bar{B}^e and \bar{B}^t are the averaged B-factors.

RF parameters	GBDT parameters
n_estimators = 1000	n_estimators = 1000
max_depth = 8	max_depth = 7
min_samples_split = 4	min_samples_split = 5
min_samples_leaf = 0.8	subsample = 0.8
	learning_rate = 0.002
	max_features = "sqrt"

Table 3: Hyperparameters of the random forest (RF) and gradient boosting decision tree (GBDT) algorithms used for the B-factor predictions.

2.3.5 Machine learning algorithms

For the blind predictions, instead of using more sophisticated methods,^{46–48} we consider two simple machine learning algorithms, namely gradient-boosting decision trees (GBDT) and random forests (RF), to highlight the proposed PSL method. The hyperparameters of these two types of algorithms are given in Table 3.

2.3.6 Machine learning results

We carried out several experiments, the first of which is a leave-one-(protein)-out prediction using the four datasets described above. We trained models five times independently with different random seeds and calculated the average Pearson correlation coefficients from the ten sets of modeling predictions. Our results are shown in Table 4, where the GBDT-based models yield better predictions than the RF-based models, as expected.

Protein set	RF	GBDT
Small	0.478	0.433
Medium	0.518	0.590
Large	0.508	0.582
Superset	0.542	0.588

Table 4: Average Pearson correlation coefficients (PCC) of leave-one (protein)-out predictions for the four B-factor datasets. The PCC results obtained with random forest (RF) and gradient boosting decision tree (GBDT) models are compared.

In our study, we additionally carried out 10-fold cross-validation at the protein level. In each fold, we use nine out of the ten subsets of the 346 proteins to train our model, while the remaining subset is reserved for testing. Specifically, features of C_{α} atoms in the training proteins are pooled together to train the models, while those in the test proteins are used for evaluation. This process is repeated across ten different splits. Table 5 shows the average PCC values for two types of machine learning models. Again, the GBDT model gives better predictions than the RF model.

Protein set	RF	GBDT
Superset	0.397	0.452

Table 5: Average Pearson correlation coefficient (PCC) from protein-level 10-fold cross validation predictions with the collected 346 proteins. The B-factor values of C_{α} atoms in each protein are predicted. The average PCC value is calculated from five independent experiments. The PCC results with random forest (RF) and gradient boosting decision tree (GBDT) modeling are compared.

We also performed an alternative C_{α} -level 10-fold cross-validation. The dataset consists of more than 74,000 C_{α} atoms from 364 proteins. In each of ten independent models, nine out of ten subsets of C_{α} atoms are used to train the models, while the remaining subset is used for testing. As shown in Table 6, GBDT modeling yields slightly better predictions than RF-based modeling.

Protein set	RF	GBDT
Superset	0.839	0.840

Table 6: Average Pearson correlation coefficient (PCC) from C_α -level 10-fold cross validation predictions with all C_α atoms in the collected 346 proteins. The average PCC value is calculated from five independent experiments. The PCC results with random forest (RF) and gradient boosting decision tree (GBDT) models are compared.

3 Methods

3.1 Persistent homology and persistent Laplacians

As one of the most abstract mathematical subjects, homology excessively simplifies complex geometry. In contrast, persistent homology balances simplification and information retrieval in data analysis and is widely used in topological data analysis.^{39,40} However, persistent homology has several drawbacks, including its insensitivity to homotopic shape evolution. To address this challenge, the persistent spectral graph, also known as persistent Laplacians, was introduced on simplicial complexes in 2019.³⁷ Since then, various persistent Laplacians, or persistent topological Laplacians, have been proposed for different topological objects, such as path complexes, directed flag complexes, hyperdigraphs, and cellular sheaves.⁴²

Given a finite set V , a simplicial complex X is a collection of subsets of V , such that if a set σ is in X , then any subset of σ is also in X . A set σ that consists of $q+1$ elements is referred to as a q -simplex. If σ is a subset of τ , then we say that σ is a face of τ and denote the face relation by $\sigma \leq \tau$. If X and Y are simplicial complexes and $X \subset Y$, then X is referred to as a subcomplex of Y . A simplicial complex X gives rise to a simplicial chain complex

$$\cdots \xrightarrow{\partial_3} C_2(X) \xrightarrow{\partial_2} C_1(X) \xrightarrow{\partial_1} C_0(X) \longrightarrow 0.$$

The real vector space $C_q(X)$ is generated by q -simplices. An element of $C_q(X)$ is called a q -chain. The boundary operator ∂_q is a linear map defined by

$$\partial_q[v_{a_0}, \dots, v_{a_q}] = \sum_i (-1)^i [v_{a_0}, \dots, \hat{v}_{a_i}, \dots, v_{a_q}],$$

where the symbol \hat{v}_{a_i} means that \hat{v}_{a_i} is deleted. The total ordering of V ensures that the boundary operator is well-defined. The q -th homology group $H_q = \ker \partial_q / \text{im} \partial_{q+1}$ is well defined since $\partial^2 = 0$. Now suppose X is a subcomplex of Y . Then we have the following diagram

$$\begin{array}{ccccccc} \cdots & \xrightarrow{\partial_{q+2}^X} & C_{q+1}(X) & \xrightarrow{\partial_{q+1}^X} & C_q(X) & \xrightarrow{\partial_q^X} & C_{q-1}(X) \xrightarrow{\partial_{q-1}^X} \cdots \\ & & \downarrow \wr & & \downarrow \wr & & \downarrow \wr \\ \cdots & \xrightarrow{\partial_{q+2}^Y} & C_{q+1}(Y) & \xrightarrow{\partial_{q+1}^Y} & C_q(Y) & \xrightarrow{\partial_q^Y} & C_{q-1}(Y) \xrightarrow{\partial_{q-1}^Y} \cdots \end{array}$$

where hooked dashed arrows represent inclusion maps $\iota : C_q(X) \hookrightarrow C_q(Y)$. The inclusion ι induces a map $\iota^\bullet : H_q(X) \rightarrow H_q(Y)$. The q -th persistent homology for the pair (X, Y) is the image

$$\iota^\bullet(H_q(X)).$$

Usually the ranks of persistent homology groups are represented by barcodes, where each bar represents a topological feature that persists in the filtration, offering a multiscale topological characterization of the input point cloud.^{39,40}

Recently, the theory of persistent Laplacians³⁷ has been proposed to extract additional information from a point cloud. A persistent Laplacian is a positive semi-definite operator whose kernel is isomorphic to the corresponding persistent homology group. The additional information provided by the non-zero eigenvalues of persistent Laplacians can be learned by machine learning algorithms. Since $C_q(X)$ is generated by q -simplices, it is equipped with a canonical inner product. Let $C_{q+1}^{X,Y} = \{c \in C_{q+1}(Y) \mid \partial_{q+1}^Y(c) \in C_q(X)\}$ and $\partial_{q+1}^{X,Y}$ the restriction of ∂_{q+1}^Y to $C_{q+1}^{X,Y}$. The q -th persistent Laplacian $\Delta_q^{X,Y}$ is defined by

$$\partial_{q+1}^{X,Y} (\partial_{q+1}^{X,Y})^\dagger + (\partial_q^X)^\dagger \partial_q^X, \quad (2)$$

where \dagger denotes the adjoint of a linear morphism. Using basic linear algebra we can prove that the kernel of $\Delta_q^{X,Y}$ is isomorphic to $\iota^\bullet(H_q(X))$. Generally speaking, any method that utilizes multiscale Laplacians to analyze data can be referred to as a persistent Laplacian method.

3.2 Cellular sheaves and persistent sheaf Laplacians

Molecular structures often contain important non-spatial information, and many applications of topological methods in analyzing molecular data require integration of non-spatial information. For example, we can use generalized distance to model the biochemical interaction between atoms or only use specific types of atoms as input to persistent homology⁴⁹ or persistent Laplacians.³⁷ An alternative approach is to integrate biological information through the construction of (co)chain complexes and extend persistent homology and persistent Laplacians to new settings. For example, one can construct a filtration of cellular sheaves and consider the persistence module of sheaf cochain complexes instead of simplicial complexes and simplicial chain complexes.⁵⁰

Roughly speaking, a cellular sheaf \mathcal{F} is a simplicial complex X with an assignment to each simplex σ of X a finite-dimensional vector space $\mathcal{S}(\sigma)$ (referred to as the stalk of \mathcal{S} over σ) and to each face relation $\sigma \leq \tau$ (i.e., $\sigma \subset \tau$) a linear morphism of vector spaces denoted by $\mathcal{S}_{\sigma \leq \tau}$ (referred to as the restriction map of the face relation $\sigma \leq \tau$), satisfying the rule

$$\rho \leq \sigma \leq \tau \Rightarrow \mathcal{S}_{\rho \leq \tau} = \mathcal{S}_{\sigma \leq \tau} \mathcal{S}_{\rho \leq \sigma}$$

and $\mathcal{S}_{\sigma \leq \sigma}$ is the identity map of $\mathcal{S}(\sigma)$. We can view stalks as information stored for each simplex, and restriction maps as the way this information interacts. A cellular sheaf gives rise to a sheaf cochain complex

$$0 \longrightarrow C^0(X; \mathcal{S}) \xrightarrow{d} C^1(X; \mathcal{S}) \xrightarrow{d} C^2(X; \mathcal{S}) \xrightarrow{d} \dots$$

The q -th sheaf cochain group $C^q(X; \mathcal{S})$ is the direct sum of stalks over q -dimensional simplices. To define coboundary maps d , we can globally orient the simplicial complex X and obtain a signed incidence relation, an assignment to each $\sigma \leq \tau$ an integer $[\sigma : \tau]$. The coboundary map $d^q : C^q(X; \mathcal{S}) \rightarrow C^{q+1}(X; \mathcal{S})$ is defined by

$$d^q|_{\mathcal{S}(\sigma)} = \sum_{\sigma \leq \tau} [\sigma : \tau] \mathcal{S}_{\sigma \leq \tau}.$$

Now suppose we have \mathcal{F} on X and \mathcal{G} on Y such that $X \subseteq Y$ and stalks and restriction maps of X are identical to those of Y . If each stalk is an inner product space then we have the following diagram

$$\begin{array}{ccc} C^{q-1}(X; \mathcal{F}) & \xrightleftharpoons[(d_{\mathcal{F}}^{q-1})^\dagger]{d_{\mathcal{F}}^{q-1}} & C^q(X; \mathcal{F}) \\ & \searrow \scriptstyle d_{\mathcal{F}, \mathcal{G}}^q & \downarrow \scriptstyle (d_{\mathcal{F}, \mathcal{G}}^q)^\dagger \\ & & \Theta_{\mathcal{F}, \mathcal{G}}^{q+1} \\ & & \swarrow \scriptstyle d_{\mathcal{G}}^q \\ C^q(Y; \mathcal{G}) & \xrightleftharpoons[(d_{\mathcal{G}}^q)^\dagger]{d_{\mathcal{G}}^q} & C^{q+1}(Y; \mathcal{G}) \end{array}$$

where $\Theta_{\mathcal{F}, \mathcal{G}}^{q+1} = \{x \in C^{q+1}(Y; \mathcal{G}) \mid (d_{\mathcal{G}}^q)^\dagger(x) \in C^q(X; \mathcal{F})\}$ and $d_{\mathcal{F}, \mathcal{G}}^q$ is the adjoint of $\pi(d_{\mathcal{G}}^q)^\dagger|_{\Theta_{\mathcal{F}, \mathcal{G}}^{q+1}} : \Theta_{\mathcal{F}, \mathcal{G}}^{q+1} \rightarrow C^q(X; \mathcal{F})$ (π is the projection map from $C^q(Y; \mathcal{G})$ to its subspace $C^q(X; \mathcal{F})$). We define the q -th persistent sheaf Laplacian $\Delta_q^{\mathcal{F}, \mathcal{G}}$ by

$$\Delta_q^{\mathcal{F}, \mathcal{G}} = (d_{\mathcal{F}, \mathcal{G}}^q)^\dagger d_{\mathcal{F}, \mathcal{G}}^q + d_{\mathcal{F}}^{q-1} (d_{\mathcal{F}}^{q-1})^\dagger.$$

When $\mathcal{F} = \mathcal{G}$, the persistent sheaf Laplacian is equal to the sheaf Laplacian of \mathcal{F} . When \mathcal{F} and \mathcal{G} are constant sheaves, persistent sheaf Laplacians coincide with persistent Laplacians. Since a sheaf cochain complex is constructed through stalks and restriction maps, we expect that persistent sheaf cohomology and persistent sheaf Laplacians contain additional information besides the underlying simplicial complex.

If a simplicial complex X is labeled (each vertex is associated with a quantity q), then a sheaf can be constructed as follows. Let $F : X \rightarrow \mathbb{R}$ be a nowhere-zero function. We let each stalk be \mathbb{R} , and for the face relation $[v_0, \dots, v_n] \leq [v_0, \dots, v_n, v_{n+1}, \dots, v_m]$ (here orientation is not relevant), the linear morphism $\mathcal{S}([v_0, \dots, v_n] \leq [v_0, \dots, v_n, v_{n+1}, \dots, v_m])$ is the scalar multiplication by

$$\frac{F([v_0, \dots, v_n]) q_{n+1} \cdots q_m}{F([v_0, \dots, v_n, v_{n+1}, \dots, v_m])}.$$

For a labeled point cloud (a point cloud where each point is associated with a quantity), if we construct a filtration of the point cloud, then for each complex in the filtration we can construct a sheaf as described above. This leads to a filtration of sheaves such as in persistent sheaf cohomology⁵¹ and persistent sheaf Laplacians.⁴¹ The harmonic

spectra of PSLs reveal the topological invariants, while the non-harmonic spectra represent geometric information on the data.^{41,42} In this work, we use sheaf Laplacians to construct features for individual C_α atoms. For a given atom A , we first pick a cutoff distance and only consider the nearby C_α atoms within the cutoff. Then we choose a radius and build an alpha complex X out of these C_α atoms. A cellular sheaf on X is constructed as follows. We denote an atom in X by v_i . We assign a label q_i to v_i , then, we let each stalk be \mathbb{R} . For face relation $v_i \leq v_j$, the restriction map is the scalar multiplication by q_j/r_{ij} , where r_{ij} is the length of $v_i v_j$. For face relation $v_i v_j \leq v_i v_j v_k$, the restriction map is the scalar multiplication by $q_k/(r_{ik}r_{jk})$. Since we want to distinguish the C_α atom A from the other atoms, we let the label of A be 0, and the labels of other nearby C_α atoms be 1. The features are then obtained from the spectra of sheaf Laplacians for this specific C_α atom A . In this manner, we can construct sheaf Laplacian features for all C_α atoms.

4 Conclusion

Protein flexibility is crucial for protein functions, and its prediction is essential for understanding protein properties, protein design, and protein engineering. However, the intrinsic complexity of proteins and their interactions present challenges in understanding protein flexibility. To address this, many effective computational approaches have been developed to predict B-factor values, which reflect protein flexibility. In the literature, a variety of techniques have been proposed, including NMA,¹⁶ GNM,^{20,21} pfFRI,²⁵ ASPH,⁵ opFRI,²⁵ and EH.⁵²

In this study, we propose a persistent sheaf Laplacian (PSL) model for protein B-factor prediction. Sheaf theory, a branch of algebraic geometry, serves as the foundation for PSL, a novel approach to topological data analysis (TDA). Unlike many global TDA tools, PSL is a localized method that captures the local topology of a point within the data. Similarly to other TDA methods, PSL also provides a multiscale analysis of the system under study.

The multiscale nature of PSL allows it to capture atomic interactions across different distance ranges, enabling a more effective analysis of protein flexibility. This characteristic makes the proposed method superior to traditional approaches, such as GNM, which fail to account for atomic interactions beyond a specific cutoff distance.

For cross-protein prediction, we further enhance the PSL by integrating additional global and local features intrinsic to protein structures and structure determination conditions. This integration enables the blind prediction of protein B-factors, which is particularly valuable for assessing protein flexibility when experimental B-factors are unavailable. The proposed PSL model has been validated using various data sets, demonstrating its effectiveness and robustness in protein flexibility analysis.

Acknowledgments

This work was supported in part by NIH grants R01AI164266 and R35GM148196, NSF grant DMS-2052983, MSU Research Foundation, and Bristol-Myers Squibb 65109.

Data and Code Availability

Code is available at https://github.com/weixiaoqimath/persistent_sheaf_Laplacians.

Data is available at https://github.com/fenghon1/MDG_bfactor.

Author Contributions

Conception and design: Guo-Wei Wei. Sample preparation and collection of data: Nicole Hayes. Algorithm implementation: Xiaoqi Wei, Hongsong Feng. Analysis and interpretation of data: Nicole Hayes, Guo-Wei Wei. Supervision: Ekaterina Merkurjev, Guo-Wei Wei. Manuscript preparation: Nicole Hayes, Xiaoqi Wei, Hongsong Feng, Ekaterina Merkurjev, Guo-Wei Wei. All authors contributed to the article and approved the submitted version.

Conflict of Interest

The authors have no conflicts to disclose.

References

- [1] Gregory A Petsko and Dagmar Ringe. *Protein structure and function*. New Science Press, 2004.
- [2] Carl Ivar Branden and John Tooze. *Introduction to protein structure*. Garland Science, 2012.
- [3] Predrag Radivojac, Zoran Obradovic, David K Smith, Guang Zhu, Slobodan Vucetic, Celeste J Brown, J David Lawson, and A Keith Dunker. Protein flexibility and intrinsic disorder. *Protein Science*, 13(1):71–80, 2004.
- [4] Zhoutong Sun, Qian Liu, Ge Qu, Yan Feng, and Manfred T. Reetz. Utility of b-factors in protein science: Interpreting rigidity, flexibility, and internal motion and engineering thermostability. *Chemical Reviews*, 119(3):1626–1665, 2019. PMID: 30698416.
- [5] David Bramer and Guo-Wei Wei. Atom-specific persistent homology and its application to protein flexibility analysis. *Computational and Mathematical Biophysics*, 8(1):1–35, 2020.
- [6] Zheng Yuan, Timothy L Bailey, and Rohan D Teasdale. Prediction of protein B-factor profiles. *Proteins: Structure, Function, and Bioinformatics*, 58(4):905–912, 2005.
- [7] Jianpeng Ma. Usefulness and limitations of normal mode analysis in modeling dynamics of biomolecular complexes. *Structure*, 13(3):373–380, 2005.
- [8] Mauno Vihinen, Esa Torkkila, and Pentti Riikonen. Accuracy of protein flexibility predictions. *Proteins: Structure, Function, and Bioinformatics*, 19(2):141–149, 1994.
- [9] Donald J Jacobs, Andrew J Rader, Leslie A Kuhn, and Michael F Thorpe. Protein flexibility predictions using graph theory. *Proteins: Structure, Function, and Bioinformatics*, 44(2):150–165, 2001.
- [10] Jordi Camps, Oliver Carrillo, Agustí Emperador, Laura Orellana, Adam Hospital, Manuel Rueda, Damjan Cicin-Sain, Marco D’Abramo, Josep Lluís Gelpí, and Modesto Orozco. Flexserv: an integrated tool for the analysis of protein flexibility. *Bioinformatics*, 25(13):1709–1710, 2009.
- [11] J. Andrew McCammon, Bruce R. Gelin, and Martin Karplus. Dynamics of folded proteins. *Nature*, 267(5612):585–590, 1977.
- [12] Robert Huber and William S. Bennett Jr. Functional significance of flexibility in proteins. *Biopolymers*, 22(1):261–279, 1983.
- [13] Kelin Xia, Kristopher Opron, and Guo-Wei Wei. Multiscale multiphysics and multidomain models—flexibility and rigidity. *The Journal of Chemical Physics*, 139(19), 2013.
- [14] Jun-Koo Park, Robert Jernigan, and Zhijun Wu. Coarse grained normal mode analysis vs. refined Gaussian network model for protein residue-level structural fluctuations. *Bulletin of Mathematical Biology*, 75:124–160, 2013.
- [15] Mitsuo Tasumi, Haruki Takeuchi, S. Ataka, Anil M. Dwivedi, and Samuel Krimm. Normal vibrations of proteins: glucagon. *Biopolymers*, 21(3):711–714, March 1982.
- [16] Bernard R. Brooks, Robert E. Bruccoleri, Barry D. Olafson, David J. States, S. Swaminathan, and Martin Karplus. CHARMM: A program for macromolecular energy, minimization, and dynamics calculations. *Journal of Computational Chemistry*, 4(2):187–217, 1983.
- [17] Nobuhiro Go, Tosiya Noguti, and Tetsuo Nishikawa. Dynamics of a small globular protein in terms of low-frequency vibrational modes. *Proceedings of the National Academy of Sciences*, 80(12):3696–3700, 1983.
- [18] Michael Levitt, Christian Sander, and Peter S. Stern. Protein normal-mode dynamics: Trypsin inhibitor, crambin, ribonuclease and lysozyme. *Journal of Molecular Biology*, 181(3):423–447, 1985.
- [19] Ali Rana Atilgan, Stewart Durell, Robert Jernigan, Melik Demirel, Özlem Keskin, and Ivet Bahar. Anisotropy of fluctuation dynamics of proteins with an elastic network model. *Biophysical Journal*, 80(1):505–515, 2001.
- [20] Ivet Bahar, Ali Rana Atilgan, Melik C. Demirel, and Burak Erman. Vibrational dynamics of folded proteins: Significance of slow and fast motions in relation to function and stability. *Physical Review Letters*, 80:2733–2736, Mar 1998.
- [21] Ivet Bahar, Ali Rana Atilgan, and Burak Erman. Direct evaluation of thermal fluctuations in proteins using a single-parameter harmonic potential. *Folding and Design*, 2(3):173–181, 1997.
- [22] Konrad Hinsén. Analysis of domain motions by approximate normal mode calculations. *Proteins: Structure, Function, and Bioinformatics*, 33(3):417–429, 1998.
- [23] Guohui Li and Qiang Cui. A coarse-grained normal mode approach for macromolecules: An efficient implementation and application to Ca²⁺-ATPase. *Biophysical Journal*, 83(5):2457–2474, 2024/07/24 2002.

- [24] Florence Tama and Yves-Henri Sanejouand. Conformational change of proteins arising from normal mode calculations. *Protein Engineering, Design and Selection*, 14(1):1–6, 01 2001.
- [25] Kristopher Opron, Kelin Xia, and Guo-Wei Wei. Fast and anisotropic flexibility-rigidity index for protein flexibility and fluctuation analysis. *The Journal of Chemical Physics*, 140(23):234105, 2014.
- [26] Kelin Xia, Kristopher Opron, and Guo-Wei Wei. Multiscale Gaussian network model (mGNM) and multiscale anisotropic network model (mANM). *The Journal of Chemical Physics*, 143(20):204106, 11 2015.
- [27] Paul J Flory. Statistical thermodynamics of random networks. *Proceedings of the Royal Society of London. A. Mathematical and Physical Sciences*, 351(1666):351–380, 1976.
- [28] Turkan Haliloglu, Ivet Bahar, and Burak Erman. Gaussian dynamics of folded proteins. *Physical Review Letters*, 79:3090–3093, Oct 1997.
- [29] Lee-Wei Yang and Choon-Peng Chng. Coarse-grained models reveal functional dynamics - i. elastic network models – theories, comparisons and perspectives. *Bioinformatics and Biology Insights*, 2:BBI.S460, 2008. PMID: 19812764.
- [30] Kristopher Opron, Kelin Xia, and Guo-Wei Wei. Communication: Capturing protein multiscale thermal fluctuations. *The Journal of Chemical Physics*, 142(21):211101, 06 2015.
- [31] Avner Schlessinger and Burkhard Rost. Protein flexibility and rigidity predicted from sequence. *Proteins: Structure, Function, and Bioinformatics*, 61(1):115–126, 2005.
- [32] Alexandre G de Brevern, Aurelie Bornot, Pierrick Craveur, Catherine Etchebest, and Jean-Christophe Gelly. Predyflexy: flexibility and local structure prediction from sequence. *Nucleic Acids Research*, 40(W1):W317–W322, 2012.
- [33] Yann Vander Meersche, Gabriel Cretin, Alexandre G de Brevern, Jean-Christophe Gelly, and Tatiana Galochkina. MEDUSA: prediction of protein flexibility from sequence. *Journal of Molecular Biology*, 433(11):166882, 2021.
- [34] Matthew R Masters, Amr H Mahmoud, Yao Wei, and Markus A Lill. Deep learning model for efficient protein–ligand docking with implicit side-chain flexibility. *Journal of Chemical Information and Modeling*, 63(6):1695–1707, 2023.
- [35] Xintao Song, Lei Bao, Chenjie Feng, Qiang Huang, Fa Zhang, Xin Gao, and Renmin Han. Accurate prediction of protein structural flexibility by deep learning integrating intricate atomic structures and cryo-em density information. *Nature Communications*, 15(1):5538, 2024.
- [36] Gang Xu, Yulu Yang, Ying Lv, Zhenwei Luo, Qinghua Wang, and Jianpeng Ma. Opus-bfactor: Predicting protein b-factor with sequence and structure information. *bioRxiv*, pages 2024–07, 2024.
- [37] Rui Wang, Duc Duy Nguyen, and Guo-Wei Wei. Persistent spectral graph. *International Journal for Numerical Methods in Biomedical Engineering*, 36(9):e3376, 2020.
- [38] Jiahui Chen, Rundong Zhao, Yiying Tong, and Guo-Wei Wei. Evolutionary de rham-hodge method. *Discrete and continuous dynamical systems. Series B*, 26(7):3785, 2021.
- [39] Gunnar Carlsson. Topology and data. *Bulletin of the American Mathematical Society*, 46(2):255–308, 2009.
- [40] Herbert Edelsbrunner, John Harer, et al. Persistent homology-a survey. *Contemporary Mathematics*, 453(26):257–282, 2008.
- [41] Xiaoqi Wei and Guo-Wei Wei. Persistent sheaf Laplacians. *Foundations of Data Science*, 7(2):446–463, 2025.
- [42] Xiaoqi Wei and Guo-Wei Wei. Persistent topological Laplacians—a survey. *Mathematics*, 13(2):208, 2025.
- [43] Frances C. Bernstein, Thomas F. Koetzle, Grahame J. B. Williams, Edgar F. Meyer Jr, Michael D. Brice, John R. Rodgers, Olga Kennard, Takehiko Shimanouchi, and Mitsuo Tasumi. The Protein Data Bank. *European Journal of Biochemistry*, 80(2):319–324, 1977.
- [44] William Humphrey, Andrew Dalke, and Klaus Schulten. VMD: Visual molecular dynamics. *Journal of Molecular Graphics*, 14(1):33–38, 1996.
- [45] Matthias Heinig and Dmitrij Frishman. STRIDE: a web server for secondary structure assignment from known atomic coordinates of proteins. *Nucleic Acids Research*, 32(Web Server issue):W500–W502, 2004.
- [46] Ekaterina Merkurjev. A fast graph-based data classification method with applications to 3d sensory data in the form of point clouds. *Pattern Recognition Letters*, 136:154–160, 2020.
- [47] Cristina Garcia-Cardona, Ekaterina Merkurjev, Andrea Bertozzi, Arjuna Flenner, and Allon Percus. Multiclass data segmentation using diffuse interface methods on graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(8):1600–1613, 2014.

- [48] Ekaterina Merkurjev, Cristina Garcia-Cardona, Andrea Bertozzi, Arjuna Flenner, and Allon Percus. Diffuse interface methods for multiclass segmentation of high-dimensional data. *Applied Mathematics Letters*, 33:29–34, 2014.
- [49] Zixuan Cang and Guo-Wei Wei. Integration of element specific persistent homology and machine learning for protein-ligand binding affinity prediction. *International Journal for Numerical Methods in Biomedical Engineering*, 34(2):e2914, 2018.
- [50] Jakob Hansen and Robert Ghrist. Learning sheaf laplacians from smooth signals. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5446–5450. IEEE, 2019.
- [51] Florian Russold. Persistent sheaf cohomology. *arXiv preprint*, 2022. doi:[10.48550/arXiv.2204.13446](https://doi.org/10.48550/arXiv.2204.13446) (accessed 2023-10-01).
- [52] Zixuan Cang, Elizabeth Munch, and Guo-Wei Wei. Evolutionary homology on coupled dynamical systems with applications to protein flexibility analysis. *Journal of Applied and Computational Topology*, 4:481–507, 2020.