

# Persistent Laplacian projected Omicron BA.4 and BA.5 to become new dominating variants

Jiahui Chen<sup>1</sup>, Yuchi Qiu<sup>1</sup>, Rui Wang<sup>1</sup>, and Guo-Wei Wei<sup>1,2,3\*</sup>

<sup>1</sup> Department of Mathematics,  
Michigan State University, MI 48824, USA.  
East Lansing, MI 48823 USA.

<sup>2</sup> Department of Electrical and Computer Engineering,  
Michigan State University, MI 48824, USA.

<sup>3</sup> Department of Biochemistry and Molecular Biology,  
Michigan State University, MI 48824, USA.

May 3, 2022

## Abstract

Due to its high transmissibility, Omicron BA.1 ousted the Delta variant to become a dominating variant in late 2021 and was replaced by more transmissible Omicron BA.2 in March 2022. An important question is which new variants will dominate in the future. Topology-based deep learning models have had tremendous success in forecasting emerging variants in the past. However, topology is insensitive to homotopic shape variations in virus-human protein-protein binding, which are crucial to viral evolution and transmission. This challenge is tackled with persistent Laplacian, which is able to capture both the topology and shape of data. Persistent Laplacian-based deep learning models are developed to systematically evaluate variant infectivity. Our comparative analysis of Alpha, Beta, Gamma, Delta, Lambda, Mu, and Omicron BA.1, BA.1.1, BA.2, BA.2.11, BA.2.12.1, BA.3, BA.4, and BA.5 unveils that Omicron BA.2.11, BA.2.12.1, BA.3, BA.4, and BA.5 are more contagious than BA.2. In particular, BA.4 and BA.5 are about 36% more infectious than BA.2 and are projected to become new dominating variants by natural selection. Moreover, the proposed models outperform the state-of-the-art methods on three major benchmark datasets for mutation-induced protein-protein binding free energy changes.

Keywords: SARS-CoV-2, evolution, infectivity, deep learning, persistent Laplacian.

---

\*Corresponding author. Email: weig@msu.edu

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Results</b>	<b>2</b>
2.1	Emerging SARS-CoV-2 variants: Infectivity . . . . .	3
2.2	The performance on the AB-Bind dataset . . . . .	4
2.3	The performance on the SKEMPI dataset . . . . .	4
2.4	The performance on the SKEMPI 2.0 dataset . . . . .	5
2.5	The performance on SARS-CoV-2 datasets . . . . .	6
<b>3</b>	<b>Theories and methods</b>	<b>7</b>
3.1	Persistent Laplacians . . . . .	7
3.1.1	Spectral graphs . . . . .	7
3.1.2	Simplicial complex . . . . .	7
3.1.3	Graph Laplacian . . . . .	8
3.1.4	Persistent spectral graphs . . . . .	9
3.2	Predictive models for mutation-induced protein-protein binding free energy changes . . . . .	10
3.2.1	Persistent Laplacian representation of PPIs . . . . .	10
3.2.2	Machine learning and deep learning algorithms . . . . .	11
3.2.3	Predictive models . . . . .	12
<b>4</b>	<b>Conclusion</b>	<b>12</b>

# 1 Introduction

The coronavirus disease, 2019 (COVID-19) caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) has lasted for more than years. The development of effective vaccines, monoclonal antibodies (mABs), and antiviral drugs has significantly improved our ability to bring COVID-19 pandemic under control. Nonetheless, the emerging SARS-CoV-2 variants become a major threat to existing vaccines, monoclonal antibodies (mABs), and antiviral drugs.

The Omicron variant has mutations on various SARS-CoV-2 proteins, such as non-structure protein 3 (NSP3), NSP4, NSP5, NSP6, NSP12, NSP14, spike (S) protein, envelope protein, membrane protein, and nucleocapsid protein. Specifically, Omicron has three main lineages, BA.1 (B.1.1.529.1), BA.2 (B.1.1.529.2), and BA.3 (B.1.1.529.3), and many sub-lineages. Many new recombinants occurred, including XD, XE, and XF. XD and XE are recombination of Delta and BA.1, while XE is basically a BA.2 Omicron lineage carrying a piece of BA.1 at the front end of its genome. The S protein of XE is still BA.2.

The research community focuses its attention on the mutations at the S protein receptor-binding domain (RBD) due to the fact that the RBD facilitates the binding between the S protein and the host angiotensin-converting enzyme 2 (ACE2), which initiates the viral entry of a host cell and infection. It turns out that the binding strength between the S protein RBD and the ACE2 is proportional to the viral infectivity [20, 25, 41, 45, 47]. An artificial intelligence (AI) study revealed that natural selection is the governing mechanism for SARS-CoV-2 evolution [10]. Specifically, viral evolution selects those mutations that are able to strengthen the RBD-ACE2 binding. This mechanism led to the occurrence of many variants, such as Alpha, Beta, Gamma, Delta, Mu, etc. Natural selection in SARS-CoV-2 mutations was conformed beyond doubt in April 2021 by the genotyping of over half a million viral genomes isolated from patients [50].

Additionally, antibodies are generated by the human immune response to infection or vaccination. A strong RBD-antibody binding would lock off RBD-ACE2 binding and directly neutralize the virus [24, 48, 60]. As such, mABs targeting the S protein, particularly the RBD, are designed to treat viral infection. It was unveiled that viral evolution also selects those mutations that are able to weaken RBD-antibody binding, leading vaccine breakthrough infections [52, 53]. Therefore, a new virus with RBD mutations that make the virus more infectious and more capable of evading the antibody protection would become the next dominating variant, which is the underlying principle for the successful forecasting of Omicron BA.2’s dominance [11].

In biophysics, the strength of protein-protein complex is measured by binding free energy (BFE). Mutation-induced BFE change  $\Delta\Delta G$  is calculated by

$$\Delta\Delta G = \Delta G_{\text{WT}} - \Delta G_{\text{MT}} \tag{1}$$

where  $\Delta G_{\text{WT}}$  and  $\Delta G_{\text{MT}}$  are the BFE of wild type and mutant. A positive (negative) BFE change indicates the strengthening (weakening) of the protein-protein binding. Protein-protein BFE changes can be carried out in a variety of ways as shown in software packages FOLDX [18], SAAMBE [38], mCSM-AB [39], mCSM-PPI2 [42], BindProfX [58], etc. AI approaches take the advantage of existing data and often outperform other methods when experimental data become available. Due to the structural complexity and high dimensionality of protein-protein interactions (PPIs), methods that are able to effectively reduce the PPI structural complexity and dimensionality have demonstrated great advantages in predicting PPI BFE changes [49]. Advanced mathematics, particularly, persistent homology [6, 15, 16, 33, 57, 61], offers tremendous abstraction of PPIs. Persistent homology is the main workhorse in popular topological data analysis (TDA) [3, 12, 13, 59]. Element-specific persistent homology (EPH) has had tremendous success in computational biology [4, 5] and worldwide competitions in computer-aided drug design [35].

Based on FPH, a topology-based network tree (TopNetTree) model was constructed from conventional neural network and decision trees for predicting PPI BFE changes [49]. In the past two years, this approach has been extended with SARS-CoV-2 related deep mutational data to predict the BFE changes RBD-ACE2

and RBD-antibody complexes up on RBD mutations [7, 8]. Initially, in early 2020, TopNetTree model was applied to successfully predict that RBD residues 452 and 501 “have high chances to mutate into significantly more infectious COVID-19 strains” [10]. These RBD mutations later appeared in all major variants, Alpha, Beta, Delta, Gamma, Delta, Epsilon, Theta, Kappa, Lambada, Mu, and Omicron L452R/Q and N501Y mutations. In April 2021, the TopNetTree model predicted a list of 31 RBD antibody-escape mutations, including W353R, I401N, Y449D, Y449S, P491R, P491L, Q493P, etc. [50]. Notably, experimental results confirmed that mutations at RBD residues Y449, E484, Q493, S494, and Y505 enable the virus to escape antibodies [2]. It was revealed that variants found in the United Kingdom and South Africa in late 2020 would strengthen virus infectivity, which is consistent with the experimental results [14]. In summer 2021, a topology-based deep neural network trained with mAbs (TopNetmAb) was developed to forecast a list of most likely vaccine-escape RBD mutations, such as S494P, Q493L, K417N, F490S, F486L, R403K, E484K, L452R, K417T, F490L, E484Q, and A475S [8], and mutations S494P, K417N, E484K/Q, and L452R were designated as the variants of concern or variants of interest denounced by the Worldwide Health Organization (WHO). The correlation between the experimental deep mutational data [26] and AI-predicted RBD-mutation-induced BFE changes for all possible 3686 RBD mutations on the RBD-ACE2 complex is 0.7 [8]. In comparison, experimental deep mutational results for the same set of RBD mutations from 2 different labs only have a correlation of 0.67 [26, 46]. TopNetmAb predictions of Omicron [9] and Omicron BA.2 [11] infectivity, vaccine breakthrough, and antibody resistance were nearly perfectly confirmed by experiments and pandemic evolution in the world. These mechanistic discovery and successful predictions may not be achievable via purely experimental means, indicate the indispensable role of AI for scientific discovery.

However, persistent homology and TDA provide only topological invariants, which may not be sufficient for representing PPI data. In particular, the shape of data arisen from a family of homotopy geometries cannot be captured by persistent homology. For example, the geometry of each drum in an acoustic drum set is designed to offer a specific sound or frequency, but persistent homology is insensitive to the change in the sizes (or shapes) in the drum set. This challenge in TDA was addressed by the introduction of persistent Laplacian, or persistent spectral graph [54]. Persistent Laplacian manifests the full set of topological invariants and the shape of data in its harmonic and non-harmonic spectra, respectively. Additional mathematical analysis [31] and a software package, i.e., HERMES [55], for persistent Laplacian have been reported in the literature. This method has been successfully applied to biological studies, including protein thermal stability [54], protein-ligand binding [32], and protein-protein binding problems [56].

In the present work, we introduce element-specific and site-specific persistent Laplacians to forecast emerging SARS-CoV-2 variants. We hypothesize that persistent Laplacians generates intrinsically low-dimensional representations of PPIs and dramatically reduce the dimensionality of PPI data, leading to a reliable high-throughput screening of emerging SARS-CoV-2 variants. To quantitatively validate this hypothesis, we integrate the harmonic and non-harmonic spectra of persistent Laplacians with efficient machine learning algorithms, i.e., gradient boosting tree (GBT) and deep neural network (Net), to predict PPI  $\Delta\Delta G$  following mutations. The resulting topological and spectral-based machine learning models are validated on three major benchmark datasets, the AB-Bind database [44], SKEMPI dataset [34] and SKEMPI v2.0 dataset [22], giving rise to the state of the art performance. Meanwhile, with additional training on SARS-CoV-2 related datasets, our models forecast emerging SARS-CoV-2 variants and recommend four Omicron subvariants, i.e., BA.2.11, BA.2.12.1, BA.4, and BA.5 for active surveillance.

## 2 Results

In this section, we first carry out the infectivity predictions on emerging SARS-CoV-2 variants. Next, three benchmark PPI datasets, i.e., the AB-Bind [44], SKEMPI [34], and SKEMPI 2.0 datasets [22] are employed to demonstrate the proposed persistent Laplacian-based AI models with ten-fold cross validations. Two

evaluation metrics, Pearson correlation  $R_p$  and the root-mean-square error (RMSE), are used to assess the quality of the present models. Lastly, we present the validation of our models on SARS-CoV-2-related datasets.

## 2.1 Emerging SARS-CoV-2 variants: Infectivity

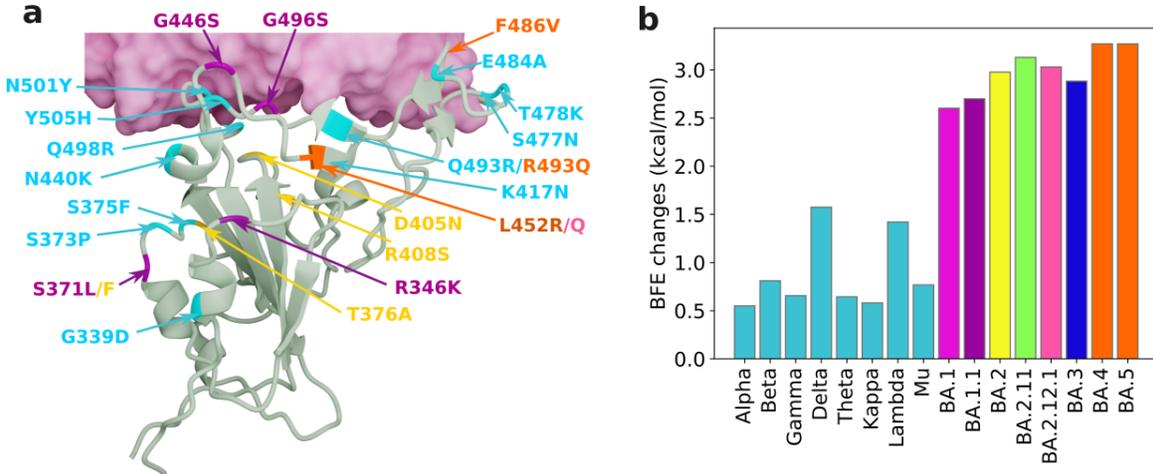


Figure 1: The RBD mutations of Omicron subvariants at the RBD-ACE2 interface and their mutation-induced BFE changes. **a** RBD mutations of Omicron subvariants at the RBD-ACE2 interface (PDB: 7T9L [30]). The shared 12 mutations are shown in cyan. BA.1 mutations are plotted with magenta. BA.2 mutations are marked in yellow. BA.4 and BA.5 mutations are labeled in orange. The rest colors can be matched from the right chart. **b** A comparison of predicted mutation-induced BFE changes for various SARS-CoV-2 variants and subvariants.

Figure 1 shows the RBD mutations of Omicron subvariants and their BFE changes of SARS-CoV-2 variants. A comparison is also given to other main SARS-CoV-2 variants Alpha, Beta, Gamma, Delta, Theta, Kappa, Lambda, and Mu variants. The Delta variant had the highest BFE change among the earlier variants and was the most infectious variant before the occurrence of the Omicron variant, which explains its dominance in 2021. Omicron BA.1, BA.2, and BA.3 have the common RBD mutations G339D, S373P, S375F, K417N, N440K, S477N, T478K, E484A, Q493R, Q498R, N501Y, and Y505H. Omicron BA.1 has three distinct RBD mutations S371L, G446S, and G496S. Four distinct mutations, S371F, T376A, D405N, and R408S, were found for Omicron BA.2. Omicron BA.3 shares three mutations either with BA.1 or BA.2: S371F, D405N, and G446S. The AI-predicted BFE changes of BA.1, BA.2, and BA.3 are 2.60, 2.98, and 2.88 kcal/mol, respectively [11]. These values are significantly higher than those of other major SARS-CoV-2 variants as shown in Figure 1. Note that Omicron BA.2 is the most infectious variant. It is about 20 and 4.2 times as infectious as the original SARS-CoV-2 and the Delta variant, respectively. The machine learning model also predicts that BA.2 is about 1.5 times as contagious as BA.1, which is highly consistent with experimental studies [1, 28]. BA.2 has been the dominating variant since late March 2022 [11].

We have also examined the other Omicron subvariants, namely, BA.1.1, BA.2.11, BA.2.12.1, BA.4, and BA.5. Compared with BA.1, BA.1.1 has one additional RBD mutation, i.e., R346K. BA.2.11 has one more RBD mutation, L452R, than BA.2 does. BA.2.12.1 has an extra RBD mutation, L452Q, compared with BA.2. BA.4 and BA.5 share the same set of RBD mutations but differ in ORF7b, nucleocapsid (N), and membrane (M) proteins. They have three additional RBD mutations, L452R, F486V, and R493Q compared with BA.2. Note that R493Q is a reversion to the wide type, Q493. It is interesting that L452R is one of Delta’s two RBD mutations. Additionally, mutations simultaneously occurred on two RBD residues, L452 and N501, which were singled out by our AI model in early 2020 [10].

Our AI-predicted BFE changes for BA.1.1, BA.2.11, BA.2.12.1, BA.4, and BA.5 are 2.70, 3.13, 3.03, 3.27, and 3.27 kcal/mol, respectively. It is noticed that BA.4 and BA.5 are predicted to be 1.36 times as

infectious as BA.2 and have high potential to become new dominating SARS-CoV-2 variants.

## 2.2 The performance on the AB-Bind dataset

The AB-Bind dataset, including 1,101 mutational data entries for experimentally determined BFE changes [44] is considered in the validation of the proposed models. Its 645 single mutations involving 29 antibody-antigen complexes are denoted as the AB-Bind S645 set. In the AB-Bind S645 set, about one-fifth of mutations strengthen the binding, while the rest are destabilizing mutations. In particular, 27 non-binders, which are mutants determined not to bind within the experimental sensitivity of the assay, are in the dataset. The mutation-induced binding free energy changes for these non-binders were set to -8 kcal/mol. For machine learning models, non-binders are outliers and can cause a very negative impact on model accuracy.

Method	$R_p$	Method	$R_p$
TopLapGBT	0.89/0.95*	mCSM-AB	0.53/0.56*
LapGBT	0.88/0.94*	Discovery Studio	0.45
TopGBT	0.88/0.95*	mCSM-PPI	0.31
TopLapNetGBT	0.87/0.93*	FoldX	0.34
LapNetGBT	0.87/0.91*	STATIUM	0.32
TopNetGBT	0.86/0.93*	DFIRE	0.31
TopNet	0.81/0.88*	bASA	0.22
TopLapNet	0.79/0.87*	dDFIRE	0.19
LapNet	0.72/0.81*	Rosetta	0.16

Table 1: Comparison of the Pearson correlation coefficients ( $R_p$ ) of various methods for the AB-bind S645 set. Except for present TopLapGBT and TopLapNet, the results of other existing methods are adopted from Ref. [39].

\*Results exclude 27 non-binders (their  $\Delta\Delta G$ s were set to -8 kcal/mol [44]).

As shown in Table 1, our TopLapGBT and LapNet models achieved the  $R_p$  of 0.89 and 0.72 for the AB-Bind S645 set. In comparison, TopNet outperforms LapNet because TopNet includes auxiliary features, while LapNet has only Laplacian features. The  $R_p$  values of our other seven models are lower than 0.89 but higher than 0.72. Note that our worst model (LapNet) still outperforms the other best model in the literature by a large margin of 36%, while our best model is about 68% better than the other best model in the literature, indicating the predictive power of our topology and Laplacian-based machine learning models. Both GBTs and Nets models are quite sensitive to system errors as the model training is based on optimizing the mean-square error of the loss function. The BFE changes of 27 non-binders (-8 kcal/mol) did not follow the distribution of the whole dataset. For the TopLapGBT model, the RMSE of AB-Bind S645set is 1.68 kcal/mol and reduces to 0.97 kcal/mol when 27 non-binder samples are excluded. In this case, the  $R_p$  of the TopLapGBT model is increased from 0.89 to 0.95. The consensus results of GBT and Net have correlations of 0.86-0.87, which are lower than that of GBT but higher than that of Net. GBT models outperform Net models in the validation, showing that GBT performs better than Net on a small dataset.

## 2.3 The performance on the SKEMPI dataset

The SKEMPI dataset [34] has 3,047 entries of BFE changes induced by mutations. This dataset is collected from the literature for protein-protein heterodimeric complexes with experimentally determined structures. It consists of single- and multi-mutations. Among them, 2,317 single mutations out of 3,047 entries are called the S2317 dataset. Recently, a subset of 1,131 non-redundant interface single-mutations is selected and denoted as the S1131 set [58]. Table 2 shows the Pearson correlation coefficients on tenfold cross-validations of various models, including topology- and Laplacian-based models. The proposed topology- and Laplacian-based models are found to be more accurate than other existing methods. One may notice that for a larger training set, the consensus predictions of GBT and Net outperform GBT methods. Additionally,

topology-based models contain topology features and auxiliary features, which include more biomolecular information than Laplacian-based models.

Method	$R_p$	Method	$R_p$
TopLapNetGBT	0.87	BindProfX	0.738
TopNetGBT	0.87	Profile-score+FoldX	0.738
TopLapNet	0.86	Profile-score	0.675
TopNet	0.86	SAAMBE	0.624
TopLapGBT	0.86	FoldX	0.457
TopGBT	0.86	BeAtMuSic	0.272
LapNetGBT	0.81	Dcomplex	0.056
LapNet	0.81		
LapGBT	0.78		

Table 2: Comparison of the Pearson correlation coefficients ( $R_p$ ) of various methods for the S1131 set in the SKEMPI dataset. The results of other methods are adopted from Ref. [58].

## 2.4 The performance on the SKEMPI 2.0 dataset

The SKEMPI 2.0 [22] database is an updated version of the original SKEMPI database with new mutations from three other databases: AB-bind [44], PROXiMATE [23], and dbMPIKT [27]. This dataset has 7,085 entries, including single-mutations and multi-mutations. To validate mCSM-PPI2, David et al. filtered only single-point mutations, selected 4169 variants in 319 different complexes, and denoted them as the S4169 set [42]. Additionally, set S8338 was derived from set S4169 by setting the BFE changes of the reverse mutations as the negative values of the original BFE changes induced by mutations. We present our tenfold cross-validation results on sets S4169 and S8338 in Table 3. For S4169, TopLapNetGBT has the most accurate result with  $R_p$  of 0.82 and RMSE of 1.06 kcal/mol. Topology-based models, aided by auxiliary features, have correlations greater than 0.80 and RMSE from 1.04 kcal/mol to 1.10 kcal/mol. Purely Laplacian-based models also performed quite well, with the Pearson correlation of 0.76, which is the same as that of the mCSM-PPI2.

S4169		S8338	
Method	$R_p$	Method	$R_p$
TopLapNetGBT	0.82	TopLapNetGBT	0.87
TopNetGBT	0.82	TopLapNet	0.87
TopLapNet	0.81	TopNetGBT	0.87
TopLapGBT	0.81	TopNet	0.86
TopNet	0.81	TopLapGBT	0.85
TopGBT	0.80	TopGBT	0.85
LapNetGBT	0.77	LapNetGBT	0.83
mCSM-PPI2	0.76	mCSM-PPI2	0.82
LapNet	0.76	LapNet	0.81
LapGBT	0.76	LapGBT	0.80

Table 3: Comparison of the Pearson correlation coefficients ( $R_p$ ) of various methods for S4169 set and S8338 set in SKEMPI 2.0. Results of mCSM-PPI2 are from Ref. [42]

For the S8338 set, TopLapNetGBT has the highest Pearson correlation  $R_p$  of 0.8702 and RMSE of 1.01 kcal/mol as shown in Table 3. TopLapNet has the most accurate results with  $R_p$  of 0.8688 and RMSE of

0.984 kcal/mol. Topology models, aided by auxiliary features, have the  $R_p$  in the range of (0.848, 0.870) and RMSE in the range of (1.070 kcal/mol, 0.984 kcal/mol). LapNet and LapGBT models have their  $R_p$  values slightly lower than that of mCSM-PPI2, but the  $R_p$  of their consensus (LapNetGBT) is higher than that of the mCSM-PPI2.

## 2.5 The performance on SARS-CoV-2 datasets

Training datasets have the utmost importance in implementing our machine learning model for SARS-CoV-2 applications. First, all the datasets mentioned above, including AB-bind, [44] PROXiMATE [23], dbMPIKT [27], SKEMPI [34], and SKEMPI 2.0 [22], are used in our model training. Additionally, SARS-CoV-2-related datasets are also employed to improve the prediction accuracy after a label transformation. These are deep mutational enrichment ratio data, including mutational scNeting data of ACE2 binding to the receptor-binding domain (RBD) of the S protein [40], mutational scNeting data of RBD binding to ACE2 [26, 46], and mutational scNeting data of RBD binding to CTC-445.2 and of CTC-445.2 binding to the RBD [26]. Note that in our validation, our training datasets exclude the test dataset, which is a mutational scNeting data of RBD binding to ACE2. Here, these datasets provide more information on SARS-CoV-2 and can be used to calibrate the models to predict the real experimental results.

Here, we present a validation of our model BFE change prediction for mutations on S protein RBD compared to the experimental deep mutational enrichment data [26]. We compare between experimental deep mutational enrichment data and BFE change predictions on SARS-CoV-2 RBD binding to ACE2 in Figure 2. Both BFE changes (Figure 2 top) and enrichment ratios (Figure 2 bottom) describe the binding affinity changes of the S protein RBD-ACE2 complex induced by mutations. It can be found that the predicted BFE changes are highly correlated to the enrichment ratio data. Pearson correlation is 0.69.

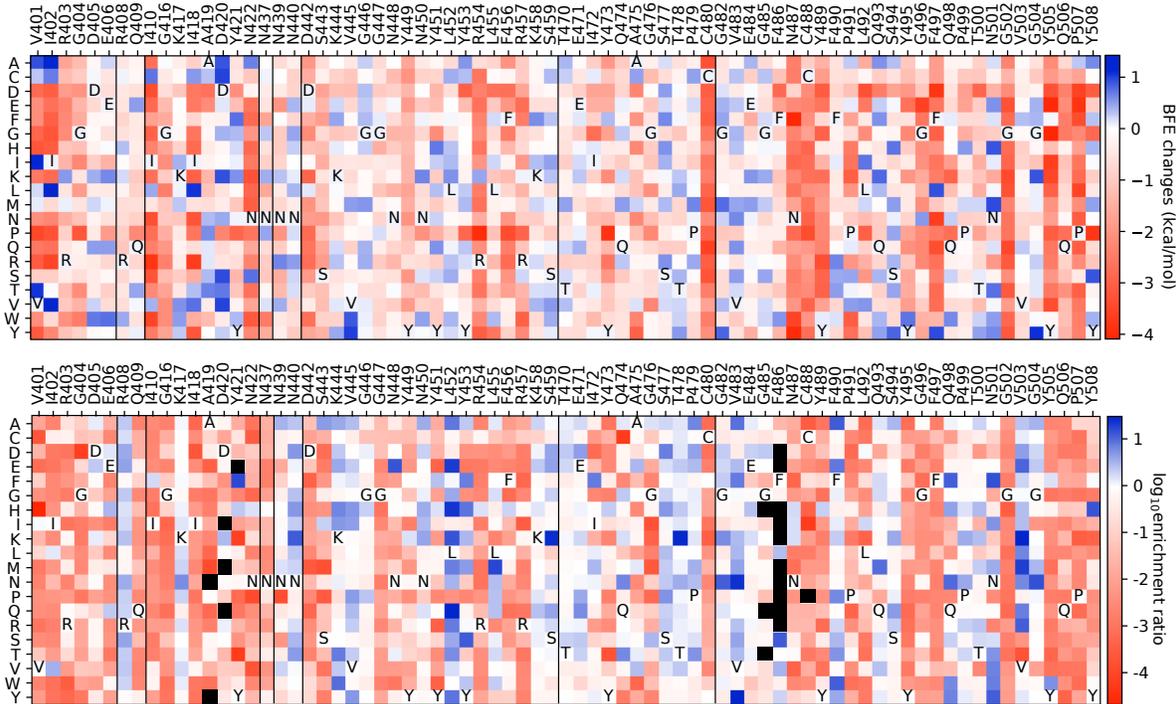


Figure 2: A comparison between experimental RBD deep mutation enrichment data and predicted BFE changes for SARS-CoV-2 RBD binding to ACE2 (6M0J) [26]. **Top:** machine learning predicted BFE changes for single-site mutants of the S protein RBD. **Bottom:** deep mutational scanning heatmap showing the average effect on the enrichment for single-site mutants of RBD when assayed by yeast display for binding to the S protein RBD [26].

### 3 Theories and methods

This section presents brief reviews of spectral graph theory, simplicial complex, and persistent Laplacian are presented. Machine learning and deep learning models are discussed in test datasets and validation settings.

#### 3.1 Persistent Laplacians

##### 3.1.1 Spectral graphs

Spectral graph theory studies the spectra of graph Laplacian matrices. It gives rise to the topological and spectral properties of underlying graphs or networks. Mathematically, a graph is an ordered pair  $G(V, E)$ , where  $V = \{v_i; i = 1, 2, \dots, N\}$  is the vertex set with size  $N$  and  $E = \{e_{ij} = (v_i, v_j); i \leq i < j \leq N\}$  is the edge set. Denote  $\deg(v)$  the degree of each vertex  $v_i \in V$ , i.e., the number of edges that connects to  $v$ . A specific Laplacian matrix  $L^G$  can be given by

$$L^G = \begin{cases} \deg(v), & \text{if } v_i = v_j, \\ -1, & \text{if } v_i \text{ and } v_j \text{ are adjacent,} \\ 0, & \text{otherwise,} \end{cases} \quad (2)$$

where ‘‘adjacent’’ is subject to a specific definition or connection rule.

Let order the eigenvalues of the graph Laplacian matrix as

$$\lambda_{\min} = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N = \lambda_{\max}. \quad (3)$$

The kernel dimension of  $L^G$  is the multiplicity of 0 eigenvalues, indicating the number of connected components of  $G(V, E)$ , which is the topological property of the graph. The non-zero eigenvalues of  $L^G$  contain the graph properties. In particular,  $\lambda_2$  is called the algebraic connectivity.

##### 3.1.2 Simplicial complex

To construct a topological description of a graph, simplicial complex is used. For a set of  $q + 1$  points,  $\{v_0, v_1, \dots, v_q\}$ , a  $q$ -plane is well defined if the  $q + 1$  points are affinely independent, i.e.,  $v_1 - v_0, v_2 - v_0, \dots, v_q - v_0$  are linearly independent. Thus, one can have at most  $n$  linearly independent vectors with at most  $n + 1$  affinely independent points in  $\mathbb{R}^n$ . An affine hull is the set of affine combinations,  $v = \sum_{i=0}^q c_i v_i$ ,  $c_i \in \mathbb{R}$ , and  $\sum_{i=0}^q c_i = 1$ . Such an affine combination is a convex combination if all  $c_i$  are non-negative. The convex hull is the set of convex combinations. A  $q$ -simplex denoted as  $\sigma_q$  is the convex hull of  $q + 1$  affinely independent points. For example, 0-, 1-, 2-, and 3-simplex are vertexes, edges, triangles, and tetrahedrons. A simplicial complex  $K$  is a collection of simplices in  $\mathbb{R}^n$  satisfying the following conditions such as the Cech complexes, Vietoris-Rips complexes, and alpha shapes. For example, the Vietoris-Rips complex of  $K$  with radius  $r$  consists of all subsets of radius  $R(\sigma)$  at most  $r$  as

$$\text{VR}(r) = \{\sigma \subseteq K | R(\sigma) \leq r\}. \quad (4)$$

For  $\sigma_q \in K$ , its face  $\sigma_{q-1}$  is also in  $K$ . The non-empty intersection of any two simplices  $\sigma_q, \sigma_p \in K$  is a face of them. The dimension of simplicial complex is defined as the maximum dimension of its simplex.

A  $q$ -chain is a finite sum of simplices as  $\sum_i c_i \sigma_i^k$  with  $\mathbb{Z}_2$  field of the coefficients  $c_i$  for the sum, and the set of all chains in a group  $C_q(K)$ . The boundary operator  $\partial_k$  maps  $C_q(K) \rightarrow C_{q-1}(K)$  defined as

$$\partial_q \sigma_q = \sum_{i=0}^q (-1)^i [v_0, \dots, \hat{v}_i, \dots, v_k] = \sum_{i=0}^q (-1)^i \sigma_{q-1}^i, \quad (5)$$

where  $\sigma_q = [v_0, v_1, \dots, v_k]$  and  $\hat{v}_i$  stands for  $v_i$  being omitted. A  $q$ -chain is called  $q$ -cycle if its boundary is zero. A chain complex is the sequence of chain groups connected by boundary operators

$$\dots \xrightarrow{\partial_{q+2}} C_{q+1}(K) \xrightarrow{\partial_{q+1}} C_q(K) \xrightarrow{\partial_q} C_{q-1}(K) \xrightarrow{\partial_{q-1}} \dots \quad (6)$$

and the  $k$ -th homology group  $H_k$  is defined by  $H_k = Z_k/B_k$  where  $Z_k = \ker \partial_k = \{c \in C_k \mid \partial_k c = 0\}$  and  $B_k = \text{im } \partial_{k+1} = \{\partial_{k+1} c \mid c \in C_{k+1}\}$ . The Betti numbers are defined by the ranks of  $k$ -th homology group  $H_k$ . This, in practice, is counting holes in  $k$ -dimension, such as  $\beta_0$  reflects the number of connected components,  $\beta_1$  gives the number of loops, and  $\beta_2$  is the number of cavities. In a nutshell, the Betti number sequence  $\{\beta_0, \beta_1, \beta_2, \dots\}$  reveals the intrinsic topological property of the system.

Recall that in graph theory, the degree of a vertex (0-simplex)  $v$  is the number of edges that are adjacent to the vertex, denoted as  $\text{deg}(v)$ . However, once we generalize this notion to  $q$ -simplex, problem aroused since  $q$ -simplex can have  $(q-1)$ -simplices and  $(q+1)$ -simplices adjacent to it at the same time. Therefore, the upper adjacency and lower adjacency are required to define the degree of a  $q$ -simplex for  $q > 0$  [29, 43].

**Definition 3.1.** Given two  $q$ -simplices  $\sigma_q^i$  and  $\sigma_q^j$  of a simplicial complex  $K$ . We say they are lower adjacent if they share a common  $(q-1)$ -face, denoted as  $\sigma_q^i \stackrel{L}{\sim} \sigma_q^j$ . The lower degree of  $q$ -simplex is the number of nonempty  $(q-1)$ -simplices in  $K$  that are faces of  $\sigma_q$ , which is denoted as  $\text{deg}_L(\sigma_q)$  and is always  $q+1$ .

**Definition 3.2.** Given two  $q$ -simplices  $\sigma_q^i$  and  $\sigma_q^j$  of a simplicial complex  $K$ . We say they are upper adjacent if they share a common  $(q+1)$ -face, denoted as  $\sigma_q^i \stackrel{U}{\sim} \sigma_q^j$ . The upper degree of  $q$ -simplex is the number of  $(q+1)$ -simplices in  $K$  of which  $\sigma_q$  is a face, which is denoted  $\text{deg}_U(\sigma_q)$ .

Then, the degree of a  $q$ -simplex ( $q > 0$ ) is defined as:

$$\text{deg}(\sigma_q) = \text{deg}_L(\sigma_q) + \text{deg}_U(\sigma_q) = \text{deg}_U(\sigma_q) + q + 1. \quad (7)$$

### 3.1.3 Graph Laplacian

The graph Laplacian was introduced to enrich topological and geometric information of simplicial complexes via a filtration process. The preliminary concepts are about the oriented simplicial complex and  $q$ -combinatorial Laplacian. More detail information can be found elsewhere [17, 19, 21, 29]. The properties of the  $q$ -combinatorial Laplacian matrix with its spectra are discussed in the following.

A  $q$ -combinatorial Laplacian is defined based on oriented simplicial complexes, and its lower- and higher-dimensional simplexes can be employed to study a specifically oriented simplicial complex. An oriented simplicial complex  $K$  is defined if all of its simplices are oriented. If  $\sigma_q^i$  and  $\sigma_q^j$  are upper adjacent with a common upper  $(q+1)$ -simplex  $\tau_{q+1}$ , they are similarly oriented if both have the same sign in  $\partial_{q+1}(\tau_{q+1})$  and dissimilarly oriented if the signs are opposite. Additionally, if  $\sigma_q^i$  and  $\sigma_q^j$  are lower adjacent with a common lower  $(q-1)$ -simplex  $\eta_{q-1}$ , they are similarly oriented if  $\eta_{q-1}$  has the same sign in  $\partial_q(\sigma_q^i)$  and  $\partial_q(\sigma_q^j)$ , and dissimilarly oriented if the signs are opposite. Similarly,  $q$ -chains can be defined on the oriented simplicial complex  $K$ , as well as  $q$ -boundary operator.

The  $q$ -combinatorial Laplacian is a linear operator  $\Delta_q : C_q(K) \rightarrow C_q(K)$  for integer  $q \geq 0$

$$\Delta_q := \partial_{q+1} \partial_{q+1}^* + \partial_q^* \partial_q \quad (8)$$

where  $\partial_q^*$  is the coboundary operator mapping  $\partial_q^* : C_{q-1}(K) \rightarrow C_q(K)$ . One property  $\partial_q \partial_{q+1} = 0$  is preserved, which implies  $\text{Im}(\partial_{q+1}) \subset \ker(\partial_q)$ . The  $q$ -combinatorial Laplacian matrix, denoted  $\mathcal{L}_q$ , is the matrix representation.

$$\mathcal{L}_q = \mathcal{B}_{q+1} \mathcal{B}_{q+1}^T + \mathcal{B}_q^T \mathcal{B}_q \quad (9)$$

of operator  $\Delta_q$ , where  $\mathcal{B}_q$  and  $\mathcal{B}_q^T$  be the matrix representation of a  $q$ -boundary operator and  $q$ -coboundary operator, respectively, with respect to the standard basis for  $C_q(K)$  and  $C_{q-1}(K)$  with some assigned orderings. Then, the number of rows in  $\mathcal{B}_q$  corresponds to the number of  $(q-1)$ -simplices and the number of

columns shows the number of  $q$ -simplices in  $K$ , respectively. In addition, the upper and lower  $q$ -combinatorial Laplacian matrices are denoted by  $\mathcal{L}_q^U = \mathcal{B}_{q+1}\mathcal{B}_{q+1}^T$  and  $\mathcal{L}_q^L = \mathcal{B}_q^T\mathcal{B}_q$ , respectively. Note that  $\partial_0$  is the zero map which leads to  $\mathcal{B}_0$  being a zero matrix. Therefore,  $\mathcal{L}_0(K) = \mathcal{B}_1\mathcal{B}_1^T + \mathcal{B}_0^T\mathcal{B}_0$ , with  $K$  the (oriented) simplicial complex of dimension 1, which is actually a simple graph. Especially, 0-combinatorial Laplacian matrix  $\mathcal{L}_0(K)$  is actually the Laplacian matrix defined in the spectral graph theory.

Given an oriented simplicial complex  $K$  with  $0 \leq q \leq \dim(K)$ , the entries of  $q$ -combinatorial Laplacian matrices are given by [17]

$$q > 0, (\mathcal{L}_q)_{ij} = \begin{cases} \deg(\sigma_q^i), & \text{if } i = j. \\ 1, & \text{if } i \neq j, \sigma_q^i \overset{U}{\approx} \sigma_q^j \text{ and } \sigma_q^i \overset{L}{\approx} \sigma_q^j \text{ with similar orientation.} \\ -1, & \text{if } i \neq j, \sigma_q^i \overset{U}{\approx} \sigma_q^j \text{ and } \sigma_q^i \overset{L}{\approx} \sigma_q^j \text{ with dissimilar orientation.} \\ 0, & \text{if } i \neq j \text{ and either } \sigma_q^i \overset{U}{\approx} \sigma_q^j \text{ or } \sigma_q^i \overset{L}{\approx} \sigma_q^j. \end{cases} \quad (10)$$

$$q = 0, (\mathcal{L}_q)_{ij} = \begin{cases} \deg(\sigma_0^i), & \text{if } i = j. \\ -1, & \text{if } \sigma_0^i \overset{U}{\approx} \sigma_0^j. \\ 0, & \text{otherwise.} \end{cases} \quad (11)$$

### 3.1.4 Persistent spectral graphs

Persistent spectral graphs were introduced by integrating graph Laplacian and multiscale filtration [54]. Both topological and geometric information (i.e. connectivity and robustness of simple graphs) can be derived from analyzing the spectra of  $q$ -combinatorial Laplacian matrix. However, this method is genuinely free of metrics or coordinates, which induced too little topological and geometric information that can be used to describe a single configuration. Therefore, persistent spectral graphs (PSG) is proposed to create a sequence of simplicial complexes induced by varying a filtration parameter, which is inspired by the idea of persistent homology and our earlier work in multiscale graphs. This section mainly introduce the construction of persistent spectral graphs.

First, a  $q$ -combinatorial Laplacian matrix is symmetric and positive semi-definite. Therefore, its eigenvalues are all real and non-negative. The multiplicity of zero spectra (also called harmonic spectra) reveals the topological information, and the geometric information will be preserved in the non-harmonic spectra. More specifically, the multiplicity of zero spectra of  $\mathcal{L}_q(K)$  is denoted by  $\beta_q$  which is actually the  $q$ -th Betti number defined in the homology:

$$\beta_q = \dim(\mathcal{L}_q(K)) - \text{rank}(\mathcal{L}_q(K)) = \text{nullity}(\mathcal{L}_q(K)) = \# \text{ of zero eigenvalues of } \mathcal{L}_q(K). \quad (12)$$

Naturally, persistent spectral theory creates a sequence of simplicial complexes induced by varying a filtration parameter [54]. A filtration of an oriented simplicial complex  $K$  is a sequence of sub-complexes  $(K_t)_{t=0}^m$  of  $K$

$$\emptyset = K_0 \subseteq K_1 \subseteq K_2 \subseteq \cdots \subseteq K_m = K. \quad (13)$$

It induces a sequence of chain complexes

$$\begin{array}{cccccccccccc}
\cdots & C_{q+1}^1 & \xrightarrow{\frac{\partial_{q+1}^1}{\partial_{q+1}^{1*}}} & C_q^1 & \xrightarrow{\frac{\partial_q^1}{\partial_q^{1*}}} & \cdots & \xrightarrow{\frac{\partial_3^1}{\partial_3^{1*}}} & C_2^1 & \xrightarrow{\frac{\partial_2^1}{\partial_2^{1*}}} & C_1^1 & \xrightarrow{\frac{\partial_1^1}{\partial_1^{1*}}} & C_0^1 & \xrightarrow{\frac{\partial_0^1}{\partial_0^{1*}}} & C_{-1}^1 \\
& \cap & & \cap & & & & \cap & & \cap & & \cap & & \\
\cdots & C_{q+1}^2 & \xrightarrow{\frac{\partial_{q+1}^2}{\partial_{q+1}^{2*}}} & C_q^2 & \xrightarrow{\frac{\partial_q^2}{\partial_q^{2*}}} & \cdots & \xrightarrow{\frac{\partial_3^2}{\partial_3^{2*}}} & C_2^2 & \xrightarrow{\frac{\partial_2^2}{\partial_2^{2*}}} & C_1^2 & \xrightarrow{\frac{\partial_1^2}{\partial_1^{2*}}} & C_0^2 & \xrightarrow{\frac{\partial_0^2}{\partial_0^{2*}}} & C_{-1}^2 \\
& \cap & & \cap & & & & \cap & & \cap & & \cap & & \\
& \vdots & & \vdots & & & & \vdots & & \vdots & & \vdots & & \\
& \cap & & \cap & & & & \cap & & \cap & & \cap & & \\
\cdots & C_{q+1}^m & \xrightarrow{\frac{\partial_{q+1}^m}{\partial_{q+1}^{m*}}} & C_q^m & \xrightarrow{\frac{\partial_q^m}{\partial_q^{m*}}} & \cdots & \xrightarrow{\frac{\partial_3^m}{\partial_3^{m*}}} & C_2^m & \xrightarrow{\frac{\partial_2^m}{\partial_2^{m*}}} & C_1^m & \xrightarrow{\frac{\partial_1^m}{\partial_1^{m*}}} & C_0^m & \xrightarrow{\frac{\partial_0^m}{\partial_0^{m*}}} & C_{-1}^m
\end{array} \tag{14}$$

For each sub-complexes  $K_t$ , we define its corresponding chain group to be  $C_q(K_t)$ , and the  $q$ -boundary operator will be denoted by  $\partial_q^t : C_q(K_t) \rightarrow C_{q-1}(K_t)$ . We say that if  $q < 0$ . then  $C_q(K_t)$  is an empty set and  $\partial_q^t$  is a zero map. If  $0 < q \leq \dim(K_t)$ , then

$$\partial_q^t(\sigma_q) = \sum_i^q (-1)^i \sigma_{q-1}^i, \sigma_q \in K_t, \tag{15}$$

with  $\sigma_q = [v_0, \dots, v_q]$  being the  $q$ -simplex, and  $\sigma_{q-1}^i = [v_0, \dots, \hat{v}_i, \dots, v_q]$  being the  $(q-1)$ -simplex for which its vertex  $v_i$  is removed. Additionally, the adjoint operator is  $\partial_q^{t*} : C_{q-1}(K_t) \rightarrow C_q(K_t)$ . The topological and spectral information of  $K_t$  can be analyzed from  $\mathcal{L}_q(K_t)$  along with the filtration parameter by diagonalizing the  $q$ -combinatorial Laplacian matrix. We call the multiplicity of zero spectra of  $\mathcal{L}_q^t$  as its persistent Betti number  $\beta_q^t$ , which counts the number of  $q$ -dimensional holes in  $K_t$ :

$$\beta_q^t = \dim(\mathcal{L}_q^t) - \text{rank}(\mathcal{L}_q^t) = \text{nullity}(\mathcal{L}_q^t) = \#\text{of harmonic spectra of } \mathcal{L}_q^t. \tag{16}$$

Specifically,  $\beta_0^t$  represents the number of connected components in  $K_t$ ,  $\beta_1^t$  reveals the number of one-dimensional loops or circles in  $K_t$ , and  $\beta_2^t$  shows the number of two-dimensional voids or cavities in  $K_t$ . Moreover, the set of spectra of  $\mathcal{L}_q^t$  is given by:

$$\text{Spectra}(\mathcal{L}_q^t) = \{(\lambda_1)_q^t, (\lambda_2)_q^t, \dots, (\lambda_N)_q^t\}, \tag{17}$$

where  $\mathcal{L}_q^t$  has dimension  $N \times N$  and spectra are arranged in ascending order. The smallest non-zero eigenvalue of  $\mathcal{L}_q^t$  is defined as  $(\tilde{\lambda}_2)_q^t$ . The  $p$ -persistent  $q$ -combinatorial Laplacian operator is defined by extending the boundary operator. Detailed descriptions can be found in Ref. [54].

### 3.2 Predictive models for mutation-induced protein-protein binding free energy changes

Since the harmonic spectra produced by the kernel of a persistent Laplacian contain exact topological information as that of persistent homology. As such, we utilize a persistent homology software, GUDHI, to generate purely topological representations of PPIs in dimensions 0, 1, and 2. Additionally, persistent Laplacian spectra, including both harmonic and non-harmonic parts, are coded in Python. Machine learning and deep learning algorithms are implemented in Pytorch [36].

#### 3.2.1 Persistent Laplacian representation of PPIs

To facilitate topological and shape analysis of PPIs via persistent Laplacians, we first composite the atoms in a protein-protein complex into various subsets.

1.  $\mathcal{A}_m$ : atoms of the mutation sites.
2.  $\mathcal{A}_{mn}(r)$ : atoms in the neighbourhood of the mutation site within a cut-off distance  $r$ .
3.  $\mathcal{A}_A(r)$ : protein A atoms within  $r$  of the binding site.
4.  $\mathcal{A}_B(r)$ : protein B atoms within  $r$  of the binding site.
5.  $\mathcal{A}_{\text{ele}}(E)$ : atoms in the system that has atoms of element type  $E$ . The distance matrix is specially designed such that it excludes the interactions between the atoms from the same set. For interactions between atoms  $a_i$  and  $a_j$  in set  $\mathcal{A}$  and/or set  $\mathcal{B}$ , the modified distance is defined as

$$D_{\text{mod}}(a_i, a_j) = \begin{cases} \infty, & \text{if } a_i, a_j \in \mathcal{A}, \text{ or } a_i, a_j \in \mathcal{B}, \\ D_e(a_i, a_j), & \text{if } a_i \in \mathcal{A} \text{ and } a_j \in \mathcal{B}, \end{cases} \quad (18)$$

where  $D_e(a_i, a_j)$  is the Euclidian distance between  $a_i$  and  $a_j$ .

Molecular atoms of different can be constructed as points presented by  $v_0, v_1, v_2, \dots, v_k$  as  $k+1$  affinely independent points in simplicial complex. Persistent spectral graph is devised to track the multiscale topological and geometrical information over different scales along a filtration [54], resulting in significant important feature vectors for the machine learning method. Features generated by binned barcode vectorization can reflect the strength of atom bonds, van der Waals interactions, and can be easily incorporated into a machine learning model, which captures and discriminates local patterns. Using the atom subsets, for example  $\mathcal{A}_A(r)$  and  $\mathcal{A}_B(r)$ , simplicial complexes are constructed by only considering the edges from  $\mathcal{A}_A(r)$  to  $\mathcal{A}_B(r)$  for Vietoris-Rips complexes. Then from the Vietoris-Rips complex filtration, barcodes generated from persistent homology are enumerated by bar lengths in certain intervals with number 0 or 1. Meanwhile, for each complexes in the filtration, eigenvalues are calculated according to the graph Laplacian analysis. The statistics of eigenvalues such as sum, maximum, minimum, mean, and standard deviation are collected to have a normalized features for machine learning methods. Another method of vectorization is to get the statistics of bar lengths, birth values, and death values, such as sum, maximum, minimum, mean, and standard deviation. This method is applied to vectorize Betti-1 ( $H_1$ ) and Betti-2 ( $H_2$ ) barcodes obtained from alpha complex filtration based on the facts that higher-dimensional barcodes are sparser than  $H_0$  barcodes.

### 3.2.2 Machine learning and deep learning algorithms

The features generated from the persistent spectral graph are tested by the gradient boosting tree (GBT) method and the deep neural network (Net) method. The validations are performed on the datasets discussed in the results section. The accurate prediction of the mutation-induced binding affinity changes of protein-protein complexes is very challenging. After effective feature-generations, a machine learning or deep learning model is also required for validations and real applications. The gradient boosting tree is a popular ensemble method for regression and classification problems. It builds a sequence of weak learners to correct training errors. By the assumption that the individual learners are likely to make different mistakes, the method combines weak learners to eliminate the overall error. Furthermore, a decision tree is added to the ensemble depending on the present prediction error on the training dataset. Therefore, this method is relatively robust against hyperparameter tuning and overfitting, especially for a moderate number of features. The GBT is shown for its robustness against overfitting, good performance for moderately small data sizes, and model interpretability. The present work uses the package provided by scikit-learn (v 0.23.0) [37]. The number of estimators and the learning rate are optimized for ensemble methods as 20000 and 0.01, respectively. For each set, ten runs (with different random seeds) were done, and the average result is reported in this work. Considering a large number of features, the maximum number of features to consider is set to the square root of the given descriptor length for GBT methods to accelerate the training process. The parameter setting shows that the performance of the average of sufficient runs is decent.

A deep neural network is a network of neurons that maps an input feature layer to an output layer. The neural network mimics the human brain to solve problems with numerous neuron units with backpropagation to update weights on each layer. To reveal the facts of input features at different levels and abstract more properties, one can construct more layers and more neurons in each layer, known as a deep neural network. Optimization methods for feedforward neural networks and dropout methods are applied to prevent overfitting. The network layers and the number of neurons in each layer are determined by grid searches based on 10-fold cross-validations. Then, the hyperparameters of stochastic gradient descent (SGD) with momentum are set up based on the network structure. The network has 7 layers with 10000 neurons in each layer. For SGD with momentum, the hyperparameters are `momentum = 0.9` and `weight_decay=0`. The learning rate is 0.002 and the epoch is 400. The Net is implemented on Pytorch [36].

### 3.2.3 Predictive models

In our previous work, topology-based deep neural network trained with mAbs (TopNetmAb) was introduced to predict mAb binding free energy changes [8]. Persistent homology is the main workhorse for TopNetmAb, but auxiliary features inherited from our earlier TopNetTree [49] are utilized.

In this work, we construct a TopNet model from TopNetmAb by excluding mAb training data. A topology-based GBT model (TopBGT) is also developed in the present work by replacing Net in the TopNet model with GBT. Both TopNet and TopGBT include a set of auxiliary features inherited from our earlier TopNetTree [49] and TopNetmAb [8] to enhance their performance.

Additionally, to evaluate the performance of persistent Laplacian (Lap) for PPIs, we construct persistent Laplacian-based GBT (LapGBT) and persistent Laplacian-based deep neural network (LapNet). Note that unlike TopNet and TopGBT, LapGBT and LapNet employ only persistent Laplacian features extracted from protein structures. Therefore, their performance depends purely on persistent Laplacian.

Moreover, TopLapGBT and TopLapNet are constructed by adding persistent Laplacian features to TopGBT and TopNet, respectively. Furthermore, the consensus of GBT and Net predictions are also used for validations, denoted as TopNetGBT and LapNetGBT, respectively. Finally, the consensus of TopLapNet and TopLapGBT is called TopLapNetGBT.

## 4 Conclusion

Due to natural selection, emerging SARS-CoV-2 variants are spreading worldwide with their increased transmissibility as a result of higher infectivity and/or stronger antibody resistance. The increase in antibody resistance also leads to vaccine breakthrough infections and jeopardizes the existing monoclonal antibody drugs. The spike protein plays the most important role in viral transmission because its receptor binding domain (RBD) binds to human ACE2 to facilitate the viral entry of host cells. Topological data analysis (TDA) of RBD-ACE2 binding free energy changes induced by RBD mutations enables the accurate forecasting of emerging SARS-CoV-2 variants [9–11, 51].

However, the earlier TDA method is not sensitive to homotopic shape evolution, which is important for protein-protein interactions (PPIs). To overcome this obstacle, persistent Laplacian, which characterizes the topology and shape of data, is introduced in this work for analyzing PPIs. Paired with advanced machine learning and deep learning algorithms, the proposed persistent Laplacian method outperforms the state-of-art approaches in validation with mutation-induced binding free energy changes of PPIs using major benchmark datasets. An important forecasting from the present work is that Omicron subvariants BA.2.11, BA.212.1, BA.4, and BA.5 have high potential to become new dominating variants in the world.

## Acknowledgment

This work was supported in part by NIH grant GM126189, NSF grants DMS-2052983, DMS-1761320, and IIS-1900473, NASA grant 80NSSC21M0023, Michigan Economic Development Corporation, MSU Foundation, Bristol-Myers Squibb 65109, and Pfizer.

## References

- [1] BA2 reinfection. <https://www.timesofisrael.com/several-cases-of-omicron-reinfection-said-detected-in-israel-with-new-ba2-strain/>.
- [2] M. Alenquer, F. Ferreira, D. Lousa, M. Valério, M. Medina-Lopes, M.-L. Bergman, J. Gonçalves, J. Demengeot, R. B. Leite, J. Lilue, et al. Signatures in sars-cov-2 spike protein conferring escape to neutralizing antibodies. *PLoS pathogens*, 17(8):e1009772, 2021.
- [3] P. Bubenik and J. A. Scott. Categorification of persistent homology. *Discrete & Computational Geometry*, 51(3):600–627, 2014.
- [4] Z. Cang, L. Mu, and G.-W. Wei. Representability of algebraic topology for biomolecules in machine learning based scoring and virtual screening. *PLoS computational biology*, 14(1):e1005929, 2018.
- [5] Z. Cang and G.-W. Wei. Topologynet: Topology based deep convolutional and multi-task neural networks for biomolecular property predictions. *PLoS computational biology*, 13(7):e1005690, 2017.
- [6] G. Carlsson. Topology and data. *Bulletin of the American Mathematical Society*, 46(2):255–308, 2009.
- [7] J. Chen, K. Gao, R. Wang, and G.-W. Wei. Prediction and mitigation of mutation threats to COVID-19 vaccines and antibody therapies. *Chemical Science*, 12(20):6929–6948, 2021.
- [8] J. Chen, K. Gao, R. Wang, and G.-W. Wei. Revealing the threat of emerging SARS-CoV-2 mutations to antibody therapies. *Journal of Molecular Biology*, 433(7744), 2021.
- [9] J. Chen, R. Wang, N. B. Gilby, and G.-W. Wei. Omicron variant (b. 1.1. 529): Infectivity, vaccine breakthrough, and antibody resistance. *J Chem Inf Model*, 62(2):412–422, 2022.
- [10] J. Chen, R. Wang, M. Wang, and G.-W. Wei. Mutations strengthened SARS-CoV-2 infectivity. *J. Mol. Biol.*, 432(19):5212–5226, 2020.
- [11] J. Chen and G.-W. Wei. Omicron ba. 2 (b. 1.1. 529.2): High potential for becoming the next dominant variant. *The Journal of Physical Chemistry Letters*, 13:3840–3849, 2022.
- [12] V. De Silva and R. Ghrist. Coverage in sensor networks via persistent homology. *Algebraic & Geometric Topology*, 7(1):339–358, 2007.
- [13] T. K. Dey, F. Fan, and Y. Wang. Computing topological persistence for simplicial maps. In *Proceedings of the thirtieth annual symposium on Computational geometry*, page 345. ACM, 2014.
- [14] L. Dupont, L. B. Snell, C. Graham, J. Seow, B. Merrick, T. Lechmere, T. J. Maguire, S. R. Hallett, S. Pickering, T. Charalampous, et al. Neutralizing antibody activity in convalescent sera from infection in humans with sars-cov-2 and variants of concern. *Nature microbiology*, pages 1–10, 2021.
- [15] H. Edelsbrunner, D. Letscher, and A. Zomorodian. Topological persistence and simplification. In *Proceedings 41st annual symposium on foundations of computer science*, pages 454–463. IEEE, 2000.
- [16] P. Frosini. Measuring shapes by size functions. In *Intelligent Robots and Computer Vision X: Algorithms and Techniques*, volume 1607, pages 122–133. International Society for Optics and Photonics, 1992.
- [17] T. E. Goldberg. Combinatorial laplacians of simplicial complexes. *Senior Thesis, Bard College*, 2002.
- [18] R. Guerois, J. E. Nielsen, and L. Serrano. Predicting changes in the stability of proteins and protein complexes: a study of more than 1000 mutations. *Journal of molecular biology*, 320(2):369–387, 2002.
- [19] D. Hernández Serrano and D. Sánchez Gómez. Higher order degree in simplicial complexes, multi combinatorial laplacian and applications of tda to complex networks. *arXiv preprint arXiv:1908.02583*, 2019.

- [20] M. Hoffmann, H. Kleine-Weber, S. Schroeder, N. Krüger, T. Herrler, S. Erichsen, T. S. Schiergens, G. Herrler, N.-H. Wu, A. Nitsche, et al. SARS-CoV-2 cell entry depends on ACE2 and TMPRSS2 and is blocked by a clinically proven protease inhibitor. *Cell*, 181(2):271–280, 2020.
- [21] D. Horak and J. Jost. Spectra of combinatorial laplace operators on simplicial complexes. *Advances in Mathematics*, 244:303–336, 2013.
- [22] J. Jankauskaitė, B. Jiménez-García, J. Dapkūnas, J. Fernández-Recio, and I. H. Moal. SKEMPI 2.0: an updated benchmark of changes in protein–protein binding energy, kinetics and thermodynamics upon mutation. *Bioinformatics*, 35(3):462–469, 2019.
- [23] S. Jemimah, K. Yugandhar, and M. Michael Gromiha. Proximate: a database of mutant protein–protein complex thermodynamics and kinetics. *Bioinformatics*, 33(17):2787–2788, 2017.
- [24] C. Li, X. Tian, X. Jia, J. Wan, L. Lu, S. Jiang, F. Lan, Y. Lu, Y. Wu, and T. Ying. The impact of receptor-binding domain natural mutations on antibody recognition of SARS-CoV-2. *Signal Transduct. Target. Ther.*, 6(1):1–3, 2021.
- [25] W. Li, Z. Shi, M. Yu, W. Ren, C. Smith, J. H. Epstein, H. Wang, G. Crameri, Z. Hu, H. Zhang, et al. Bats are natural reservoirs of SARS-like coronaviruses. *Science*, 310(5748):676–679, 2005.
- [26] T. W. Linsky, R. Vergara, N. Codina, J. W. Nelson, M. J. Walker, W. Su, C. O. Barnes, T.-Y. Hsiang, K. Esser-Nobis, K. Yu, et al. De novo design of potent and resilient hACE2 decoys to neutralize SARS-CoV-2. *Science*, 370(6521):1208–1214, 2020.
- [27] Q. Liu, P. Chen, B. Wang, J. Zhang, and J. Li. dbmpikt: a database of kinetic and thermodynamic mutant protein interactions. *Bmc Bioinformatics*, 19(1):1–7, 2018.
- [28] F. P. Lyngse, C. T. Kirkeby, M. Denwood, L. E. Christiansen, K. Mølbak, C. H. Møller, R. L. Skov, T. G. Krause, M. Rasmussen, R. N. Sieber, et al. Transmission of sars-cov-2 omicron voc subvariants ba. 1 and ba. 2: Evidence from danish households. *medRxiv*, 2022.
- [29] S. Maletić and M. Rajković. Consensus formation on a simplicial complex of opinions. *Physica A: Statistical Mechanics and its Applications*, 397(March):111–120, 2014.
- [30] D. Mannar, J. W. Saville, X. Zhu, S. S. Srivastava, A. M. Berezuk, K. Tuttle, C. Marquez, I. Sekirov, and S. Subramaniam. Sars-cov-2 omicron variant: Ace2 binding, cryo-em structure of spike protein-ace2 complex and antibody evasion. *BioRxiv*, 2021.
- [31] F. Mémoli, Z. Wan, and Y. Wang. Persistent laplacians: Properties, algorithms and implications. *arXiv preprint arXiv:2012.02808*, 2020.
- [32] Z. Meng and K. Xia. Persistent spectral–based machine learning (perspect ml) for protein-ligand binding affinity prediction. *Science Advances*, 7(19):eabc5329, 2021.
- [33] K. Mischaikow and V. Nanda. Morse theory for filtrations and efficient computation of persistent homology. *Discrete & Computational Geometry*, 50(2):330–353, 2013.
- [34] I. H. Moal and J. Fernández-Recio. SKEMPI: a structural kinetic and energetic database of mutant protein interactions and its use in empirical models. *Bioinformatics*, 28(20):2600–2607, 2012.
- [35] D. D. Nguyen, Z. Cang, K. Wu, M. Wang, Y. Cao, and G.-W. Wei. Mathematical deep learning for pose and binding affinity prediction and ranking in D3R Grand Challenges. *Journal of Computer-aided Molecular Design*, 33(1):71–82, 2019.
- [36] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.

- [37] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.
- [38] M. Petukh, L. Dai, and E. Alexov. Saambe: webserver to predict the charge of binding free energy caused by amino acids mutations. *International journal of molecular sciences*, 17(4):547, 2016.
- [39] D. E. Pires and D. B. Ascher. mcsm-ab: a web server for predicting antibody–antigen affinity changes upon mutation with graph-based signatures. *Nucleic acids research*, 44(W1):W469–W473, 2016.
- [40] E. Procko. The sequence of human ace2 is suboptimal for binding the s spike protein of sars coronavirus 2. *BioRxiv*, 2020.
- [41] X.-X. Qu, P. Hao, X.-J. Song, S.-M. Jiang, Y.-X. Liu, P.-G. Wang, X. Rao, H.-D. Song, S.-Y. Wang, Y. Zuo, et al. Identification of two critical amino acid residues of the severe acute respiratory syndrome coronavirus spike protein for its variation in zoonotic tropism transition via a double substitution strategy. *Journal of Biological Chemistry*, 280(33):29588–29595, 2005.
- [42] C. H. Rodrigues, Y. Myung, D. E. Pires, and D. B. Ascher. mcsm-ppi2: predicting the effects of mutations on protein–protein interactions. *Nucleic acids research*, 47(W1):W338–W344, 2019.
- [43] D. H. Serrano and D. S. Gómez. Centrality measures in simplicial complexes: applications of tda to network science. *arXiv preprint arXiv:1908.02967*, 2019.
- [44] S. Sirin, J. R. Apgar, E. M. Bennett, and A. E. Keating. AB-Bind: antibody binding mutational database for computational affinity predictions. *Protein Science*, 25(2):393–409, 2016.
- [45] H.-D. Song, C.-C. Tu, G.-W. Zhang, S.-Y. Wang, K. Zheng, L.-C. Lei, Q.-X. Chen, Y.-W. Gao, H.-Q. Zhou, H. Xiang, et al. Cross-host evolution of severe acute respiratory syndrome coronavirus in palm civet and human. *Proceedings of the National Academy of Sciences*, 102(7):2430–2435, 2005.
- [46] T. N. Starr, A. J. Greaney, S. K. Hilton, D. Ellis, K. H. Crawford, A. S. Dingens, M. J. Navarro, J. E. Bowen, M. A. Tortorici, A. C. Walls, et al. Deep mutational scanning of SARS-CoV-2 receptor binding domain reveals constraints on folding and ACE2 binding. *Cell*, 182(5):1295–1310, 2020.
- [47] A. C. Walls, Y.-J. Park, M. A. Tortorici, A. Wall, A. T. McGuire, and D. Veelsler. Structure, function, and antigenicity of the SARS-CoV-2 spike glycoprotein. *Cell*, 2020.
- [48] C. Wang, W. Li, D. Drabek, N. M. Okba, R. van Haperen, A. D. Osterhaus, F. J. van Kuppeveld, B. L. Haagmans, F. Grosveld, and B.-J. Bosch. A human monoclonal antibody blocking SARS-CoV-2 infection. *Nature communications*, 11(1):1–6, 2020.
- [49] M. Wang, Z. Cang, and G.-W. Wei. A topology-based network tree for the prediction of protein–protein binding affinity changes following mutation. *Nature Machine Intelligence*, 2(2):116–123, 2020.
- [50] R. Wang, J. Chen, K. Gao, and G.-W. Wei. Vaccine-escape and fast-growing mutations in the United Kingdom, the United States, Singapore, Spain, India, and other COVID-19-devastated countries. *Genomics*, 113(4):2158–2170, 2021.
- [51] R. Wang, J. Chen, Y. Hozumi, C. Yin, and G.-W. Wei. Emerging vaccine-breakthrough SARS-CoV-2 variants. *ACS Infect. Dis.*, 8(3):546–556, 2021.
- [52] R. Wang, J. Chen, Y. Hozumi, C. Yin, and G.-W. Wei. Emerging vaccine-breakthrough SARS-CoV-2 variants. *ACS Infect. Dis.*, 8(3):546–556, 2022.
- [53] R. Wang, J. Chen, and G.-W. Wei. Mechanisms of sars-cov-2 evolution revealing vaccine-resistant mutations in europe and america. *The Journal of Physical Chemistry Letters*, 12:11850–11857, 2021.

- [54] R. Wang, D. D. Nguyen, and G.-W. Wei. Persistent spectral graph. *International journal for numerical methods in biomedical engineering*, 36(9):e3376, 2020.
- [55] R. Wang, R. Zhao, E. Ribando-Gros, J. Chen, Y. Tong, and G.-W. Wei. Hermes: Persistent spectral graph software. *Foundations of data science (Springfield, Mo.)*, 3(1):67, 2021.
- [56] J. Wee and K. Xia. Persistent spectral based ensemble learning (perspect-el) for protein–protein binding affinity prediction. *Briefings in Bioinformatics*, 23(2), 2022.
- [57] K. L. Xia and G. W. Wei. Persistent homology analysis of protein structure, flexibility and folding. *International Journal for Numerical Methods in Biomedical Engineering*, 30:814–844, 2014.
- [58] P. Xiong, C. Zhang, W. Zheng, and Y. Zhang. Bindprofx: assessing mutation-induced binding affinity change by protein interface profiles with pseudo-counts. *Journal of molecular biology*, 429(3):426–434, 2017.
- [59] Y. Yao, J. Sun, X. H. Huang, G. R. Bowman, G. Singh, M. Lesnick, L. J. Guibas, V. S. Pande, and G. Carlsson. Topological methods for exploring low-density states in biomolecular folding pathways. *The Journal of Chemical Physics*, 130:144115, 2009.
- [60] F. Yu, R. Xiang, X. Deng, L. Wang, Z. Yu, S. Tian, R. Liang, Y. Li, T. Ying, and S. Jiang. Receptor-binding domain-specific human neutralizing monoclonal antibodies against SARS-CoV and SARS-CoV-2. *Signal Transduction and Targeted Therapy*, 5(1):1–12, 2020.
- [61] A. Zomorodian and G. Carlsson. Computing persistent homology. *Discrete & Computational Geometry*, 33(2):249–274, 2005.