

EquiPocket: an E(3)-Equivariant Geometric Graph Neural Network for Ligand Binding Site Prediction

Yang Zhang^{a,c}, Wenbing Huang^{a*}, Zhewei Wei^{a*}, Ye Yuan^b, Zhaohan Ding^c

^aRenmin University of China, ^bBeijing Institute of Technology, ^cDP Technology
China

fengyuewuya@ruc.edu.cn, hwenbing@126.com, zhewei@ruc.edu.cn, yuan-ye@bit.edu.cn, dingzh@dp.tech

ABSTRACT

Predicting the binding sites of the target proteins plays a fundamental role in drug discovery. Most existing deep-learning methods consider a protein as a 3D image by spatially clustering its atoms into voxels and then feed the voxelized protein into a 3D CNN for prediction. However, the CNN-based methods encounter several critical issues: 1) defective in representing irregular protein structures; 2) sensitive to rotations; 3) insufficient to characterize the protein surface; 4) unaware of data distribution shift. To address the above issues, this work proposes EquiPocket, an E(3)-equivariant Graph Neural Network (GNN) for binding site prediction. In particular, EquiPocket consists of three modules: the first one to extract local geometric information for each surface atom, the second one to model both the chemical and spatial structure of the protein, and the last one to capture the geometry of the surface via equivariant message passing over the surface atoms. We further propose a dense attention output layer to better alleviate the data distribution shift effect incurred by the variable protein size. Extensive experiments on several representative benchmarks demonstrate the superiority of our framework to the state-of-the-art methods.

CCS CONCEPTS

• Applied computing → Bioinformatics; • Computing methodologies → Artificial intelligence.

KEYWORDS

Binding Site Prediction, Graph Neural Network, Drug Discovery

ACM Reference Format:

Yang Zhang^{a,c}, Wenbing Huang^{a*}, Zhewei Wei^{a*}, Ye Yuan^b, Zhaohan Ding^c. 2018. EquiPocket: an E(3)-Equivariant Geometric Graph Neural Network for Ligand Binding Site Prediction . In *Proceedings of (Conference acronym 'XX)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/XXXXXX.XXXXXXXX>

1 INTRODUCTION

Nearly all biological and pharmacological processes in living systems involve interactions between receptors (*i.e.* target proteins)

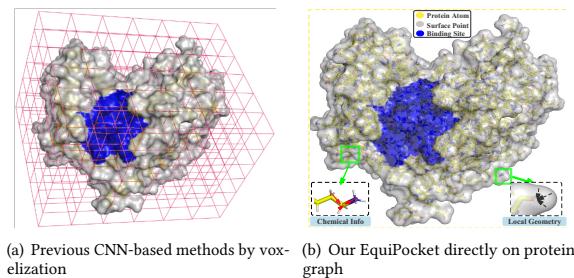
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/XXXXXX.XXXXXXXX>



(a) Previous CNN-based methods by voxelization (b) Our EquiPocket directly on protein graph

Figure 1: Illustrative comparison between previous CNN-based methods and our EquiPocket.

and ligands (*i.e.* small molecules or other proteins) [34]. The places where such interactions occur are known as the binding sites/pockets of the ligands on the target protein structures, which are essential to determine whether or not the ligands are druggable and functionally relevant. Moreover, the knowledge of the binding site is able to facilitate many downstream tasks, such as docking [50], the design of drug molecules [47], etc. Therefore, predicting the binding sites of the target proteins via in-silico algorithms forms an indispensable and even the first step in drug discovery.

Through the past years, plenty of computational methods have been proposed to detect binding sites, which can be roughly classified into three categories [30]: the geometry-based [4, 9, 13, 27, 28, 48], the probe-energy-based [10, 25, 26, 32] and the template-based [3, 44]. The computational methods exploit hand-crafted algorithmic procedures guided by domain knowledge or external templates, leading to insufficient expressivity in representing complicated proteins. Recently, with the accumulation of labeled data and the development of machine learning techniques, the learning-based approaches have been developed [24], which manage to analyze and extract the underlying patterns of the input data that eventually align with the assign labels through the iterative process of learning. Although the learning-based methods have exhibited clear superiority over the classical computational counterparts, their performance and flexibility are still limited by the hand-crafted input features and the insufficiently-expressive models they used [24].

More recently, motivated by the breakthrough of deep learning in a variety of fields, Convolutional Neural Networks (CNNs) have been applied successfully for the prediction of binding sites [18]. Typical works include DeepSite [17], DeepPocket [1], DeepSurf [31], etc. The CNN-based methods consider a protein as a 3D image by spatially clustering its atoms into the nearest voxels, and then model

the binding site prediction as a object detection problem or a semantic segmentation task on 3D grids. Thanks to the multi-layer and end-to-end learning paradigm, the CNN-based methods are observed to outperform traditional learning-based approaches and generally achieve the best performance on several public benchmarks [31].

In spite of the impressive progress, existing CNN-based models still encounter several issues as below:

Issue 1. Defective in leveraging regular voxels to model the proteins of irregular shape. First, a considerable number of voxels probably contain no atom since the protein atoms are unevenly distributed in space, which yields unnecessary redundancy in computation and memory. Moreover, the voxelization is usually constrained within a fixed-size space (e.g. $70\text{\AA} \times 70\text{\AA} \times 70\text{\AA}$) [17, 41]. The atoms beyond this size will be directly discarded, resulting in incomplete and inaccurate modeling particularly for large proteins. Besides, although the voxelization process is able to encode certain spatial structure of the protein, it overlooks the irregular chemical interactions (*i.e.* the chemical bonds) between atoms and the topological structure upon that is also useful for binding site detection.

Issue 2. Sensitive to rotations. To discretize the protein into 3D grids, the CNN methods fix the three bases of the coordinates beforehand. When rotating the protein, the voxelization results could be distinct, and the predicted binding sites will change, which, however, conflicts with the fact that any rotation of the protein keeps the binding sites invariant. While it can be alleviated by the local grid [31] or augmenting training data with random rotations [33, 41], which yet is data-dependent and unable to guarantee rotation invariance in theory.

Issue 3. Insufficient to characterize the geometry of the protein surface. The surface atoms comprise the major part of the binding pocket, which should be elaborately modeled. In the CNN-based methods, the surface atoms are located in the voxels that are surrounded by empty voxels, which somehow encodes the surface geometry. Nevertheless, such information is coarse to depict how the surface atoms interact and what their local geometry is. Indeed, the description of the surface atoms is purely driven by the geometric shape of the solvent-accessible surface of the protein [35] (Figure 1(b)), which, unfortunately, is less explored in current works.

Issue 4. Unaware of data distribution shift. In practical scenarios, the size of the proteins varies greatly across different datasets. It

requires the deep learning model we apply to be well generalizable and adaptive, so that it is able to overcome the distribution shift incurred by the variable protein size. However, this point is not seriously discussed previously.

In this paper, to address the above issues, we propose to apply Graph Neural Networks (GNNs) [5, 19, 39] instead of CNNs to represent proteins. By considering atoms as nodes, interactions as edges, GNNs are able to encode the irregular protein structures by multi-layer message passing. More importantly, a recent line of researches [12, 14, 39] has enhanced GNNs by encapsulating E(3) equivariance/invariance with respect to translations/rotations; in this way, equivariant GNNs yield outputs that are independent of the choice of the coordinate systems, leading to improved generalization ability. That being said, trivially applying equivariant GNNs for the binding site prediction task is still incapable of providing desirable performance, and even achieves worse accuracy than the

CNN-based counterparts. By looking into the design of the architecture, equivariant GNNs naturally cope with the first two issues as mentioned above, yet leave the other two unsolved. To this end, we make the contributions as follows:

- To the best of our knowledge, we are the first to apply an E(3)-equivariant GNN for ligand binding site prediction, which is dubbed **EquiPocket**. In contrast to conventional CNN-based methods, EquiPocket is free of the voxelization process, able to model irregular protein structures by nature, and insensitive to any Euclidean transformation, thereby addressing Issue 1 and 2.
- EquiPocket consists of three modules: the first one to extract local geometric information for each surface atom with the help of solvent-accessible surface [35], the second one to model both the chemical and spatial structure of the protein, and the last one to capture the comprehensive geometry of the surface via equivariant message passing over the surface atoms. The first and last module are proposed to tackle Issue 3, while the second module attempts to involve both the chemical and spatial interactions, as presented in Issue 1.
- To resolve Issue 4, namely, alleviating the effect by data distribution shift, we further propose a novel output layer called *dense attention output layer* in Equiocket, which enables us to adaptively balance the scope of the receptive field for each atom in accordance to the density distribution of the neighbor atoms.
- Extensive experiments on severral representative benchmarks demonstrate the superiority of our framework to the state-of-the-art methods in prediction accuracy. The design of our model is sufficiently ablated as well.

2 RELATED WORK

2.1 Binding Site Prediction

Computational Methods. The computational methods for binding site prediction include geometry-based [4, 9, 13, 27, 28, 48], probe- and energy-based [25, 26, 32] and template-based [3, 44] methods: 1) Since most ligand binding sites occur on the 3D structure, geometry-based methods (POCKET [28], CriticalFinder [9], LigSite [13], Fpocket [27], etc.) are designed to identify these hollow spaces and then rank them using the expert design geometric features. 2) Probe-based methods (SURFNET [25], Q-SiteFinder [26], etc. [10]), also known as energy-based methods, calculate the energy resulting from the interaction between protein atoms and a small-molecule probe, whose value dictates the existence of binding sites. 3)Template-based methods (FINDSITE [3], LIBRA [44], etc.) are mainly to compare the required query protein with the published protein structure database to identify the binding sites.

Traditional Learning-based Methods. PRANK [23] is a learning-based method that employs the traditional machine learning algorithm random forest(RF) [2]. Based on the pocket points and chemical properties from Fpocket [27] and Concavity [4], this method measures the "ligandibility" as the binding ability of a candidate pocket using the RF model. However, those methods require the manual extraction of numerous features with limit upgrading.

CNN-based Methods. Over the last few years, deep learning has surpassed far more traditional ML methods in many domains.

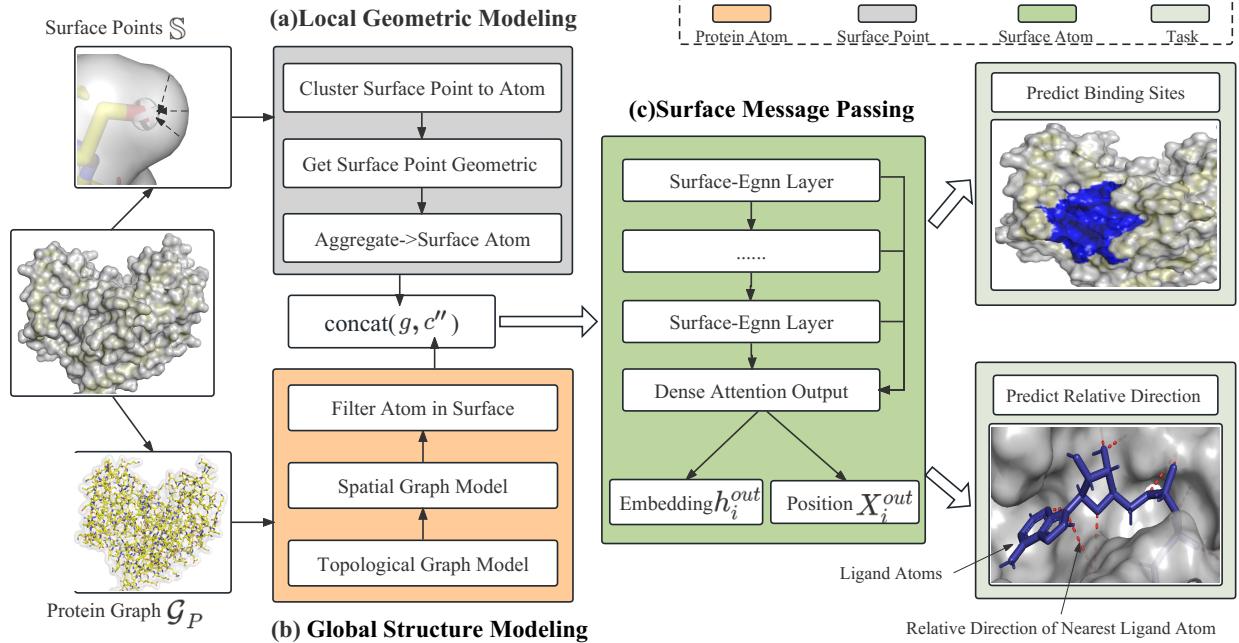


Figure 2: An illustration of the scheme of our EquiPocket framework.

For binding site prediction task, many researchers [1, 17, 18, 31, 41] regard a protein as a 3D image, and model this task as a computer vision problem. DeepSite [17] is the first attempt to employ the CNN architecture for binding site prediction, which like P2Rank [24] treats this task as a binary classification problem and converts a protein to 3D voxelized grids. The methods FRSite [16] and Kalasany [41] adhere to the principle of deepsite, but the former regards this task as an object detection problem, and the latter regards this task as a semantic segmentation task. Deppocket [1] is a method similar to p2rank, but implements a CNN-based segmentation model as the scoring function in order to more precisely locate the the binding sites. The recent CNN-based method DeepSurf [31] constructs a local 3D grid and updates the 3D-CNN architecture to mitigate the detrimental effects of protein rotation.

2.2 Graph Neural Networks for Molecule Modeling

There are multi-level information in molecules including atom info, chemical bonds, spatial structure, physical constraints, etc. Numerous researchers view molecules as topological structures and apply topological-based GNN models (like graph2vec [11], GAT [45], GCN [19], GCN2 [5], GIN [49] and etc. [42]) to extract the chemical info, which achieve positive outcomes. With the accumulation of structure data for molecules, spatial-based graph models (DimeNet [22], DimeNet++ [21], SphereNet [29], SchNet [40], Egnn [39], [14], [12] and etc.) are proposed for molecule task which aggregates spatial and topological information. However, these models may not be adequate for macro-molecules due to their high calculation and resource requirements.

3 NOTATIONS AND DEFINITIONS

Protein Graph. A protein such as the example in Figure 1(b) is denoted as a graph $\mathcal{G}_P = (\mathcal{V}_P, \mathcal{E}_C, \mathcal{E}_D)$, where $\mathcal{V}_P = \{v_0, \dots, v_N\}$ forms the set of N atoms, \mathcal{E}_C represents the chemical-bond edges, and \mathcal{E}_D collects the spatial edges between any two atoms if their spatial distance is less than a cutoff $\theta > 0$. In particular, each node (*i.e.* atom) is associated with a feature (x_i, c_i) , where $x_i \in \mathbb{R}^3$ denotes the 3D coordinates and $c_i \in \mathbb{R}^5$ is the chemical feature.

Surface Point Set. The surface geometry of a protein is of crucial interest for binding site detection. Here we define the set of surface points, by $\mathbb{S} = \{s_0, \dots, s_M\}$, $M \gg N$. Each surface point s_i is NOT necessarily an atom of the protein, and it corresponds to (x_i, p_i) , where $x_i \in \mathbb{R}^3$ represents the 3D coordinates of s_i and $p_i \in \mathcal{V}_P$ indicates the index of the nearest protein atom in \mathcal{V}_P to s_i . We employ the open source MSMS [37] to derive surface points.

Protein Surface Graph. By referring to the surface points defined above, we collect all the nearest protein atoms p_i of the surface points, giving rise to the surface graph $\mathcal{G}_S = (\mathcal{V}_S, \mathcal{E}_S)$, and clearly $\mathcal{G}_S \subseteq \mathcal{G}_P$. We call the atoms in the surface graph as *surface atoms*, which are distinguished from surface points defined in the last paragraph. Notably, the edges of the surface graph, *i.e.*, \mathcal{E}_S is only composed of spatial edges from \mathcal{E}_D , since those chemical edges are mostly broken among the extracted atoms.

Equivariance and Invariance. In 3D space, the symmetry of the physical laws requires the detection model to be equivariant with respect to arbitrary coordinate systems [12]. In form, suppose X to be 3D geometric vectors (positions, velocities, etc) that are steerable by $E(3)$ group (rotations/translations/reflections), and \mathbf{h} non-steerable

features. The function f is $E(3)$ -equivariant, if for any transformation $g \in E(3)$, $f(g \cdot X, \mathbf{h}) = g \cdot f(X, \mathbf{h})$, $\forall X \in \mathbb{R}^{3 \times m}, \mathbf{h} \in \mathbb{R}^d$. Similarly, f is invariant if $f(g \cdot X, \mathbf{h}) = f(X, \mathbf{h})$. The group action \cdot is instantiated as $g \cdot X := X + \mathbf{b}$ for translation $\mathbf{b} \in \mathbb{R}^3$ and $g \cdot X(t) := OX$ for rotation/reflection $O \in \mathbb{R}^{3 \times 3}$.

Problem Statement. Given a protein \mathcal{G}_P and its surface points \mathbb{S} , as well as the constructed surface graph \mathcal{G}_S , our goal is to learn an $E(3)$ -invariant model $f(\mathcal{G}_P, \mathbb{S}, \mathcal{G}_S)$ to predict the atoms of the binding site: $\mathcal{V}_B \subseteq \mathcal{V}_P$.

4 THE PROPOSED METHODOLOGY

Figure 2 illustrates the overall framework of our EquiPocket, which consists of three modules: the *local geometric modeling module* § 4.1 that focuses on extracting the geometric information of each surface atom, the *global structure modeling module* § 4.2 to characterize both the chemical and spatial structures of the protein, and the *surface message passing module* § 4.3 which concentrates on capturing the entire surface geometry based on the extracted information by the two former modules. The training losses are also presented. We defer the pseudo codes of EquiPocket to Appendix 1.

4.1 Local Geometric Modeling Module

This subsection presents how to extract the local geometric information of the protein surface \mathcal{G}_S , with the help of surface points \mathbb{S} . The local geometry of each protein atom closely determines if the region nearby is appropriate or not to become part of binding sites. We adopt the surrounding surface points of each protein surface atom to describe the local geometry.

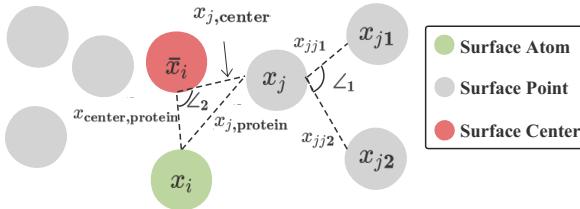


Figure 3: An illustration of local geometric features.

To be specific, for every surface atom $i \in \mathcal{V}_S$, its surrounding surface points are returned by a subset of \mathbb{S} , namely, $\mathbb{S}_i = \{s_j = (x_j, p_j) \in \mathbb{S} \mid p_j = i\}$, where p_j , as defined before, indicates the nearest protein atom. We now construct the geometric information based on \mathbb{S}_i . We denote the center/mean of all 3D coordinates in \mathbb{S}_i as \bar{x}_i . For each surrounding surface point $s_j \in \mathbb{S}_i$, we first search its two nearest surface points from \mathbb{S} as s_{j_1} and s_{j_2} , and then calculate the following relative position vectors:

$$\begin{cases} x_{jj_1} = x_j - x_{j_1}, \\ x_{jj_2} = x_j - x_{j_2}, \\ x_{j,center} = x_j - \bar{x}_i, \\ x_{j,protein} = x_j - x_i, \\ x_{center,protein} = \bar{x}_i - x_i. \end{cases} \quad (1)$$

We further derive the following scalars upon Eq. 1:

$$\begin{aligned} \mathbf{g}(s_j) &:= [\|x_{jj_1}\|_2, \|x_{jj_2}\|_2, \angle_1, \\ &\quad \|x_{j,center}\|_2, \|x_{j,protein}\|_2, \|x_{protein,center}\|_2, \angle_2], \end{aligned} \quad (2)$$

where the angels are computed by $\angle_1 = \frac{\mathbf{x}_{jj_1} \cdot \mathbf{x}_{jj_2}}{\|\mathbf{x}_{jj_1}\|_2 \|\mathbf{x}_{jj_2}\|_2}$ and $\angle_2 = \frac{\mathbf{x}_{j,center} \cdot \mathbf{x}_{center,protein}}{\|\mathbf{x}_{j,center}\|_2 \|\mathbf{x}_{center,protein}\|_2}$; here the operator \cdot defines the inner-product between two vectors. Basically, as displayed in Figure 3, the first three quantities in $\mathbf{g}(s_j)$ depict how the nearby surface points are arranged around s_j , and the last four ones describe where s_j is located within the global region of \mathbb{S}_i .

We aggregate the geometric information $\mathbf{g}(s_j)$ over all surface points in \mathbb{S}_i and obtain a readout descriptor for surface atom i as

$$\begin{aligned} \mathbf{g}_i &= [\text{Pooling}(\{\mathbf{g}(s_j)\}_{s_j \in \mathbb{S}_i}), \\ &\quad \text{MLP}(\text{Pooling}(\{\mathbf{g}(s_j)\}_{s_j \in \mathbb{S}_i}))] \end{aligned} \quad (3)$$

Here, MLP denotes multi-layer perceptron, and the function Pooling is implemented as a concatenation of mean pooling and max pooling throughout our experiments. The front part in Eq. 3 is used to gather local geometric features, while the latter part attempts to compute the global size of surrounding surface points. Notably, the geometric descriptor \mathbf{g}_i is $E(3)$ -invariant.

4.2 Global Structure Modeling Module

This module aims at processing the information of the whole protein \mathcal{G}_P , including atom type, chemical bonds, relevant spatial positions, etc. Although the binding pocket is majorly comprised of surface atoms, the global structure of the protein in general influences how the ligand is interacted with and how the pocket is formulated, which should be modeled. We fulfil this purpose via two concatenated processes: chemical-graph modeling and spatial-graph modeling.

The chemical-graph modeling process copes with the chemical features $\{\mathbf{c}_i\}_{i \in \mathcal{V}_P}$ and the chemical interactions \mathcal{E}_C of the protein graph. For each atom in the protein, its chemical type, the numbers of electrons around, and the chemical bonds connected to other atoms are important clues to identify the interaction between the protein and the ligand [50]. We employ typical GNNs [20, 36, 45] to distill this type of information. Formally, we proceed:

$$\mathbf{c}'_i \}_{i \in \mathcal{V}_P} = \text{GNN}(\{\mathbf{c}_i\}_{i \in \mathcal{V}_P}, \mathcal{E}_C), \quad (4)$$

where \mathbf{c}'_i is the updated chemical feature for atom v_i . While various GNNs can be used in Eq. 4, here we implement GAT [45] given its desirable performance observed in our experiments.

The spatial-graph modeling process further involves the 3D coordinates $\{\mathbf{x}_i\}_{i \in \mathcal{V}_P}$ to better depict the spatial interactions \mathcal{E}_D within the protein. Different from chemical features \mathbf{c}'_i , the 3D coordinates provide the spatial position of each atom and reflect the pair-wise distances in 3D space, which is helpful for physical interaction modeling. We leverage EGNN [39] as it conforms to $E(3)$ equivariance/invariance and achieves promising performance on modeling spatial graphs. Specifically, we process EGNN as follows:

$$\{\mathbf{c}''_i\}_{i \in \mathcal{V}_P} = \text{EGNN}(\{\mathbf{x}_i, \mathbf{c}'_i\}_{i \in \mathcal{V}_P}, \mathcal{E}_D). \quad (5)$$

Here, we only reserve the invariant output (*i.e.*, \mathbf{c}''_i) and have discarded the equivariant output (*e.g.* updated 3D coordinates) of EGNN, since the goal of this module is to provide invariant features.

We select the updated features of the surface atoms \mathcal{V}_S , which will be fed into the module in § 4.3.

4.3 Surface Message Passing Module.

Given the local geometric features $\{g_i\}_{i \in \mathcal{V}_S}$ from § 4.1, and the globally-encoded features of the surface atoms $\{c''_i\}_{i \in \mathcal{V}_S}$ from § 4.2, the module in this subsection carries out equivariant message passing on the surface graph \mathcal{G}_S to renew the entire features of the protein surface. We mainly focus on the surface atoms here, because firstly the surface atoms are more relevant to the binding sites than the interior atoms, and secondly the features $\{c''_i\}_{i \in \mathcal{V}_S}$ that are considered as the input have somehow encoded the information of the interior structure via the processes in 4.2.

Surface-EGNN. During the l -th layer message passing, each node is associated with an invariant feature $\mathbf{h}_i^{(l)} \in \mathbb{R}^{m_l}$ and an equivariant double-channel matrix $X_i^{(l)} \in \mathbb{R}^{3 \times 2}$. We first concatenate c''_i with g_i as the initial invariant feature:

$$\mathbf{h}_i^{(0)} = [c''_i, g_i]. \quad (6)$$

The equivariant matrix $X_i^{(0)}$ is initialized by the 3D coordinates of the atom and the center of its surrounding surface points, that is,

$$X_i^{(0)} = [\mathbf{x}_i, \bar{\mathbf{x}}_i]. \quad (7)$$

We update $\mathbf{h}_i^{(l)} \in \mathbb{R}^{d_l}$ and $X_i^{(l)} \in \mathbb{R}^{3 \times 2}$ synchronously to unveil both the topological and geometrical patterns. Inspired from EGNN [39] and its multi-channel version GMN [14], we formulate the l -th layer for each surface atom $i \in \mathcal{V}_S$ as:

$$\mathbf{m}_{ij} = \phi_m(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)}, f_x(X_i^{(l)}, X_j^{(l)}), e_{ij}), \quad (8)$$

$$\mathbf{h}_i^{(l+1)} = \phi_h(\mathbf{h}_i^{(l)}, \sum_{j \in \mathcal{N}(i)} \mathbf{m}_{ij}), \quad (9)$$

$$\mathbf{X}_i^{(l+1)} = \mathbf{X}_i^{(l)} + \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} (\mathbf{x}_i^{(l)} - \mathbf{x}_j^{(l)}) \phi_x(\mathbf{m}_{ij}), \quad (10)$$

where the functions ϕ_m, ϕ_h, ϕ_x are all MLPs, $\mathcal{N}(i)$ denotes the neighbors of node i in terms of the spatial edges \mathcal{E}_d , $|\cdot|$ counts the size of the input set, and the invariant message \mathbf{m}_{ij} from node j to i is employed to update the invariant feature $\mathbf{h}_i^{(l+1)}$ via ϕ_h and the equivariant matrix $X_i^{(l+1)}$ via the aggregation of the relative position $\mathbf{x}_i^{(l)} - \mathbf{x}_j^{(l)}$ multiplied with ϕ_x .

As a core operator in the message passing above, the function $f_x(\mathbf{X}_i, \mathbf{X}_j)$ is defined as follows:

$$f_x(\mathbf{X}_i, \mathbf{X}_j) := [\|\mathbf{x}_{ij}\|_2, \|\mathbf{x}_{ci}\|_2, \|\mathbf{x}_{cj}\|_2, \angle_{ci,ij}, \angle_{cj,ij}, \angle_{ci,cj}], \quad (11)$$

where, the relative positions are given by $\mathbf{x}_{ij} = \mathbf{x}_i - \mathbf{x}_j$, $\mathbf{x}_{ci} = \bar{\mathbf{x}}_i - \mathbf{x}_i$ and $\mathbf{x}_{cj} = \bar{\mathbf{x}}_j - \mathbf{x}_j$; the angles $\angle_{ci,ij}, \angle_{cj,ij}, \angle_{ci,cj}$ are defined as the inner-products of the corresponding vectors denoted in the subscripts, e.g., $\angle_{ci,ij} = \frac{\mathbf{x}_{ci} \cdot \mathbf{x}_{ij}}{\|\mathbf{x}_{ci}\|_2 \|\mathbf{x}_{ij}\|_2}$. Through the design in Eq. 11, $f_x(\mathbf{X}_i, \mathbf{X}_j)$ elaborates the critical information (including relative distances and angles) around the four points: $\mathbf{x}_i, \bar{\mathbf{x}}_i, \mathbf{x}_j, \bar{\mathbf{x}}_j$, which largely characterizes the geometrical interaction between the two input matrices. Nicely, $f_x(\mathbf{X}_i, \mathbf{X}_j)$ is invariant, ensuring the equivariance of the proposed Surface-EGNN.

Dense Attention Output Layer. Conventionally, we can apply the output of the final layer, i.e., $(\mathbf{h}_i^{(L)}, X_i^{(L)})$ to estimate the binding

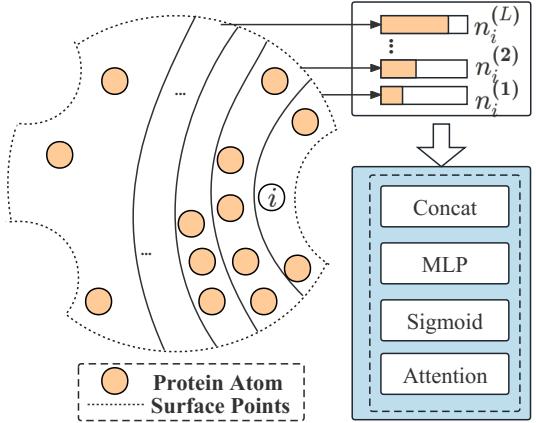


Figure 4: An illustration of Dense Attention in a Protein.

site. Nevertheless, such flat output overlooks the discrepancy of size and shape between different proteins. As showed in Figure 5(b), for small or densely-connected proteins, the receptive field of each node will easily cover most nodes after a small number of message-passing layers, and excessive message passing will lead to over-smoothing [15] that will incur performance detriment. For large or sparsely-connected proteins, on the contrary, insufficient message passing can hardly attain the receptive field with a desirable scope, which will also decrease the performance. It thus requires us to develop an adaptive mechanism to balance the message passing scope between different proteins. We propose the *dense attention output layer* to achieve this goal.

Intuitively, for each target atom, the spatial distribution of the neighbors is able to reflect the density of the spatial connections around. This motivates us to calculate the proportion of the atoms with different distance ranges. As θ is the cutoff to create the spatial graph, we use it as the distance unit. We compute by:

$$\mathbf{n}_i^{(l)} = \frac{|\{j \in \mathcal{V}_P \mid 0 \leq \|\mathbf{x}_i - \mathbf{x}_j\|_2 < l\theta\}|}{N_P}, \quad (12)$$

where, the proportion is evaluated within the distance range $[0, l\theta]$, $N_P = |\mathcal{V}_P|$, and the neighbor hop $l \in \mathbb{Z}^+$. We collect the proportions of all hops from 0 to L , yielding the proportion vector $\mathbf{n}_i = [n_i^{(0)}, n_i^{(1)}, \dots, n_i^{(L)}, N_P] \in \mathbb{R}^{L+2}$ with N_P plus to emphasize the total number of the protein atoms. Clearly, \mathbf{n}_i contains rich information of the spatial density, and we apply it to determine the importance of different layers, by producing the attention as:

$$\mathbf{a}_i = \text{Sigmoid}(\phi_a(\mathbf{n}_i)). \quad (13)$$

Here, ϕ_a is an MLP with the number of output channels as $L+1$, the Sigmoid function¹ is applied for each channel, implying that $\mathbf{a}_i \in (0, 1)^{L+1}$. We then multiply the hidden feature of the corresponding layer with each channel of the attention vector, and concatenate them into a vector:

$$\mathbf{h}_i^{\text{out}} = \text{Concat}(a_{i0} \mathbf{h}_i^{(0)}, \dots, a_{iL} \mathbf{h}_i^{(L)}), \quad (14)$$

¹Note that the sum of all channels of \mathbf{a}_i is unnecessarily equal to 1, since the Sigmoid function instead of the previously-used SoftMax function is applied here.

where a_{il} is the l -th channel of \mathbf{a}_i . By making use of Eq. 14, the learnable attentions enable the model to adaptively balance the importance of different layers for different input proteins. We will illustrate the benefit of the proposed strategy in our experiments. As for the coordinates, we simply compute the mean of all layers to retain translation equivariance:

$$\mathbf{X}_i^{\text{out}} = \frac{1}{L+1} \sum_{l=0}^L \mathbf{X}_i^{(l)}. \quad (15)$$

4.4 Optimization Objective

We set $y_i = 1$ if a surface atom i is within 4Å to any ligand atom [31]. We predict the probability \hat{y}_i of being a part of binding site according its dense embedding \mathbf{h}_i .

$$\hat{y}_i = \text{Sigmoid}(\text{MLP}(\mathbf{h}_i^{\text{out}})). \quad (16)$$

Following [17, 18], Dice loss is used:

$$\mathcal{L}_b = 1 - \frac{2 \cdot \sum(\hat{y}_i \cdot y_i)}{\sum(\hat{y}_i) + \sum(y_i) + \epsilon}, \quad (17)$$

where $\epsilon > 0$ is a small value to maintain numeric stability.

Predict the relative direction of nearest ligand atom. Beyond the CNN-based methods, our EquiPocket is an E(3)-equivariant model, which can not only output the embedding $\mathbf{h}_i^{\text{out}}$ but also the coordinate matrix $\mathbf{X}_i^{\text{out}}$ (with initial position vector \mathbf{x}_i). We further leverage the position vector to predict the relative direction \mathbf{d}_i of its nearest ligand atom (with position vector \mathbf{m}_i), in order to enhance our framework to gather local geometric features.

$$\mathbf{d}_i = \frac{\mathbf{m}_i - \mathbf{x}_i}{\|\mathbf{m}_i - \mathbf{x}_i\|_2}, \quad \hat{\mathbf{d}}_i = \frac{\mathbf{x}_i^{\text{out}} - \mathbf{x}_i}{\|\mathbf{x}_i^{\text{out}} - \mathbf{x}_i\|_2}. \quad (18)$$

The cosine loss function is used for the direction loss \mathcal{L}_d :

$$\mathcal{L}_d = \sum (1 - \cos(\hat{\mathbf{d}}_i, \mathbf{d}_i)). \quad (19)$$

The eventual loss is $\mathcal{L} = \mathcal{L}_b + \mathcal{L}_d$. We train the parameters of all the three modules end to end.

5 EXPERIMENTS

In this section, we will conduct experiments on multiple datasets to evaluate the performance of our framework in comparison to baseline methods and investigate the following tasks:

- **Task 1.** Can our framework's performance match or surpass that of existing methods?
- **Task 2.** Can different modules of our framework bring significant improvement for binding site prediction?
- **Task 3.** Can our framework mitigate the detrimental effects of the data distribution shift?
- **Task 4.** How do the hyperparameters (the cutoff θ and depth of surface-egnn) affect the performance and computational cost?

5.1 Experimental Settings

5.1.1 *Dataset.* We conduct experiments based on the following datasets:

- **scPDB** [8] is the famous dataset for binding site prediction, which contains the protein structure, ligand structure, and 3D cavity structure generated by VolSite [7]. The 2017 release of scPDB is used for training and cross-validation of our

Table 1: Summary of Dataset

DataSet	Average			
	Atom Num	Atom in Surface	Surface Points	Target Atoms
scPDB	4205	2317	24010	47
COACH420	2123	1217	12325	58
HOLO4k	3845	2052	20023	106
PDBbind	3104	1677	17357	37

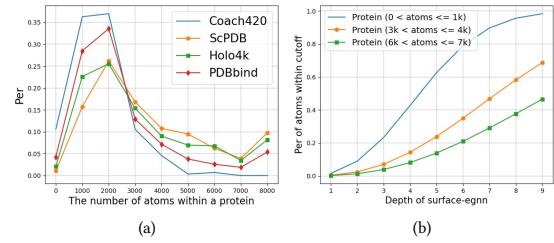


Figure 5: The protein distribution of datasets and spatial distribution of atom neighbors

framework, which contains 17,594 structures, 16,034 entries, 4,782 proteins, and 6,326 ligands.

- **PDBbind** [46] is a well-known and commonly used dataset for the research of protein-ligand complex. It contains the 3D structures of proteins, ligands, binding sites, and accurate binding affinity results determined in the laboratory. We use the release of v2020, which consists of two parts: general set (14, 127 complexes) and refined set (5,316 complexes). The general set contains all protein-ligand complexes. The refined set contains better-quality compounds selected from the general set, which is used for the test in our experiments.
- **COACH 420 and HOLO4K** are two test datasets for the binding site prediction, which are first introduced by [24]. Consistent with [1, 24, 31], we use the mlig subsets of each dataset for evaluation, which contain the relevant ligands for binding site prediction.

Data Distribution Shift. As depicted in the Figure 5(a) and Table 1 that after data processing, there is a significant gap in protein size and protein distribution between the training dataset (scPDB) and the test dataset (COACH420, HOLO4k, PDBbind). The number of atoms within a protein ranges from hundreds to tens of thousands. As for protein distribution in datasets, scPDB has the longest average structure, followed by HOLO4k and PDBbind, with COACH420 having the shortest average protein structure. This fact will hurt model learning and generalization, as discussed in § 5.2.3.

5.1.2 *Target of Binding Sites.* The CNN-based methods [1, 17, 41] label the subgrid as positive if its geometric center is closer than 4Å to the binding sites geometric center. In our experiment, consistent with [31], we set the protein atoms within 4Å of any ligand atom as positive and negative otherwise. After obtaining the probability that an atom is a candidate binding site, we use the mean-shift algorithm [6] to predict the binding site center, which can determine the number of clusters on its own (details in Appendix A.2.4).

5.1.3 *Data Preparation.* We perform the following four processing steps: i) Cluster the structures in scPDB by their Uniprot IDs, and

Table 2: Experimental and ablation results of baseline models and our framework.^a

Methods	Type	Param (M)	Failure Rate ↓	COACH420		HOLO4K		PDBbind2020	
				DCC↑	DCA↑	DCC↑	DCA↑	DCC↑	DCA↑
Fpocket ^b	Geometric-based	\	0.000	0.228	0.444	0.192	0.457	0.253	0.371
DeepSite ^b		1.00	\	\	0.564	\	0.456	\	\
Kalasanty ^b	3D-CNN	70.64	0.120	0.335	0.636	0.244	0.515	0.416	0.625
DeepSurf ^b		33.06	0.054	0.386	0.658	0.289	0.635	0.510	0.708
GAT		0.03	0.11	0.039(0.005)	0.130(0.009)	0.036(0.003)	0.110(0.010)	0.032(0.001)	0.088(0.011)
GCN	Topological Graph	0.06	0.163	0.049(0.001)	0.139(0.010)	0.044(0.003)	0.174(0.003)	0.018(0.001)	0.070(0.002)
GAT + GCN		0.08	0.31	0.036(0.009)	0.131(0.021)	0.042(0.003)	0.152(0.020)	0.022(0.008)	0.074(0.007)
GCN2		0.11	0.466	0.042(0.098)	0.131(0.017)	0.051(0.004)	0.163(0.008)	0.023(0.007)	0.089(0.013)
SchNet	Spatial	0.49	0.14	0.168(0.019)	0.444(0.020)	0.192(0.005)	0.501(0.004)	0.263(0.003)	0.457(0.004)
Eggn	Graph	0.41	0.270	0.156(0.017)	0.361(0.020)	0.127(0.005)	0.406(0.004)	0.143(0.007)	0.302(0.006)
EquiPocket-L		0.15	0.552	0.070(0.009)	0.171(0.008)	0.044(0.004)	0.138(0.006)	0.051(0.003)	0.132(0.009)
EquiPocket-G	Ours	0.42	0.292	0.159(0.016)	0.373(0.021)	0.129(0.005)	0.411(0.005)	0.145(0.007)	0.311(0.007)
EquiPocket-LG		0.50	0.220	0.212(0.016)	0.443(0.011)	0.183(0.004)	0.502(0.008)	0.274(0.004)	0.462(0.005)
EquiPocket		1.70	0.051	0.423(0.014)	0.656(0.007)	0.337(0.006)	0.662(0.007)	0.545(0.010)	0.721(0.004)

^a The standard deviation of each index is indicated in brackets. The result of 5-fold for EquiPocket is shown in Appendix A.2.5.

^b We use their published pre-train models or published result, details in Appendix A.2.3.

select the longest sequenced protein structures from every cluster as the train data [18]. Finally, 5,372 structures are selected out. ii) Split proteins and ligands for the structures in COACH420 and HOLO4k, according to the research [24]. iii) Clean protein by removing the solvent, hydrogens atoms. Using MSMS [37] to generate the solvent-accessible surface of a protein. iv) Read the protein file by RDKIT [43], and extract the atom and chemical bond features. Remove the error structures.

5.1.4 Evaluation Metrics. DCC is the distance between the predicted binding site center and the true binding site center. DCA is the shortest distance between the predicted binding site center and any atom of the ligand. The samples with DCC(DCA) less than the threshold are considered successful. The samples without any binding site center are considered failures. Consistent with [1, 17, 31, 41], threshold is set to 4 Å. We use **Success Rate** and **Failure Rate** to evaluate experimental performance.

$$\begin{aligned} \text{Success Rate(DCC)} &= \frac{1(\{\text{Predicted sites} | \text{DCC} < \text{threshold}\})}{1(\{\text{True sites}\})}, \\ \text{Success Rate(DCA)} &= \frac{1(\{\text{Predicted sites} | \text{DCA} < \text{threshold}\})}{1(\{\text{True sites}\})}, \\ \text{Failure Rate} &= \frac{1(\{\text{Protein} | 1(\text{predicted binding center}) = 0\})}{1(\{\text{Protein}\})}, \end{aligned} \quad (20)$$

where $1(\cdot)$ represents the cardinality of a set. After ranking the predicted binding sites, we take the same number with the true binding sites to calculate the success rate.

5.1.5 EquiPocket Framework. We implement our EquiPocket framework based on (GAT [45]+EGNN [39]) as our global structure modeling module. The cutoff θ and depth in our surface-egnn model are set to 6 and 4.

To indicate the EquiPocket Framework with different modules, we adopt the following symbol as follows: i) **EquiPocket-L**: Only contain the local geometric modeling module. ii) **EquiPocket-G**:

Only contain the global structure modeling module. iii) **EquiPocket-LG**: Only contain both the local geometric and global structure modeling modules. iii) **EquiPocket**: Contain all the modules.

5.1.6 Baseline Models. We compare our framework with the following models: 1) geometric-based method(Fpocket [27]), 2) CNN-based methods (DeepSite [17], Kalasanty [41] and DeepSurf [31]), 3) topological graph-based models (GAT [45], GCN [19] and GCN2 [5]), 4) spatial graph-based models (SchNet [40], EGNN [38]).

5.1.7 Environment and Parameter. We implement our EquiPocket framework in PyTorch Geometric, all the experiments are conducted on a machine with an NVIDIA A100 GPU (80GB memory). We take 5-fold cross validation on training data scPDB and use valid loss to save checkpoint. The batch size is set to 8. For baseline models and models in global structure modeling module of EquiPocket, we use them suggested settings to get optimal performance. More details and related resources for our experiments can be found in **Appendix A.2.2**.

5.2 Result analysis

5.2.1 Model Comparison. In Table 2, we compared our EquiPocket framework with baseline methods mentioned above. As can be observed, the performance of the computational method Fpocket is inferior, with no failure rate, since it simply employs the geometric feature of a protein. The performance of CNN-based methods is much superior to that of the conventional method, with DCC and DCA metrics improving by more than 50 percent but requiring enormous parameter values and computing resources. However, these two early methods DeepSite and Kalasanty are hampered by data distribution shift (Issue 4) and their inability to process big proteins, which may fail prediction. The recently proposed method DeepSurf employs the local-grid concept to handle any size of proteins, although CNN architecture also still results in inevitable failures.

For graph models, the poor performance of topological-graph models (GCN, GAT, GCN2) is primarily due to the fact that they only

consider atom attributes and chemical bond information, ignoring the spatial structure in a protein. The performance of spatial-graph models is generally better than that of topological-graph models. EGNN model utilizes not only the properties of atoms but also their relative and absolute spatial positions, resulting in a better effect. SchNet merely updates the information of atoms based on the relative distance of atoms. We attempt to execute the Dimenet++ [21], which uses the angle info between atoms and atoms, but it requires too many computing resources, resulting in an OOM (Out Of Memory) error. However, the performance of the spatial-graph model is worse than that of the CNN-based and geometric-based methods because the former cannot obtain enough geometric features (Issue 3) and cannot address the data distribution shift (Issue 4).

As the above results indicate, geometric info of protein surface and multi-level structure info in a protein is essential for binding site prediction. In addition, it reflects the limitations of the current GNN models, where it is difficult to collect sufficient geometric information from the protein surface or the calculation resources are too large to apply to macromolecular systems like proteins. Consequently, our EquiPocket framework is not only able to update chemical and spatial information from an atomic perspective but also able to effectively collect geometric information without excessive computing expense, resulting in a 10-20% increase in effect over previous results. Case study based on different methods is showed in Appendix 5.2.5.

5.2.2 Ablation Study. As shown in Table 2, we conduct ablation experiments on our EquiPocket framework with different modules.

Local Geometric Modeling Module. This module is used to extract the geometric features of protein atoms from their nearest surface points. EquiPocket-G consists solely of this module, and the performance is negligible. There are two primary causes for this result. First, geometric information can only determine part of the binding sites. Second, it can only reflect the geometric features over a relatively small distance and cannot cover an expansive area.

Global Structure Modeling Module. The primary purpose of this module is to extract information about the whole protein, such as atom type, chemical bonds, relevant spatial positions, etc. We implement EquiPocket-G based on (GAT + EGNN) models, which is E(3) equivariance/invariance and has a better effect than its predecessor, EquiPocket-L. In comparison, the value of DCC increased by about 10%, and DCA increased by about 20%. This demonstrates that structure information of the whole protein is necessary for binding site prediction. In addition, when the two modules are combined as the EquiPocket-LG, the prediction effect is significantly improved, proving the complementarity of surface geometric information and global structure information.

Surface Message Passing Module. In the previous model, EquiPocket-LG, information was extracted solely from atoms and their closest surface points. Nonetheless, the binding site is determined not only by the information of a single atom but also by the atoms surrounding it. Therefore, the surface message passing module is proposed to collect and update the atom's features from its neighbors. After adding this module, the performance of EquiPocket has been significantly enhanced, DCC and DCA have increased by approximately 20% on average, and the failure rate has been significantly reduced. Through the addition of multiple modules, we address the Issue 3

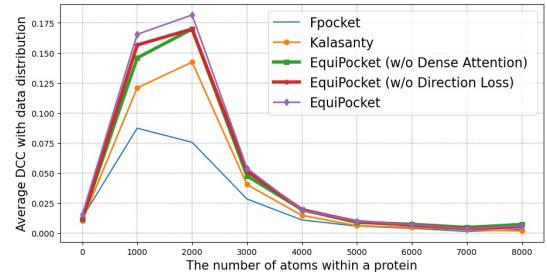


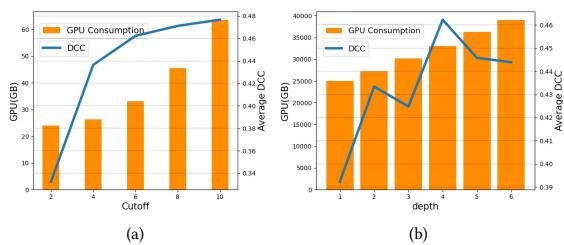
Figure 6: The performance of EquiPocket and baseline methods for proteins of various sizes.

and the performance of our framework eventually surpasses that of the existing SOTA method, demonstrating the efficacy of our framework design.

5.2.3 Data Distribution Shift. As shown in Figure 6, we calculate the average DCC with the distribution of various sizes proteins. The geometric-based method Fpocket only utilizes the geometric features of a protein surface. Therefore, its performance is superior to that of most other methods for proteins with fewer than 1,000 atoms, but its prediction effect decreases significantly as the size of the protein increases. Kalasanty is a CNN-based and learn-based method. As the number of atoms in the protein varies, the prediction effect exhibits an increasing and then a decreasing trend, which is not only influenced by the size of the protein but also has a significant correlation with the dataset's distribution. According to the train data (scPDB), the majority of proteins contain fewer than 2,000 protein atoms (as depicted in Figure 5(a)). Consequently, the model's parameters will be biased toward this protein size. In addition, for proteins with more than 8000 atoms, the prediction effect is not even as good as the geometric-based method. This is due to the fact that CNN methods typically restrict the protein space to $70\text{\AA} \times 70\text{\AA} \times 70\text{\AA}$, and for proteins larger than this size, the prediction frequently fails. For our EquiPocket framework, we do not need to cut the protein into grids, and we utilize both geometric information from the surface points and global structure information from the whole protein, so the performance for proteins of varying sizes is significantly superior to that of other methods.

Dense Attention. The Dense Attention is introduced in § 4.3 to reduce the negative impact caused by the data distribution shift (Issue 4). As shown in 6, when the number of atoms contained in a protein is less than 3000, the result of the EquiPocket (w/o Dense Attention) is weaker than that of the original EquiPocket, whereas when the protein is larger, there is no significant distinction between the two models. It simply reflects the role of Dense Attention, which, by weighting the surface-egnn layer at different depths, mitigates the detrimental effect of the data distribution shift.

Direction Loss. Direction loss is a novel task designed to improve the extraction of local geometric features. The result of the EquiPocket (w/o Direction Loss) in Figure 6 demonstrates conclusively that the prediction performance of small proteins with fewer than 3,000 atoms is diminished in the absence of this task, which reveals the importance of the task.

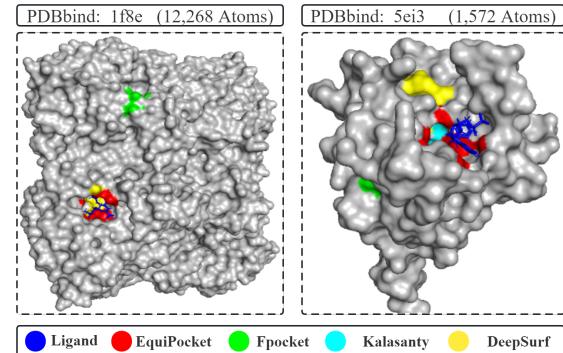
**Figure 7: The influence of cutoff θ and depth for experiment**

5.2.4 Hyperparameters Analysis. In our EquiPocket framework, the cutoff θ and depth of surface-egnn are two crucial parameters that can impact performance and computational efficiency.

Cutoff θ . We set the depth of surface-egnn to 4 and implement various cutoff values (2, 4, 6, 8, 10). Figure 7(a) indicates that when the cutoff is set to 2, the average DCC of our framework is poor, and GPU memory consumption is relatively low (22GB). This is due to the fact that when the cutoff is small, the surface-egnn can only observe a tiny receptive field. As the cutoff increases, the performance and GPU memory continue to rise until the DCC reaches a bottleneck when the cutoff is 10, and the GPU memory reaches 62GB. Therefore, when selecting parameters for our framework, we must strike a balance between performance and efficiency.

Depth. The depth of surface-egnn has an immediate influences on the performance and computation cost. We set the cutoff to 6 and implement various depth (1, 2, 3, 4, 5, 6). Figure 7(b) demonstrates that as depth increases, performance steadily improves and becomes stable as GPU memory continues to expand. Because the prediction of binding sites is highly influenced by their surrounding atoms, therefore, an excessively large receptive field may not offer any benefits but will necessitate additional computing resources.

5.2.5 Case Study. We also display two examples of our EquiPocket and other methods in Figure 8. We take two proteins, 1f8e (with 12,268 atoms) and 5ei3 (with 1,572 atoms), from the test dataset PDBbind. As can be seen from Figure 8: The binding sites predicted by the geometry-based method Fpocket are extremely distant from the actual binding sites. This is due to the fact that this method prioritizes local geometric information and disregards the multi-level structure information of proteins, resulting in limited scope and weak performance. The CNN-based method Kalasanty did not provide any predicted binding site for protein 1f8e. We conjecture that this method restricts the protein within a specific space size which is highly susceptible to failure with large proteins. The recently-proposed CNN-based method DeepSurf takes local grids on the protein surface, which can address the issue of fixed space size. However, the prediction of binding sites in protein 5ei3 by DeepSurf is far from the ground truth because the CNN-based methods are defective in obtaining geometric and chemical features. Our EquiPocket framework is unaffected by the shortcomings of the aforementioned methods, allowing it to achieve superior outcomes for both large and small proteins.

**Figure 8: Case Study.**

6 CONCLUSION

In this paper, concentrating on the ligand binding site prediction, we propose a novel E(3)-Equivariant geometric graph framework called EquiPocket, which contains the local geometric modeling module, global structure modeling module, and surface passing module to gather the surface geometric and multi-level structure features in a protein. Experiments demonstrate that our framework is highly generalizable and beneficial, and achieves superior prediction accuracy and computational efficiency compared with the existing methods.

6.1 Future Work

6.1.1 Protein Surface. As demonstrated by our experiments, the geometric information derived from the protein surface plays a significant role in the prediction of binding sites. In this work, we use MSMS [37] to generate the protein surface, which may have an uncertain impact on the prediction results. Therefore, in the future, we will be able to establish a more efficient surface generation method or gather the geometric information of a protein without a fixed protein surface.

6.1.2 Global Structure Model. In this work, we take the existing graph-based models (GAT + EGNN) to gather the multi-level structure information in a protein. However, their capabilities are limited because these models are not entirely tailored to complex structures such as proteins. In the future, we will be able to develop more effective models to gather information from proteins in order to improve the prediction performance of binding sites.

6.1.3 Computing Resources. As the experimental results show, both our EquiPocket and CNN-based methods require a significant amount of computing resources to analyze protein data, which will have a negative impact on their actual deployment. Consequently, when applying our method to real-world problems, it is crucial to consider how to compress algorithm parameters, increase computational efficiency, and decrease resource consumption.

REFERENCES

- [1] Rishal Aggarwal, Akash Gupta, Vineeth Chelur, CV Jawahar, and U Deva Priyakumar. 2021. Deepocket: ligand binding site detection and segmentation using 3d

- convolutional neural networks. *Journal of Chemical Information and Modeling* (2021).
- [2] Mariana Belgiu and Lucian Drăguț. 2016. Random forest in remote sensing: A review of applications and future directions. *ISPRS journal of photogrammetry and remote sensing* 114 (2016), 24–31.
- [3] Michał Brylinski and Jeffrey Skolnick. 2008. A threading-based method (FIND-SITE) for ligand-binding site prediction and functional annotation. *Proceedings of the National Academy of sciences* 105, 1 (2008), 129–134.
- [4] Ke Chen, Marcin J Mizianty, Jianzhao Gao, and Lukasz Kurgan. 2011. A critical comparative assessment of predictions of protein-binding sites for biologically relevant organic compounds. *Structure* 19, 5 (2011), 613–621.
- [5] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. 2020. Simple and Deep Graph Convolutional Networks. arXiv:2007.02133 [cs.LG]
- [6] Dorin Comaniciu and Peter Meer. 2002. Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on pattern analysis and machine intelligence* 24, 5 (2002), 603–619.
- [7] Franck Da Silva, Jeremy Desaphy, and Didier Rognan. 2018. ICHEM: a versatile toolkit for detecting, comparing, and predicting protein-ligand interactions. *ChemMedChem* 13, 6 (2018), 507–510.
- [8] Jérémie Desaphy, Guillaume Bret, Didier Rognan, and Esther Kellenberger. 2015. sc-PDB: a 3D-database of ligandable binding sites—10 years on. *Nucleic acids research* 43, D1 (2015), D399–D404.
- [9] Sérgio ED Dias, Quoc T Nguyen, Joaquim A Jorge, and Abel JP Gomes. 2017. Multi-GPU-based detection of protein cavities using critical points. *Future Generation Computer Systems* 67 (2017), 430–440.
- [10] Christina E Faller, E Prabhu Raman, Alexander D MacKerell, and Olgun Guvench. 2015. Site Identification by Ligand Competitive Saturation (SILCS) simulations for fragment-based drug design. In *Fragment-Based Methods in Drug Discovery*. Springer, 75–87.
- [11] A. Gonczarek, J. M. Tomczak, S. Zaręba, J. Kaczmar, P. Dąbrowski, and M. J. Walczak. 2016. Learning Deep Architectures for Interaction Prediction in Structure-based Virtual Screening. (2016).
- [12] J. Han, Y. Rong, T. Xu, and W. Huang. 2022. Geometrically Equivariant Graph Neural Networks: A Survey. (2022).
- [13] Manfred Hendlich, Friedrich Rippmann, and Gerhard Barnickel. 1997. LIGSITE: automatic and efficient detection of potential small molecule-binding sites in proteins. *Journal of Molecular Graphics and Modelling* 15, 6 (1997), 359–363.
- [14] W. Huang, J. Han, Y. Rong, T. Xu, F. Sun, and J. Huang. 2022. Equivariant Graph Mechanics Networks with Constraints. *arXiv e-prints* (2022).
- [15] W. Huang, Y. Rong, T. Xu, F. Sun, and J. Huang. 2020. Tackling Over-Smoothing for General Graph Convolutional Networks. (2020).
- [16] Mingjian Jiang, Zhiqiang Wei, Shugang Zhang, Shuang Wang, Xiaofeng Wang, and Zhen Li. 2019. Frsite: protein drug binding site prediction based on faster R-CNN. *Journal of Molecular Graphics and Modelling* 93 (2019), 107454.
- [17] José Jiménez, Stefan Doerr, Gerard Martínez-Rosell, Alexander S Rose, and Gianni De Fabritiis. 2017. DeepSite: protein-binding site predictor using 3D convolutional neural networks. *Bioinformatics* 33, 19 (2017), 3036–3042.
- [18] Jeevan Kandel, Hilal Tayara, and Kil To Chong. 2021. PUResNet: prediction of protein-ligand binding sites using deep residual neural network. *Journal of cheminformatics* 13, 1 (2021), 1–14.
- [19] T. N. Kip and M. Welling. 2016. Semi-Supervised Classification with Graph Convolutional Networks. (2016).
- [20] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
- [21] J. Klicpera, S. Giri, J. T. Margraf, and S. Günnemann. 2020. Fast and Uncertainty-Aware Directional Message Passing for Non-Equilibrium Molecules. (2020).
- [22] J. Klicpera, J. Gro, and S. Günnemann. 2020. Directional Message Passing for Molecular Graphs. (2020).
- [23] Radoslav Krivák and David Hoksza. 2015. Improving protein-ligand binding site prediction accuracy by classification of inner pocket points using local features. *Journal of cheminformatics* 7, 1 (2015), 1–13.
- [24] Radoslav Krivák and David Hoksza. 2018. P2Rank: machine learning based tool for rapid and accurate prediction of ligand binding sites from protein structure. *Journal of cheminformatics* 10, 1 (2018), 1–12.
- [25] Roman A Laskowski. 1995. SURFNET: a program for visualizing molecular surfaces, cavities, and intermolecular interactions. *Journal of molecular graphics* 13, 5 (1995), 323–330.
- [26] Alasdair TR Laurie and Richard M Jackson. 2005. Q-SiteFinder: an energy-based method for the prediction of protein-ligand binding sites. *Bioinformatics* 21, 9 (2005), 1908–1916.
- [27] Vincent Le Guilloux, Peter Schmidke, and Pierre Tuffery. 2009. Fpocket: an open source platform for ligand pocket detection. *BMC bioinformatics* 10, 1 (2009), 1–11.
- [28] David G Levitt and Leonard J Banaszak. 1992. POCKET: a computer graphics method for identifying and displaying protein cavities and their surrounding amino acids. *Journal of molecular graphics* 10, 4 (1992), 229–234.
- [29] Yi Liu, Limei Wang, Meng Liu, Xuan Zhang, Bora Oztekin, and Shuiwang Ji. 2021. Spherical Message Passing for 3D Graph Networks. arXiv:2102.05013 [cs.LG]
- [30] Gabriele Macari, Daniele Toti, and Fabio Polticelli. 2019. Computational methods and tools for binding site recognition between proteins and small molecules: from classical geometrical approaches to modern machine learning strategies. *Journal of computer-aided molecular design* 33, 10 (2019), 887–903.
- [31] Stelios K Mylonas, Apostolos Axenopoulos, and Petros Daras. 2021. DeepSurf: a surface-based deep learning approach for the prediction of ligand binding sites on proteins. *Bioinformatics* 37, 12 (2021), 1681–1690.
- [32] Chi-Ho Ngan, David R Hall, Brandon Zerbe, Laurie E Grove, Dima Kozakov, and Sandor Vajda. 2012. FTsite: high accuracy detection of ligand binding sites on unbound protein structures. *Bioinformatics* 28, 2 (2012), 286–287.
- [33] Matthew Ragoza, Joshua Hochuli, Elisa Idrobo, Jocelyn Sunseri, and David Ryan Koes. 2017. Protein-ligand scoring with convolutional neural networks. *Journal of chemical information and modeling* 57, 4 (2017), 942–957.
- [34] HP Rang. 2006. The receptor concept: pharmacology's big idea. *British journal of pharmacology* 147, S1 (2006), S9–S16.
- [35] Timothy J Richmond. 1984. Solvent accessible surface area and excluded volume in proteins: Analytical equations for overlapping spheres and implications for the hydrophobic effect. *Journal of molecular biology* 178, 1 (1984), 63–89.
- [36] Sabrina, Jaeger, Simone, Fulle, Samo, and Turk. 2018. Mol2vec: Unsupervised Machine Learning Approach with Chemical Intuition. *Journal of Chemical Information & Modeling* (2018).
- [37] M. F. Samner, A. J. Olson, and J. C. Spehner. 1996. Reduced SURFACE: An efficient way to compute molecular surfaces. *Peptide Science* 38, 3 (1996), 305–320.
- [38] Victor Garcia Satorras, Emiel Hoogeboom, and Max Welling. 2021. E (n) equivariant graph neural networks. In *International conference on machine learning*. PMLR, 9323–9332.
- [39] V. G. Satorras, E. Hoogeboom, and M. Welling. 2021. E(n) Equivariant Graph Neural Networks. (2021).
- [40] Kristof T. Schütt, Pieter Jan Kindermans, Huziel E. Sauceda, Stefan Chmiela, Alexandre Tkatchenko, and Klaus-Robert Müller. 2017. SchNet: A continuous-filter convolutional neural network for modeling quantum interactions. In *Advances in Neural Information Processing Systems*.
- [41] Marta M Stepniewska-Dziubinska, Piotr Zielenkiewicz, and Paweł Siedlecki. 2020. Improving detection of protein-ligand binding sites with 3D segmentation. *Scientific reports* 10, 1 (2020), 1–9.
- [42] M. Sun, S. Zhao, G. Coryandar, E. Olivier, J. Zhou, and F. Wang. 2019. Graph convolutional networks for computational drug development and discovery. *Briefings in Bioinformatics* (2019).
- [43] Paolo Tosco, Nikolaus Stiefl, and Gregory Landrum. 2014. Bringing the MMFF force field to the RDKit: implementation and validation. *Journal of Cheminformatics* 6, 1 (2014-07-12), 1, 1 (2014), 37.
- [44] Daniele Toti, Le Viet Hung, Valentina Tortosa, Valentina Brandi, and Fabio Polticelli. 2018. LIBRA-WA: a web application for ligand binding site detection and protein function recognition. *Bioinformatics* 34, 5 (2018), 878–880.
- [45] Petar Veliković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. 2017. Graph Attention Networks. (2017).
- [46] Renxiao Wang, Xueliang Fang, Yipin Lu, and Shaomeng Wang. 2004. The PDBbind database: collection of binding affinities for protein-ligand complexes with known three-dimensional structures. *Journal of Medicinal Chemistry* 47, 12 (2004), 2977–80.
- [47] Renxiao Wang, Ying Gao, and Luhua Lai. 2000. LigBuilder: a multi-purpose program for structure-based drug design. *Molecular modeling annual* 6 (2000), 498–516.
- [48] Martin Weisel, Ewgenij Proschak, and Gisbert Schneider. 2007. PocketPicker: analysis of ligand binding-sites with shape descriptors. *Chemistry Central Journal* 1, 1 (2007), 1–17.
- [49] K. Xu, W. Hu, J. Leskovec, and S. Jegelka. 2018. How Powerful are Graph Neural Networks? (2018).
- [50] Yang Zhang, Gengmo Zhou, Zhewei Wei, and Hongteng Xu. 2022. Predicting Protein-Ligand Binding Affinity via Joint Global-Local Interaction Modeling. *arXiv preprint arXiv:2209.13014* (2022).

A APPENDIX

A.1 The pseudo-code of our EquiPocket framework

Algorithm 1: EquiPocket

Input: Protein structure \mathcal{G}_P
Output: Candidate Binding sites and their ligandability score

- 1: Clean Structure by removing solvent, hydrogens atoms
- 2: Create the solvent accessible surface of the protein \mathbb{S} use MSMS
- 3: **for** every s_i in \mathbb{S} **do**
- 4: Get its closed protein atom p_i
- 5: **end for**
- 6: Get the surface atom \mathcal{V}_S according to the surface points's closed protein atom
- 7: **for** every surface atom $i \in \mathcal{V}_S$ **do**
- 8: Get their surrounding surface points set \mathbb{S}_i
- 9: Get the geometric embedding \mathbf{g}_i
- 10: **end for**
- 11: Get the global structure embedding \mathbf{c}''_i of the protein
- 12: **for** every surface atom $i \in \mathcal{V}_S$ **do**
- 13: Get its invariant feature $\mathbf{h}_i^{(0)} = [\mathbf{c}''_i, \mathbf{g}_i]$ and equivariant position matrix $\mathbf{X}_i^{(0)} = [\mathbf{x}_i, \bar{\mathbf{x}}_i]$
- 14: Get the updated embedding $\mathbf{h}_i^{(l)}$ and updated coordinates $\mathbf{X}_i^{(l)}$ based on our surface-egnn model
- 15: Get the dense embedding \mathbf{h}_i and position \mathbf{X}_i according to its dense attention \mathbf{a}_i
- 16: predict the probability $\hat{\mathbf{y}}_i$ as ligandability score and the nearest ligand atom direction \mathbf{d}_i
- 17: **end for**
- 18: Discard protein atoms with probability less than T (T=0.5 in our experiments);
- 19: Cluster the remaining protein atoms;
- 20: Form binding sites and get the average ligandability score for each cluster;
- 21: Rank the predicted binding sites by their ligandability score;
- 22: **return** The candidate binding sites and ligandability score;

A.2 Experiment Details

A.2.1 Cross-validation. We shuffled the training data and divided the data into 5 parts, taking one of them at a time as the validation set. We use 5-fold cross-validation and report the mean and standard deviation.

A.2.2 Parameter Settings. For geometric-based method Fpocket, we use its published tool. For CNN-based methods kalasanty and DeepSurf, we use their published pre-train models. For GNN-based models, the number of layers is set to 3 except GAT. For GAT, we set the number to 1. For GIN, we set the initial ϵ to 0 and make it trainable. For GCN2, we set the strength of the initial residual connection α to 0.5 and the strength of the identity mapping β to 1. For SchNet, EGNN, DimeNet++, SphereNet as baseline models, we set the cutoff distance to 5. For our EquiPocket, we use Adam

Table 3: Sources of baseline codes and pre-train models.

Methods	URL
fpointer	https://github.com/Discrengine/fpocket
kalasanty	https://gitlab.com/cheminfIBB/kalasanty
DeepSurf	https://github.com/stemylonas/DeepSurf
GAT	https://github.com/pyg-team/pytorch_geometric
GCN	https://github.com/pyg-team/pytorch_geometric
GCN2	https://github.com/chennnM/GCNII
SchNet	https://github.com/pyg-team/pytorch_geometric
DimeNet++	https://github.com/pyg-team/pytorch_geometric
EGNN	https://github.com/vgsatorras/egnn/

optimizer for model training with a learning rate of 0.0001 and set the batch size as 8. The basic dimensions of node and edge embeddings are both set to 128. The dropout rate is set to 0.1. The probe radius in MSMS to generate solvent-accessible surface of a protein is set to 1.5.

A.2.3 Baseline Codes. The result of DeepSite comes from [31], because they did not provide a pre-train model. Table 3 describes sources of baseline codes.

A.2.4 Binding sites center. The CNN-based methods [1, 17, 41] consider a protein as a 3D image, convert it to a voxel representation by discretizing it into grids and calculate the geometric center of binding site $center_{cnn}$ according to the grid of the cavity or ligand. They label the **grid** as positive if its geometric center is closer than 4Å to the binding sites geometric center. Therefore, the prediction objects of these models actually contain the grid of **ligand atoms**. The predicted binding site center $center_{cnn}$ of CNN-based methods is calculated according to the positive grid. For our EquiPocket, we label the **protein atoms** within 4Å of any ligand atom as positive and negative otherwise. Therefore, there is a natural gap in the prediction object between our framework and CNN-based methods, which also lead to the natural gap for the center of predicted binding site. In order to reduce the metric difference caused by the different prediction objects, we get the predicted binding site center $center_{equipocket}$ as follow: We use $pos_i \in \mathbb{R}^3$ to represent the position of protein atom v_i , $center_i \in \mathbb{R}^3$ to represent the nearest surface point center, $\hat{pos}_i^L \in \mathbb{R}^3$ to represent the predicted position of nearest ligand atom from the protein atom v_i . The \hat{pos}_i^L is used to calculate the geometric center of binding site.

$$\hat{pos}_i^L = pos_i + threshold \cdot \frac{(center_i - pos_i)}{|center_i - pos_i|}, \quad (21)$$

Where $threshold$ is set to 4, because we label the protein atoms within 4Å of any ligand atom as positive and negative otherwise.

A.2.5 The 5-fold results for EquiPocket.

Table 4: The 5-fold results for EquiPocket.

Methods	Fold	Param (M)	failure Rate ↓	COACH420		HOLO4K		PDBbind2020	
				DCC↑	DCA↑	DCC↑	DCA↑	DCC↑	DCA↑
EquiPocket-L	0	0.15	0.598	0.083	0.160	0.038	0.128	0.049	0.124
EquiPocket-L	1	0.15	0.557	0.064	0.165	0.046	0.138	0.055	0.142
EquiPocket-L	2	0.15	0.571	0.074	0.177	0.045	0.139	0.052	0.122
EquiPocket-L	3	0.15	0.462	0.059	0.173	0.042	0.138	0.052	0.129
EquiPocket-L	4	0.15	0.472	0.072	0.180	0.048	0.146	0.049	0.143
EquiPocket-G	0	0.42	0.305	0.135	0.330	0.122	0.400	0.142	0.302
EquiPocket-G	1	0.42	0.291	0.175	0.385	0.128	0.405	0.145	0.302
EquiPocket-G	2	0.42	0.295	0.145	0.357	0.121	0.407	0.145	0.305
EquiPocket-G	3	0.42	0.278	0.169	0.367	0.127	0.406	0.133	0.292
EquiPocket-G	4	0.42	0.292	0.152	0.367	0.133	0.411	0.151	0.308
EquiPocket-LG	0	0.50	0.235	0.225	0.442	0.183	0.498	0.273	0.463
EquiPocket-LG	1	0.50	0.207	0.220	0.460	0.189	0.509	0.280	0.468
EquiPocket-LG	2	0.50	0.203	0.184	0.440	0.180	0.510	0.269	0.459
EquiPocket-LG	3	0.50	0.224	0.215	0.448	0.186	0.500	0.275	0.465
EquiPocket-LG	4	0.50	0.231	0.213	0.431	0.179	0.492	0.272	0.456
EquiPocket	0	1.70	0.054	0.423	0.656	0.341	0.665	0.558	0.715
EquiPocket	1	1.70	0.053	0.431	0.660	0.329	0.668	0.538	0.725
EquiPocket	2	1.70	0.041	0.443	0.664	0.336	0.660	0.550	0.724
EquiPocket	3	1.70	0.051	0.411	0.646	0.338	0.668	0.532	0.723
EquiPocket	4	1.70	0.053	0.407	0.654	0.345	0.652	0.546	0.719