VN-EGNN: E(3)- AND SE(3)-EQUIVARIANT GRAPH NEURAL NETWORKS WITH VIRTUAL NODES ENHANCE PROTEIN BINDING SITE IDENTIFICATION

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

020

021

024

025

026

027 028 029 Paper under double-blind review

ABSTRACT

Being able to identify regions within or around proteins, to which ligands can potentially bind, is an essential step in developing new drugs. Binding site identification methods can now profit from the availability of large amounts of 3D structures in protein structure databases or from AlphaFold predictions. Current binding site identification methods heavily rely on graph neural networks (GNNs), usually designed to output E(3)-equivariant predictions. Such methods turned out to be very beneficial for physics-related tasks like binding energy or motion trajectory prediction. However, the performance of GNNs at binding site identification is still limited potentially due to a lack of expressiveness capable of modeling higher-order geometric entities, such as binding pockets. In this work, we extend E(n)-equivariant graph neural networks (EGNNs) by adding virtual nodes and applying an extended message passing scheme. The virtual nodes in these graphs are dedicated entities to learn representations of binding sites, which leads to improved predictive performance. In our experiments, we show that our proposed method, VN-EGNN, sets a new state-of-the-art at locating binding site centers on COACH420, HOLO4K and PDBbind2020.

1 INTRODUCTION

031 Binding site identification remains a central computational problem in drug discovery. With the advent of AlphaFold (Jumper et al., 2021; Abramson et al., 2024), millions of 3D structures of 033 proteins have been unlocked for further investigation by the scientific community (Tunyasuvunakool 034 et al., 2021; Cheng et al., 2023). The 3D structure of a protein can provide crucial information about its function, and drug discovery is one of the most important fields that profits from these 035 3D structures (Ren et al., 2023; Sadybekov and Katritch, 2023). It has been envisioned that the 036 availability of 3D structures will allow to purposefully design drugs that alter protein function in a 037 desired way. However, to enable structure-based drug design, further computational approaches, such as *docking* or *binding site identification* methods, have to be employed (Lengauer and Rarey, 1996; Cheng et al., 2007; Halgren, 2009). While docking approaches predict the location of a specific small 040 molecule, called a ligand, within a protein's active site upon binding, binding site identification aims 041 at finding regions on the protein likely to form a binding pocket and interact with unknown ligands 042 (Schmidtke and Barril, 2010). Note that docking and binding site identification are fundamentally 043 different tasks in structure-based drug design: for the vast majority of proteins no ligand is known, 044 and binding site identification methods can provide valuable information for understanding protein function, guiding rational drug design or identifying a protein as a potential drug target. For both approaches, deep learning methods, and specifically geometric deep learning have brought significant 046 advances (Gainza et al., 2020; Sverrisson et al., 2021; Méndez-Lucio et al., 2021; Ganea et al., 2022; 047 Stärk et al., 2022; Lu et al., 2022; Corso et al., 2023). 048

Methods for binding site identification. The identification of binding sites relies on the successful combination of physical, chemical and geometric information. Initially, machine learning methods for binding site prediction were based on carefully designed input features due to their tabular processing structure. For instance, FPocket (Le Guilloux et al., 2009) relies on Voronoi tessellation and alpha spheres (Liang et al., 1998) and additionally takes an electronegativity criterion into account. P2Rank, a random-forest-based method, makes use of the protein surface (Krivák and Hoksza, 2018). With



Figure 1: Overview of binding site identification methods. Top Left: Traditional methods, based on segmentation of a voxel grid, in which the pocket center is calculated as the geometric center of the positively labeled voxels. Bottom Left: Geometric Deep Learning approaches, such as EGNN, in which the pocket center is calculated as the geometric center of the positively labeled nodes. Right:
VN-EGNN approach (ours): the predicted binding site center is the position of the virtual node after *L* message passing layers.

076

069

077 the advent of end-to-end deep learning and especially with the breakthrough of convolutional (Lecun 078 et al., 1998) and graph neural networks (GNNs) (Scarselli et al., 2009; Defferrard et al., 2016; Kipf 079 and Welling, 2017; Gilmer et al., 2017; Satorras et al., 2021), the construction of input features 080 can be learned which helped to advance predictive quality. For instance, DeepSite (Jiménez et al., 081 2017) is a voxel-based 3D convolutional neural network for binding site prediction. Convolutional operations on the 3D space are, however, computationally very demanding and so quickly other 083 approaches to tackle binding site identification were developed, e.g., DeepSurf (Mylonas et al., 2021) or PointSite (Yan et al., 2022). DeepSurf operates on surface-based representations and places several 084 voxelized grids on the protein's surface, while PointSite is based on a form of sparse convolutions 085 to reduce the computational overhead and keep sparse regions in the 3D space sparse. Typical 086 convolutional networks, however, do not perform well at binding site identification, likely because 087 of the irregularity of protein structures and due to the fact that proteins may be arbitrarily rotated 088 and shifted in space (Zhang et al., 2023b). Thus, geometric deep learning approaches, most notably 089 (graph-based) group-equivariant architectures, such as EquiPocket (Zhang et al., 2023b), which are 090 equivariant to the group of Euclidean transformations in 3D space (E(3)), are powerful methods for 091 binding site identification. 092

E(3)-equivariant graph neural networks. We use graph neural networks (GNNs) that are robust 093 to transformations of the Euclidean group, i.e., rotations, reflections, and translations, as well as to 094 permutations. From a technical point of view, equivariance of a function f to certain transformations means that for any transformation parameter g and all inputs x we have $f(T_q(x)) = S_q(f(x))$, where 096 T_q and S_q denote transformations on the input and output domain of f, respectively (see App. D for further information). Equivariant operators applied to molecular graphs allow to preserve the 098 geometric structure of the molecule. We build on E(n)-equivariant GNNs (EGNNs) of Satorras et al. 099 (2021) applied to three dimensional space and the problem of binding site identification. In contrast to methods such as MACE (Batatia et al., 2022), Nequip (Batzner et al., 2022), or SEGNN (Brandstetter 100 et al., 2021), EGNNs operate on scalar features, e.g., distances, and use scale operations for coordinate 101 updates. Thus, EGNNs operate efficiently (Villar et al., 2021) without resorting to compute-expensive 102 higher order features, and, most importantly, allow for efficient coordinate update of virtual nodes. 103

Limitations of GNNs and a mitigation strategy. Graph neural networks can suffer from limited
 expressiveness (Morris et al., 2019; Xu et al., 2019), oversmoothing (Li et al., 2018; Rusch et al., 2023), or oversquashing (Alon and Yahav, 2021; Topping et al., 2022), which can lead to unfavorable
 learning dynamics or weak predictive performance. To improve the learning dynamics of GNNs, several works have introduced virtual nodes, sometimes called super-nodes or supersource-nodes,

108 that are introduced into a message passing scheme and connected to all other nodes. In a benchmark 109 setting, Hu et al. (2020) and Rosenbluth et al. (2024) showed that adding virtual nodes tends to 110 increase the predictive performance. Hwang et al. (2022) provide a theoretical analysis of the benefits 111 of virtual nodes in terms of expressiveness, demonstrate the increased expressiveness of GNNs with 112 virtual nodes and also hint at the fact that such nodes can decrease oversmoothing. Alon and Yahay (2021) mention that virtual nodes might be used as a technique to overcome oversquashing effects. 113 Cai et al. (2023) and Cai (2023) show that an MPNN with one virtual node, connected to all nodes, 114 can approximate a Transformer layer. Low rank global attention (Puny et al., 2020) can be seen 115 as one virtual node, which improves expressiveness. Practically, virtual nodes have already been 116 suggested in the original work by Gilmer et al. (2017), and they were even mentioned earlier in 117 Scarselli et al. (2009) and used in application areas such as drug discovery (Li et al., 2017; Pham 118 et al., 2017; Ishiguro et al., 2019). Joshi et al. (2023) investigated the expressive power of EGNNs 119 in greater detail and argue that these networks can suffer from oversquashing. In order to alleviate 120 the oversquashing problem of EGNNs for binding site identification, we suggest to extend EGNNs 121 with virtual nodes and introduce an adapted message passing scheme. We refer to this method as 122 Virtual-Node Equivariant GNN (VN-EGNN). A related method (MEAN) to ours, which is in the 123 context of EGNNs and which uses global nodes, that are connected to many other graph nodes (i.e., all within components), was suggested by Kong et al. (2023) for conditional antibody design. MEAN 124 (applied to components) can potentially be considered to be the first EGNN-like architecture using 125 virtual nodes. 126

127 EGNNs with virtual nodes for binding site identification. In accordance with previous approaches, 128 we consider binding site identification as a segmentation task. While other methods are, e.g., based 129 on voxel grids, our method is based on EGNNs with virtual nodes, where all atoms or residues of the protein (the physical entities) are represented by physical, i.e., non-virtual, nodes in the graph 130 (see Fig. 1, left). The objective is to correctly classify whether a node is within a certain radius of a 131 region to which potential ligands can bind. Therefore, binding site identification can be considered 132 as a node-level binary classification task and thus a semantic segmentation task. For this task, the 133 ground truth is whether an atom was within a certain radius to experimentally observed protein 134 binding ligands. In addition to node features, EGNNs act on coordinate features associated with each 135 node, and both feature types are updated during message passing. While it appears straightforward to 136 associate physical nodes with the protein's atoms, it is a-priori unclear if the coordinate embeddings 137 of virtual nodes are useful for the task at hand. In an initial experiment, we trained VN-EGNN to learn 138 a semantic segmentation task using multiple virtual nodes to which we assigned random coordinates. 139 In an analysis of the results, we could empirically observe that coordinates of the virtual nodes 140 converged towards the actual physical binding positions of ligands on the protein (see App. H.1 for further details and Fig. 1, right, for a visualization). The results of this initial experiment gave rise to 141 the assumption that virtual nodes enable VN-EGNNs to form useful neural representations of binding 142 sites and especially allow the prediction of locations of binding site centers. This, however, further 143 implies that the binding site center itself may be a useful optimization target to train VN-EGNNs. 144 Thus, we extended our objective from only predicting whether physical nodes are close to binding 145 regions, to also directly taking the distance between observed and predicted binding site centers into 146 account. With this multi-modal objective, the coordinate embeddings of virtual nodes are trained 147 to predict the locations of binding site centers. The remaining features of the virtual nodes are 148 considered to form an abstract neural representation of a protein binding site. 149

Contributions. In this work, we aim at improving binding site identification through geometric deep learning methods. Here, we follow the approach of using EGNNs (Satorras et al., 2021; Zhang et al., 2023b) for identifying binding pockets. Although EGNNs are prime candidates for this task, traditional EGNNs exhibit poor performance at binding site identification (Zhang et al., 2023b), which might be due to a) their lack of dedicated nodes that can learn representations of binding sites, and b) oversquashing effects which hamper learning (Alon and Yahav, 2021; Topping et al., 2022; Joshi et al., 2023). We aim to alleviate both problems by using EGNNs with virtual nodes. In this work, we contribute the following:

157 158 159

• We propose to adapt E(3)-equivariant GNNs towards the identification of binding sites of proteins.

• We demonstrate that the virtual nodes in the message passing scheme learn useful representations and accurate locations of binding pockets. • We assess the performance of other methods, baselines and our method on benchmarking datasets.

167

162

2 E(3)-EQUIVARIANT GRAPH NEURAL NETWORKS WITH VIRTUAL NODES

168 2.1 NOTATIONAL PRELIMINARIES

169 We give an overview on variable and symbol notation in App. B, and a more detailed description and 170 discussion on how we represent proteins and binding sites in Apps. C.1 and C.2, respectively. To 171 quickly summarize, the coordinates of the *i*-th physical node, e.g., the location of an atom of a protein, are denoted as $\mathbf{x}_i \in \mathbb{R}^3$, and its other node features as $h_i \in \mathbb{R}^D$. l added as an index to symbols 172 173 will indicate neural network layers, but might be omitted for simplicity sometimes if it is clear from 174 the context. We will consider virtual nodes and the k-th virtual node coordinates will be denoted as \mathbf{z}_k , while the other virtual node features will be denoted as v_k . We use upper-case bold letters to 175 denote the matrices collecting the coordinates and features of the N physical and the K virtual nodes, 176 respectively: $H^{l} := (h_{1}^{l}, \dots, h_{N}^{l}), \mathbf{X}^{l} := (\mathbf{x}_{1}^{l}, \dots, \mathbf{x}_{N}^{l}), V^{l} := (v_{1}^{l}, \dots, v_{K}^{l}), \mathbf{Z}^{l} := (\mathbf{z}_{1}^{l}, \dots, \mathbf{z}_{K}^{l}).$ 177 178 We denote the graph, which VN-EGNN works upon as \mathcal{G} and A as the associated adjacency matrix. 179 $\mathcal{N}(i)$ indicates neighbouring nodes to node i within \mathcal{G} and edge features between nodes i and j as a_{ij} (for two non-virtual nodes) and d_{ij} (between physical and virtual nodes). At training time, we have access to node-level labels $y_i \in \{0, 1\}$ and a set of M center coordinates $\{\mathbf{y}_m\}_{m=1}^M$ with $\mathbf{y}_m \in \mathbb{R}^3$, 181 which the model should predict. We denote predicted node-level labels by $\hat{y}_i \in [0,1]$ and the set of 182 K predicted center coordinates from the model by $\{\hat{\mathbf{y}}_k\}_{k=1}^K$ with $\hat{\mathbf{y}}_k \in \mathbb{R}^3$.

183 184 185

2.2 EGNNS AND THEIR APPLICATION TO BINDING SITE IDENTIFICATION

EGNNs are straightforward to apply to proteins, when they are represented by a neighborhood graph \mathcal{P} , in which each node represents an atom and edges between two atoms represent spatially close atoms (distance between the atoms in the protein is below some threshold). To apply EGNNs, we first set $\mathcal{G} = \mathcal{P}$. For binding pocket identification, one could predict node labels y_i , which indicate whether the atom belongs to a binding pocket or not.

The physical nodes represent atoms and their initial coordinate features are set to the location of the atoms \mathbf{x}_i^0 , and the initial node features h_i^0 to, e.g., the atom or residue type. Then, we apply the layer-wise message passing scheme $(\mathbf{X}^{l+1}, \mathbf{H}^{l+1}) = \text{EGNN}(\mathbf{X}^l, \mathbf{H}^l, \mathbf{A})$ (Eqs. (1) to (4)) as given by Satorras et al. (2021):

$$\boldsymbol{m}_{ij} = \boldsymbol{\phi}_e(\boldsymbol{h}_i^l, \boldsymbol{h}_j^l, \|\mathbf{x}_i^l - \mathbf{x}_j^l\|^2, a_{ij})$$
(1)

$$\boldsymbol{m}_i = \sum_{j \in \mathcal{N}(i)} \boldsymbol{m}_{ij} \tag{2}$$

199 200

196

$$\mathbf{x}_{i}^{l+1} = \mathbf{x}_{i}^{l} + \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} \frac{\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}}{\|\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}\|} \phi_{x}(\boldsymbol{m}_{ij})$$
(3)

201

$$\boldsymbol{h}_{i}^{l+1} = \boldsymbol{\phi}_{h}\left(\boldsymbol{h}_{i}^{l}, \boldsymbol{m}_{i}\right), \tag{4}$$

where ϕ_e , ϕ_x and ϕ_h denote multilayer-perceptrons (MLPs). To identify binding pockets, we can extract predictions \hat{y}_i for each atom *i* by a read-out function applied to the output of the last message passing step *L*, i.e., $\hat{y}_i = \sigma(\boldsymbol{w}^\top \boldsymbol{h}_i^L)$ with an activation function σ and parameters \boldsymbol{w} . Our model does not incorporate edge features, symbolized by a_{ij} . Hence, we will exclude these from the subsequent EGNN formulations and their related derivations.

209 210

211

2.3 VN-EGNN: EXTENSION OF EGNN WITH VIRTUAL NODES

1

We now extend \mathcal{G} (which is set to the protein neighborhood graph \mathcal{P} for the task of protein binding site identification) with a set of K virtual nodes, which exhibit edges to all other nodes, which will allow us to learn representations of hidden geometric entities, such as binding sites, and simultaneously ameliorate oversquashing. To be able to process this extended graph, we modify EGNNs by locating the virtual nodes at coordinates $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_K) \in \mathbb{R}^3$ and associating them with a set of properties $V = (v_1, ..., v_K) \in \mathbb{R}^D$. The new message passing scheme ($X^{l+1}, H^{l+1}, Z^{l+1}, V^{l+1}$) = VN-EGNN (X^l, H^l, Z^l, V^l) of a single VN-EGNN layer consists of three phases (Eqs. (7) to (10), Eqs. (11) to (14), and, Eqs. (15) to (18)), in which the feature and coordinate embeddings of the physical nodes are updated twice:

$$\boldsymbol{h}_{i}^{l} \rightarrow \boldsymbol{h}_{i}^{l+1/2} \rightarrow \boldsymbol{h}_{i}^{l+1}, \quad \mathbf{x}_{i}^{l} \rightarrow \mathbf{x}_{i}^{l+1/2} \rightarrow \mathbf{x}_{i}^{l+1} \quad \forall i$$
 (5)

while virtual node embeddings are only updated once per message passing layer

$$\mathbf{v}_k^l \to \mathbf{v}_k^{l+1}, \quad \mathbf{z}_k^l \to \mathbf{z}_k^{l+1} \quad \forall k.$$
 (6)

Message passing phase I between physical nodes (analogous to EGNN):

$$\boldsymbol{m}_{ij}^{(aa)} = \boldsymbol{\phi}_{e^{(aa)}}(\boldsymbol{h}_{i}^{l}, \boldsymbol{h}_{j}^{l}, \|\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}\|)$$
(7)

$$\boldsymbol{m}_{i}^{(aa)} = \sum_{j \in \mathcal{N}(i)} \boldsymbol{m}_{ij}^{(aa)} \tag{8}$$

$$\mathbf{x}_{i}^{l+1/2} = \mathbf{x}_{i}^{l} + \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} \frac{\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}}{\|\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}\|} \phi_{\mathbf{x}^{aa}}(\boldsymbol{m}_{ij}^{(aa)})$$
(9)

$$\boldsymbol{h}_{i}^{l+1/2} = \boldsymbol{h}_{i}^{l} + \boldsymbol{\phi}_{h^{(aa)}} \left(\boldsymbol{h}_{i}^{l}, \boldsymbol{m}_{i}^{(aa)} \right).$$
(10)

Message passing phase II from physical nodes to virtual nodes:

$$\boldsymbol{m}_{ij}^{(av)} = \boldsymbol{\phi}_{e^{(av)}}(\boldsymbol{v}_i^l, \boldsymbol{h}_j^{l+1/2}, \|\mathbf{z}_i^l - \mathbf{x}_j^{l+1/2}\|)$$
(11)

$$\boldsymbol{m}_{i}^{(av)} = \frac{1}{N} \sum_{j=1}^{N} \boldsymbol{m}_{ij}^{(av)}$$
(12)

$$\mathbf{z}_{i}^{l+1} = \mathbf{z}_{i}^{l} + \frac{1}{N} \sum_{j=1}^{N} \frac{\mathbf{z}_{i}^{l} - \mathbf{x}_{j}^{l+1/2}}{\|\mathbf{z}_{i}^{l} - \mathbf{x}_{j}^{l+1/2}\|} \phi_{\mathbf{x}^{av}}(\boldsymbol{m}_{ij}^{(av)})$$
(13)

$$\boldsymbol{v}_{i}^{l+1} = \boldsymbol{v}_{i}^{l} + \boldsymbol{\phi}_{h^{(av)}} \left(\boldsymbol{v}_{i}^{l}, \boldsymbol{m}_{i}^{(av)} \right)$$
(14)

Message passing phase III from virtual nodes to physical nodes:

$$\boldsymbol{m}_{ij}^{(va)} = \boldsymbol{\phi}_{e^{(va)}}(\boldsymbol{h}_i^{l+1/2}, \boldsymbol{v}_j^{l+1}, \|\mathbf{x}_i^{l+1/2} - \mathbf{z}_j^{l+1}\|)$$
(15)

$$m_i^{(va)} = \sum_{j=1}^K m_{ij}^{(va)}$$
 (16)

$$\mathbf{x}_{i}^{l+1} = \mathbf{x}_{i}^{l+1/2} + \frac{1}{K} \sum_{j=1}^{K} \frac{\mathbf{x}_{i}^{l+1/2} - \mathbf{z}_{j}^{l+1}}{\|\mathbf{x}_{i}^{l+1/2} - \mathbf{z}_{j}^{l+1}\|} \phi_{\mathbf{x}^{va}}(\boldsymbol{m}_{ij}^{(va)})$$
(17)

$$\boldsymbol{h}_{j=1}^{l+1} = \boldsymbol{h}_{i}^{l+1/2} + \boldsymbol{\phi}_{i}(v_{a}) \left(\boldsymbol{h}_{i}^{l+1/2}, \boldsymbol{m}_{i}^{(v_{a})} \right)$$
(18)

$$\boldsymbol{h}_{i}^{l+1} = \boldsymbol{h}_{i}^{l+1/2} + \boldsymbol{\phi}_{h^{(va)}} \left(\boldsymbol{h}_{i}^{l+1/2}, \boldsymbol{m}_{i}^{(va)} \right)$$
(18)

Here, $\phi_{e^{(aa)}}, \ldots, \phi_{h^{(va)}}$ are again MLPs. The MLPs ϕ_l are layer-specific, i.e. ϕ_l^l and currently do not consider edge features. To keep the notation uncluttered, we skipped the layer index l for the MLPs in the formulae above. We call this message passing scheme *heterogeneous* because the different types of messages are generated in subsequent phases.

2.4 INITIALIZATION OF VIRTUAL NODES IN VN-EGNN.

It is possibly to equivariantly initialize virtual nodes. To do so, we can initialize the coordinates of the virtual nodes z_k^0 at the center of mass, concretely the average coordinates of the physical nodes (Zhang et al., 2024; Kaba et al., 2023), while the initial features v_k^0 are learned feature vectors that are pairwise distinct. These choices lead to maintenance of the equivariance properties, but guarantee that the virtual nodes are updated differently during message passing. Note, that the center of mass is
not necessarily a very meaningful value for binding site identification, although it ensures invariant
binding site predictions. In practice, this might therefore restrict the architecture quite a lot. In line
with new developments such as Abramson et al. (2024), we relaxed the initialization procedure to
specifically exploit our prior knowledge, that ligands might tend to bind to surface areas of proteins.

In detail, the relaxed initialization procedure distributes the K virtual nodes evenly across a sphere using a Fibonacci grid (Swinbank and James Purser, 2006), of which the radius is defined as the distance between the protein center and its most distant atom. The virtual node properties v_k^0 are initialized by averaging over the initial features h_i^0 . This procedure is simple and efficient.

Data augmentation To prevent virtual nodes from starting at identical locations during different epochs of training, we randomly rotate the sphere with the initial locations of virtual nodes.

2.5 PROPERTIES OF VN-EGNN

The following proposition shows, that analogously to EGNNs, VN-EGNNs are equivariant with respect to roto-translations and reflections by construction.

Proposition 1. Equivariant graph neural networks with virtual nodes as defined in Eqs. (7) to (18) are equivariant with respect to roto-translations and reflections of the input and virtual node coordinates.

Proof. See App. E.

280

281 282

283 284

285

286

287

288 289

290

291 Virtual nodes ameliorate oversquashing by bounding the maximal shortest-path distance 292 between nodes and required message passing steps. Several works (Alon and Yahav, 2021; 293 Topping et al., 2022; Di Giovanni et al., 2023) have investigated the relation between oversquashing 294 and characteristics of the MPNN layers and the adjacency matrix. According to Topping et al. (2022) oversquashing is defined as $\frac{\partial h_i^{r+1}}{\partial h_j^0}$, which is the effect that one node with index j has on a node with 295 296 index i during learning, where the nodes are at a shortest-path distance of r + 1. Critically, this 297 quantity can be bounded by the model parameters of the involved MLPs and the topology of the 298 graph (Di Giovanni et al., 2023), concretely the normalized adjacency matrix. We use Topping et al. 299 (2022, Lemma 1), which states that $\left|\frac{\partial h_i^{r+1}}{\partial h_j^0}\right| \leq (\alpha\beta)^{r+1} (\hat{A}^{r+1})_{ij}$ where α and β are bounds on the 300 element-wise gradients of the MLPs of the message passing network, h_i^{r+1} is one component of the 301 302 node representation of node i in message passing layer r + 1. The quantity r + 1 is both the number 303 of message passing layers and the shortest-path distance between nodes i and j in the graph, and \hat{A} 304 is the normalized adjacency matrix, for which the diagonal values of the original matrix are set to 1. The normalized adjacency matrix \hat{A} is a symmetric positive matrix that has a leading eigenvalue at 305 306 1 (Perron, 1907; Frobenius, 1912), such that all eigenvalues of all other eigenvectors of \hat{A}^r decay 307 exponentially with r. Depending on the weights and activation functions of the MLPs, $|\alpha\beta|^{r+1}$ either grows or vanishes exponentially with r, which might lead to either exploding or vanishing gradients, 308 respectively. Thus, learning can only be stabilized via keeping r stable, which virtual nodes that are 309 connected to all other nodes can provide since they bound both the maximal path distance and the 310 necessary number of message passing steps by r + 1 = 2. 311 312

Expressiveness of VN-EGNN. The expressive power of GNNs is linked to their ability to distinguish non-isomorphic graphs. While a minimum of k layers of an EGNN is required to distinguish two k-hop distinct graphs, one layer of VN-EGNN is presumed to be sufficient, as can be shown by the application of the Geometric Weisfeiler-Leman test, which serves as an upper bound on the expressiveness of EGNNs. Experimental findings on k-chain geometric graphs support this proposition and demonstrate the increased expressive power of VN-EGNN compared to EGNNs without virtual nodes. For further details and an empirical study see App. K.

319

320 2.6 Adjustments of VN-EGNN for binding site identification.321

SE(3) equivariance through feature encoding. We break the equivariance property for mirroring through feature encoding. Since each node has an initial feature that codes for the amino acid, either one-hot encoding or ESM embeddings (Lin et al., 2023), we naturally encode L- and D-amino

324 acids differently, which leads to different initial features of the initial nodes after mirroring, and 325 consequently breaks E(3) symmetry to SE(3). 326

2.7 TRAINING VN-EGNNS

Objective. Previous methods, which consider binding site identification as a node-level prediction task (see Section 2.2) with $\hat{y}_n = \sigma(\boldsymbol{w}^{\top} \boldsymbol{h}_n^{\perp})$, use a type of segmentation loss. The segmentation loss can either be the cross-entropy loss CE:

331 332

327

328

330

335

336 337 338

347

348 349

351

353

361

362

364 365

366

or the Dice loss, that is based on the continuous Dice coefficient (Shamir et al., 2019), with $\epsilon = 1$:

$$\mathcal{L}_{\text{dice}} \coloneqq 1 - \frac{2 \sum_{n=1}^{N} y_n \, \hat{y}_n + \epsilon}{\sum_{n=1}^{N} y_n + \sum_{n=1}^{N} \hat{y}_n + \epsilon}$$

 $\mathcal{L}_{\text{segm}} = \frac{1}{N} \sum_{n=1}^{N} \text{CE}(y_n, \hat{y}_n)$

339 The introduction of virtual nodes with coordinates allows to directly tackle the more challenging 340 problem of predicting binding site center points and extracting predictions for these points as outputs 341 of the last EGNN layer. For each protein in the training set, we know the geometric centers of its 342 annotated binding sites, which we denote as $\{\mathbf{y}_1, \dots, \mathbf{y}_M\}$. The read-out $\hat{\mathbf{y}}_k$ for each virtual node 343 $\hat{\mathbf{y}}_k := \mathbf{z}_k^L$ $(1 \le k \le K)$ corresponds to its coordinate embedding \mathbf{z}_k^L in the last layer L. Each known 344 binding site center should be detected by at least one virtual node, via its read-out, which leads to the 345 following objective

$$\mathcal{L}_{\rm bsc} = \frac{1}{M} \sum_{m=1}^{M} \min_{k \in 1, \dots, K} \|\mathbf{y}_m - \hat{\mathbf{y}}_k\|^2.$$
(19)

The full objective of VN-EGNN for binding site identification is 350

 $\mathcal{L} = \mathcal{L}_{\rm bsc} + \mathcal{L}_{\rm dice},$

352 in which the two terms could also be balanced against each other through a hyperparameter, which we found was not necessary.

354 Self-confidence module. We employ a self-confidence module (Jumper et al., 2021; Zhang et al., 355 2023a), to assess the quality of predicted binding sites, by equipping each prediction with a confidence 356 score. This allows a ranking of the predictions, similar to Krivák and Hoksza (2018). The confidence 357 value, indicated by \hat{c}_k , is computed through $\hat{c}_k = \psi(\boldsymbol{v}_k)$, with ψ implemented as an MLP. During training, the target values for the confidence prediction are generated on-the-fly from the predicted positions $\hat{\mathbf{y}}_k$ and the closest known binding pocket center \mathbf{y}_m , in analogy with confidence scores for 360 object detection methods in computer vision.

The confidence label for the k-th virtual node is obtained by (Zhang et al., 2023a):

$$c_{k} = \begin{cases} 1 - \frac{1}{2\gamma} \cdot \|\mathbf{y}_{m} - \hat{\mathbf{y}}_{k}\| & \text{if } \|\mathbf{y}_{m} - \hat{\mathbf{y}}_{k}\| \le \gamma, \\ c_{0} & \text{otherwise} \end{cases},$$
(20)

with $c_0 = 0.001$. To align with the commonly accepted threshold value of 4Å for the DCC/DCA success rates, we choose $\gamma = 4$. The loss on the confidence score is a mean squared error loss:

$$\mathcal{L}_{\text{confidence}} = \frac{1}{K} \sum_{k=1}^{K} (c_k - \hat{c}_k)^2.$$
(21)

EXPERIMENTS 3

372 373 374

375

3.1 DATA

We use the benchmarking setting of Zhang et al. (2023b) to perform experiments on four datasets for 376 binding site identification: scPDB (Desaphy et al., 2015), PDBbind (Wang et al., 2004), COACH420 377 and HOLO4K. For details, see App. G.1.

Table 1: Performance at binding site identification in terms of DCC and DCA success rates.^a The first column provides the method, the second the number of parameters of the model, the fourth and the fifth column the performance on the COACH420 dataset, the sixth and seventh column the performance on the HOLO4K dataset, and the remaining columns the performance on PDBbind2020. The best performing method(s) per column are marked bold. The second best in italics.

Methods		Param COACH420		HOLO4K ^d		PDBbind2020	
methods	(M)	DCC↑	DCA↑	DCC↑	DCA↑	DCC↑	DCA↑
Fpocket (Le Guilloux et al., 2009) ^b P2Rank (Krivák and Hoksza, 2018) ^c	\ \	0.228 0.464	0.444 0.728	0.192 0.474	0.457 0.787	0.253 0.653	0.371 0.826
DeepSite (Jiménez et al., 2017) ^b Kalasanty (Stepniewska-Dziubinska et al., 2020) ^b	1.00 70.64	0.335	0.564	\ 0.244	0.456	\ 0.416	\ 0.625
DeepSurf (Mylonas et al., 2021) ^b DeepPocket (Aggarwal et al., 2022b) ^c	33.06	0.386 0.399	0.658 0.645	0.289 0.456	0.635 0.734	0.510 0.644	0.708 0.813
GAT (Veličković et al., 2018) ^b GCN (Kipf and Welling, 2017) ^b GAT + GCN ^b GCN2 (Chen et al., 2020) ^b	0.03 0.06 0.08 0.11	0.039(0.005) 0.049(0.001) 0.036(0.009) 0.042(0.098)	0.130(0.009) 0.139(0.010) 0.131(0.021) 0.131(0.017)	0.036(0.003) 0.044(0.003) 0.042(0.003) 0.051(0.004)	0.110(0.010) 0.174(0.003) 0.152(0.020) 0.163(0.008)	0.032(0.001) 0.018(0.001) 0.022(0.008) 0.023(0.007)	0.088(0.01 0.070(0.00) 0.074(0.00) 0.089(0.01)
SchNet (Schütt et al., 2017) ^b EGNN (Satorras et al., 2021) ^b	0.49 0.41	0.168(0.019) 0.156(0.017)	0.444(0.020) 0.361(0.020)	0.192(0.005) 0.127(0.005)	0.501(0.004) 0.406(0.004)	0.263(0.003) 0.143(0.007)	0.457(0.004
EquiPocket (Zhang et al., 2023b) ^b	1.70	0.423(0.014)	0.656(0.007)	0.337(0.006)	0.662(0.007)	0.545(0.010)	0.721(0.00
VN-EGNN (ours)	1.20	0.605(0.009)	0.750(0.008)	0.532(0.021)	0.659(0.026)	0.669(0.015)	0.820(0.01

^a The standard deviation across training re-runs is indicated in parentheses.

^b Results from Zhang et al. (2023b).

^c Uses different training set and, thus, limited comparability

^d This dataset represents a strong domain shift from the training data for all methods (except for P2Rank). Details on the domain shift in App. J.

Table 2: Ablation study. The main components of the VN-EGNN architecture are ablated and tested for their performance on the benchmarking datasets. The first column reports the variant of the ablated method, the second column whether the method contains virtual nodes (VN), the third column whether the method applies heterogenous message passing, and the fourth column whether ESM embeddings were used. The remaining columns are analogous to Table 1.

M-d-d-	VN	heterog.	EGM	COAC	CH420	HOL	.04K	PDBbi	nd2020
Methods	VIN	MP	ESM	DCC↑	DCA↑	DCC↑	DCA↑	DCC↑	DCA↑
EGNN+VN (Satorras et al., 2021) ^b	X	X	x	0.156(0.017)	0.361(0.020)	0.127(0.005)	0.406(0.004)	0.143(0.007)	0.302(0.006)
VN-EGNN (VN only)	~	×	X	0.497(0.014)	0.700(0.013)	0.414(0.023)	0.618(0.024)	0.502(0.029)	0.717(0.025)
VN-EGNN (residue emb.)	\checkmark	\checkmark	X	0.503(0.022)	0.684(0.016)	0.438(0.019)	0.605(0.013)	0.551(0.017)	0.751(0.009)
VN-EGNN (homog.)	\checkmark	×	~	0.575(0.008)	0.708(0.009)	0.479(0.012)	0.595(0.010)	0.649(0.010)	0.805(0.006)
VN-EGNN (full)	\checkmark	\checkmark	\checkmark	0.605(0.009)	0.750(0.008)	0.532(0.021)	0.659(0.026)	0.669(0.015)	0.820(0.010)

a The standard deviation across training re-runs is indicated in parentheses

^b Results from Zhang et al. (2023b)

405 406 407

408

404

397

3.2 EVALUATION

Methods compared. We compare the following binding site identification methods from different categories: *Geometry-based*: Fpocket (Le Guilloux et al., 2009) and P2Rank (Krivák and Hoksza, 2018). *CNN-based*: DeepSite (Jiménez et al., 2017), Kalasanty (Stepniewska-Dziubinska et al., 2020), and DeepSurf (Mylonas et al., 2021). *Topological graph-based*: GAT (Veličković et al., 2018), GCN (Kipf and Welling, 2017), and GCN2 (Chen et al., 2020). *Spatial graph-based*: SchNet (Schütt et al., 2017), EGNN (Satorras et al., 2021), EquiPocket (Zhang et al., 2023b), and our proposed VN-EGNN.

Evaluation metrics. We used the DCC/DCA success rate, which are well-established metrics for 416 binding site identification (see e.g., Chen et al., 2011). DCC is defined as the distance between the 417 predicted and known binding site centers, whereas DCA is defined as the shortest distance between 418 the predicted center and any atom of the ligand. Following Stepniewska-Dziubinska et al. (2020) 419 and Zhang et al. (2023b), predictions within a certain threshold of DCC and DCA, are considered as 420 successful, which is commonly referred to as DCC/DCA success rate. Adhering to these works, we 421 maintained a threshold of 4Å throughout our experiments (for other thresholds, see Fig. 2). In line 422 with Chen et al. (2011); Zhang et al. (2023b); Stepniewska-Dziubinska et al. (2020) for each protein 423 only M predicted binding sites with the highest self-confidence scores \hat{c}_k are considered, where M is the number of known binding sites of the protein. Subsequently, each predicted binding site was 424 aligned with the closest real binding site and DCC/DCA success rate was calculated. 425

426 427

428

3.3 IMPLEMENTATION DETAILS.

We used an AdamW (Loshchilov and Hutter, 2019) optimizer for 1500 epochs, selecting the best
checkpoint based on the validation dataset. We used 5 VN-EGNN layers, where each layer consists
of the three-step message passing scheme described in Section 2.3, and the feature and message size
was set to 100, in all layers. Due to the possibility of different virtual nodes converging to identical



Figure 2: Left: Model prediction showing initial positions of the virtual nodes (yellow spheres), ground truth ligand (violet), annotated binding site (violet protein regions), and node position changes (arrows). Violet spheres show clustered virtual node predictions with self-confidence scores. For better visualization, only a subset of initial node positions is shown. (PDB: 1MXI-A) **Right:** DCC success rates at varying thresholds for the distance between predicted and known binding pocket centers in Å.

451 locations, we employed the Mean Shift Algorithm (Comaniciu and Meer, 2002) at inference time, to 452 cluster virtual nodes that are in close spatial proximity. By averaging their self-confidence scores and 453 positions, we treated these clustered nodes as a single instance. Because of the large complexes in 454 HOLO4K, we ran VN-EGNN for each chain and merged the predicted pocket centers. For the initial 455 residue node features we used pre-trained ESM-2 (Lin et al., 2023) protein embeddings following 456 Corso et al. (2023); Pei et al. (2023). For virtual nodes, we derived their features by averaging the 457 residue node features across the entire protein. We used the position of the α -carbons as residue node locations. Virtual nodes are connected to all residue nodes, but not to each other. A linear 458 layer was used to map these initial features to the required dimensions (h_n^0, v_k^0) of the model. Layer 459 normalization and Dropout (Srivastava et al., 2014) was applied in each message passing layer. SiLU 460 Hendrycks and Gimpel (2016) activation was used across all layers. Analogous to Pei et al. (2023) we 461 applied normalization (divided by 5) and unnormalization (multiplied by 5) on the coordinates and 462 used the Huber loss (Huber, 1964) for the coordinates, which empirically proved to be slightly more 463 effective. The learning rate was set to 10^{-3} , after 100 epochs we reduced the learning rate by factor 464 of 10^{-1} if the model did not improve for 10 epochs. For training we used 4x NVIDIA A100 40GB 465 GPUs with a batch size of 64 on each GPU. The training time was about 8 hours. Hyperparameters 466 were selected based on a validation dataset which consisted of a 10%-split of the training data (see 467 Table G1).

468 469

470

445

446

447

448

449

450

3.4 Results

Our experimental results demonstrate that our method, VN-EGNN, surpasses all prior approaches 471 in terms of the DCC metric on COACH420, PDBbind2020 and even on the challenging HOLO4K 472 dataset, see Table 1. On COACH420, VN-EGNN exhibits the best DCA score and on PDBbind it 473 yields the same DCA score as P2Rank. Note that there is limited comparability with P2rank since 474 this method uses a different training set that might be closer to HOLO4K. HOLO4K contains many 475 complexes of symmetric proteins (see App. G.2), which should be considered as a strong domain 476 shift to the training data and thus pose a problem for all methods, except P2rank. For a more detailed 477 discussion and a visual analysis, we refer to App. J. Visualizations of the predictions of our model 478 are shown in Fig. 2 and Fig. 11. We evaluate predictions of our model with respect to the Dice Loss 479 and to IoU in App. H.3. Memory utilization is shown in Fig. M1 (see App. M).

480 481 482

3.5 ABLATION STUDY

Our proposed method comprises three main components as compared to typical other methods: (a)
 virtual nodes, (b) heterogenous message passing, and (c) pre-trained protein embeddings as node
 representations. We ablate these three components in a set of experiments (see Table 2). (a) Removing
 virtual nodes. We compared our model to a standard EGNN framework to determine the added value

486 of virtual nodes. In this study, we analysed how a standard EGNN performs compared to our method. 487 Table 2, shows that the standard EGNN architecture did not perform well. (b) Homogeneous message 488 passing. Our approach to message passing, which is applied in a sequential manner, was contrasted 489 with the traditional method where updates across nodes occur in parallel or homogeneously. This 490 evaluation was further enriched by employing identical MLPs for both graph and virtual nodes across all layers, providing a direct comparison of the impact of our message passing strategy. (c) Atom 491 type embeddings. We evaluated the impact of the type of embeddings, as outlined in Section 2.2. 492 Table 2 illustrates that, regardless of the initial embeddings used, our model surpasses all preceding 493 approaches, except P2Rank, in achieving higher DCC success rate across the COACH420 and 494 PDBbind2020 datasets. This was accomplished by adopting a one-hot encoding scheme solely for 495 the amino acid types, complemented by an additional category for the virtual nodes. 496

We provide additional insights on the usage of a different number of virtual nodes and a different number of layers in App. L.

497 498 499

DISCUSSION AND CONCLUSIONS. 4

500 501

502 Main findings. We have introduced a novel method that extends EGNNs (Satorras et al., 2021) with 503 virtual nodes and a heterogeneous message passing scheme. These new assets improve the learning 504 dynamics by ameliorating the oversquashing problem and enable the learning of representations of hidden geometric entities. Concretely, we have developed this method for binding site identification, 505 for which our experiments show that VN-EGNN exhibits high predictive performance in terms of 506 DCC and DCA and sets a new state-of-the-art on COACH420, HOLO4K and PDBbind2020. We 507 attribute our improvement to the direct prediction of binding site centers, rather than inferring them 508 from the geometric center of segmented areas, a common practice in previous methods. Relying 509 on segmentation can lead to inaccuracies, especially if a single erroneous prediction impacts the 510 calculated center. Overall, VN-EGNN yields highly accurate predictions of binding site centers. 511

Comparison with previous work. In contrast to previous methods (Mylonas et al., 2021; Zhang 512 et al., 2023b), which primarily utilized surface information, based on Sanner et al. (1996); Eisenhaber 513 et al. (1995), or methods that operated on atom-level information (Jiménez et al., 2017; Stepniewska-514 Dziubinska et al., 2020; Aggarwal et al., 2022b), our approach exclusively relies on residue-level 515 information, specifically using α -carbons as physical nodes. This strategy significantly enhances 516 computational efficiency during both training and inference due to the reduced size of the input graphs. 517 Our results support the finding by Jumper et al. (2021), that residue-level information inherently 518 includes all relevant side-chain conformations. 519

Limitations. Our method is currently limited to predicting binding pockets of proteins similar to 520 those in PDB. We expect that VN-EGNNs can also be applied to other physical or geometric problems 521 with hidden geometric entities, such as particle flows, however, their performance in these fields 522 remains to be shown. We have developed VN-EGNN with having the application of binding site 523 identification in mind. Usually in this field, there is a very limited number of training data points 524 and therefore methods taking more prior knowledge into account (e.g., in the design of the network 525 architecture, etc.) could prove beneficial over methods not relying much on this knowledge. With 526 more data points available for training the advantage to take prior knowledge into account may however diminish. Note that our method is not a docking method and thus cannot be used to dock 527 ligands to protein structure. However, our predicted binding sites can be used as proposal regions for 528 other methods, which could lead to improved performance and efficiency for docking methods. 529

530 531

532

References

533 Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., Ronneberger, O., Willmore, L., Ballard, A. J., Bambrick, J., Bodenstein, S. W., Evans, D. A., Hung, C.-C., O'Neill, M., Reiman, 534 D., Tunyasuvunakool, K., Wu, Z., Žemgulytė, A., Arvaniti, E., Beattie, C., Bertolli, O., Bridgland, A., Cherepanov, A., Congreve, M., Cowen-Rivers, A. I., Cowie, A., Figurnov, M., Fuchs, F. B., 536 Gladman, H., Jain, R., Khan, Y. A., Low, C. M. R., Perlin, K., Potapenko, A., Savy, P., Singh, S., Stecula, A., Thillaisundaram, A., Tong, C., Yakneen, S., Zhong, E. D., Zielinski, M., Žídek, A., 538 Bapst, V., Kohli, P., Jaderberg, M., Hassabis, D., and Jumper, J. M. (2024). Accurate structure prediction of biomolecular interactions with AlphaFold 3. Nature, 630(8016):493-500.

540 541 542	Aggarwal, R., Gupta, A., Chelur, V., Jawahar, C., and Priyakumar, U. D. (2022a). DeepPocket: Ligand Binding Site Detection and Segmentation using 3D Convolutional Neural Networks. <i>Journal of Chemical Information and Modeling</i> , 62(21):5069–5079.
543 544 545 546	Aggarwal, R., Gupta, A., Chelur, V., Jawahar, C. V., and Priyakumar, U. D. (2022b). Deeppocket: Ligand binding site detection and segmentation using 3d convolutional neural networks. <i>Journal</i> <i>of Chemical Information and Modeling</i> , 62(21):5069–5079.
547 548 549	Alon, U. and Yahav, E. (2021). On the Bottleneck of Graph Neural Networks and its Practical Implications. In <i>International Conference on Learning Representations</i> .
550 551 552	Batatia, I., Kovacs, D. P., Simm, G., Ortner, C., and Csányi, G. (2022). MACE: Higher order equivariant message passing neural networks for fast and accurate force fields. In <i>Advances in Neural Information Processing Systems</i> .
553 554 555	Batzner, S., Musaelian, A., Sun, L., Geiger, M., Mailoa, J. P., Kornbluth, M., Molinari, N., Smidt, T. E., and Kozinsky, B. (2022). E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. <i>Nature communications</i> , 13(1):2453.
557 558 559	Brandstetter, J., Hesselink, R., van der Pol, E., Bekkers, E. J., and Welling, M. (2021). Geometric and Physical Quantities improve E(3) Equivariant Message Passing. In <i>International Conference on Learning Representations</i> .
560 561 562	Cai, C. (2023). Local-to-global Perspectives on Graph Neural Networks. <i>arXiv preprint</i> arXiv:2306.06547.
563 564	Cai, C., Hy, T. S., Yu, R., and Wang, Y. (2023). On the Connection Between MPNN and Graph Transformer. In <i>International Conference on Machine Learning</i> .
566 567 568	Chen, K., Mizianty, M. J., Gao, J., and Kurgan, L. (2011). A Critical Comparative Assessment of Predictions of Protein-Binding Sites for Biologically Relevant Organic Compounds. <i>Structure</i> , 19(5):613–621.
569 570	Chen, M., Wei, Z., Huang, Z., Ding, B., and Li, Y. (2020). Simple and Deep Graph Convolutional Networks. In <i>International Conference on Machine Learning</i> .
571 572 573 574	Cheng, A. C., Coleman, R. G., Smyth, K. T., Cao, Q., Soulard, P., Caffrey, D. R., Salzberg, A. C., and Huang, E. S. (2007). Structure-based maximal affinity model predicts small-molecule druggability. <i>Nature Biotechnology</i> , 25(1):71–75.
575 576 577 578	Cheng, J., Novati, G., Pan, J., Bycroft, C., Žemgulytė, A., Applebaum, T., Pritzel, A., Wong, L. H., Zielinski, M., Sargeant, T., Schneider, R. G., Senior, A. W., Jumper, J., Hassabis, D., Kohli, P., and Avsec, Ž. (2023). Accurate proteome-wide missense variant effect prediction with AlphaMissense. <i>Science</i> , 381(6664):eadg7492.
579 580 581	Comaniciu, D. and Meer, P. (2002). Mean shift: a robust approach toward feature space analysis. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 24(5):603–619.
582 583 584	Corso, G., Stärk, H., Jing, B., Barzilay, R., and Jaakkola, T. S. (2023). DiffDock: Diffusion Steps, Twists, and Turns for Molecular Docking. In <i>International Conference on Learning Representations</i> .
586 587 588	Defferrard, M., Bresson, X., and Vandergheynst, P. (2016). Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. In <i>Advances in Neural Information Processing Systems</i> .
589 590 591	Desaphy, J., Bret, G., Rognan, D., and Kellenberger, E. (2015). sc-PDB: a 3D-database of ligandable binding sites—10 years on. <i>Nucleic Acids Research</i> , 43(D1):D399–D404.
592 593	Di Giovanni, F., Giusti, L., Barbero, F., Luise, G., Lio, P., and Bronstein, M. M. (2023). On over- squashing in message passing neural networks: The impact of width, depth, and topology. In <i>International Conference on Machine Learning</i> .

594 595 596	Eisenhaber, F., Lijnzaad, P., Argos, P., Sander, C., and Scharf, M. (1995). The double cubic lattice method: Efficient approaches to numerical integration of surface area and volume and to dot surface contouring of molecular assemblies. <i>Journal of Computational Chemistry</i> , 16(3):273–284.
597 598 599	Frobenius, F. G. (1912). Über Matrizen aus nicht negativen Elementen. Königliche Akademie der Wissenschaften Berlin.
600 601 602	Gainza, P., Sverrisson, F., Monti, F., Rodolà, E., Boscaini, D., Bronstein, M. M., and Correia, B. E. (2020). Deciphering interaction fingerprints from protein molecular surfaces using geometric deep learning. <i>Nature Methods</i> , 17(2):184–192.
603 604 605 606	Ganea, OE., Huang, X., Bunne, C., Bian, Y., Barzilay, R., Jaakkola, T. S., and Krause, A. (2022). Independent SE(3)-Equivariant Models for End-to-End Rigid Protein Docking. In <i>International</i> <i>Conference on Learning Representations</i> .
607 608 609	Gaulton, A., Bellis, L. J., Bento, A. P., Chambers, J., Davies, M., Hersey, A., Light, Y., McGlinchey, S., Michalovich, D., Al-Lazikani, B., and Overington, J. P. (2011). ChEMBL: a large-scale bioactivity database for drug discovery. <i>Nucleic Acids Research</i> , 40(D1):D1100–D1107.
610 611 612	Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. (2017). Neural Message Passing for Quantum Chemistry. In <i>International Conference on Machine Learning</i> .
613 614	Halgren, T. A. (2009). Identifying and Characterizing Binding Sites and Assessing Druggability. <i>Journal of Chemical Information and Modeling</i> , 49(2):377–389.
615 616 617	Hendrycks, D. and Gimpel, K. (2016). Gaussian Error Linear Units (GELUs). arXiv preprint arXiv:1606.08415.
618 619 620	Hu, W., Fey, M., Zitnik, M., Dong, Y., Ren, H., Liu, B., Catasta, M., and Leskovec, J. (2020). Open Graph Benchmark: Datasets for Machine Learning on Graphs. In <i>Advances in Neural Information Processing Systems</i> .
621 622 623	Huber, P. J. (1964). Robust Estimation of a Location Parameter. <i>Annals of Mathematical Statistics</i> , 35:492–518.
624 625	Hwang, E., Thost, V., Dasgupta, S. S., and Ma, T. (2022). An Analysis of Virtual Nodes in Graph Neural Networks for Link Prediction. In <i>Learning on Graphs Conference</i> .
626 627 628	Ishiguro, K., Maeda, Si., and Koyama, M. (2019). Graph Warp Module: an Auxiliary Module for Boosting the Power of Graph Neural Networks in Molecular Graph Analysis. <i>arXiv preprint arXiv:1902.01020</i> .
629 630 631	Jiménez, J., Doerr, S., Martínez-Rosell, G., Rose, A. S., and de Fabritiis, G. (2017). DeepSite: protein- binding site predictor using 3D-convolutional neural networks. <i>Bioinformatics</i> , 33(19):3036–3042.
632 633	Joshi, C. K., Bodnar, C., Mathis, S. V., Cohen, T., and Lio, P. (2023). On the Expressive Power of Geometric Graph Neural Networks. In <i>International Conference on Machine Learning</i> .
634 635 636 637 638 639 640	Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek, A., Potapenko, A., Bridgland, A., Meyer, C., Kohl, S. A. A., Ballard, A. J., Cowie, A., Romera-Paredes, B., Nikolov, S., Jain, R., Adler, J., Back, T., Petersen, S., Reiman, D., Clancy, E., Zielinski, M., Steinegger, M., Pacholska, M., Berghammer, T., Bodenstein, S., Silver, D., Vinyals, O., Senior, A. W., Kavukcuoglu, K., Kohli, P., and Hassabis, D. (2021). Highly accurate protein structure prediction with AlphaFold. <i>Nature</i> , 596(7873):583–589.
641 642 643	Kaba, SO., Mondal, A. K., Zhang, Y., Bengio, Y., and Ravanbakhsh, S. (2023). Equivariance with learned canonicalization functions. In <i>International Conference on Machine Learning</i> , pages 15546–15566. PMLR.
644 645 646	Kandel, J., Tayara, H., and Chong, K. T. (2021). PUResNet: prediction of protein-ligand binding sites using deep residual neural network. <i>Journal of Cheminformatics</i> , 13(1):1–14.
	Kinf T. N. and Walling M. (2017) Sami Supervised Classification with Graph Convolutional

647 Kipf, T. N. and Welling, M. (2017). Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations*.

648 649 650	Kong, X., Huang, W., and Liu, Y. (2023). Conditional antibody design as 3d equivariant graph translation. In <i>The Eleventh International Conference on Learning Representations</i> .
651 652	Krivák, R. and Hoksza, D. (2018). P2Rank: machine learning based tool for rapid and accurate prediction of ligand binding sites from protein structure. <i>Journal of Cheminformatics</i> , 10(1):1–12.
653 654 655	Le Guilloux, V., Schmidtke, P., and Tuffery, P. (2009). Fpocket: An open source platform for ligand pocket detection. <i>BMC Bioinformatics</i> , 10:168.
656 657	Lecun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-Based Learning Applied to Document Recognition. <i>Proceedings of the IEEE</i> , 86(11):2278–2324.
659 660	Lengauer, T. and Rarey, M. (1996). Computational methods for biomolecular docking. <i>Current Opinion in Structural Biology</i> , 6(3):402–406.
661 662 663	Li, J., Cai, D., and He, X. (2017). Learning Graph-Level Representation for Drug Discovery. <i>arXiv</i> preprint arXiv:1709.03741.
664 665	Li, Q., Han, Z., and Wu, XM. (2018). Deeper Insights Into Graph Convolutional Networks for Semi-Supervised Learning. In AAAI Conference on Artificial Intelligence.
6667 668 669	Liang, J., Edelsbrunner, H., and Woodward, C. (1998). Anatomy of protein pockets and cavities: Measurement of binding site geometry and implications for ligand design. <i>Protein Science</i> , 7(9):1884–1897.
670 671 672	Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W., Smetanin, N., Verkuil, R., Kabeli, O., Shmueli, Y., et al. (2023). Evolutionary-scale prediction of atomic-level protein structure with a language model. <i>Science</i> , 379(6637):1123–1130.
673 674 675	Loshchilov, I. and Hutter, F. (2019). Decoupled Weight Decay Regularization. In International Conference on Learning Representations.
676 677 678	Lu, W., Wu, Q., Zhang, J., Rao, J., Li, C., and Zheng, S. (2022). TANKBind: Trigonometry- Aware Neural NetworKs for Drug-Protein Binding Structure Prediction. In <i>Advances in Neural</i> <i>Information Processing Systems</i> .
680 681 682	Méndez-Lucio, O., Ahmad, M., del Rio-Chanona, E. A., and Wegner, J. K. (2021). A geometric deep learning approach to predict binding conformations of bioactive molecules. <i>Nature Machine Intelligence</i> , 3(12):1033–1039.
683 684 685	Morris, C., Ritzert, M., Fey, M., Hamilton, W. L., Lenssen, J. E., Rattan, G., and Grohe, M. (2019). Weisfeiler and Leman Go Neural: Higher-Order Graph Neural Networks. In AAAI Conference on Artificial Intelligence.
687 688	Mylonas, S. K., Axenopoulos, A., and Daras, P. (2021). DeepSurf: a surface-based deep learning approach for the prediction of ligand binding sites on proteins. <i>Bioinformatics</i> , 37(12):1681–1690.
689 690 691 692	Pei, Q., Gao, K., Wu, L., Zhu, J., Xia, Y., Xie, S., Qin, T., He, K., Liu, TY., and Yan, R. (2023). FABind: Fast and Accurate Protein-Ligand Binding. In <i>Advances in Neural Information Processing Systems</i> .
693	Perron, O. (1907). Zur Theorie der Matrices. Mathematische Annalen, 64:248-263.
694 695 696	Pham, T., Tran, T., Dam, H., and Venkatesh, S. (2017). Graph Classification via Deep Learning with Virtual Nodes. <i>arXiv preprint arXiv:1708.04357</i> .
697 698 699	Puny, O., Atzmon, M., Smith, E. J., Misra, I., Grover, A., Ben-Hamu, H., and Lipman, Y. (2022). Frame averaging for invariant and equivariant network design. In <i>International Conference on Learning Representations</i> .
700	Puny, O., Ben-Hamu, H., and Lipman, Y. (2020). Global Attention Improves Graph Networks Generalization. <i>arXiv preprint arXiv:2006.07846</i> .

702 703 704 705 706 707	Ren, F., Ding, X., Zheng, M., Korzinkin, M., Cai, X., Zhu, W., Mantsyzov, A., Aliper, A., Aladinskiy, V., Cao, Z., Kong, S., Long, X., Man Liu, B. H., Liu, Y., Naumov, V., Shneyderman, A., Ozerov, I. V., Wang, J., Pun, F. W., Polykovskiy, D. A., Sun, C., Levitt, M., Aspuru-Guzik, A., and Zhavoronkov, A. (2023). AlphaFold accelerates artificial intelligence powered drug discovery: efficient discovery of a novel CDK20 small molecule inhibitor. <i>Chemical Science</i> , 14(6):1443–1452.
708 709 710	Rosenbluth, E., Tönshoff, J., Ritzert, M., Kisin, B., and Grohe, M. (2024). Distinguished in uniform: Self attention vs. virtual nodes.
711 712	Rusch, T. K., Bronstein, M. M., and Mishra, S. (2023). A Survey on Oversmoothing in Graph Neural Networks. <i>arXiv preprint arXiv:2303.10993</i> .
713 714 715	Sadybekov, A. V. and Katritch, V. (2023). Computational approaches streamlining drug discovery. <i>Nature</i> , 616(7958):673–685.
716 717	Sanner, M. F., Olson, A. J., and Spehner, J. C. (1996). Reduced surface: An efficient way to compute molecular surfaces. <i>Biopolymers</i> , 38(3):305–320.
718 719 720	Satorras, V. G., Hoogeboom, E., and Welling, M. (2021). E(n) Equivariant Graph Neural Networks. In <i>International Conference on Machine Learning</i> .
721 722 723	Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., and Monfardini, G. (2009). The Graph Neural Network Model. <i>IEEE Transactions on Neural Networks</i> , 20(1):61–80.
724 725	Schmidtke, P. and Barril, X. (2010). Understanding and Predicting Druggability. A High-Throughput Method for Detection of Drug Binding Sites. <i>Journal of Medicinal Chemistry</i> , 53(15):5858–5867.
726 727 728 729	Schütt, K., Kindermans, PJ., Sauceda Felix, H. E., Chmiela, S., Tkatchenko, A., and Müller, KR. (2017). SchNet: A continuous-filter convolutional neural network for modeling quantum interactions. In Advances in Neural Information Processing Systems.
730 731	Shamir, R. R., Duchin, Y., Kim, J., Sapiro, G., and Harel, N. (2019). Continuous Dice Coefficient: a Method for Evaluating Probabilistic Segmentations. <i>arXiv preprint arXiv:1906.11031</i> .
732 733 734 735	Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. <i>Journal of Machine Learning Research</i> , 15(1):1929–1958.
736 737 738	Stärk, H., Ganea, O., Pattanaik, L., Barzilay, R., and Jaakkola, T. (2022). EquiBind: Geometric Deep Learning for Drug Binding Structure Prediction. In <i>International Conference on Machine Learning</i> .
739 740 741	Stepniewska-Dziubinska, M. M., Zielenkiewicz, P., and Siedlecki, P. (2020). Improving detection of protein-ligand binding sites with 3D segmentation. <i>Scientific Reports</i> , 10(1):5035.
742 743 744	Sverrisson, F., Feydy, J., Correia, B. E., and Bronstein, M. M. (2021). Fast End-to-End Learning on Protein Surfaces. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> .
745 746	Swinbank, R. and James Purser, R. (2006). Fibonacci grids: A novel approach to global modelling. <i>Quarterly Journal of the Royal Meteorological Society</i> , 132(619):1769–1793.
747 748 749 750	Topping, J., Giovanni, F. D., Chamberlain, B. P., Dong, X., and Bronstein, M. M. (2022). Under- standing over-squashing and bottlenecks on graphs via curvature. In <i>International Conference on</i> <i>Learning Representations</i> .
751 752 753 754 755	 Tunyasuvunakool, K., Adler, J., Wu, Z., Green, T., Zielinski, M., Žídek, A., Bridgland, A., Cowie, A., Meyer, C., Laydon, A., Velankar, S., Kleywegt, G. J., Bateman, A., Evans, R., Pritzel, A., Figurnov, M., Ronneberger, O., Bates, R., Kohl, S. A. A., Potapenko, A., Ballard, A. J., Romera-Paredes, B., Nikolov, S., Jain, R., Clancy, E., Reiman, D., Petersen, S., Senior, A. W., Kavukcuoglu, K., Birney, E., Kohli, P., Jumper, J., and Hassabis, D. (2021). Highly accurate protein structure prediction for the human proteome. <i>Nature</i>, 596(7873):590–596

756 757 758	Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., and Bengio, Y. (2018). Graph Attention Networks. In <i>International Conference on Learning Representations</i> .
759 760 761	Villar, S., Hogg, D. W., Storey-Fisher, K., Yao, W., and Blum-Smith, B. (2021). Scalars are universal: Equivariant machine learning, structured like classical physics. <i>Advances in Neural Information Processing Systems</i> .
762 763 764	Wang, R., Fang, X., Lu, Y., and Wang, S. (2004). The PDBbind Database: Collection of Binding Affinities for Protein Ligand Complexes with Known Three-Dimensional Structures. <i>Journal of</i> <i>Medicinal Chemistry</i> , 47(12):2977–80.
765 766 767	Weisfeiler, B. and Leman, A. (1968). The reduction of a graph to canonical form and the algebra which appears therein. <i>NTI, Series</i> , 2(9):12–16.
768 769 770 771	Widdowson, D. and Kurlin, V. (2023). Recognizing rigid patterns of unlabeled point clouds by complete and continuous isometry invariants with no false negatives and no false positives. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 1275–1284.
772 773 774	Xu, K., Hu, W., Leskovec, J., and Jegelka, S. (2019). How Powerful are Graph Neural Networks? In <i>International Conference on Learning Representations</i> .
775 776 777	Yan, X., Lu, Y., Li, Z., Wei, Q., Gao, X., Wang, S., Wu, S., and Cui, S. (2022). PointSite: A Point Cloud Segmentation Tool for Identification of Protein Ligand Binding Atoms. <i>Journal of Chemical</i> <i>Information and Modeling</i> , 62(11):2835–2845.
778 779 780	Zhang, Y., Cai, H., Shi, C., Zhong, B., and Tang, J. (2023a). E3Bind: An End-to-End Equivariant Network for Protein-Ligand Docking. In <i>International Conference on Learning Representations</i> .
781 782 783	Zhang, Y., Cen, J., Han, J., Zhang, Z., Zhou, J., and Huang, W. (2024). Improving equivariant graph neural networks on large geometric graphs via virtual nodes learning. In <i>Forty-first International Conference on Machine Learning</i> .
784 785 786 787 788 789 790 791	Zhang, Y., Huang, W., Wei, Z., Yuan, Y., and Ding, Z. (2023b). EquiPocket: an E(3)- Equivariant Geometric Graph Neural Network for Ligand Binding Site Prediction. <i>arXiv preprint</i> <i>arXiv:2302.12177</i> .
792 793 794	
795 796 797	
798 799 800	
802 803	
804 805 806	
807 808	

APPENDIX **Contents A** Implications and future work **B** Notation Overview C Problem Statement C.1 Representation of proteins. C.2 Representation of binding sites. C.3 Objective. Utilized Loss Functions. C.4 **D** Background on Group Theory and Equivariance E **Equivariance of VN-EGNN** F **Properties of virtual node initialization** Invariance with respect to the initial coordinates of the virtual nodes. E.1 A potential alternative strategy for initialization virtual node coordinates. F.2 **G** Experimental Settings **H** Additional Insights H.1 Initial experiment Virtual node initialization strategies H.2 H.3 Segmentation Loss Evaluation **Visualizations** Ι Domain shift of the HOLO4K dataset J **Expressiveness of VN-EGNN** Κ Additional hyperparameter evaluation L **M** Memory Utilization

864 A IMPLICATIONS AND FUTURE WORK

A.1 IMPLICATIONS

We expect that our empirical results and our new method re-new the interest in theoretically investigating the effect of virtual nodes on the expressivity and the oversquashing problem of GNNs. Practically,
we envision that VN-EGNN will be a useful tool in molecular biology and structure-based drug
design, which is regularly used to analyze proteins for potential binding pockets and their druggability.
On the long run, this could make the drug development process more time- and cost-efficient.

874 A.2 FUTURE WORK

We aim at using VN-EGNN to annotate all proteins in PDB with binding sites, and potentially a subset of the 200 million structures AlphaFold DB, with predicted binding pockets and release this annotated dataset.

918 B NOTATION OVERVIEW

Definition	Symbol/Notation	Туре
number of physical nodes	Ν	N
number of virtual nodes	K	\mathbb{N}_0
number of known binding pockets	M	\mathbb{N}_0
dimension of node features	D	\mathbb{N}
dimension of messages	E	\mathbb{N}
number of message passing layers/steps	L	\mathbb{N}
node indices	i,j,k,n	$\{1,, K\}$ or $\{1,, K\}$
binding pocket index	m	$\{1,, M\}$
layer/step index	l	$\{1,, L\}$
index set of 10 nearest neighbor atoms	$\mathcal{N}(i)$	$\{1,, N\}^{10}$
physical node coordinates	\mathbf{x}_{i}^{l}	\mathbb{R}^3
virtual node coordinates	\mathbf{z}_{i}^{l}	\mathbb{R}^3
physical node feature representation	$oldsymbol{h}_{i}^{'l}$	\mathbb{R}^{D}
virtual node feature representation	$oldsymbol{v}_{i}^{l}$	\mathbb{R}^{D}
edge feature between atoms	a_{ii}	\mathbb{R}
edge feature between atom and virtual node	d_{ij}^{j}	\mathbb{R}
ground-truth atom label	y_n	{0,1}
predicted atom label	\hat{y}_n	[0, 1]
ground-truth binding site center	\mathbf{y}_m	\mathbb{R}^3
prediction of binding site center	$\hat{\mathbf{y}}_k$	\mathbb{R}^3
messages*	$\boldsymbol{m}_{\cdots}^{(aa)} \; \boldsymbol{m}_{\cdots}^{(av)} \; \boldsymbol{m}_{\cdots}^{(va)}$	\mathbb{R}^{E}
neural networks for message passing*:	\dots_{ij} , \dots_{ij} , \dots_{ij}	10
message calculation	$\phi_{(aa)}, \phi_{(av)}, \phi_{(va)}$	$\mathbb{R}^D \times \mathbb{R}^D \times \mathbb{R} \times \mathbb{R}$
coordinate update	$\varphi_{e(uu)}, \varphi_{e(uv)}, \varphi_{e(vu)}$ $\phi_{e(uu)}, \phi_{e(vu)}, \phi_{e(vu)}$	$\mathbb{R}^E \to \mathbb{R}$
feature undate	$\phi_{1}(aa)$, $\phi_{1}(av)$, $\phi_{2}(va)$	$\mathbb{R}^D \times \mathbb{R}^E \to \mathbb{R}^D$
	$\varphi_{h^{(uv)}}, \varphi_{h^{(uv)}}, \varphi_{h^{(vu)}}$	
segmentation loss	$\mathcal{L}_{ ext{segm}}$	$\mathbb{K}'' \times \mathbb{K}'' \to \mathbb{K}$
binding site center loss	$\mathcal{L}_{ ext{bsc}}$	$\mathbb{K}^{3 \times M} \times \mathbb{K}^{3 \times K} -$

* The superscripts (aa), (av) and (va) represent the message passing direction (atom to atom, atom to virtual node, virtual node to atom).

972 C PROBLEM STATEMENT 973

974 C.1 REPRESENTATION OF PROTEINS. 975

The 3D structure of a protein is usually given by some measurement of its atoms that form the 976 primary amino acid sequence of the protein and the absolute coordinates for the atoms are given 977 as 3D points $\mathbf{x} \in \mathbb{R}^3$. The atoms themselves as well as the amino acids are characterized by their 978 physical, chemical and biological properties. We assume that these properties are summarized by 979 feature vectors $h \in \mathbb{R}^D$, which are located at the atom centers (either of all the atoms or only 980 the ones forming the protein backbone). We formally represent proteins by a neighborhood graph 981 $\mathcal{P} = (\mathcal{P}_N, \mathcal{P}_E)$ with N atom-property pairs, i.e. $\mathcal{P}_N = \{(\mathbf{x}_n, \mathbf{h}_n)\}_{n=1}^N$ with $\mathbf{x}_n \in \mathbb{R}^3$ and $\mathbf{h}_n \in \mathbb{R}^D$ 982 and a set of directed edges \mathcal{P}_E which consist of atom-property pairs (i, j). Each node i has incoming 983 edges from the 10 nearest nodes i that are closer than 10Å according to the Euclidean distance 984 $\|\mathbf{x}_i - \mathbf{x}_j\|.$ 985

C.2 REPRESENTATION OF BINDING SITES.

Binding sites are regions around or within proteins, to which ligands can potentially bind. Basically,
one can either describe binding sites *explicitly* or *implicitly*. In their explicit representation binding
sites would be directly described by the location specifics of the regions, where ligands are located,
especially by a region center point. In their implicit representation, binding sites would be described
by the atoms of the protein, which surround the ligand. Atoms close to the ligand would be marked as
binding site atoms. It might be worth mentioning, that several binding sites per protein are possible.

Formally, for the explicit representation, we describe the (experimentally observed) binding site center points of M distinct binding sites by $\mathbf{y}_m \in \mathbb{R}^3$ with $1 \le m \le M$. For the implicit representation, we assign to each protein atom n a label $y_n \in \{0, 1\}$, which is set to 1 if the atom center is within the threshold distance of observed binding ligands, and 0 otherwise.

C.3 OBJECTIVE.

From an abstract point of view, we want to learn a predictive machine learning model \mathcal{F} , parameterized by ω , which maps proteins characterized by the positions of their atoms together with their properties to a binary prediction per atom, whether it might form a binding site and to K 3D coordinates representing binding site region center points:

 $\mathcal{F}_{\omega} : \underset{n=1}{\overset{N}{\underset{\text{protein 3D atom}\\ \text{coords with}}{\sum}}} \mapsto \underbrace{[0,1]^{N}}_{\text{sem. segm.}} \times \underset{\text{protein atoms}}{\overset{K=1}{\underset{\text{prot. pred.}\\ \text{virt. nodes}}{\sum}} = \mathcal{F}_{\omega}^{\text{segm}}$

998

999

986

987

1006

1009

1010

1012

1013

1015

1024 1025 $\mathcal{F}_{\omega}\left((\mathbf{x}_{1}, \boldsymbol{h}_{1}), \dots, (\mathbf{x}_{N}, \boldsymbol{h}_{N})\right) = \left((\hat{y}_{1}, \dots, \hat{y}_{N}), (\hat{\mathbf{y}}_{1}, \dots, \hat{\mathbf{y}}_{K})\right)$ $\mathcal{F}_{\omega}^{\text{segm}}\left((\mathbf{x}_{1}, \boldsymbol{h}_{1}), \dots, (\mathbf{x}_{N}, \boldsymbol{h}_{N})\right) \coloneqq \text{proj}_{1}\mathcal{F}_{\omega}\left((\mathbf{x}_{1}, \boldsymbol{h}_{1}), \dots, (\mathbf{x}_{N}, \boldsymbol{h}_{N})\right)$

(C.1)

 $\mathcal{F}^{\mathrm{bsc}}_{\omega}\left((\mathbf{x}_1, \boldsymbol{h}_1), \ldots, (\mathbf{x}_N, \boldsymbol{h}_N)\right) \coloneqq \mathrm{proj}_2 \mathcal{F}_{\omega}\left((\mathbf{x}_1, \boldsymbol{h}_1), \ldots, (\mathbf{x}_N, \boldsymbol{h}_N)\right),$

where proj_i is a projection, that gives the *i*-th component (i.e., prediction of the semantic segmenation part or coordinate predictions or virtual nodes). Note, that for our predictive model, we use a fixed number K of binding point centers, while indeed M might have been observed for a specific protein.

1020 C.4 UTILIZED LOSS FUNCTIONS.

Segmentation loss. For semantic segmentation (i.e., the prediction of $\mathcal{F}_{\omega}^{\text{segm}}$), we use a Dice loss, that is based on the continuous Dice coefficient (Shamir et al., 2019), with $\epsilon = 1$:

$$\mathcal{L}_{\text{dice}} = \text{Dice}\left((y_1, \dots, y_N), (\hat{y}_1, \dots, \hat{y}_N)\right) \coloneqq 1 - \frac{2\sum_{n=1}^N y_n \, \hat{y}_n + \epsilon}{\sum_{n=1}^N y_n + \sum_{n=1}^N \hat{y}_n + \epsilon} \tag{C.2}$$

1026 Perfect predictions lead to a Dice loss of 0, while perfectly wrong predictions would lead to a Dice of 1027 1 (in case $\epsilon = 0$ and the denominator is > 0). 1028

Binding site center loss. For prediction of the binding site region center points (i.e., the prediction 1029 of $\mathcal{F}_{\omega}^{\text{bsc}}$), we use the (squared) Euclidean distance between the set of predicted points and the set of 1030 observed ones. More specifically, we assume to be given M observed center points $\{\mathbf{y}_1, \ldots, \mathbf{y}_M\}$. 1031 Each of the binding site center points should be detected by at least one of the K outputs from $\mathcal{F}_{\omega}^{bsc}$, 1032 which translates to using the minimum squared distance to any predicted center point for any of the 1033 observed center points: 1034

 $\mathcal{L}_{\text{bsc}} = \text{Dist}\left(\{\mathbf{y}_1, \dots, \mathbf{y}_M\}, \{\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_K\}\right) \coloneqq \frac{1}{M} \sum_{m=1}^M \min_{k \in 1, \dots, K} \|\mathbf{y}_m - \hat{\mathbf{y}}_k\|^2.$ (C.3)

Our optimization objective is then the sum of the Dice and the Dist loss:

$$\alpha \operatorname{Dice}\left((y_1,\ldots,y_N),(\hat{y}_1,\ldots,\hat{y}_N)\right) + \operatorname{Dist}\left(\{\mathbf{y}_1,\ldots,\mathbf{y}_M\},\{\hat{\mathbf{y}}_1,\ldots,\hat{\mathbf{y}}_K\}\right)$$
(C.4)

with the hyperparameter $\alpha = 1$.

¹⁰⁸⁰ D BACKGROUND ON GROUP THEORY AND EQUIVARIANCE

A group in the mathematical sense is a set G along with a binary operation $\circ : G \times G \rightarrow G$ with the following properties:

- Associativity: The group operation is associative, i.e. $(g \circ h) \circ k = g \circ (h \circ k)$ for all $g, h, k \in G$.
- *Identity:* There exists a unique identity element $e \in G$, such that $e \circ g = g \circ e = g$ for all $g \in G$.
 - *Inverse:* For each $g \in G$ there is a unique inverse element $g^{-1} \in G$, such that $g \circ g^{-1} = g^{-1} \circ g = e$.
- *Closure:* For each $g, h \in G$ their combination $g \circ h$ is also an element of G.

1094 A group action of group G on a set X is defined as a set of mappings $T_g: X \to X$ which associate each element $g \in G$ with a transformation on X, whereby the identity element $e \in G$ leaves X unchanged $(T_e(x) = x \quad \forall x \in X)$.

An example is the group of translations \mathbb{T} on \mathbb{R}^n with group action $T_t(x) = \mathbf{x} + \mathbf{t} \quad \forall \mathbf{x}, \mathbf{t} \in \mathbb{R}^n$, which shifts points in \mathbb{R}^n by a vector \mathbf{t} .

1099 A function $f: X \to Y$ is equivariant to group G with group action T_g if there exists an equivalent 1100 group action $S_g: Y \to Y$ on G such that

$$f(T_g(x)) = S_g(f(x)) \quad \forall x \in X, g \in G.$$

For example, a function $f : \mathbb{R}^n \to \mathbb{R}^n$ is translation-equivariant if a translation of an input vector $\mathbf{x} \in \mathbb{R}^n$ by $\mathbf{t} \in \mathbb{R}^n$ leads to the same transformation of the output vector $f(x) \in \mathbb{R}^n$, i.e. $f(\mathbf{x} + \mathbf{t}) = f(\mathbf{x}) + \mathbf{t}$.

Equivariant graph neural networks (EGNNs) ψ as defined by Satorras et al. (2021) exhibit three types of equivariances:

- 1. *Translation equivariance:* EGNNs are equivariant to column-wise addition of a vector $\mathbf{t} \in \mathbb{R}^n$ to all points in a point cloud $\mathbf{X} \in \mathbb{R}^{n \times N}$: $\psi(\mathbf{X} + \mathbf{t}) = \psi(\mathbf{X}) + \mathbf{t}$.
- 2. *Rotation and reflection equivariance*: Rotation or reflection of all points in the point cloud by multiplication with an orthogonal matrix $\mathbf{R} \in \mathbb{R}^{n \times n}$ leads to an equivalent rotation of the output coordinates: $\psi(\mathbf{RX}) = \mathbf{R}\psi(\mathbf{X})$.

The group spanning all translations, rotations and reflections in \mathbb{R}^n is called Euclidean group, denoted E(n), as it preserves Euclidean distances. A proof for E(n)-equivariance of VN-EGNN can be found in App. E.

3. *Permutation equivariance:* The numbering of elements in a point cloud or graph nodes does not influence the output, i.e. multiplication with a permutation matrix $P \in \mathbb{R}^{N \times N}$ leads to the same permutation of output nodes: $\psi(\mathbf{X}P) = \psi(\mathbf{X})P$. This property holds for message passing graph neural networks in general, as they aggregate and update node information based on local neighborhood structure, regardless of the order in which nodes are presented.

¹¹³⁴ E EQUIVARIANCE OF VN-EGNN

1135

1146 1147

1154

1164

1136 In this section we show that the equivariance property of EGNN (Satorras et al., 2021) extends to 1137 VN-EGNN, i.e., that rotation and reflection by an orthogonal matrix $\mathbf{R} \in \mathbb{R}^{3\times3}$, and translation 1138 by a vector $\mathbf{t} \in \mathbb{R}^3$ of atom and virtual node coordinates leads to an equivalent transformation of 1139 output coordinates while leaving node features invariant when applying the message passing steps of 1140 VN-EGNN.

1141 **Proposition 1.** (more formal) E(3) equivariant graph neural networks with virtual nodes (VN-EGNN) as defined by the message passing scheme $(\mathbf{X}^{l+1}, \mathbf{H}^{l+1}, \mathbf{Z}^{l+1}, \mathbf{V}^{l+1}) =$ 1143 VN-EGNN $(\mathbf{X}^l, \mathbf{H}^l, \mathbf{Z}^l, \mathbf{V}^l, \mathbf{A})$ in Eqs. (7) to (18) are equivariant with respect to reflections and 1144 roto-translations of the input and virtual node coordinates, i.e., the following holds (equivariance to 1145 reflections and roto-translations):

$$\left(\mathbf{R}\mathbf{X}^{l+1} + \mathbf{t}, \mathbf{H}^{l+1}, \mathbf{R}\mathbf{Z}^{l+1} + \mathbf{t}, \mathbf{V}^{l+1}\right) = \text{VN-EGNN}\left(\mathbf{R}\mathbf{X}^{l} + \mathbf{t}, \mathbf{H}^{l}, \mathbf{R}\mathbf{Z}^{l} + \mathbf{t}, \mathbf{V}^{l}\right), \quad (E.1)$$

where the addition $\mathbf{X}^{l} + \mathbf{t}$ is defined as column-wise addition of the vector \mathbf{t} to the matrix \mathbf{X} .

Proof. We use the notation from Section 2.1 and proceed by tracking the propagation of node roto-translations through the VN-EGNN network. First, we want to show invariance in Eq. (7) in phase I of message passing, equivalently to Satorras et al. (2021), i.e.:

$$\boldsymbol{m}_{ij}^{(aa)} = \boldsymbol{\phi}_{e^{(aa)}}(\boldsymbol{h}_{i}^{l}, \boldsymbol{h}_{j}^{l}, \|\mathbf{R}\mathbf{x}_{i}^{l} + \mathbf{t} - [\mathbf{R}\mathbf{x}_{j}^{l} + \mathbf{t}]\|^{2}) = \boldsymbol{\phi}_{e^{(aa)}}(\boldsymbol{h}_{i}^{l}, \boldsymbol{h}_{j}^{l}, \|\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}\|^{2})$$
(E.2)

1155 1156 Assuming the initial node features h_i^0 do not encode information about the original coordinates \mathbf{x}_i^0 , it 1157 remains to be shown that the Euclidean distance between two nodes is also invariant to translation 1158 and rotation:

$$\|\mathbf{R}\mathbf{x}_{i}^{l} + \mathbf{t} - [\mathbf{R}\mathbf{x}_{j}^{l} + \mathbf{t}]\|^{2} = \|\mathbf{R}\mathbf{x}_{i}^{l} - \mathbf{R}\mathbf{x}_{j}^{l}\|^{2}$$
$$= (\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l})^{\top}\mathbf{R}^{\top}\mathbf{R}(\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l})$$
$$= (\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l})^{\top}\mathbf{I}(\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l})$$
$$= \|\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}\|^{2}$$
(E.3)

1165 Consequently, the sum over messages (Eq. (8)) and the feature update function (Eq. (10)), which only 1166 uses the summed messages and previous node features as input, are invariant as well, leaving the 1167 intermediate output feature representations $h_i^{l+1/2}$ independent of coordinate transformations.

For the remaining equation (Eq. (9)) of phase I the equivariance property can be shown as follows, where Eq. (E.3) is used in the first equality:

$$\begin{split} \mathbf{R}\mathbf{x}_{i}^{l} + \mathbf{t} + \frac{1}{\mathcal{N}(i)} \sum_{j \in \mathcal{N}(i)} \frac{\mathbf{R}\mathbf{x}_{i}^{l} + \mathbf{t} - [\mathbf{R}\mathbf{x}_{j}^{l} + \mathbf{t}]]}{\|\mathbf{R}\mathbf{x}_{i}^{l} + \mathbf{t} - [\mathbf{R}\mathbf{x}_{j}^{l} + \mathbf{t}]\|} \phi_{\mathbf{x}^{aa}}(\boldsymbol{m}_{ij}^{(aa)}) \\ &= \mathbf{R}\mathbf{x}_{i}^{l} + \mathbf{t} + \frac{1}{\mathcal{N}(i)} \sum_{j \in \mathcal{N}(i)} \frac{\mathbf{R}\mathbf{x}_{i}^{l} + \mathbf{t} - [\mathbf{R}\mathbf{x}_{j}^{l} + \mathbf{t}]}{\|\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}\|} \phi_{\mathbf{x}^{aa}}(\boldsymbol{m}_{ij}^{(aa)}) \\ &= \mathbf{R}\mathbf{x}_{i}^{l} + \mathbf{t} + \frac{1}{|\mathcal{N}(i)|} \mathbf{R} \sum_{j \in \mathcal{N}(i), j \neq i} \frac{\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}}{\|\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}\|} \phi_{\mathbf{x}^{aa}}(\boldsymbol{m}_{ij}^{(aa)}) \\ &= \mathbf{R} \left(\mathbf{x}_{i}^{l} + \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i), j \neq i} \frac{\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}}{\|\mathbf{x}_{i}^{l} - \mathbf{x}_{j}^{l}\|} \phi_{\mathbf{x}^{aa}}(\boldsymbol{m}_{ij}^{(aa)}) \right) + \mathbf{t} \\ &= \mathbf{R} \mathbf{x}_{i}^{l+1/2} + \mathbf{t} \end{split}$$

1184 1185

In phase II of message passing, we input the updated physical node coordinates $\mathbf{Rx}_{j}^{l+1/2} + \mathbf{t}$ from Eq. (E.4) together with virtual node coordinates $\mathbf{Rz}_{i}^{l} + \mathbf{t}$, both subjected to identical rotation and

translation. Invariance of Eqs. (11), (12) and (14) can be deduced similarly to above, using the invariance properties of node features v_i^l and $h_j^{l+1/2}$, and Euclidean distance (Eq. (E.3)):

$$\begin{split} \boldsymbol{m}_{ij}^{(av)} &= \phi_{e^{(av)}}(\boldsymbol{v}_i^l, \boldsymbol{h}_j^{l+1/2}, \|\mathbf{R}\mathbf{z}_i^l + \mathbf{t} - [\mathbf{R}\mathbf{x}_j^{l+1/2} + \mathbf{t}]\|^2) \\ &= \phi_{e^{(av)}}(\boldsymbol{v}_i^l, \boldsymbol{h}_j^{l+1/2}, \|\mathbf{z}_i^l - \mathbf{x}_j^{l+1/2}\|^2) \end{split}$$

Thus, the output virtual node features v_i^{l+1} are invariant to roto-translations of node coordinates. Note that reflections are also covered by Eq. (E.4) since the distance of two points does not change under reflection.

Equivariance of output virtual node coordinates \mathbf{z}_i^{l+1} follows analogously to Eq. (E.4):

$$\mathbf{R}\mathbf{z}_{i}^{l} + \mathbf{t} + \frac{1}{N} \sum_{j=1}^{N} \frac{\mathbf{R}\mathbf{z}_{i}^{l} + \mathbf{t} - [\mathbf{R}\mathbf{x}_{j}^{l+1/2} + \mathbf{t}]}{\|\mathbf{R}\mathbf{z}_{i}^{l} + \mathbf{t} - [\mathbf{R}\mathbf{x}_{j}^{l+1/2} + \mathbf{t}]\|} \phi_{\mathbf{x}^{av}}(\boldsymbol{m}_{ij}^{(av)}) = \mathbf{R}\mathbf{z}_{i}^{l+1} + \mathbf{t}$$
(E.6)

The same derivations of message invariance

and coordinate equivariance

$$\mathbf{R}\mathbf{x}_{i}^{l+1/2} + \mathbf{t} + \frac{1}{K} \sum_{j=1}^{K} \frac{\mathbf{R}\mathbf{x}_{i}^{l+1/2} + \mathbf{t} - [\mathbf{R}\mathbf{z}_{j}^{l+1} + \mathbf{t}]}{\|\mathbf{R}\mathbf{x}_{i}^{l+1/2} + \mathbf{t} - [\mathbf{R}\mathbf{z}_{j}^{l+1} + \mathbf{t}]\|} \phi_{\mathbf{x}^{va}}(\boldsymbol{m}_{ij}^{(va)}) = \mathbf{R}\mathbf{x}_{i}^{l+1} + \mathbf{t}$$
(E.8)

 $\boldsymbol{m}_{ij}^{(va)} = \boldsymbol{\phi}_{e^{(va)}}(\boldsymbol{h}_{i}^{l+1/2}, \boldsymbol{v}_{j}^{l+1}, \|\mathbf{R} + \mathbf{x}_{i}^{l+1/2}\mathbf{t} - [\mathbf{R}\mathbf{z}_{j}^{l+1} + \mathbf{t}]\|^{2})$

 $= \boldsymbol{\phi}_{e^{(va)}}(\boldsymbol{h}_{i}^{l+1/2}, \boldsymbol{v}_{j}^{l+1}, \|\mathbf{x}_{i}^{l+1/2} - \mathbf{z}_{j}^{l+1}\|^{2})$

can be applied to phase III (Eqs. (15) to (18)), proving that invariance of feature representations h_i^{l+1} and equivariance of coordinates \mathbf{x}_i^{l+1} holds true for physical nodes as well, thus, proving Proposition 1.

_		

(E.7)

(E.5)

¹²⁴² F PROPERTIES OF VIRTUAL NODE INITIALIZATION

1244 F.1 INVARIANCE WITH RESPECT TO THE INITIAL COORDINATES OF THE VIRTUAL NODES.

1246 Note that Proposition 1 aims for equivariance with respect to rotations of the physical protein nodes \mathbf{X} and arbitrary, but fixed initialized virtual nodes \mathbf{Z} . We further want to have predictions, which 1247 are approximately invariant to differently chosen initial virtual node coordinates \mathbf{Z}^0 . This ultimately 1248 leads to predictions that are approximately equivariant with respect to E(3)-transformations of the 1249 physical protein nodes. In practice, we distribute the initial virtual node coordinates evenly on 1250 a sphere according to an algorithm, which constructs a spherical Fibonacci grid (Swinbank and 1251 James Purser, 2006). The algorithm provides spherically distributed grid points, which are fixed at 1252 certain locations in the 3D space. In order to achieve invariance with respect to differently chosen 1253 initial virtual node coordinates, we randomly rotate this grid of initial virtual node coordinates for 1254 each sample in every epoch, i.e. there is variation in the relative alignment of the Fibonacci grid 1255 points, that represent the virtual node positions, to physical protein nodes. Empirically, we observe 1256 that this training strategy leads to approximate invariance to different initializations of the virtual 1257 node coordinates (see Table F1).

F.2 A POTENTIAL ALTERNATIVE STRATEGY FOR INITIALIZATION VIRTUAL NODE COORDINATES.

An idea to avoid random alignments between physical node coordinates and initial virtual node 1262 coordinates, would be, to change initial coordinates in an equivariant way with respect to E(3) group 1263 transformations of the protein physical nodes. This could be achieved, e.g., by defining frames (Puny 1264 et al., 2022) based on Principal Component Analysis of physical protein node coordinates and by 1265 aligning the Fibonacci grid relative to these frames. Consequently, we would achieve that binding 1266 pocket predictions would change equivariantly with E(3)- transformations of the protein. Thereby 1267 the definition of such frames via Principal Component Analysis (PCA) is possible up to certain 1268 degenerate cases, that occur with probability zero for proteins. Since the orientation of axes might 1269 still not be unique, a strategy might be to compute properties such as the overall molecular weight 1270 for each octant in the coordinate system spanned by PCA eigenvectors. The orientation can then be set, such that for the octant with the maximum overall molecular weight, all coordinates get positive 1271 values. 1272

1273

1277 1278

1279

1280

1291

1293 1294 1295

1258

Table F1: Mean Performance at binding site identification in terms of DCC and DCA success rates together with their standard deviations (in parentheses). Means and standard deviations are across different random rotations of the Fibonacci grid.

Dataset	DCC	DCA
COACH420 HOLO4K PDBbind2020	$\begin{array}{c} 0.612(0.005) \\ 0.524(0.002) \\ 0.702(0.001) \end{array}$	0.741(0.006) 0.632(0.002) 0.833(0.002)
PDB0ind2020	0.702(0.001)	0.833(0.002)

¹²⁹⁶ G EXPERIMENTAL SETTINGS

1298 G.1 DATASETS

scPDB (Desaphy et al., 2015) is a frequently utilized dataset for binding site prediction (Kandel et al., 2021; Stepniewska-Dziubinska et al., 2020), encompassing both protein and ligand structures.
We employed the 2017 release of scPDB in the training and validation. This release comprises 17,594 structures, 16,034 entries, 4,782 proteins, and 6,326 ligands. Structures were clustered based on their Uniprot IDs. From each cluster, protein structures with the longest sequences were selected, in alignment to the strategies used in Kandel et al. (2021) and Zhang et al. (2023b).

1306 (Source: https://github.com/jivankandel/PUResNet/blob/main/scpdb_ 1307 subset.zip)

PDBbind (Wang et al., 2004) is a widely recognized dataset integral to the study of protein-ligand interactions. This dataset provides detailed 3D structural information of proteins, ligands, and their respective binding sites, complemented by rigorously determined binding affinity values derived from laboratory evaluations. For our work, we draw upon the v2020 edition, which is divided into two sets: the general set (comprising 14,127 complexes) and the refined set (containing 5,316 complexes).
While the general set encompasses all protein-ligand interactions, only the refined set, curated for its superior quality from the general collection, is used in our experiments.

1317 COACH420 and HOLO4K are benchmark datasets utilized for the prediction of binding sites, as originally detailed by Krivák and Hoksza (2018). Following the methodologies of Krivák and Hoksza (2018); Mylonas et al. (2021); Aggarwal et al. (2022a), we adopt the so-called mlig subsets from each of these datasets, which encompass the significant ligands pertinent to binding site prediction. Note that the HOLO4K contains many multi-chain systems and complexes with multiple copies of the protein (see Section App. J), such that this dataset's distribution is strongly differs from the other datasets.

1324
(Source: https://github.com/rdk/p2rank-datasets)

For comprehensive data preparation across all datasets, solvent atoms were excluded and erroneousstructures were removed.

1328 1329

1330

1331 1332

G.2 HYPERPARAMETERS AND HYPERPARAMETER SELECTION

Table G1 shows the evaluated hyperparameters. Bold indicates the parameters used in final model.

1333	Table G1: A list of considered and s	Table G1: A list of considered and selected hyperparameters.			
1334	hyperparmater	considered and selected values			
1335		considered and selected values			
1336	optimizer	{AdamW, Adam }			
1337	learning rate	$\{0.001, 0.0001\}$			
1338	activation function	{ SiLU , ReLU }			
1339	dimension of node features D	$\{20, 30, 100\}$			
1340	dimension of the messages P	$\{40, 50, 100\}$			
1341	number of message passing layers/steps L	$\{2, 3, 4, 5\}$			
12/12	number of virtual nodes K	$\{4, 8\}$			
1040	Huber loss δ	$\{1\}$			
1343					
1344					
1345					
1346					
1347					
1348					
1349					

^{1315 (}Source: http://www.pdbbind.org.cn/download/PDBbind_v2020_refined.tar. gz)

1350 Η **ADDITIONAL INSIGHTS** 1351

H.1 INITIAL EXPERIMENT

1352

1353

1354 Fig. H1 shows the training curves for a VN-EGNN during the development phase. VN-EGNN were 1355 only trained to minimize the segmentation loss \mathcal{L}_{segm} . Even in the absence of a the binding site 1356 center loss \mathcal{L}_{bsc} , the virtual nodes tend to converge towards the actual binding site center. This 1357 finding inspired us to further refine the position of the virtual nodes by including it directly to the 1358 optimization objective, which further improved the results significantly. 1359



1373 Figure H1: Validation curves of a VN-EGNN during development. Despite only being trained to 1374 minimize the segmentation loss, the virtual nodes converged towards the known binding sites. Left: DCC success rate during training. Middle: DCA success rate during training. Right: Segmentation 1376 loss during training. 1377

1379 H.2 VIRTUAL NODE INITIALIZATION STRATEGIES 1380

Fig. H2 shows learning curves for DCC and DCA on the validation set for a strict equivariant initialization strategy (i.e., center of mass (CoM) based initialization (Zhang et al., 2024; Kaba et al., 2023)) and for a relaxed version of this initialization strategy (i.e, Fibonacci grid based initialization, where the grid is rotated during training). For CoM based initialization the node embeddings are learned, whereas for the Fibonacci grid based initialization the virtual node features are initialized by the average of the physical node ESM embeddings of the protein. Our findings suggest that relaxing strict equivariance could be beneficial for binding site identification.

1389 0.80 1390 1391 0.75 1392 0.70 1393 DCA DCD 0.65 1394 0.4 1395 0.60 0.3 1396 0.55 0.2 0.50 1398 ò 50 100 150 200 250 300 350 400 0 50 100 150 200 250 300 Epoch Epoch 1399 1400



350 400

1375

1378

1404 H.3 SEGMENTATION LOSS EVALUATION

In Table H1, we evaluate predictions for the Dice loss, which is a part of our objective and for the
Intersection over Union (IoU, see Eq. (H.1)), which is related to Dice loss.

1408 1409 1410

1411 1412

$$\operatorname{IoU}\left((y_1, \dots, y_N), (\hat{y}_1, \dots, \hat{y}_N)\right) \coloneqq \frac{\sum_{n=1}^N y_n \, \mathbf{1}_{[\hat{y}_n > 0.5]}}{\sum_{n=1}^N y_n + \sum_{n=1}^N \mathbf{1}_{[\hat{y}_n > 0.5]} - \sum_{n=1}^N y_n \, \mathbf{1}_{[\hat{y}_n > 0.5]}} \quad (\mathbf{H.1})$$

	COACH420	HOLO4K	PDBbind2020
Dice loss	0.397(0.015)	0.584(0.031)	0.357(0.010)
IoU	0.437(0.005)	0.263(0.025)	0.477(0.003)

Table H1: Dice and Intersection over Union (IoU) Loss



Ι VISUALIZATIONS





classes are shown, remaining proteins are colored in grey.

¹⁵¹² J DOMAIN SHIFT OF THE HOLO4K DATASET

The HOLO4K benchmark comprises a large set of protein complexes and their annotated binding sites. HOLO4K has often been used as a benchmarking dataset (Krivák and Hoksza, 2018), while it exhibits different characteristics than other datasets, such as scPDB, COACH420 and PDBBind2020.
The number of chains per sample, i.e. PDB file, is larger than in these datasets (see Fig. J1), and also the number of binding sites per entry is higher (see Fig. J1). HOLO4K contains many symmetric units of large complexes which lead to these statistics. Thus, for machine learning methods trained on scPDB, the HOLO4K dataset represents a difficult test case due to the mentioned domain shifts.



Figure J1: Histogram of the number of protein chains (left) and of the number of binding sites frequency (right) per sample for each of the datasets used in this study. Note that the HOLO4K dataset has highly different characteristics compared to the other datasets and thus represents a strong domain shift for methods trained on scPDB.

¹⁵⁶⁶ K EXPRESSIVENESS OF VN-EGNN

1567 1568

The expressive power of GNNs is often described in terms of their ability to distinguish nonisomorphic graphs. The Weisfeiler-Leman (WL) (Weisfeiler and Leman, 1968) test, an iterative method to determine whether two attributed graphs are isomorphic, provides an upper bound to the expressiveness of GNNs. To extend the applicability of this framework to geometric graphs, Joshi et al. (2023) introduced the Geometric Weisfeiler-Leman test (GWL) which assesses whether two graphs are *geometrically isomorphic*.

Definitions (Joshi et al., 2023): Two graphs \mathcal{G}_1 and \mathcal{G}_2 with node features $h_i^{\mathcal{G}_j}$ and coordinates $\mathbf{x}_i^{\mathcal{G}_j}$ for $j \in \{1, 2\}$ are called *geometrically isomorphic* if there exists an edge-preserving bijection $b : \mathcal{V}(\mathcal{G}_1) \to \mathcal{V}(\mathcal{G}_2)$ between their corresponding node indices $\mathcal{V}(\mathcal{G}_j)$, such that their geometric features are equivalent up to E(n) group actions, i.e. global rotations/reflections \mathbf{R} and translations t:

1580

1581

1582 1583

 $\left(h_{b(i)}^{\mathcal{G}_2}, \mathbf{x}_{b(i)}^{\mathcal{G}_2}\right) = \left(h_i^{\mathcal{G}_1}, \mathbf{R}\mathbf{x}_i^{\mathcal{G}_1} + \mathbf{t}\right) \quad \forall i \in \mathcal{V}(\mathcal{G}_1).$ (K.1)

Two graphs \mathcal{G}_1 and \mathcal{G}_2 are called *k*-hop distinct if for all graph isomorphisms *b*, there is some node $i \in \mathcal{V}(\mathcal{G}_1), b(i) \in \mathcal{V}(\mathcal{G}_2)$ such that the corresponding *k*-hop neighborhood subgraphs $\mathcal{N}_i^{(\mathcal{G}_1,k)}$ and $\mathcal{N}_{b(i)}^{(\mathcal{G}_2,k)}$ are distinct. Otherwise, if $\mathcal{N}_i^{(\mathcal{G}_1,k)}$ and $\mathcal{N}_{b(i)}^{(\mathcal{G}_2,k)}$ are identical up to group actions for all $i \in \mathcal{V}(\mathcal{G}_1)$, we say \mathcal{G}_1 and \mathcal{G}_2 are *k*-hop identical.

In addition to iteratively updating node colors depending on node features in the local neighborhood analogously to the WL test, GWL keeps track of E(n)-equivariant hash values of each node's local geometry, i.e., distances to and angles between neighboring nodes. Thus, k iterations of GWL are necessary and sufficient to distinguish any k-hop distinct, (k - 1)-hop identical geometric graphs (Joshi et al., 2023).

Proposition 2. Any two geometrically distinct graphs G_1 and G_2 , where the underlying attributed graphs are isomorphic, can be distinguished with one iteration of GWL by adding one virtual node connected to all other nodes.

1597 1598

Proof. For 1-hop distinct graphs one iteration of GWL suffices to distinguish them even without virtual nodes and, thus, the proposition holds.

Now, we assume that \mathcal{G}_1 and \mathcal{G}_2 are k-hop distinct and (k-1)-hop identical for any $k > 1 \in \mathbb{N}$ and place one virtual node connected to all other nodes in an equivalent position in both graphs.

Note that finding equivalent virtual node positions is not trivial, as there is no straightforward way to spatially align the two graphs. Since the graphs are (k - 1)-hop identical there is at least one bijective mapping between (k-1)-hop sub-graphs of \mathcal{G}_1 and \mathcal{G}_2 , such that the neighborhood structure between matching sub-graphs is preserved.¹ For each such mapping between sub-graphs, we align the two graphs in space by overlaying one matching pair of (k - 1)-hop sub-graphs (consisting of more than two nodes that are not arranged in a straight line) and position the virtual node in the same coordinates in both aligned graphs.

1610 Since the virtual node is connected to each node in the graph, its 1-hop neighborhood and therefore 1611 the receptive field of the first GWL iteration contains the entire graph. Due to the k-hop distinctness 1612 of the graphs, there exists at least one node for which the geometric orientation relative to the matched 1613 subgraph deviates between \mathcal{G}_1 and \mathcal{G}_2 . Thus, the hash values corresponding to the virtual nodes' 1614 geometric information differ, and the graphs can be distinguished by only one iteration of GWL. 1615 Note that in case there are multiple possible mappings between (k-1)-hop subgraphs, \mathcal{G}_1 and \mathcal{G}_2 have to be distinguishable by one GWL iteration for each such mapping, in order to be classified as 1616 geometrically distinct. \square 1617

¹In practice, finding such a mapping is non-trivial and can be computationally expensive (Widdowson and Kurlin, 2023).

As k iterations of GWL act as an upper bound on the expressiveness of a k-layer geometric GNN, we propose that one layer of VN-EGNN is sufficient to distinguish two k-hop distinct graphs while without virtual nodes k EGNN layers are necessary to complete the same task.

We demonstrate this on the example of *n*-chain geometric graphs, where each pair of graphs comprises nodes arranged in a line and two end points with distinct orientations (Fig. K1). These graphs are $(\lfloor \frac{n}{2} \rfloor + 1)$ -hop distinct and should therefore be distinguishable by $(\lfloor \frac{n}{2} \rfloor + 1)$ EGNN layers or $(\lfloor \frac{n}{2} \rfloor + 1)$ iterations of GWL.



Figure K1: A pair of *n*-chain geometric graphs consisting of *n* nodes arranged in a line and two end points with opposite orientations. Without the addition of a virtual node, these graphs are $(\lfloor \frac{n}{2} \rfloor + 1)$ -hop distinct.

We trained EGNNs with an increasing number of layers to classify 4-chain graphs, both with and 1641 without the addition of a virtual node. We tested two different strategies for the virtual node's initial 1642 1643 position. Firstly, the virtual node was placed at the same relative position in both graphs, such that when the first (n + 1) nodes of the two *n*-chain graphs are overlaid, the virtual nodes coincide. 1644 The exact position was randomly selected on a sphere centered at the midpoint of the n-chain. In 1645 the second setting, the virtual node was initialized at the center of mass (Zhang et al., 2024; Kaba 1646 et al., 2023) of each graph, ensuring equivariance with respect to the initialization. Note that for the 1647 experiments including virtual nodes, we did not use the heterogeneous message passing scheme as 1648 described in Section 2.3, but apply the EGNN to the entire graph, including virtual nodes, at once. 1649

The results shown in Table K1 demonstrate that, as expected, 3 layers of EGNN are necessary to distinguish the 4-chain graphs while after adding a virtual node, one iteration is sufficient for correct classification, indicating the increased expressiveness of VN-EGNN. Although we used the setting of Joshi et al. (2023), we could not reproduce their finding that 6 EGNN layers are necessary to solve this task, which they explained with possible oversmoothing or oversquashing effects. The differences might arise from the use of different features dimensions, which is why we include results for 5 different feature dimensions.

1656 1657

1623

1628 1629 1630

1633

Table K1: Classification accuracy of EGNNs with and without virtual nodes and increasing node embedding dimensions on 4-chain geometric graphs. The standard deviation across 100 training re-runs is indicated with \pm and column "Dim." indicates the used node feature dimension. Note that VN-EGNN can distinguish these graphs already with one message passing layer (see columns "1 Layer" and "2 Layers").

Layer and 2 Layers).												
	Dim.	1 Layer	2 Layers	3 Layers	4 Layers	5 Layers	6 Layers	7 Layers	8 Layers			
EGNN	8	50.0 ± 0.0	50.0 ± 0.0	50.0 ± 0.0	98.0 ± 9.8	94.0 ± 16.2	93.0 ± 17.3	99.5 ± 5.0	99.5 ± 5.0			
	16	50.0 ± 0.0	50.0 ± 0.0	86.0 ± 22.4	97.5 ± 10.9	99.5 ± 5.0	99.5 ± 5.0	99.5 ± 5.0	100.0 ± 0.0			
	32	50.0 ± 0.0	50.0 ± 0.0	56.5 ± 16.8	50.0 ± 0.0	50.0 ± 0.0	96.5 ± 12.8	99.0 ± 7.0	93.5 ± 16.8			
	64	50.0 ± 0.0	50.0 ± 0.0	100.0 ± 0.0	99.0 ± 7.0	100.0 ± 0.0	99.0 ± 7.0	100.0 ± 0.0	100.0 ± 0.0			
	128	50.0 ± 0.0	50.0 ± 0.0	96.5 ± 12.8	98.5 ± 8.5	95.0 ± 15.0	99.5 ± 5.0	99.5 ± 5.0	99.5 ± 5.0			
VN-EGNN (initialization on a sphere)	8	65.5 ± 23.1	50.0 ± 0.0	84.5 ± 23.1	92.5 ± 17.9	64.0 ± 22.4	97.0 ± 11.9	86.5 ± 23.3	97.5 ± 10.9			
	16	86.0 ± 23.5	95.0 ± 15.0	98.5 ± 8.5	99.5 ± 5.0	99.5 ± 5.0	98.0 ± 9.8	99.5 ± 5.0	100.0 ± 0.0			
	32	95.0 ± 15.0	100.0 ± 0.0	99.5 ± 5.0	99.5 ± 5.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0			
	64	97.5 ± 10.9	100.0 ± 0.0	99.5 ± 5.0	99.5 ± 5.0	99.0 ± 7.0	100.0 ± 0.0	100.0 ± 0.0	99.5 ± 5.0			
	128	99.0 ± 7.0	99.5 ± 5.0	99.5 ± 5.0	99.0 ± 7.0	99.5 ± 5.0	99.5 ± 5.0	99.0 ± 7.0	99.0 ± 7.0			
	8	95.0 ± 15.0	89.0 ± 20.7	99.5 ± 5.0	97.5 ± 10.9	100.0 ± 0.0	99.5 ± 5.0	98.5 ± 8.5	100.0 ± 0.0			
VN-EGNN	16	98.0 ± 9.8	97.0 ± 11.9	99.5 ± 5.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0			
(with center of mass initialization)	32	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	99.5 ± 5.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0			
	64	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0			
	128	99.0 ± 7.0	99.5 ± 5.0	100.0 ± 0.0	99.5 ± 5.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0			

¹⁶⁷⁴ L ADDITIONAL HYPERPARAMETER EVALUATION

In Table L1 we benchmark VN-EGNN for different numbers of virtual nodes and different numbers of layers than in the main part of this paper. In many cases we still obtain decent results.

	COAC	CH420	HOL	.O4K	PDBbind2020		
	DCC	DCA	DCC	DCA	DCC	DCA	
10 virtual nodes	0.591(0.010)	0.736(0.010)	0.530(0.011)	0.649(0.010)	0.677(0.013)	0.813(0.012)	
12 virtual nodes	0.609(0.011)	0.738(0.020)	0.521(0.019)	0.642(0.017)	0.677(0.010)	0.825(0.010)	
4 layer	0.595(0.008)	0.740(0.012)	0.530(0.019)	0.646(0.023)	0.685(0.013)	0.827(0.008)	
6 layer	0.598(0.017)	0.731(0.001)	0.507(0.016)	0.614(0.029)	0.673(0.024)	0.816(0.013)	
Default	0.605(0.009)	0.750(0.008)	0.532(0.021)	0.659(0.026)	0.669(0.015)	0.820(0.010)	

 Table L1: DCC/DCA performance for varying numbers of virtual nodes and message passing layers.

 All other parameters were maintained as specified in the original paper.

1728 M MEMORY UTILIZATION

Since the number of virtual nodes are connected to all other nodes, we show how much memory utilization increases when increasing the number of virtual nodes in Fig. M1. Thereby we consider a range for the number of virtual nodes, which seems to be practically relevant for binding site identification (i.e., up to 16 virtual nodes).



Figure M1: Memory utilization for varying numbers of virtual nodes was assessed as a percentage of total capacity on an NVIDIA A100 GPU (80GB).