

# Graph Attention Site Prediction (GrASP): Identifying Druggable Binding Sites Using Graph Neural Networks with Attention

Zachary Smith,<sup>1</sup> Michael Strobel,<sup>1</sup> Bodhi P. Vani, and Pratyush Tiwary\*



Cite This: *J. Chem. Inf. Model.* 2024, 64, 2637–2644



Read Online

ACCESS |



Metrics & More

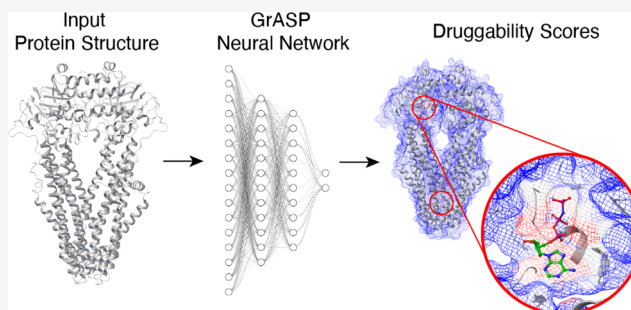


Article Recommendations



Supporting Information

**ABSTRACT:** Identifying and discovering druggable protein binding sites is an important early step in computer-aided drug discovery, but it remains a difficult task where most campaigns rely on *a priori* knowledge of binding sites from experiments. Here, we present a binding site prediction method called Graph Attention Site Prediction (GrASP) and re-evaluate assumptions in nearly every step in the site prediction workflow from data set preparation to model evaluation. GrASP is able to achieve state-of-the-art performance at recovering binding sites in PDB structures while maintaining a high degree of precision which will minimize wasted computation in downstream tasks such as docking and free energy perturbation.



## INTRODUCTION

A critical early step in computer-aided drug discovery is identifying druggable binding sites or those that can bind ligands likely to alter activity. Virtual screening of ligands with docking methods is often done for a specific binding site which requires *a priori* knowledge of where ligands are likely to bind.<sup>1–4</sup> Recently, modern structure prediction methods such as AlphaFold2<sup>5,6</sup> and RoseTTAFold<sup>7</sup> have greatly expanded the number of predicted structures for the human proteome,<sup>8</sup> while enhanced sampling methods for molecular dynamics have revealed conformations with cryptic pockets inaccessible in the protein's crystal structure.<sup>9–11</sup> The combination of advances in these two areas has led to a deluge of protein conformations that have not been probed for binding sites in experiments. For drug discovery to keep pace with structure discovery, accurate high-throughput binding site identification methods must be developed.

Initially, binding site prediction methods used human-designed representations of proteins based on geometry,<sup>12–18</sup> sequence conservation,<sup>19,20</sup> interactions with probe molecules,<sup>21,22</sup> or a combination of these features.<sup>2,23</sup> Recent methods, however, have leveraged machine learning combined with binding-site databases<sup>24,25</sup> to learn how to predict binding sites.<sup>26–33</sup> Despite the existence of large databases and modern machine learning architectures, one of the most popular and successful methods in this area is P2Rank, a random forest classifier trained on 251 protein structures.<sup>27</sup> It is striking that this model is able to outperform a Convolutional Neural Network (CNN) trained on thousands of structures.<sup>26</sup> The reason behind P2Rank's success might be the use of better representations such as an accessible surface area mesh with a

rotationally invariant model or the use of a smaller but more carefully curated data set.

One more recently developed class of machine learning architectures that employs a natural representation for molecules is Graph Neural Networks (GNNs)<sup>34,35</sup> which represent inputs as graphs and pass messages between connected nodes. GNNs have been shown to excel at closely related tasks such as binding affinity prediction,<sup>36,37</sup> docking,<sup>38</sup> predicting which sites will open midsimulation,<sup>39</sup> predicting the type of molecule that binds to a known site,<sup>40</sup> and even predicting protein–protein interactions.<sup>41</sup> Like P2Rank, GNNs also have rotational invariance, guaranteeing that the orientation of an input molecule does not affect the internal representation.

With this motivation, we developed a GNN-based method called Graph Attention Site Prediction (GrASP). GrASP is designed with the representational advantages of P2Rank in mind and performs a rotationally invariant featurization of solvent-accessible atoms. As a deeper model, GrASP requires a larger data set for training, and to achieve this goal, we have created a new publicly available version of the sc-PDB database containing 26,196 binding sites across 16,889 protein structures. GrASP is able to recover a higher number of ground truth binding sites when evaluated on P2Rank's test

**Special Issue:** Machine Learning in Bio-cheminformatics

**Received:** October 20, 2023

**Revised:** February 22, 2024

**Accepted:** February 23, 2024

**Published:** March 7, 2024



sets but has the important advantage that over 70% of its output binding sites correspond to real binding sites, whereas under 30% of P2Rank sites correspond to real sites.

## METHODS

In this section, we introduce Graph Neural Networks and show each step of the site prediction pipeline including data set creation, protein representation, and the model architecture.

**Graph Neural Networks (GNNs).** For the sake of better motivating the architecture underlying GrASP, we start with a brief pedagogical overview. Graph Neural Networks (GNNs) are a family of architectures that operate on a graph structure to represent the features of individual nodes and the relational structure between them. In this work, we represent proteins as graphs in which nodes represent heavy atoms, and edges are drawn between all pairs of atoms within 5 Å of each other. Node features include both atomic features such as formal charge and residue features such as residue name. Edges also have features of the inverse distance and bond order. A full list of features can be found in the [Supporting Information \(SI\)](#). GNNs represent nodes using message-passing layers, which perform the following three operations:

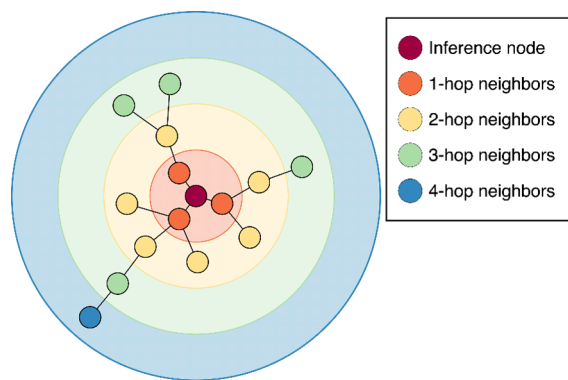
1. Message: Neighboring nodes send information to one another about their current state.
2. Aggregate: Each node collects the messages from its neighbors and aggregates them by applying an aggregation function.
3. Update: Each node incorporates the aggregated information with its own representation to generate a new latent representation of itself.

This process can be formalized as the following:<sup>42</sup>

$$x'_i = f_{\Theta}(x_i, \text{Aggregate}(\{x_j \mid j \in N(i)\})) \quad (1)$$

Here  $x_i$  is the current representation of node  $i$ ,  $x'_i$  is the updated representation of node  $i$ ,  $N(i)$  denotes the set of neighbors connected to node  $i$ , and  $f_{\Theta}$  denotes a parametrized update function.

This process can be repeated with multiple GNN layers for a node's representation to incorporate information from a larger region of the graph. Since each message includes information about a node's immediate neighbors, each GNN layer allows the node to access information influenced by nodes one hop further than the previous layer.<sup>43</sup> This can be seen in [Figure 1](#) where the inference node's hidden representation would



**Figure 1.**  $k$ -hop neighborhoods for a given inference node in the input graph. The  $k$ -th GNN layer representation is affected by neighbors up to  $k$  hops away.

include information about  $k$ -hop neighbors after passing through  $k$  GNN layers. These repeated GNN layers are commonly used within an encoder-processor-decoder framework implemented through multilayer perceptrons (MLPs) before and after a set of GNN layers.<sup>44</sup>

Repeated aggregation comes at the cost of oversmoothing, a phenomenon where deeper GNNs cause node representations to become increasingly similar.<sup>45</sup> A number of methods have been developed to encourage diverse latent representations and to allow for deeper GNN architectures. Three of these are used in this work: ResNet skip connections,<sup>46</sup> jumping knowledge skip connections,<sup>47</sup> and Noisy Nodes.<sup>45</sup> Both ResNet and jumping knowledge skip connections preserve information from earlier GNN layers (equivalently  $k$ -hop neighborhoods) by combining their latent representations with those of later layers. ResNet skip connections do so locally by adding the input and output of each GNN layer, while jumping knowledge skip connections feed the latent representations of multiple GNN layers into the decoder. In contrast, Noisy Nodes is a regularization procedure where noise is added to the input features, and an additional decoder head that attempts to reconstruct the denoised inputs is added after the processor layers, forcing the intermediate processor layer's latent representations to maintain enough diversity to reconstruct inputs.

**Graph Attention Networks (GATs).** Graph attention networks (GATs) are GNNs that use attention to learn weights for each neighbor and perform a weighted average aggregation.<sup>48</sup> A GAT layer is shown in [eq 2](#) where  $\Theta_s$  and  $\Theta_t$  are linear layers, and  $\alpha_{ij}$  represents the attention coefficient for messages from node  $j$  to node  $i$ .

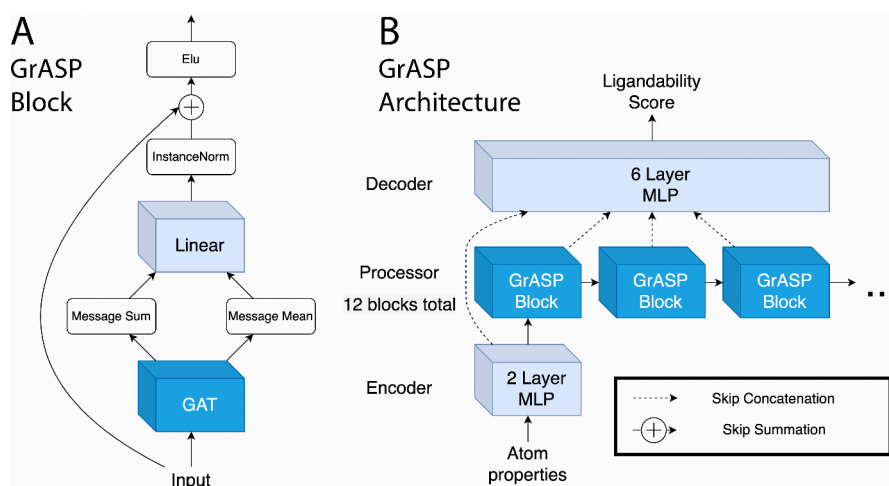
$$x'_i = \alpha_{i,i} \Theta_s x_i + \sum_{j \in N(i)} \alpha_{i,j} \Theta_t x_j \quad (2)$$

We use the attention function from GATv2 which calculates weights with the softmax of an MLP over a concatenation of both node and edge features.<sup>42</sup> This function is shown in [eq 3](#) where  $\parallel$  represents concatenation,  $e_{ij}$  are edge features, and the linear layers  $\Theta$  (composed of  $\Theta_s$ ,  $\Theta_v$ , and  $\Theta_e$  operating on  $x_i$ ,  $x_j$ , and  $e_{ij}$  respectively) and  $a^T$  form the MLP.

$$\alpha_{i,j} = \frac{\exp(a^T \text{LeakyReLU}(\Theta[x_i \parallel x_j \parallel e_{ij}]))}{\sum_{k \in N(i) \cup \{i\}} \exp(a^T \text{LeakyReLU}(\Theta[x_i \parallel x_k \parallel e_{i,k}]))} \quad (3)$$

**Graph Attention Site Prediction (GrASP).** GrASP is a GAT-based model for binding site prediction. GrASP first employs the GAT-based model to perform semantic segmentation on all protein surface atoms, scoring which atoms are likely part of a binding site. These atomic scores are then aggregated into binding sites using average linkage clustering<sup>49</sup> and ranked as a function of their constituent atoms' scores. This overall workflow performs an instance segmentation task (binding site prediction) by postprocessing the semantic segmentation predictions (atomic binding scores).

**Preprocessing.** The first issue we address is the definition of a binding site, for which there is no consensus definition in the literature. Definitions range from atoms within 2.5 Å<sup>50</sup> of the ligand to residues within 6.5 Å<sup>24</sup> and choose to include different combinations of empty space, surface atoms (or surface meshes), and buried atoms. This wide range of representations has two implications. The first implication is



**Figure 2.** Diagram of the GrASP model. A) The GrASP blocks used to represent each atom's local chemical environment. B) The full architecture combining GrASP blocks in an encoder-processor-decoder framework. Layers that do not consider neighbors are light blue, while layers that consider neighbors are blue.

that we cannot perform an unbiased comparison with metrics based on a specific definition because we would artificially skew success rates toward methods trained with a similar definition. For example, one metric we can not use is the volume overlap between the “true” and predicted binding sites. We focus on a metric that directly compares predictions to the ligand instead of a prescribed area around it: the distance from the predicted site center to any ligand-heavy atom. This is not the only metric that fits this criterion, but we choose to use it for fair comparison because P2Rank was also tuned using this metric. The second implication of not having a consensus binding site definition is that we can tune the definition used during training to maximize the model's performance on our chosen metrics. Since these metrics do not rely on the site definition, we can tune this hyperparameter without affecting the evaluation of other methods. To achieve this goal, we assign a continuous target score to each surface atom using a sigmoid function on the distance between the ligand and the protein atom. This representational choice, for which we provide details in the SI, makes it so that GrASP is penalized more for incorrectly characterizing atoms near ligands instead of treating all atoms within a cutoff distance as the same.

The second issue we address is the definition of the protein graph. We do this using the same inductive bias as for the binding site definition: only surface atoms can be considered binding sites. This means that we will only score surface atoms, but we wish to characterize the local chemical environment of these atoms using their neighbors. We construct a near-surface graph consisting of both surface atoms, defined using a solvent-accessible surface area, and buried atoms within 5 Å of surface atoms. In other words, we use the induced subgraph consisting of the surface atoms' one-hop neighborhood. More precise details about the implementation of this representation are available in the SI. This representational choice gives GrASP the inductive bias that only surface atoms are accessible and allows it to learn druggability without first learning which atoms a ligand can reach.

**Architecture.** It has been shown that there is no best aggregator for graphs with continuous features.<sup>51</sup> This led to the development of GNNs using multiple aggregators. This multiaggregation strategy is the inspiration for the GrASP block shown in Figure 2A. This block consists of a GAT layer

with four attention heads that pass both summed and averaged messages through a linear layer, an InstanceNorm,<sup>52</sup> a residual skip connection,<sup>46</sup> and an Elu activation.<sup>53</sup> The linear layer after the multiaggregation allows the model to decide how much weight to give the sum and mean for each feature.

These GrASP blocks are combined with an MLP encoder and MLP decoder to make the full GrASP model shown in Figure 2B. The output of each hybrid block is concatenated using jumping knowledge skip connections<sup>47</sup> as an input for the decoder. During training, GrASP also receives inputs with Gaussian noise added and uses a second Noisy Nodes<sup>45</sup> head to reconstruct denoised inputs. This denoising head operates on outputs from the last GrASP block and aims to reduce oversmoothing, as oversmoothed outputs cannot be used to reconstruct nodes with different features.

**Postprocessing.** The neural network architecture outlined so far scores the likelihood for any given heavy atom to be a part of a binding site, as shown in Figure 3. For applications to drug discovery and model evaluation, it is necessary to aggregate predicted binding site atoms into discrete binding sites. We accomplish this by using average linkage clustering<sup>49</sup> on all heavy atoms with a predicted binding likelihood above 0.3. The output clusters are then ranked using the same scoring function as P2Rank except replacing surface points with atoms,  $S_s = \sum S_a^2$ , where  $S_s$  is the score for a binding site, and  $S_a$  is the score for an individual atom.<sup>27</sup> We then obtain the center for each binding site by computing the convex hull of the atom cluster and calculating its center.<sup>54,55</sup>

**Relationship to P2Rank.** P2Rank is one of the most popular and successful methods for the binding site prediction. This method applies a random forest to score points on the protein's solvent-accessible surface and then aggregates these surface points into sites using single linkage clustering.<sup>27</sup> While P2Rank uses a different class of model and operates on surface points instead of atoms, P2Rank and GrASP share significant representational similarities. Each surface point in P2Rank describes its local chemical environment using a distance-weighted average of nearby atom properties (up to 6 Å away) with weights,  $w(d) = 1 - \frac{d}{6}$ .<sup>27</sup> This average can be written as a message passing layer shown in eq 4 describing a bipartite graph where surface points  $x_i$  receive messages from nearby



**Table 1. GrASP Validation Performance Averaged Across 10 Models Corresponding to Each Cross-Validation Fold in the Modified sc-PDB Set**

	sc-PDB Cross-validation				
	DCA Recall Top $N$ ( $\uparrow$ )	DCA Recall Top $N + 2$ ( $\uparrow$ )	DCA Precision Top 3 ( $\uparrow$ )	DCA Precision Top 5 ( $\uparrow$ )	DCA Precision All Sites ( $\uparrow$ )
GrASP	85.3	91.4	69.7	66.4	65.0

atoms  $x_i$  with distance-based weights shown in eq 5. Here, we see that P2Rank parametrizes the local chemical environment with a single pass through a hand-designed message-passing function. GrASP generalizes this featurization process by learning these aggregation weights through attention and applying multiple message-passing steps.

$$x'_i = \sum_{j \in N(i)} \alpha_{ij} x_j \quad (4)$$

$$\alpha_{ij} = \frac{w(d_{ij})}{\sum_{k \in N(i)} w(d_{ik})} \quad (5)$$

**Data Sets.** Our training and validation were performed using a modified version of the sc-PDB (v.2017) database.<sup>24</sup> sc-PDB is a curated database designed for small ligand docking which contains nonrepeating protein–ligand pairs. The crystal structures for these pairs are split into mol2 files, which contain the ligand, the binding site (all residues within 6.5 Å), the binding cavity (empty space around the ligand), the full protein, and other structures useful for docking. This database provides 17,594 binding sites and is commonly used to train binding site prediction models but has the shortcoming of unique protein–ligand pairs, which means that a large number of binding sites are not labeled. To address this shortcoming, we modified sc-PDB to contain binding sites corresponding to protein–ligand pairs that are already labeled once (for example, labeling sites on both chains in a symmetric dimer).

We first modify the sc-PDB database by combining entries with the same PDB ID and with protein mol2 files that can be aligned exactly. We then identify unlabeled buried ligands that have the chemical composition as ligands already labeled for any entry with the same PDB ID. We found almost 9,000 additional ligands that fit our criteria which led to a total of 26,196 binding sites across 16,889 protein structures in our final modified data set. This procedure converts the single-site entries of sc-PDB into multisite entries more suitable for binding site prediction methods. The resulting modified data set is available at <https://github.com/tiwarylab/GrASP>, and additional details on data set preparation are available in the SI.

We train and validate our model on the modified data set with the 10-fold cross-validation splits of sc-PDB from ref.<sup>31</sup> which are made to prevent data leakage with respect to UniProt IDs as well as binding site similarity.

We also modify the test sets used to evaluate P2Rank<sup>27</sup> to ensure that all ligands are both bound and biologically or pharmacologically relevant. The main preparation of the COACH420 and HOLO4K sets used (i) geometric criteria to ensure the ligand is interacting with the protein and (ii) simple name filters to avoid the inclusion of water, salt, or sugar as ligands. The P2Rank authors also propose an alternative preparation of these data sets referred to as Mlig sets which use the Binding MOAD database to check that ligands are either biologically or pharmacologically relevant but do not employ previous geometric criteria. We apply both sets of criteria to these sets to ensure both bound and relevant

ligands and title the new sets COACH420(Mlig+) and HOLO4K(Mlig+). We also found that HOLO4K contains many multimers with repetitions of the same binding mode. In a real-world setting, multimers would only be considered when they are known to occur *in vivo*, and their interface is suspected to be druggable. To reflect this setting, we consider each ligand bound to all proteins within 4 Å and connect all chains that share an interfacial ligand. We then split all systems into subsystems consisting of single chains without interfacial ligands and connected subsystems with interfacial ligands. This processing should more closely reflect the workflow used in practice, avoiding evaluation on homomultimers while preserving evaluation on interfacial binding. The consideration of chains and interfaces does not affect COACH420(Mlig+) as this set consists of only single chains.

## RESULTS

Here, we introduce a new metric to evaluate binding site prediction based on standard metrics in semantic segmentation and compare GrASP to P2Rank on updated versions of the original P2Rank data sets.

**Metrics.** A commonly used metric to evaluate binding site performance is the distance from the predicted site center to any ligand-heavy atom (DCA). A binding site prediction is considered successful if this distance is below 4 Å, and DCA is reported as the percentage of successful predictions over the total number of “ground truth” binding sites (or equivalently bound ligands), usually subject to the constraint that only the top  $N$  or top  $N + 2$  ranked predictions are considered for each system where  $N$  is the number of binding sites in the ground truth. This metric can be seen as a constrained analogy to recall a metric commonly used for classification problems defined as  $\frac{TP}{TP + FN}$ , where TP is the number of true positives, and FN is the number of false negatives. This ratio can be equivalently defined as the total number of correct predictions divided by the total number of members of the class being predicted. Because DCA refers to both the success criteria and the metric, we will distinguish these two by calling the criteria DCA and the metric DCA recall.

DCA recall evaluates the number of correct predictions among the top  $N$  binding sites, but in a discovery setting, the number of binding sites is not known *a priori*. This means that in a real setting any predictions beyond  $N$  can waste computational resources in downstream tasks even if ranked correctly, and likely a fixed maximum number of sites would be considered for each system to stay within a computational budget. To reflect this cost, we propose a constrained analog to the precision metric called DCA precision. DCA precision is the ratio of correctly predicted sites over the total number of predicted sites. This can be computed over all predictions or among the top  $M$  sites, where  $M$  is a constant that reflects a more realistic cap on the number of sites a user is willing to study per system. DCA precision and DCA recall can be used similarly to the standard precision and recall metrics from

**Table 2. Comparison between P2Rank and GrASP Performance on the COACH420(Mlig+) Test Set<sup>a</sup>**

	COACH420(Mlig+)				
	DCA Recall Top N (↑)	DCA Recall Top N + 2 (↑)	DCA Precision Top 3 (↑)	DCA Precision Top 5 (↑)	DCA Precision All Sites (↑)
P2Rank	74.9	79.4	41.0	33.2	28.3
GrASP	77.5	<b>80.6</b>	<b>71.2</b>	<b>71.0</b>	<b>71.0</b>

<sup>a</sup>Arrows denote whether each metric increases or decreases with higher performance, and the highest performance is shown in bold for each metric.

**Table 3. Comparison between P2Rank and GrASP Performance on the HOLO4K(Mlig+) Test Set<sup>a</sup>**

	HOLO4K(Mlig+)				
	DCA Recall Top N (↑)	DCA Recall Top N + 2 (↑)	DCA Precision Top 3 (↑)	DCA Precision Top 5 (↑)	DCA Precision All Sites (↑)
P2Rank	81.2	<b>86.5</b>	45.9	35.4	25.5
GrASP	<b>81.3</b>	84.3	<b>72.8</b>	<b>71.6</b>	<b>71.4</b>

<sup>a</sup>Arrows denote whether each metric increases or decreases with higher performance, and the highest performance is shown in bold for each metric.

machine learning, which are always shown together to evaluate the trade-off between false negative and false positive errors.

**Validation Set Results.** To evaluate and tune our model, we performed 10-fold cross-validation on our augmented sc-PDB database.<sup>31</sup> The averaged binding site metrics across the 10-fold scale are shown in Table 1 with GrASP crossing 90% recall in the top  $N + 2$  category. Hyperparameter and model architecture choices were made to maximize the top  $N$  DCA recall in this setting.

**Test Set Results.** We evaluated both GrASP and P2Rank on our new versions of the COACH420 and HOLO4K sets previously used by P2Rank. COACH420(Mlig+) contains 256 single-chain systems, with 315 ligands bound across these systems. This set represents the setting where a small number of predictions are needed and where interfacial binding sites are not considered. Table 2 contains the DCA precision and recall metrics for both methods and shows that GrASP has gained 2.6% recall in the top  $N$  category as well as 30% or greater precision in all categories. To assess the significance of the difference in recall, we used McNemar's test<sup>56</sup> comparing which binding sites each method succeeded on. We found that in both the  $N$  and  $N + 2$  categories, the difference in recall was not significant. We also assessed the difference in the total number of binding sites returned by each method using the Wilcoxon signed-rank test.<sup>57</sup> We found the difference in site quantity significant with a  $p$ -value less than 0.001 when running three comparisons: comparing the total number of sites, the number in the top 3, and the number in the top 5. This difference explains the contrast in precision between the methods, with P2Rank consistently returning more sites. GrASP's precision is invariant with respect to the number of sites considered in this set, while P2Rank's precision falls as more sites are considered. This difference with respect to the number of sites considered is a consequence of reliance on ranking as there will be many sites returned outside of the top  $N$ . This shows the necessity of using a maximum number of binding sites or a site score threshold when using ranking-based methods in production.

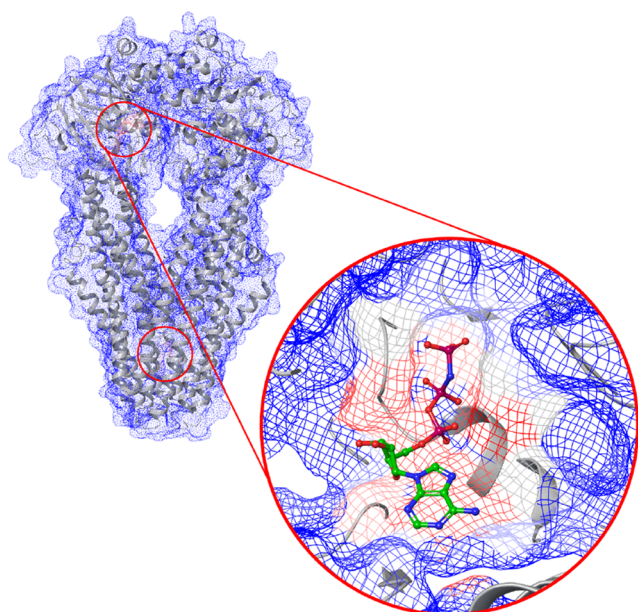
HOLO4K(Mlig+) contains a mix of single-chain and multichain systems with 6,368 ligands across 4,514 systems. Like COACH420(Mlig+), these systems primarily have one ligand bound but occasionally contain up to 12 ligands. We show in Table 3 that GrASP has a similar recall to P2Rank, is even outperformed by 2.2% in the top  $N + 2$  recall, but still outperforms P2Rank in precision by a wide margin. We again assessed the significance of these differences using McNemar's

test and the Wilcoxon signed-rank test. The difference in the Top  $N$  recall was not significant, but the difference in the top  $N + 2$  recall was significant with a  $p$ -value below 0.001. Similarly, the difference in the number of binding sites was significant with a  $p$ -value below 0.001 whether considering all sites, the top 3, or the top 5. As before, GrASP's precision falls by a much smaller amount as more sites are considered, highlighting that ranking too many sites without constraints is insufficient for real-world applications.

While computing the contingency tables for McNemar's test, we saw that many of the binding sites that were failure cases for one method were successes for the other. This prompted us to calculate the percentage of binding sites where either GrASP or P2Rank are successful. For  $N$  and  $N + 2$  on COACH420(Mlig+), either method succeeded on 84.13% and 86.35% of sites, respectively. For HOLO4K(Mlig+), either method succeeded on 90.33% for the top  $N$  and 92.73% for the top  $N + 2$ . Using predictions from both models provides a significant increase in binding site coverage and may be beneficial in studies where precision is not valued.

We also compute DCA recall with varying success thresholds for both test sets in Figure 4. Interestingly with less strict DCA success thresholds, P2Rank outperforms GrASP on both the top  $N$  and  $N + 2$  on HOLO4K(Mlig+), but GrASP's top  $N$  recall improves so significantly on COACH420(Mlig+) that it outperforms P2Rank's top  $N + 2$  recall.

**Sequence Identity Generalization.** The UniProt splitting criterion commonly used to prevent leakage between train and test sets is insufficient to assess a model's ability to generalize to novel proteins. While this approach mirrors the original P2Rank approach, we can quantify generalization more carefully by analyzing success rates as a function of sequence identity between the training and test sets. We used MMseqs2<sup>58</sup> to find the most similar entry in the training set for each system in the test set and assigned this sequence identity to all labeled binding sites in the test system. We then assigned each test binding site into histogram bins with 10% intervals in sequence identity (including the lower bound but not the upper bound). We recalculated the top  $N$  and  $N + 2$  DCA recall for each sequence identity bin individually to assess GrASP's performance with respect to the novelty of the test system's sequence. We show in Figure 5 that GrASP's DCA recall has a very small variance with respect to sequence identity for all bins with sufficient data (above 20% identity). Notably, GrASP is still able to maintain the same success rate

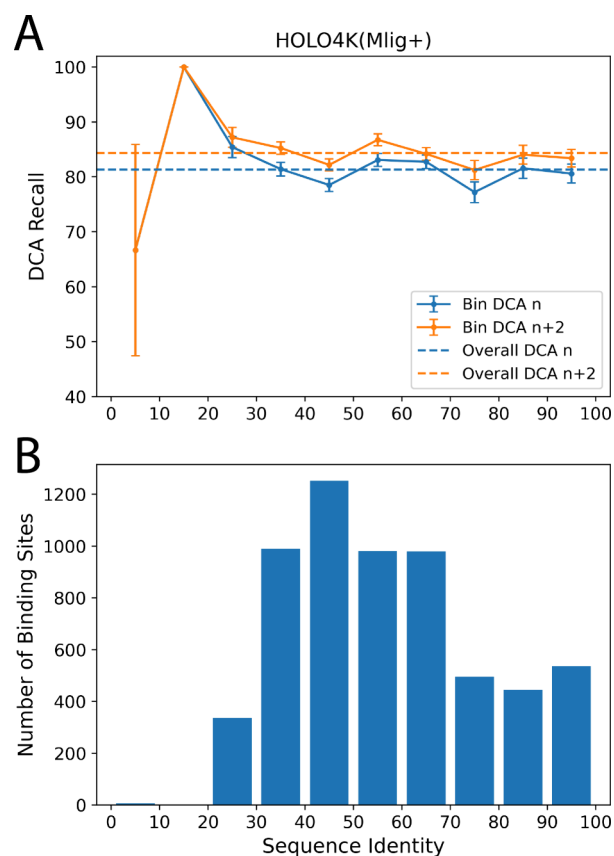


**Figure 3.** An example of GrASP atom druggability scores ranging from 0 (blue) to 1 (red) for PDB 4Q4A: an ABC transporter that does not have its UniProt ID in GrASP's training data. High-scoring regions are highlighted with red circles, and the scores around the ligand in this structure are shown.

for the 20–30% range where proteins are much less likely to be homologous. Here, we show this analysis for GrASP on HOLO4K(Mlig+) because the size of the test set allows for small standard error, but we show this analysis for both GrASP and P2Rank on both test sets in the SI. P2Rank's performance is also similar in all well-sampled bins but with higher variance than GrASP.

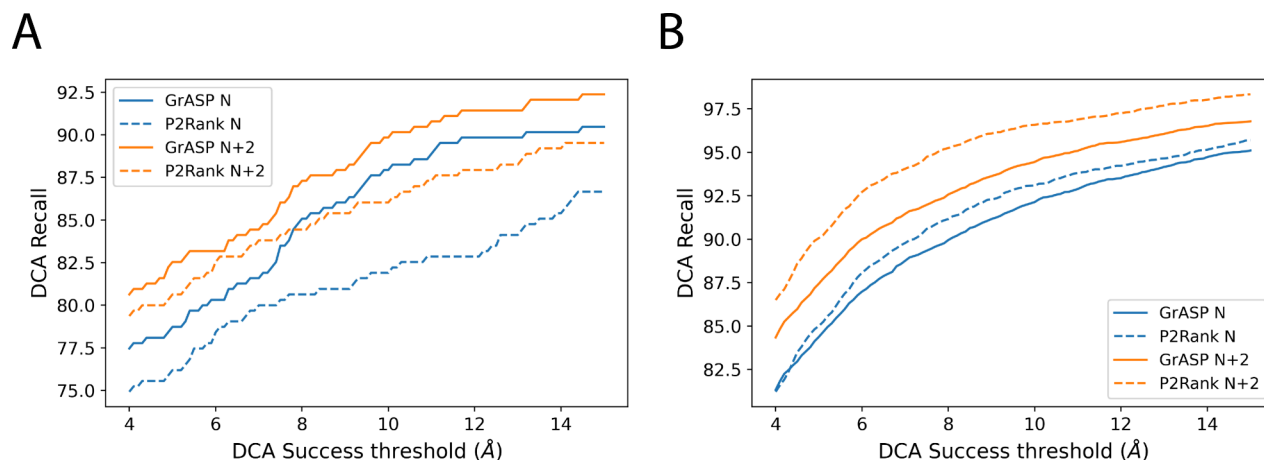
## DISCUSSION

In this work, we have developed a new method called Graph Attention Site Prediction which reaches state-of-the-art performance in binding site recall and does so with high precision, a metric that has not yet been reported for binding site prediction but affects the computational cost to use predicted binding sites for other tasks. Precision analysis in the



**Figure 5.** GrASP's performance on the HOLO4K(Mlig+) set as a function of sequence similarity between train and test sets. A) GrASP's performance on samples in each sequence similarity bin with standard error is displayed as bars, and the performance on the full set is shown as dashed lines. B) Histogram of sequence similarity between GrASP's training data and HOLO4K(Mlig+). Note that the 0–20% range has insufficient data to draw meaningful conclusions.

setting where the number of binding sites is unknown shows the weakness of ranking-based methods. If the true number of sites is not known, there is not a clear stopping point when using a ranked list, and downstream tasks may be frequently performed on poor predictions. We predict that coupling a ranked binding site list with a site score threshold to discard



**Figure 4.** Comparison of the DCA recall for GrASP and P2Rank with varying DCA success thresholds for A) COACH420(Mlig+) and B) HOLO4K(Mlig+).



poor predictions would improve precision and, in turn, reduce waste in downstream tasks for drug discovery. We recommend future methods aim to optimize such thresholds and report both precision and recall for DCA or other metrics of their choice.

Currently, binding site prediction methods either rank binding sites generated with geometric criteria or perform semantic segmentation and then cluster the segmentation mask. Future methods should treat binding site prediction as an instance segmentation task, where the model predicts which atoms (or surface points) are part of a binding site and which binding site they belong to. The current clustering-based instance segmentation is not end-to-end differentiable and lags behind the methodology used in image segmentation.<sup>59</sup> Given this suboptimal step in current methods, we recommend that small-scale projects use the raw semantic segmentation scores on surface atoms and handpick where to dock ligands. We also recommend that the community increases focus on treating the task as instance segmentation instead of perfecting methods for semantic segmentation, because clustering quality may set a cap on performance.

## ■ ASSOCIATED CONTENT

### Data Availability Statement

The trained GrASP model together with code, an easy-to-use web interface through Google Colab, and associated data sets to retrain the model are available at <https://github.com/tiwarylab/GrASP>.

### ■ Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.jcim.3c01698>.

Detailed description of methods and additional model evaluation (PDF)

## ■ AUTHOR INFORMATION

### Corresponding Author

**Pratyush Tiwary** – Institute for Physical Science and Technology, University of Maryland, College Park 20742, United States; Department of Chemistry and Biochemistry, University of Maryland, College Park 20742, United States; [orcid.org/0000-0002-2412-6922](https://orcid.org/0000-0002-2412-6922); Email: [ptiwary@umd.edu](mailto:ptiwary@umd.edu)

### Authors

**Zachary Smith** – Institute for Physical Science and Technology, University of Maryland, College Park 20742, United States; Biophysics Program, University of Maryland, College Park 20742, United States

**Michael Strobel** – Department of Computer Science, University of Maryland, College Park 20742, United States

**Bodhi P. Vani** – Institute for Physical Science and Technology, University of Maryland, College Park 20742, United States; [orcid.org/0000-0002-7747-279X](https://orcid.org/0000-0002-7747-279X)

Complete contact information is available at: <https://pubs.acs.org/doi/10.1021/acs.jcim.3c01698>

### Author Contributions

<sup>1</sup>Z.S. and M.S. contributed equally.

### Notes

The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

The authors declare the following competing financial interest(s): P.T. is a consultant to Schrodinger, Inc. and is on their Scientific Advisory Board.

## ■ ACKNOWLEDGMENTS

The research reported in this publication was supported by the National Institute of General Medical Sciences of the National Institutes of Health under Award Number R35GM142719. We are grateful to NSF ACCESS Bridges2 (project CHE180053) and University of Maryland Zaratán High-Performance Computing cluster for enabling the work performed here. The authors thank the P2Rank team for discussing their data sets, Michelle Girvan for suggesting average linkage clustering, Schrödinger and Nimbus Therapeutics scientists for discussing test set preparation, and Pavan Ravindra for discussing GNNs.

## ■ REFERENCES

- (1) McInnes, C. *Curr. Opin. Chem. Biol.* **2007**, *11*, 494–502.
- (2) Zhang, Z.; Li, Y.; Lin, B.; Schroeder, M.; Huang, B. *Bioinformatics* **2011**, *27*, 2083–2088.
- (3) Sherman, W.; Day, T.; Jacobson, M. P.; Friesner, R. A.; Farid, R. *J. Med. Chem.* **2006**, *49*, 534–553.
- (4) Clark, A. J.; Tiwary, P.; Borrelli, K.; Feng, S.; Miller, E. B.; Abel, R.; Friesner, R. A.; Berne, B. J. *J. Chem. Theory Comput.* **2016**, *12*, 2990–2998.
- (5) Jumper, J.; et al. *Nature* **2021**, *596*, 583–589.
- (6) Vani, B. P.; Aranganathan, A.; Wang, D.; Tiwary, P. AlphaFold2-RAVE: From Sequence to Boltzmann Ranking. *J. Chem. Theory Comput.* **2023**, *19*, 4351.
- (7) Baek, M.; et al. *Science* **2021**, *373*, 871–876.
- (8) Tunyasuvunakool, K.; et al. *Nature* **2021**, *596*, 590–596.
- (9) Kuzmanic, A.; Bowman, G. R.; Juarez-Jimenez, J.; Michel, J.; Gervasio, F. L. *Acc. Chem. Res.* **2020**, *53*, 654–661.
- (10) Benabderrahmane, M.; Bureau, R.; Voisin-Chiret, A. S.; Santos, J. S.-d. O. *J. Chem. Inf. Model.* **2021**, *61*, 5581–5588.
- (11) Oleinikovas, V.; Saladino, G.; Cossins, B. P.; Gervasio, F. L. *J. Am. Chem. Soc.* **2016**, *138*, 14257–14263.
- (12) Le Guilloux, V.; Schmidtke, P.; Tuffery, P. Fpocket: An open source platform for ligand pocket detection. *BMC Bioinf* **2009**, *10*, 168.
- (13) Levitt, D. G.; Banaszak, L. J. *J. Mol. Graphics* **1992**, *10*, 229–234.
- (14) Hendlich, M.; Rippmann, F.; Barnickel, G. *J. Mol. Graphics Modell.* **1997**, *15*, 359–363.
- (15) Weisel, M.; Proschak, E.; Schneider, G. PocketPicker: analysis of ligand binding-sites with shape descriptors. *Chem. Cent. J.* **2007**, *1*, 7.
- (16) Laskowski, R. A. *J. Mol. Graphics* **1995**, *13*, 323–330.
- (17) Brady, G. P.; Stouten, P. F. *J. Comput.-Aided Mol. Des.* **2000**, *14*, 383–401.
- (18) Tan, K. P.; Nguyen, T. B.; Patel, S.; Varadarajan, R.; Madhusudhan, M. S. *Nucleic Acids Res.* **2013**, *41*, W314–W321.
- (19) Capra, J. A.; Laskowski, R. A.; Thornton, J. M.; Singh, M.; Funkhouser, T. A. *PLoS Comput. Biol.* **2009**, *5*, e1000585.
- (20) Brylinski, M.; Skolnick, J. *Proc. Natl. Acad. Sci. U. S. A.* **2008**, *105*, 129–134.
- (21) Ngan, C.-H.; Hall, D. R.; Zerbe, B.; Grove, L. E.; Kozakov, D.; Vajda, S. *Bioinformatics* **2012**, *28*, 286–287.
- (22) Hernandez, M.; Ghersi, D.; Sanchez, R. *Nucleic Acids Res.* **2009**, *37*, W413–W416.
- (23) Huang, B.; Schroeder, M. LIGSITE(csc): predicting ligand binding sites using the Connolly surface and degree of conservation. *BMC Struct. Biol.* **2006**, *6*, 19.
- (24) Desaphy, J.; Bret, G.; Rognan, D.; Kellenberger, E. *Nucleic Acids Res.* **2015**, *43*, D399–D404.

- (25) Govindaraj, R. G.; Brylinski, M. Comparative assessment of strategies to identify similar ligand-binding pockets in proteins. *BMC Bioinf* **2018**, *19*, 91.
- (26) Jiménez, J.; Doerr, S.; Martínez-Rosell, G.; Rose, A. S.; De Fabritiis, G. *Bioinformatics* **2017**, *33*, 3036–3042.
- (27) Krivák, R.; Hoksza, D. P2Rank: machine learning based tool for rapid and accurate prediction of ligand binding sites from protein structure. *J. Cheminf.* **2018**, *10*, 39.
- (28) Jendele, L.; Krivák, R.; Skoda, P.; Novotny, M.; Hoksza, D. *Nucleic Acids Res.* **2019**, *47*, W345–W349.
- (29) Jakubec, D.; Skoda, P.; Krivák, R.; Novotny, M.; Hoksza, D. *Nucleic Acids Res.* **2022**, *50*, W593–W597.
- (30) Aggarwal, R.; Gupta, A.; Chelur, V.; Jawahar, C. V.; Priyakumar, U. D. DeepPocket: Ligand Binding Site Detection and Segmentation using 3D Convolutional Neural Networks. *J. Chem. Inf. Model.* **2022**, *62*, S069.
- (31) Stepniewska-Dziubinska, M. M.; Zielenkiewicz, P.; Siedlecki, P. Improving detection of protein-ligand binding sites with 3D segmentation. *Sci. Rep.* **2020**, *10*, S035.
- (32) Simonovsky, M.; Meyers, J. *J. Chem. Inf. Model.* **2020**, *60*, 2356–2366.
- (33) Pu, L.; Govindaraj, R. G.; Lemoine, J. M.; Wu, H.-C.; Brylinski, M. DeepDrug3D: Classification of ligand-binding pockets in proteins with a convolutional neural network. *PLoS Comput. Biol.* **2019**, *15*, e1006718.
- (34) Kipf, T. N.; Welling, M. Semi-Supervised Classification with Graph Convolutional Networks. arXiv:1609.02907. *arXiv Preprint*. 2016. <https://arxiv.org/abs/1609.02907> (accessed 2024-02-29).
- (35) Fey, M.; Lenssen, J. E. Fast Graph Representation Learning with PyTorch Geometric. arXiv:1903.02428. *arXiv Preprint*. 2019. <https://arxiv.org/abs/1903.02428> (accessed 2024-02-29).
- (36) Feinberg, E. N.; Sur, D.; Wu, Z.; Husic, B. E.; Mai, H.; Li, Y.; Sun, S.; Yang, J.; Ramsundar, B.; Pande, V. S. *ACS Cent. Sci.* **2018**, *4*, 1520–1530.
- (37) Townshend, R. J. L.; Vögele, M.; Suriana, P.; Derry, A.; Powers, A.; Laloudakis, Y.; Balachandar, S.; Jing, B.; Anderson, B.; Eismann, S.; Kondor, R.; Altman, R. B.; Dror, R. O. ATOM3D: Tasks On Molecules in Three Dimensions. arXiv:2012.04035. *arXiv Preprint*. 2020. <https://arxiv.org/abs/2012.04035> (accessed 2024-02-29).
- (38) Corso, G.; Stärk, H.; Jing, B.; Barzilay, R.; Jaakkola, T. DiffDock: Diffusion Steps, Twists, and Turns for Molecular Docking. arXiv:2210.01776. *arXiv Preprint*. 2022. <https://arxiv.org/abs/2210.01776> (accessed 2024-02-29).
- (39) Meller, A.; Ward, M.; Borowsky, J.; Lotthammer, J. M.; Kshirsagar, M.; Oveido, F.; Lavista Ferres, J.; Bowman, G. R. Predicting the locations of cryptic pockets from single protein structures using the PocketMiner graph neural network. *bioRxiv* **2022**, DOI: 10.1101/2022.06.28.497399.
- (40) Shi, W.; Singha, M.; Pu, L.; Ramanujam, J. R.; Brylinski, M. Graphsite: Ligand-binding site classification using Deep Graph Neural Network. *bioRxiv* **2021**, DOI: 10.1101/2021.12.06.471420.
- (41) Sverrisson, F.; Feydy, J.; Correia, B. E.; Bronstein, M. M. Fast End-to-End Learning on Protein Surfaces. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*; 2021; pp 15267–15276.
- (42) Brody, S.; Alon, U.; Yahav, E. How Attentive are Graph Attention Networks?. arXiv:2105.14491. *arXiv Preprint*. 2021. <https://arxiv.org/abs/2105.14491> (accessed 2024-02-29).
- (43) Xu, K.; Hu, W.; Leskovec, J.; Jegelka, S. How Powerful are Graph Neural Networks? arXiv:1810.00826. *arXiv Preprint*. 2018. <https://arxiv.org/abs/1810.00826> (accessed 2024-02-29).
- (44) Battaglia, P. W. et al. *CoRR*; 2018; abs/1806.01261.
- (45) Godwin, J.; Schaarschmidt, M.; Gaunt, A. L.; Sanchez-Gonzalez, A.; Rubanova, Y.; Veličković, P.; Kirkpatrick, J.; Battaglia, P. Simple GNN Regularisation for 3D Molecular Property Prediction and Beyond. *International Conference on Learning Representations*; 2022.
- (46) He, K.; Zhang, X.; Ren, S.; Sun, J. Deep Residual Learning for Image Recognition. *Proceedings of the IEEE conference on computer vision and pattern recognition*; 2016; pp 770–778.
- (47) Xu, K.; Li, C.; Tian, Y.; Sonobe, T.; Kawarabayashi, K.-i.; Jegelka, S. Representation Learning on Graphs with Jumping Knowledge Networks. *Proceedings of the 35th International Conference on Machine Learning*; 2018; pp 5453–5462.
- (48) Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Lio, P.; Bengio, Y. Graph Attention Networks. arXiv:1710.10903. *arXiv Preprint*. 2017. <https://arxiv.org/abs/1710.10903> (accessed 2024-02-29).
- (49) Sokal, R. R.; Michener, C. D. A statistical method for evaluating systematic relationships [J]. *Univ. Kans. Sci. Bull.* **1958**, *38*, 1409–1438.
- (50) Krivák, R.; Hoksza, D. Improving protein-ligand binding site prediction accuracy by classification of inner pocket points using local features. *J. Cheminf.* **2015**, *7*, 12.
- (51) Corso, G.; Cavalleri, L.; Beaini, D.; Liò, P.; Veličković, P. Principal Neighbourhood Aggregation for Graph Nets. *Advances in Neural Information Processing Systems* **2020**, 13260–13271.
- (52) Ulyanov, D.; Vedaldi, A.; Lempitsky, V. Instance Normalization: The Missing Ingredient for Fast Stylization. arXiv:1607.08022. *arXiv Preprint*. 2016. <https://arxiv.org/abs/1607.08022> (accessed 2024-02-29).
- (53) Clevert, D.-A.; Unterthiner, T.; Hochreiter, S. Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs). arXiv:1511.07289. *arXiv Preprint*. 2015. <https://arxiv.org/abs/1511.07289> (accessed 2024-02-29).
- (54) Rockafellar, R.T. *Convex Analysis*; 1970.
- (55) Preparata, F. P.; Shamos, M. I. *Computational Geometry: An Introduction*; Springer Science & Business Media: 2012.
- (56) McNemar, Q. *Psychometrika* **1947**, *12*, 153–157.
- (57) Conover, W. J. *Practical Nonparametric Statistics*; John Wiley & Sons: 1999; Vol. 350.
- (58) Steinegger, M.; Söding, J. *Nat. Biotechnol.* **2017**, *35*, 1026–1028.
- (59) Kirillov, A.; Mintun, E.; Ravi, N.; Mao, H.; Rolland, C.; Gustafson, L.; Xiao, T.; Whitehead, S.; Berg, A. C.; Lo, W.-Y.; Dollár, P.; Girshick, R. Segment Anything. arXiv:2304.02643. *arXiv Preprint*. 2023. <https://arxiv.org/abs/2304.02643> (accessed 2024-02-29).