

RESEARCH

Open Access



Comparative evaluation of methods for the prediction of protein–ligand binding sites

Javier S. Utgés¹ and Geoffrey J. Barton^{1*}

Abstract

The accurate identification of protein–ligand binding sites is of critical importance in understanding and modulating protein function. Accordingly, ligand binding site prediction has remained a research focus for over three decades with over 50 methods developed and a change of paradigm from geometry-based to machine learning. In this work, we collate 13 ligand binding site predictors, spanning 30 years, focusing on the latest machine learning-based methods such as VN-EGNN, IF-SitePred, GrASP, PURESNet, and DeepPocket and compare them to the established P2Rank, PRANK and fpocket and earlier methods like PocketFinder, Ligsite and Surfnet. We benchmark the methods against the human subset of our new curated reference dataset, LIGYSIS. LIGYSIS is a comprehensive protein–ligand complex dataset comprising 30,000 proteins with bound ligands which aggregates biologically relevant unique protein–ligand interfaces across biological units of multiple structures from the same protein. LIGYSIS is an improvement for testing methods over earlier datasets like sc-PDB, PDBbind, binding MOAD, COACH420 and HOLO4K which either include 1:1 protein–ligand complexes or consider asymmetric units. Re-scoring of fpocket predictions by PRANK and DeepPocket display the highest recall (60%) whilst IF-SitePred presents the lowest recall (39%). We demonstrate the detrimental effect that redundant prediction of binding sites has on performance as well as the beneficial impact of stronger pocket scoring schemes, with improvements up to 14% in recall (IF-SitePred) and 30% in precision (Surfnet). Finally, we propose top- $N+2$ recall as the universal benchmark metric for ligand binding site prediction and urge authors to share not only the source code of their methods, but also of their benchmark.

Scientific contributions

This study conducts the largest benchmark of ligand binding site prediction methods to date, comparing 13 original methods and 15 variants using 10 informative metrics. The LIGYSIS dataset is introduced, which aggregates biologically relevant protein–ligand interfaces across multiple structures of the same protein. The study highlights the detrimental effect of redundant binding site prediction and demonstrates significant improvement in recall and precision through stronger scoring schemes. Finally, top- $N+2$ recall is proposed as a universal benchmark metric for ligand binding site prediction, with a recommendation for open-source sharing of both methods and benchmarks.

Keywords Ligand binding site prediction, Binding pocket, Benchmark, Reference dataset, Machine learning, Drug discovery

Introduction

Identifying where ligands can bind to proteins is of critical importance in understanding and modulating protein function. While X-ray crystallography remains the gold-standard to identify and characterise binding sites

*Correspondence:

Geoffrey J. Barton
gjbarton@dunee.ac.uk

¹ Division of Computational Biology, School of Life Sciences, University of Dundee, Dow Street, Dundee DD1 5EH, Scotland, UK



© The Author(s) 2024. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

[1–5], over the last three decades, significant effort has been made to develop computational methods that predict binding sites from an apo three-dimensional protein structure [6].

Prediction methods exploit a variety of different techniques to suggest binding sites. Geometry-based techniques like fpocket [7], Ligsite [8] and Surfnets [9] identify cavities by analysing the geometry of the molecular surface of a protein and usually rely on the use of a grid, gaps, spheres, or tessellation [7–14]. Energy-based methods such as PocketFinder [15] rely on the calculation of interaction energies between the protein and chemical group or probe to identify cavities [15–20]. Conservation-based methods make use of sequence evolutionary conservation information to find patterns in multiple sequence alignments and identify conserved key residues for ligand site identification [21–23]. Template-based methods rely on structural information from homologues and the assumption that structurally conserved proteins might bind ligands at a similar location [24–29]. Combined approaches or meta-predictors combine multiple methods, or the use of multiple types of data to infer ligand binding sites, e.g., geometric features with sequence conservation [30–36]. Finally, machine learning methods utilise a wide range of machine learning techniques including random forest, as well as deep, graph, residual, or convolutional neural networks [37–57]. Machine learning methods comprise the bulk of the methods reviewed in this analysis and are exemplified by PRANK [37], P2Rank [38, 40], DeepPocket [48], PURESNet [44, 57], GrASP [54], IF-SitePred [55] and VN-EGNN [56]. Open source, peer-reviewed and easy to install methods were prioritised (Supplementary Table 1). This set of methods represents the most complete and relevant set of ligand binding site prediction tools benchmarked to date and is representative of the state-of-the-art within the field.

Table 1 summarises the methods evaluated in this work, which were executed with their standard settings. VN-EGNN [56] combines virtual nodes with equivariant graph neural networks. Virtual nodes, represented by ESM-2 embeddings [58] are passed through a series of message-passing layers until they reach their final coordinates, which represent the centroid of predicted pockets. Pocket residues are not reported. IF-SitePred [55] represents protein residues with ESM-IF1 embeddings [59] and employs 40 different LightGBM models [60] to classify residues as ligand-binding if all forty models return a $p > 0.5$. It later utilises PyMOL [61] to place a series of cloud points which are clustered using DBSCAN [62] and a threshold of 1.7 Å. Pocket centroids are obtained by averaging the clustered points' coordinates, scored and ranked based on the number of

cloud points. Like VN-EGNN, no pocket residues are defined. GrASP [54] employs graph attention networks to perform semantic segmentation on all surface protein atoms, represented by 17 atom, residue and bond-level features, scoring which are likely part of a binding site. Atoms with a score > 0.3 are clustered into binding sites using average linkage and a threshold of 15 Å. Pocket scores are calculated as the sum of squares of binding site atom scores. PURESNet [44] combines deep residual and convolutional neural networks to predict ligand binding sites using an 18-element vector of atom-level features and one-hot encoding to represent grid voxels. Voxels with a score > 0.34 are clustered into binding sites using DBSCAN and a threshold of 5.5 Å [57]. Pockets are represented by their residues, but neither pocket centroid, nor score or ranking are reported. Similarly to PURESNet, DeepPocket [48] exploits convolutional neural networks on grid voxels represented by 14 atom-level features to re-score (DeepPocket_{RESC}) and additionally extract new pocket shapes (DeepPocket_{SEG}) from fpocket candidates. P2Rank [40] relies on solvent accessible surface (SAS) points placed over the protein surface, represented by 35 atom and residue-level features, and a random forest classifier to score them based on their likelihood of binding a ligand. SAS points with a score > 0.35 are clustered into sites using single linkage and a threshold of 3 Å. P2Rank_{CONS} [63] works in the same manner but considers an extra feature: amino acid conservation as measured by Jensen–Shannon divergence [64]. Both report residue and pocket level scores, as well as pocket centroids and rank. PRANK [37] uses the P2Rank algorithm to score/re-score pockets predicted by other methods. In this work, fpocket predictions are re-scored with PRANK (fpocket_{PRANK}). PocketFinder [15] uses the Lennard–Jones [65] transformation on a 1 Å grid surrounding the protein surface to predict protein cavities. PocketFinder does not report pocket centroid, score or rank. Finally, geometry-based methods: fpocket [7], Ligsite [8] and Surfnets [9] rely on the geometry of the molecular surface to find cavities. fpocket is the only one of these three methods that reports pocket centroid, score, rank and residues.

In this work, we have gathered thirteen ligand binding site prediction tools, spanning three decades of research and compared them against our reference dataset LIGYSIS. LIGYSIS extends the work by Utgés et al. [5] that identifies human protein–ligand binding sites for biologically relevant ligands, defined by BioLiP [66], from protein structures determined by X-ray crystallography. We compared the thirteen methods to each other and to the LIGYSIS reference dataset according to a range of metrics including the number of ligand sites, their size, shape, proximity, overlap and redundancy. We have

Table 1 Summary of ligand binding site prediction methods analysed in this study

Method	Approach	Features	# Features	P centroid	P residues	P score	P ranking	R score	R threshold	Cluster	Algorithm	Threshold (Å)
VN-EGNN	EGNN + VN	ESM-2 embeddings	1280	✓	×	✓	✓	×	–	–	–	–
IF-SitePred	LightGBM	ESM-IF1 embeddings	512	✓	×	✓	✓	×	0.5 (ALL 40)	Cloud points	DBSCAN	1.7
GrASP	GAT-GNN	Atom, residue, bond...	17	✓	×	✓	✓	✓	0.3	Atoms	Average	15
PUResNet	DRN + 3D-CNN	Atom + one-hot encoding	18	×	✓	×	×	×	0.34	Atoms	DBSCAN	5.5
DeepPocket	fpocket + 3D-CNN	Atom	14	✓	✓	✓	✓	×	–	–	–	–
P2Rank _{CONS}	Random Forest	Atom and residue	36	✓	✓	✓	✓	✓	0.35	SAS points	Single	3
P2Rank	Random Forest	Atom and residue	35	✓	✓	✓	✓	✓	0.35	SAS points	Single	3
fpocket _{P-RANK}	fpocket + Random Forest	Atom and residue	34	✓	✓	✓	✓	✓	–	–	–	–
fpocket	α-spheres	–	–	×	✓	✓	✓	×	–	α-spheres	Multiple	1.7
PocketFinder ⁺	LJ potential	–	–	×	×	×	×	✓	–	–	–	–
Ligsite ⁺	Cubic grid	–	–	×	×	×	×	✓	–	–	–	–
Surfnet ⁺	Gap regions	–	–	×	×	×	×	✓	–	–	–	–

All these methods were used with their default settings. Check marks (✓) indicate that a method provides a given output and crosses (×) the contrary. Dashes (–) indicate a field is not applicable for a given method, e.g., features for non-machine learning-based methods. Approach: the techniques applied by the method; Features/#Features: the features and their number if the method is machine learning-based; P centroid/P residues/P score/P ranking/R score: whether the method reports the pocket centroid, pocket residues, pocket score, pocket ranking and residue *ligandability* score. Information about their clustering strategies is also relevant: whether the method uses a residue *ligandability* threshold (R threshold), the instances they cluster (Cluster) to define the distinct pockets, the clustering algorithm used (Algorithm) and threshold employed (Threshold). For example, P2Rank uses a random forest classifier on SAS points represented by 35 atom and residue features. Points with a score > 0.35 are later clustered into binding sites using single linkage and a threshold of 3 Å. DeepPocket and fpocket_{P-RANK} use fpocket predictions as a starting point and later employ different technologies to re-score or re-define pockets. EGNN + VN: equivariant graph neural network + virtual nodes; LightGBM: light gradient boosting machine; GAT: graph attention network; GNN: graph neural network; DRN: deep residual network; 3D-CNN: three-dimensional convolutional neural network; LJ potential: Lennard–Jones potential

explored in detail and demonstrated the importance of a robust pocket scoring scheme and highlight how some of the methods can improve their performance significantly by re-scoring their predictions. This analysis identifies the strengths and weaknesses of prediction assessment metrics and leads to guidance for developing ligand binding site prediction tools or using these methods to understand protein function and in drug development. This work represents the first independent ligand site prediction benchmark for over a decade, since Chen et al. [67] and the largest to date in terms of dataset size (2775), methods compared (13 original + 15 variants) and metrics employed (>10).

Results

The LIGYSIS dataset

The ligand binding site analysis, or LIGYSIS, dataset comprises protein–ligand complexes for 3448 human proteins. For each protein, biologically relevant protein–ligand interactions, in accordance with BioLiP [66], are considered across the PISA-defined [68] biological assemblies of the multiple entries deposited in the PDB [69]. Ligands are clustered using their protein interaction fingerprint to identify ligand binding sites as described by Utgés et al. [5, 70] (see “Methods”). The full LIGYSIS dataset includes $\approx 30,000$ proteins with known ligand-bound complexes. Here, we focus on the *human subset* of LIGYSIS as a manageable set to run all prediction methods on and refer to this as *LIGYSIS* for brevity.

The LIGYSIS dataset differs from previous train and test sets for ligand binding sites by considering biological units, aggregating multiple structures of the same protein, and removing redundant protein–ligand interfaces. The asymmetric unit is the smallest portion of a crystal structure that can reproduce the complete unit cell through a series of symmetry operations. The asymmetric unit often does not correspond to the biological assembly or unit and relying on it can lead to artificial crystal contacts or redundant protein–ligand interfaces. The biological unit is the biologically relevant and functional macromolecular assembly for a given structure and might be formed by one, multiple copies or a portion of the asymmetric unit [71]. LIGYSIS consistently

considers biological units, which is key in any analysis that delves into molecular interactions at residue or atomistic level. An example of this illustrated in Fig. 1A is PDB: 1JQY [72], present in the HOLO4K dataset, where the asymmetric unit is formed by three copies of a homopentamer, whereas the biological unit comprises a single pentamer. In the asymmetric unit of PDB: 1JQY 14 molecules of BMSC-0010 (A32) interact with 14 copies of *Escherichia coli* heat-labile enterotoxin B chain (P32890). This protein–ligand interface is the same repeated 14 times (Fig. 1A). Protein–ligand interface redundancy can also be an issue when the asymmetric unit equals the biological assembly (Fig. 1B). In PDB: 1PPR [73], also in HOLO4K, molecules of chlorophyll A (CLA), peridinin (PID) and digalactosyl diacyl glycerol (DGD) bind to the three copies of a peridinin-chlorophyll a-binding protein 1, chloroplastic, PCP, (P80484) trimer, resulting in a redundancy of 3 \times . To account for this, LIGYSIS considers unique non-redundant protein–ligand interfaces by retrieving the UniProt residue numbers of the residues the ligands interact with, so 1/14 interfaces are considered for PDB: 1JQY and 12/36 for PDB: 1PPR. Finally, unique ligand interactions are aggregated across different structures for the same protein and ligand sites defined.

Figure 2 shows the comparison between PDB: 4GQQ [74], present in the PDBbind dataset, and the LIGYSIS entry for human pancreatic alpha-amylase (P04746), which representative structure is also PDB: 4GQQ. The entry in PDBbind represents a single protein–ligand complex, whereas LIGYSIS makes use of 51 structures, 195 ligands to define 13 different ligand binding sites. LIGYSIS aggregates all unique biologically relevant protein–ligand interactions for a protein in a non-redundant manner, thus representing the most complete and integrative protein–ligand binding dataset up to date. For this reason, we propose LIGYSIS as a new benchmark dataset for the prediction of ligand binding sites and use it in this work to evaluate a set of thirteen ligand binding site prediction and cavity identification tools.

Comparison of datasets

To assess the scope and limitations of the predictive methods surveyed in this work, their training and test

(See figure on next page.)

Fig. 1 Redundancy in protein–ligand interfaces. Two examples of how the type of macromolecular assembly and difference between the asymmetric and biological unit of a protein–ligand complex results in redundant protein–ligand interfaces. **A** For PDB: 1JQY [72], the asymmetric unit comprises three copies of a homopentamer, whereas the biologically functional assembly is a single pentamer. A BMSC-0010 ligand molecule binds to each copy, except for one, of each of the three pentamers. This results in the same protein–ligand interface repeated 14 times, i.e., 14 \times redundancy; **B** for PDB: 1PPR [73] both the asymmetric and biological units are a homo-trimer. Different molecules of the same ligands are binding to the same interfaces across the three copies of the trimer, i.e., 3 \times redundancy. Dashed rectangles indicate the asymmetric/biological units

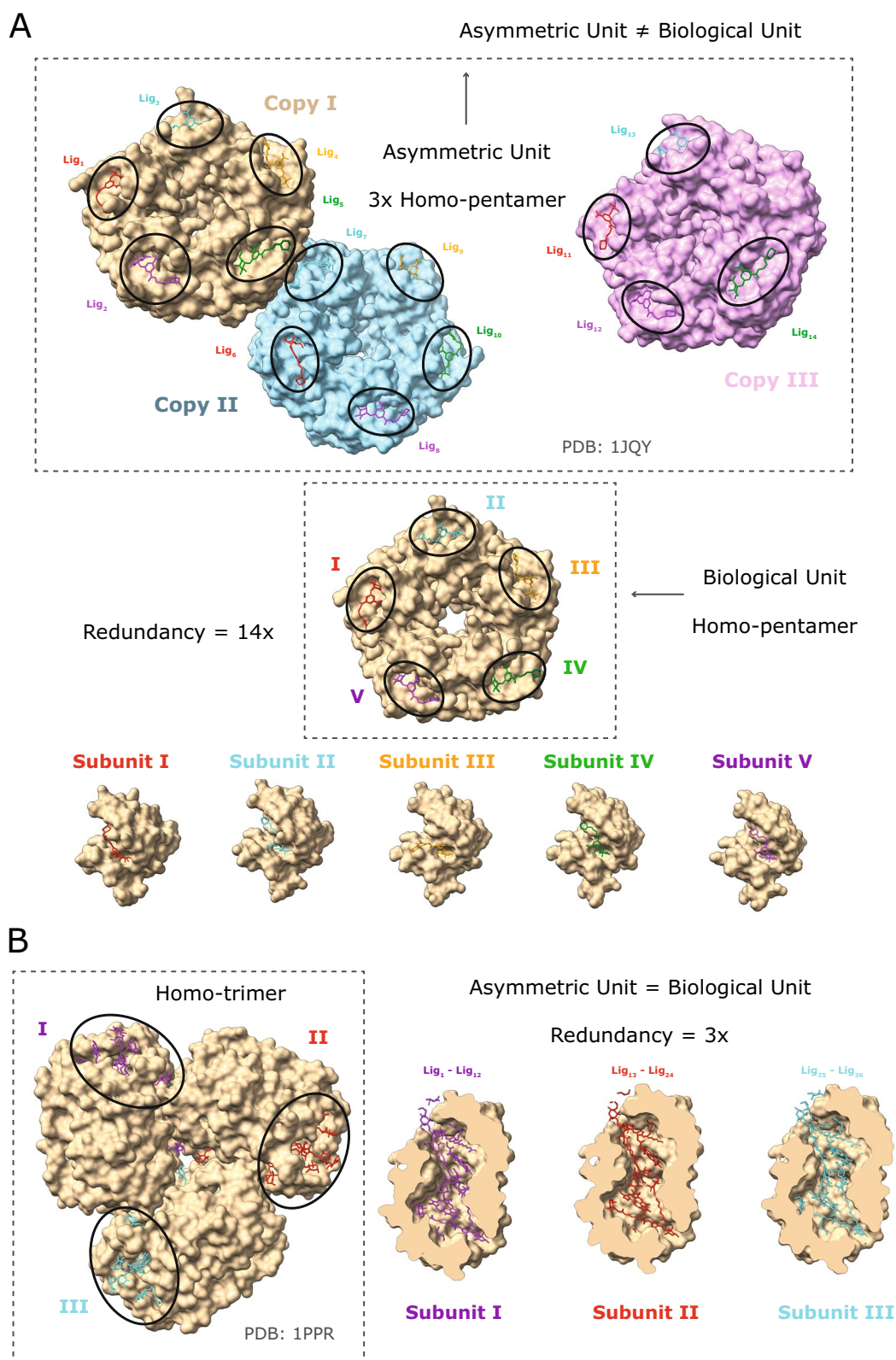


Fig. 1 (See legend on previous page.)

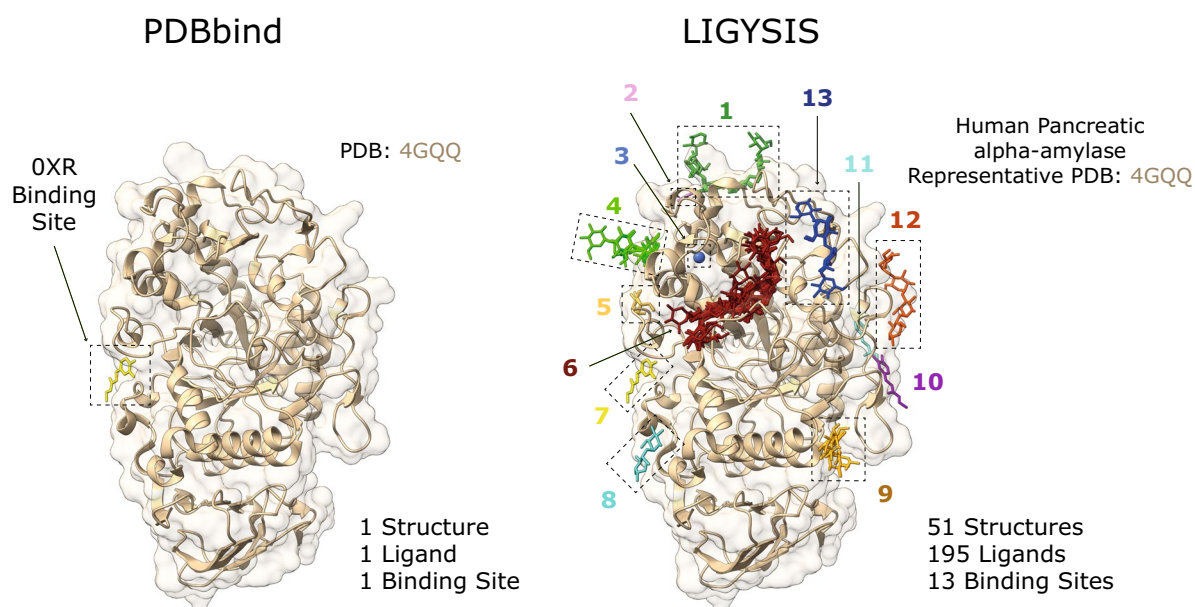


Fig. 2 Comparison of PDBbind and LIGYSIS. PDBbind is comprised by complexes between a protein and the most biologically relevant ligand in a structure. For PDB: 4GQQ [74], this is ethyl caffeate (OXR). LIGYSIS considers all unique biologically relevant protein–ligand interactions across all the structures for a given protein. For human pancreatic alpha-amylase (P04746), which representative structure is PDB: 4GQQ, 13 ligand binding sites are defined from 195 ligands across 51 structures. LIGYSIS provides a better representation of the ligand-binding capabilities of a protein than a single protein–ligand complex and constitutes therefore a better benchmark for ligand binding site prediction tools

sets were compared with LIGYSIS by number of sites per protein, number of interacting protein chains per ligand site, ligand size, ligand site size, and ligand composition. $sc\text{-PDB}_{\text{FULL}}$ represents the full $sc\text{-PDB}$ dataset used for training by DeepPocket, $b\text{MOAD}_{\text{SUB}}$ the subset of binding MOAD used for training by IF-SitePred and $\text{PDBbind}_{\text{REF}}$ the reference subset of PDBbind which VN-EGNN uses for testing. Table 2 summarises the size of the datasets, which methods employ them and their overlap with the LIGYSIS set. LIGYSIS differs from all other datasets since biologically relevant ions are considered, comprising $\approx 40\%$ of the ligand sites. For additional reference, $\text{LIGYSIS}_{\text{NI}}$, a subset of LIGYSIS without ions, is also included in this analysis. Figure 3A shows the number of binding sites per entry across datasets. $sc\text{-PDB}_{\text{FULL}}$, $\text{PDBbind}_{\text{REF}}$ and SC6K only consider the most relevant ligand for each entry. COACH420 and JOINED mostly present single-ligand entries ($\approx 70\%$). $b\text{MOAD}_{\text{SUB}}$ (46%) and CHEN11 (58%) present more similar distributions to LIGYSIS, where 54% of the protein chains present more than one binding site. This percentage decreases for $\text{LIGYSIS}_{\text{NI}}$ (38%) as ion sites are removed. HOLO4K presents the highest proportion (62%) of multi-ligand entries. Both HOLO4K and COACH420 are based on asymmetric units, and not biological assemblies. For HOLO4K, 1811 (40%) of structures present different numbers of chains between the

asymmetric and biological units. This is even more frequent in COACH420: 234 (56%). Moreover, multimeric complexes might present the same protein–ligand interface repeated across the copies of the complex (Fig. 1). Considering predictions of these interfaces as independent can lead to an overestimate of the performance of a predictor. Regarding the number of chains interacting with a given ligand (Fig. 3B), CHEN11, COACH420 and JOINED present the smaller fraction of multimeric protein–ligand interactions: 3, 6 and 8%, respectively, whereas SC6K presents the highest (44%). The rest of the methods range between 20 and 30%. There are no striking differences regarding the size of the interacting protein chains, represented by the number of residues (Fig. 3C).

Figure 3D represents the ligand type composition of the datasets. Non-polymer ligands dominate all datasets ($>66\%$), and the proportion of peptides and nucleic acids differ across datasets, with JOINED and LIGYSIS presenting fewer ligands of these types (0.9% and 1.6%). $sc\text{-PDB}_{\text{FULL}}$ and SC6K are depleted in saccharides ($<1\%$). Figure 3E depicts the difference in the number of atoms of the ligands in each dataset. LIGYSIS is of course different, due to its ion ligand content, however, there is no difference between $\text{LIGYSIS}_{\text{NI}}$ and the other datasets. Figure 3F conveys the difference in the number of ligand-interacting residues. LIGYSIS has the largest proportion

Table 2 Summary statistics of the different datasets analysed in this study

Dataset	Type	# Structures	# Sites	# Ligands	Overlap (%)	Methods
LIGYSIS	NEW	3448	8244	65,116⁺	–	–
LIGYSIS _{NI}	NEW	2275	4572	38,595	–	–
sc-PDB _{FULL}	TRAIN	17,594⁺	17,594⁺	17,594	801[–] (9.7)	VN-EGNN, GrASP, P2ResNet, DeepPocket
bMOAD _{SUB}	TRAIN	5899	11,184	11,184	606 (7.6)	IF-SitePred
CHEN11	TRAIN	244[–]	479[–]	479[–]	40⁺ (0.5)	PRANK, P2Rank
PDBbind _{REF}	TEST	5316	5316	5316	310 (3.8)	VN-EGNN
SC6K	TEST	6147	6147	6147	259 (3.1)	DeepPocket
HOLO4K	TEST	4009	10,175	10,175	207 (2.5)	<i>ALL</i> *
COACH420	TEST	413	624	624	41 (0.5)	VN-EGNN, GrASP, DeepPocket, P2Rank, P2ResNet
JOINED	TEST	557	752	752	110 (1.3)	PRANK

LIGYSIS is our reference dataset, LIGYSIS_{NI} is a subset with no ion (NI) ligand binding sites, sc-PDB_{FULL}, bMOAD_{SUB} and CHEN11 constitute the training datasets, whereas PDBbind_{REF}, SC6K, HOLO4K, COACH420 and JOINED represent test sets. # Structures, # Sites and # Ligands represent the number of PDB structures, ligand sites and total number of ligands for each dataset. Note that for LIGYSIS and LIGYSIS_{NI}, 3448 and 2775, are the number of human structural segments considered, each represented by a single chain. For each segment, all biologically relevant ligand-binding structures were considered: $N = 23,321$ (LIGYSIS) and $N = 19,012$ (LIGYSIS_{NI}). The number of ligands, or protein–ligand complexes, is not equal to the number of sites for LIGYSIS, as data from multiple structures of the same protein are aggregated into unique sites, i.e., a LIGYSIS site often includes multiple ligands. Overlap is the number of LIGYSIS binding sites represented by at least one protein–ligand complex for a given dataset. Percentage relative to LIGYSIS also reported. Methods represents the ligand site predictors that use these datasets for training or test. Only the original version of each dataset is considered in the analysis, e.g., HOLO4K is analysed, but not HOLO4K_{Mlig}, nor HOLO4K_{Mlig+}, HAP, or HAP-small. The same goes for Mlig, Mlig+ versions of COACH420, sc-PDB_{SUB} and sc-PDB_{RICH}. *ALL** represents all the methods compared in this work except for PRANK, fpocket, PocketFinder⁺, Ligsite⁺ and Surfnet⁺. For # Structures, # Sites and # Ligands, highest values are indicated with “⁺” bold superscript and lowest with “[–]”. This is the other way around for Overlap

of small sites, 1–10 residues (56%). This is directly related to the prominent ion component, and the frequency decreases when ions are removed (LIGYSIS_{NI}: 36%). CHEN11, COACH420, JOINED and HOLO4K are more similar to LIGYSIS_{NI}, whereas sc-PDB_{FULL}, PDBbind_{REF} and SC6K are clearly different and present almost exclusively large sites, larger than 20 amino acids (>90%). Figure 3G explores the ligand diversity on each dataset by showing the top-5 most frequent ligands per dataset, the proportion of the total number of ligands they represent, as well as Shannon’s entropy H' . Shannon’s entropy is a measure of diversity. Larger numbers indicate a more

evenly spread distribution of a larger number of different molecules, whereas small numbers indicate higher frequency of a few ligands. Whilst four out of the top-5 ligands of LIGYSIS are ions (Zn^{+2} , Ca^{+2} , Mg^{+2} , Mn^{+2}) and represent 19.2% of all ligands, its diverse composition is comparable to that of PDBbind_{REF}. Without ions, LIGYSIS_{NI} is the most diverse dataset with $H' = 8.8$ and its top-5 ligands only covering 5.3% of all ligands in the set. SC6K is the least diverse with its top-5 most frequent ligands covering 33% of all ligands. All datasets except for LIGYSIS, LIGYSIS_{NI} and PDBbind_{REF} are dominated by co-factor ligands, such as flavin-adenine dinucleotide

(See figure on next page.)

Fig. 3 Comparison of datasets. **A, B, C, E** and **F** plot the frequencies of binned intervals of a discrete variable, coloured using the *cividis* palette. Interval ranges were selected to facilitate data interpretation. **A** Number of ligand binding sites per dataset entry; **B** number of ligand-interacting protein chains. This represents whether the ligand interacts with a single protein chain or more; **C** length of ligand-interacting protein chains (number of amino acids); **D** ligand molecule type frequency as described in the CCD; **E** number of ligand atoms; **F** binding site size, i.e., number of ligand-interacting protein residues. Dashed lines drawn at $p = 0.5$; **G** five most frequent ligands per dataset. Error bars represent 95% confidence interval of the proportion [75]. Ligands of similar type are coloured in shades of the same colour: greens for ions, reds for co-factors, blues for energy-carrier molecules, yellows for sugars, grey for peptides and white for other non-polymeric ligands. Above the bars, Shannon’s Entropy and the proportion of all ligands in each set covered by these top-5 can be found. Both are measures of ligand diversity within each dataset. A subset of LIGYSIS with no ions (NI), LIGYSIS_{NI}, is included in this analysis, as most training and test datasets do not consider ions. *ADE* adenine, *CLR* cholesterol, *GAI* guanidine, *GSH* glutathione, *MAN* mannose, *FUC* fucose, *NAP* nicotinamide–adenine–dinucleotide phosphate, *SAH* S-adenosyl-L-homocysteine, *FMN* flavin mononucleotide, *GAL* galactose, *N-mer* protein peptides of N amino acids, *017* darunavir, *PLP* vitamin B6 phosphate, *BGC* glucose, *XYP* xylose

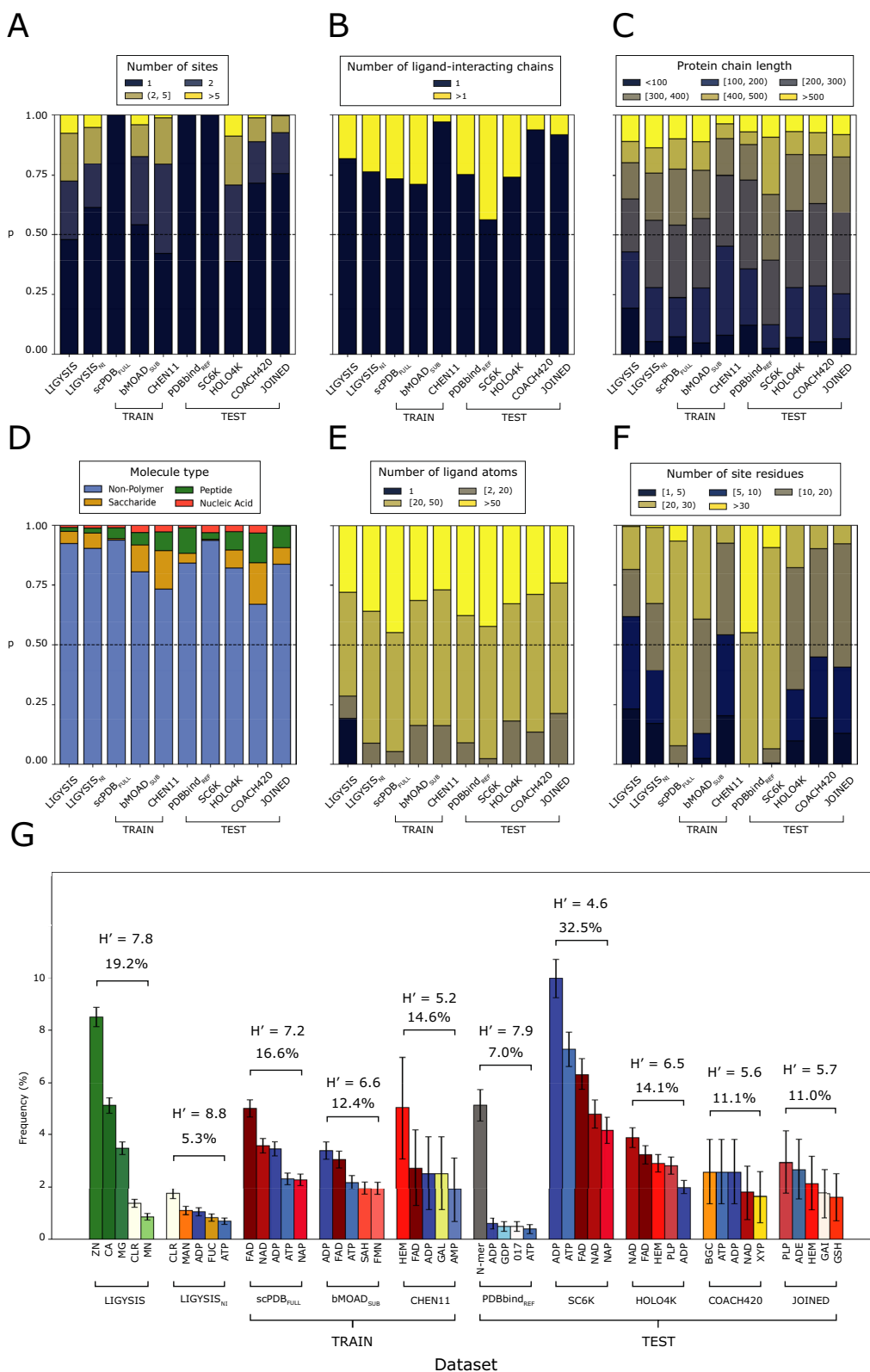


Fig. 3 (See legend on previous page.)

(FAD), nicotinamide-adenine dinucleotide (NAD), and haem (HEM), or energy carrier molecules such as adenine tri-, di- and monophosphate (ATP, ADP, AMP). Short peptides (<10 aas) are the most common ligands in PDBbind_{REF} (5%), and energy carriers represent <2% of the top-5 ligands. Cholesterol, mannose and fucose are some of the most common ligands in LIGYSIS_{NI}. For all further analysis, including characterisation of binding pockets and performance evaluation of the prediction methods, LIGYSIS, i.e., *including ions*, was utilised. This was done to have a larger benchmark set and to challenge the methods as they have not been trained on such sites.

Characterisation of binding pockets

After removing representative chains with missing residue mappings to UniProt, the final LIGYSIS set which was employed for the benchmark of the methods comprises 2775 protein chains. All methods were executed with default parameters. P2Rank_{CONS} represents predictions by P2Rank with an extra feature of amino acid conservation and DeepPocket_{SEG} are the pockets extracted by the segmentation module of DeepPocket (see “Methods”). Not all methods predict pockets on all the chains. VN-EGNN, GrASP, fpocket, PocketFinder⁺, Ligsite⁺, and Surfnet⁺ predict in >99% of the chains, P2Rank_{CONS} on 93%, followed by P2Rank on 86%, PURESNet and DeepPocket_{SEG} (85%), and finally IF-SitePred only predicts pockets on 75% of the chains. PURESNet, DeepPocket, P2Rank_{CONS} and P2Rank often don't predict on smaller proteins (<100 amino acids) as well as non-globular or elongated proteins, representing 60–80% of proteins with no predicted pockets. However, for IF-SitePred larger globular proteins represent ≈50% of all proteins where this method fails to predict a pocket (Supplementary Figure 1). Predicted residue ligandability scores for P2Rank_{CONS}, P2Rank and IF-SitePred (which we derived in this work), were examined for the proteins with no predicted pockets. Figure 4 illustrates 8 examples of proteins where residues with IF-SitePred (Eq. 13) high ligandability scores cluster in space into clear binding sites that are not reported as predictions by this method. This suggests that IF-SitePred is too strict in selecting only those residues predicted as ligand-binding by all 40 models. The cloud point selection and clustering approach or threshold in this method may also play a role in this.

Table 3 summarises the ligand site characterisation analysis. fpocket predicts the most sites out of all the methods, with 57,859, followed by IF-SitePred (44,948), DeepPocket_{SEG} (21,718), VN-EGNN (13,582), P2Rank (12,412), P2Rank_{CONS} (10,180), Surfnet⁺ (9043), PocketFinder⁺ (8913), Ligsite⁺ (6903), GrASP (4694) and PURESNet, which predicts fewest sites (2621). LIGYSIS defines 6882 binding sites from experimental data.

Relative to LIGYSIS, the prediction methods have ratios of predicted/defined sites ranging from 8.4 (fpocket) to 0.4 (PURESNet) with P2Rank_{CONS} in the middle, predicting 1.5 pockets per observed reference site. IF-SitePred, DeepPocket_{SEG} as well as P2Rank and fpocket predict more pockets on larger protein chains, whereas the rest of methods do not (Supplementary Note 1 and Supplementary Figure 2). This effect is most clear with fpocket, which predicts 350 pockets for chain A of PDB: 7SUD [85], a structure of the DNA-dependent protein kinase catalytic subunit, DNPK1, (P78527) with 3736 amino acid residues. In contrast, VN-EGNN, which initially places $K=8$ virtual nodes, results in a maximum of 8 predicted pockets, regardless of protein chain size, and PURESNet predicts a single pocket in 90% of the protein chains.

Figure 5 represents how the eleven unique sets of ligand site predictions compare to each other as well to LIGYSIS, which *defines* ligand sites from experimentally determined biologically relevant protein–ligand complexes. There are eleven unique sets of predictions since DeepPocket_{RESC} and fpocket_{PRANK} do not predict, but re-score and re-rank fpocket predictions. However, DeepPocket_{SEG} predictions are different as new pocket shapes are extracted by the CNN segmentation module. Figure 5A shows how VN-EGNN, GrASP and PURESNet differ from the other methods with a maximum of 7, 12 and 4 predicted pockets, respectively. PocketFinder⁺, Ligsite⁺ and Surfnet⁺ present narrow distributions like LIGYSIS and with medians of 1–3 pockets per protein. P2Rank_{CONS} and P2Rank also present a median of 3 pockets per protein but display wider distributions as they can predict up to 60, and 80 pockets per protein, respectively. Overall, P2Rank_{CONS} predicts fewer pockets than P2Rank. DeepPocket_{SEG}, fpocket, and IF-SitePred follow with a median of 6, 17, and 20 pockets. The difference in number of pockets between DeepPocket_{SEG} and DeepPocket_{RESC} or fpocket is because 60% of fpocket candidates are not extracted by the CNN segmentation module implemented in DeepPocket.

Figure 5B shows the distribution of pocket radius of gyration, R_g . VN-EGNN and IF-SitePred differ from the rest of methods with narrow distributions and medians around 6 Å. These two methods do not report pocket residues. Instead, they were obtained using a distance threshold of 6 Å from the centroid, for VN-EGNN, and cloud points, for IF-SitePred. This is reflected by examining the percentage of pockets with $R_g > 10$ Å which is 0% and 0.1% for VN-EGNN and IF-SitePred. This is a striking difference compared to LIGYSIS and the other methods: 1.8% (fpocket), 4.8% (P2Rank), 5.7% (DeepPocket_{SEG}), 6.4% (GrASP), 6.5% (P2Rank_{CONS}), 11.6% (PURESNet), 12.6% (LIGYSIS), 16% (Surfnet⁺), 21.4% (PocketFinder⁺) and 33.5% (Ligsite⁺). The latter three predict the sites

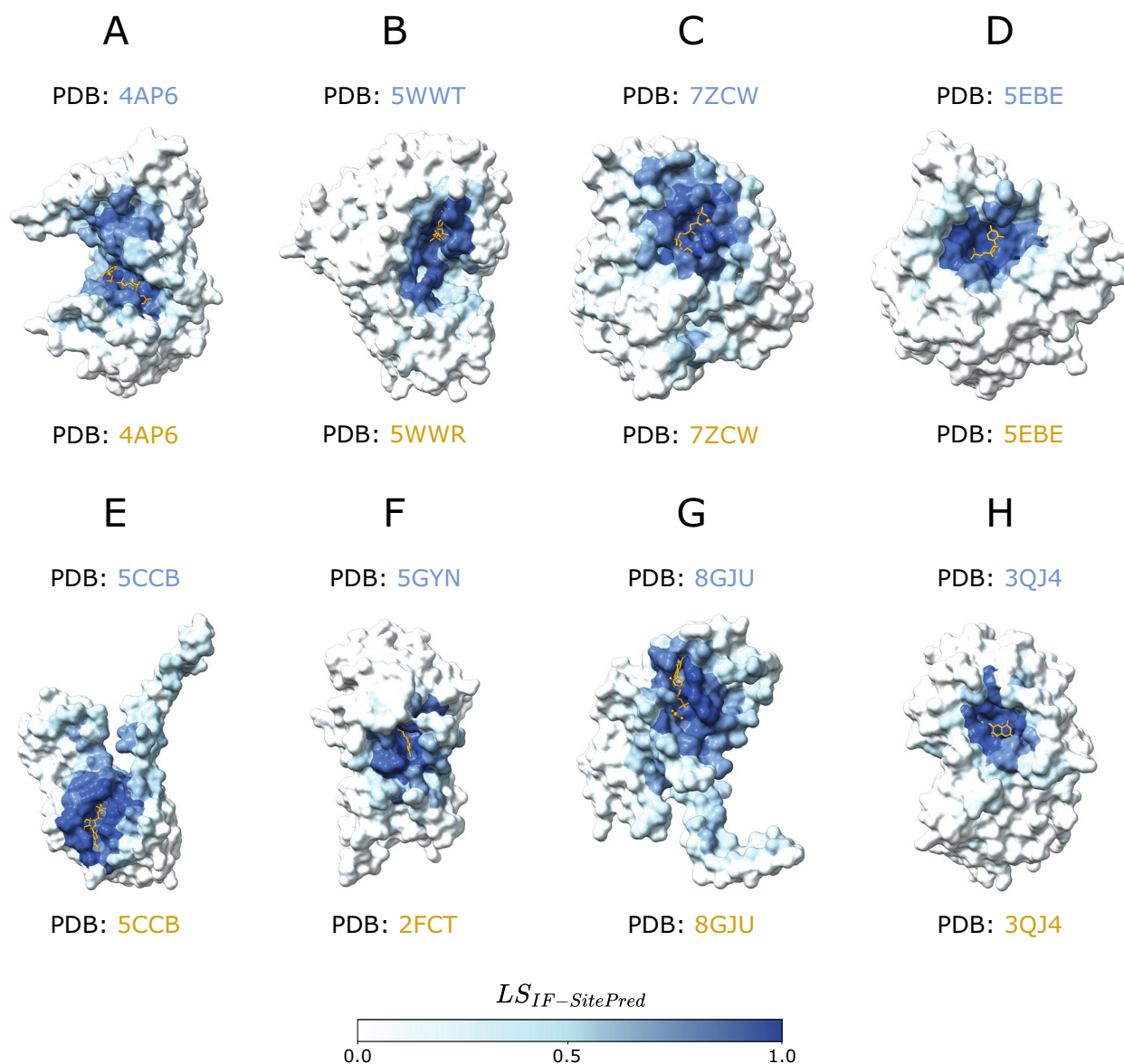


Fig. 4 IF-SitePred "missed" predictions. Eight examples of human protein chains where IF-SitePred does not report any predicted ligand binding sites. Predictions are made on ligand-stripped chains. Ligand molecules, in orange, are superposed to illustrate how the ligandability scores recapitulate the observed binding site. These are protein representative chains and ligand molecules might not be observed in the same entry. **A** GDP-fucose protein *O*-fucosyltransferase 2, Q9Y2G5, with GDP-L-fucose (GFB) superimposed (PDB: 4AP6) [76]; **B** tRNA (cytosine(72)-C (5))-methyltransferase NSUN6, Q8TEA1, (PDB: 5WWT) [77] with superposed sinefungin (SFG) (PDB: 5WWT); **C** tubulin beta-2B chain, Q9BVA1, with phosphomethylphosphonic acid guanylate ester (G2P) (PDB: 7ZCW) [78]; **D** cyclic GMP-AMP phosphodiesterase SMPDL3A, Q92484, with cytidine-5'-monophosphate (CSP) (PDB: 5EBE) [79]; **E** tRNA (adenine(58)-N(1))-methyltransferase catalytic subunit, Q96FX7, with SAH (PDB: 5CCB) [80]; **F** chronophin, Q96GD0, (PDB: 5GYN) [81] with PLP (PDB: 2FCT) [82]. **G** Mitochondrial methylmalonic aciduria type A protein, Q8IVH4, with GDP (PDB: 8GJU) [83]; **H** renalase, Q5VYX0, (PDB: 3QJ4) with FAD [84]. Residues are coloured based on the ligandability score calculated by averaging the probabilities predicted by each of the 40 IF-SitePred prediction models. This is a score (Eq. 13) ranging 0–1 which is indicative of the likelihood of a given residue binding a ligand (see "Methods"). Clear pockets can be observed formed by residues with high ligandability scores (darker blue colour), which agree with the sites where ligands bind. These "missed" pockets contribute to the lower recall of IF-SitePred and strongly suggest that a more permissive ligand-binding residue selection threshold (currently *all* 40 models) or a different clustering scheme might be able to capture these predictions and increase the recall of this method

Table 3 Ligand site characterisation summary

Method	Coverage (%)	# Total pockets	# Pockets per protein	R_g (Å)	MCD (Å)	MRO	Redundancy (%)
LIGYSIS (reference)	2775	6882	1, 1, 27	5.9	14.1	0	2.3
(d) VN-EGNN	2764 (99.6)	13,582 (×2.0)	1, 5, 7	5.9	1.1*	0.85*	66.8*
(d) IF-SitePred	2075 (74.8*)	44,948 (×6.5)	1, 20* , 129	5.9	3.4	0.55	49.5
(d) GrASP	2771 (99.9)	4694 (×0.7)	1, 1, 12	7.9	21.4	0	0.0
(d) PURESnet	2360 (85.1)	2621 (×0.4)	1, 1, 4	8.1	27.0	0	0.0
(d) DeepPocket _{SEG}	2349 (84.7)	21,718 (×3.2)	1, 6, 196	7.7	4.6	0.4	31.1
(d) P2Rank _{CONS}	2759 (92.9)	12,412(×1.8)	1, 3, 57	7.1	13.9	0.05	0.7
(d) P2Rank	2402 (86.6)	10,180 (×1.5)	1, 3, 85	7.1	13.8	0.05	0.6
(d) fpocket	2759 (99.4)	57,859 (×8.4*)	1, 17, 349*	6.3	9.7	0.15	0.7
(d) PocketFinder ⁺	2775 (100)	8913 (×1.3)	1, 3, 23	8.6	18.7	0.05	0.0
(d) Ligsite ⁺	2775 (100)	6903 (×1.0)	1, 2, 12	9.1*	16.7	0.09	0.0
(d) Surfnet ⁺	2775 (100)	9043 (×1.3)	1, 3, 40	8.4	17.2	0.07	0.0

LIGYSIS is not a ligand binding site predictor, but a reference dataset curated from experimentally determined structures of biologically relevant protein–ligand complexes. These predictions result from the default prediction of the methods, indicated by **(d)** preceding method names. Coverage represents the number of protein chains for which the different methods return at least one prediction. Percentage is relative to LIGYSIS protein chains. VN-EGNN failed with an error for PDB: 6BCU chain: A [86]. The rest of the methods ran successfully for all protein chains; # total pockets and ratio of predicted pockets by reference pockets in parenthesis, e.g., for each LIGYSIS site, fpocket predicts on average 8.4 pockets; minimum, median and maximum number of pockets per protein; Median pocket radius of gyration R_g (Å); median minimum centroid distance (MCD) (Å) for all pockets. For proteins where multiple pockets are predicted, MCD represents the distance to the closest pocket centroid for each of the different predicted pockets within a protein. This is a measure of how close predicted pockets are to each other; maximum residue overlap (MRO). For a given pocket, MRO is the maximum residue overlap with other pockets' residues within a protein. MRO is a measure of how similar, in terms of shared residues, the predicted pockets are (see "Methods" for detailed explanation). For example, the median overlap between VN-EGNN predicted pockets is 85%. Redundancy represents the percentage of predicted pockets that are redundant, i.e., the closest pocket centroid is within 5 Å, or overlap is at least 3/4 ($\geq 75\%$) residues. This is the case for 67% of VN-EGNN pockets and 0% for GrASP pockets. Bold font and "*" indicate the most extreme values within each column

with largest median $R_g \approx 9$ Å. VN-EGNN, GrASP, PURESnet and DeepPocket_{SEG} predict sites with $R_g = 0$ Å. This is rather infrequent (7.8% GrASP) and <3% for the other three. These examples correspond to *singletons*, i.e., pockets formed by only one amino acid.

Figure 5C illustrates how close predicted sites are to each other within a protein chain. Pairwise distances between the centroids of all ligand site pairs for a protein are calculated, and for each site, the minimum distance is taken. Sites predicted by VN-EGNN, IF-SitePred and DeepPocket_{SEG} are very close to each other, with median distances (\tilde{d}) of 1.1, 3.4 and 4.6 Å, respectively. fpocket follows with $\tilde{d} = 9.7$ Å. The rest of the methods, and LIGYSIS, present median distances ranging between 13 and 18 Å (LIGYSIS, P2Rank_{CONS}, P2Rank, PocketFinder⁺, Ligsite⁺, Surfnet⁺), and finally GrASP ($\tilde{d} = 21.7$ Å) and

PURESnet ($\tilde{d} = 27$ Å). Both versions of P2Rank present the most similar distribution to what is observed on LIGYSIS.

Figure 5D depicts the overlap existing between residues that form the predicted pockets within a protein. All pairwise overlaps, i.e., Jaccard Index, are calculated between pockets in a chain, and for each pocket, the maximum is taken. This is a measure of how much predicted pockets overlap. This is directly related to how close pockets are, and so VN-EGNN, IF-SitePred and DeepPocket_{SEG} present very high overlaps $\tilde{o} = 0.85$, $\tilde{o} = 0.55$ and $\tilde{o} = 0.40$, respectively. fpocket follows with $\tilde{o} = 0.15$, Ligsite⁺ ($\tilde{o} = 0.09$), Surfnet⁺ ($\tilde{o} = 0.07$), P2Rank_{CONS}, P2Rank and PocketFinder⁺ ($\tilde{o} = 0.05$), and finally LIGYSIS, GrASP and PURESnet with $\tilde{o} = 0.0$. GrASP is the only out of the thirteen methods presented here that clusters atoms directly,

(See figure on next page.)

Fig. 5 Binding pocket characterisation. Ligand binding site predictions by eleven methods with default parameters **(d)** are compared to the LIGYSIS reference dataset. A black dashed line separates LIGYSIS from the eleven ligand binding site predictors. Predictions by DeepPocket_{RES} and fpocket_{PRANK} are not included in this analysis as they are re-scored and re-ranked fpocket predictions and their plotting would be redundant. Data points farther than four standard deviations (SD) from the mean are considered outliers. The limit of the Y axis is the maximum *non-outlier* value plus a buffer value. This way, only the most extreme outliers are hidden, which maximises visual interpretation of the data whilst minimising the number of data points not shown. Within the violin plots are box plots representing the underlying distribution. Line represents the median, box contains the interquartile range (IQR), and whiskers extend to $1.5 \times \text{IQR}$. **A** Number of pockets per protein; **B** pocket radius of gyration R_g (Å); **C** minimum inter-pocket centroid distance (MCD) (Å). This is a measure of how close predicted pockets are to each other within a protein; **D** Maximum inter-pocket residue overlap (MRO). Residue overlap was calculated as Jaccard Index. This is a measure of how much pockets overlap in terms of binding residues

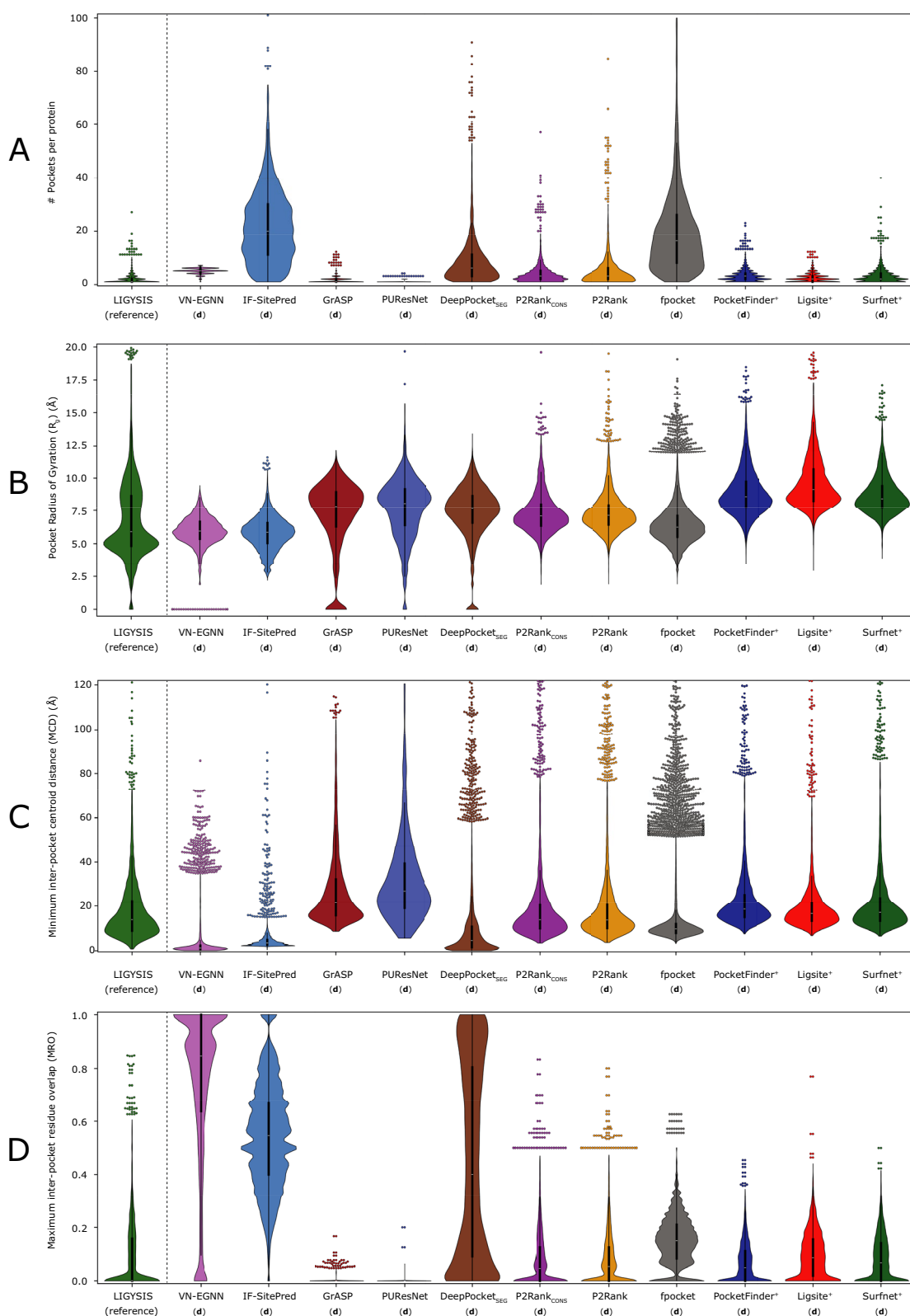


Fig. 5 (See legend on previous page.)

and as a result, overlap between pockets is minimised. Other methods cluster cloud points (IF-SitePred), SAS points (P2Ranks) voxels (PUNet, DeepPocket_{SEG}), alpha spheres (fpocket), or grid points (PocketFinder⁺, Ligsite⁺, Surfnet⁺) but not residues, resulting consequently in higher overlapping.

Proximity in space between predicted sites as well as residue overlap are indicators of redundant ligand binding site prediction, i.e., duplicate predictions of a unique observed ligand site. This is the case for VN-EGNN, IF-SitePred and DeepPocket_{SEG}. This phenomenon can negatively impact the precision and recall of these methods. Accordingly, correcting for redundancy should have a significant impact on the performance of these methods. In contrast, GrASP and PUNet which predict a small number of pockets show low proximity and overlap of predicted sites and so redundancy is not an issue.

Evaluation of predictive performance

Pocket level evaluation

The ideal ligand binding site predictor would have a high precision, i.e., most of the predictions it makes are correct, whilst maintaining a high recall, i.e., recapitulating most of the observed sites. Moreover, the ideal predictor returns predictions that are non-redundant, i.e., it does not predict the same pocket multiple times. Additionally, pockets are ranked in a systematic manner according to a strong and meaningful pocket scoring scheme which captures well the nature of existing ligand binding sites and therefore ranks the predicted pockets from more likely (high score, top) to least likely (low score, bottom). A good predictor would also perform well at the residue level. This means it is able of capturing the likelihood for a residue to bind a ligand. This can be done by means of a residue ligandability score, which additionally might highlight key residues, the more “ligandable” within a binding site. Ligand site prediction methods were benchmarked with these criteria in mind. The following results concern the original version of the predictors. The non-redundant and re-scored variants of the methods will be discussed in a later section.

Figure 6A illustrates the recall curve for top- $N+2$ pockets for each method, where N is the number of observed sites for a target protein. Reported recall is obtained using $DCC=12 \text{ \AA}$. Re-scored fpocket predictions by PRANK (fpocket_{PRANK}) and DeepPocket (DeepPocket_{RESC}) yield the highest recall with 60.4% and 58.1%, respectively, closely followed by, P2Rank_{CONS} (53.9%) and P2Rank (51.9%). The rest of the methods present recall <50% with PUNet, VN-EGNN and IF-SitePred with lowest recall of 41.1%, 40.9% and 25.7%, respectively (Table 4). Figure 6B, shows the recall curve considering different top- $N+X$ predictions. Most methods reach a plateau

by top- $N+5$, as they do not predict that many pockets. However, methods that predict more pockets per protein, such as IF-SitePred, or fpocket, fpocket_{PRANK}, DeepPocket_{SEG} and DeepPocket_{RESC}, which take fpocket predictions as a base, increase their recall as more predictions are considered. fpocket, fpocket_{PRANK} and DeepPocket_{RESC} reach a maximum recall of $\approx 90\%$ when *all* predictions are considered, regardless of rank. Most other methods present a maximum recall of $\approx 50\text{--}60\%$ and PUNet presents the lowest maximum recall of 41%. Figure 6C illustrates the top- $N+2$ recall curve if residue overlap, I_{rel} , was used instead of DCC as a criterion. In this case, Ligsite⁺, PocketFinder⁺, and Surfnet⁺ come on top with recall $\approx 45\%$ at $I_{rel} \geq 0.5$. This is explained by their prediction of massive cavities, that while often fully contain or overlap with the observed pocket, do not meet the DCC criterion, as their centroids are farther than 12 \AA from the observed site.

Figure 6D represents the cumulative number of TP against FP when predictions across the proteins in the reference dataset are sorted by score. This shows how effective the scoring scheme of each method is in ranking their predictions to reflect the nature of ligand binding sites. At 100 FP the #TP fall into three different blocks and one outlier: Ligsite⁺, fpocket, Surfnet⁺ and PocketFinder⁺ are at the bottom with # TP $\in (60, 120)$. Secondly, fpocket_{PRANK}, DeepPocket_{SEG}, DeepPocket_{RESC}, P2Rank, and PUNet follow with # TP $\in (530, 670)$. Re-scoring fpocket predictions with PRANK or DeepPocket results in up to +500 TP. GrASP, IF-SitePred and P2Rank_{CONS}, present a high # TP ranging 900–1000 at 100 FP. Finally, VN-EGNN sits at the top with 1301 TP. However, this number is not representative, as the # TP are inflated due to the redundancy in the predictions of VN-EGNN. This is the same for IF-SitePred and DeepPocket_{SEG}. Redundant correct predictions of the same pocket will count as multiple TPs, whereas they should only count as 1 TP. Newer methods, e.g., GrASP, P2Rank_{CONS}, and VN-EGNN, and IF-SitePred, despite redundancy in their prediction for the latter two, are better at ranking their predicted pockets, presenting up to 1000 more TP for 100 FP than earlier methods. This means their scoring schemes are significantly better at capturing the essence of a ligand binding site. Including evolutionary conservation in P2Rank (P2Rank_{CONS}) predictions results in an increase of +346 TP relative to default P2Rank, indicating that the fewer predicted pockets, and their scores are a more faithful representation of the observed LIGYSIS dataset.

Figure 6E provides insight into the precision of the methods by examining how this metric changes as more predictions are considered. In the same manner as for Fig. 7D, predictions across proteins in the

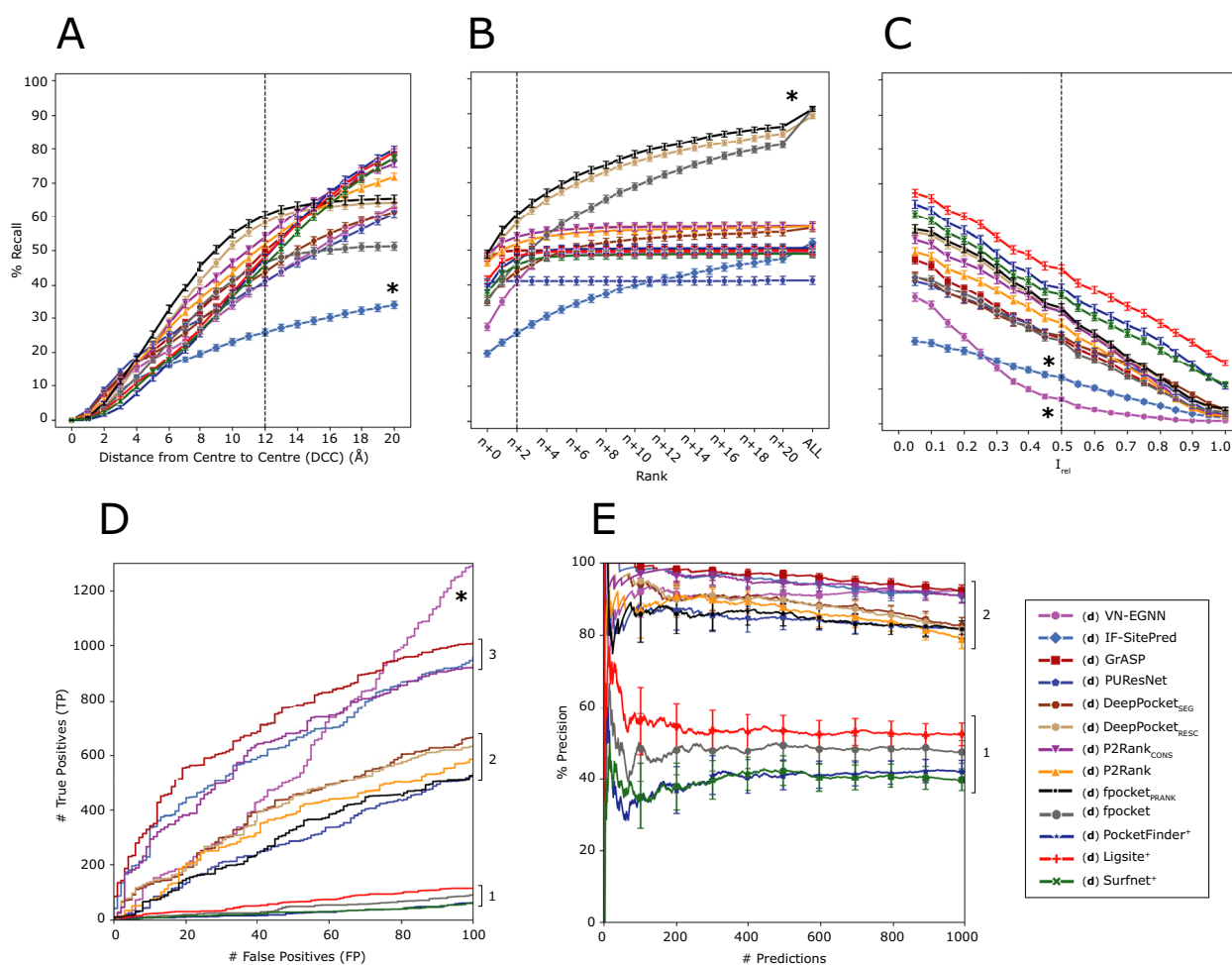


Fig. 6 Ligand binding site prediction benchmark at the pocket level. These curves correspond to the default predictions of the thirteen methods, indicated by (d) preceding their names. **A** Recall, percentage of observed sites that are correctly predicted by a method within the top- $N+2$ predictions; **B** recall using DCC = 12 Å but considering increasing rank thresholds, i.e., top- N , $N+1$, $N+2$, etc. ALL represents the maximum recall of a method, obtained by considering all predictions, regardless of their rank or score; **C** recall curve of top- $N+2$ predictions using I_{rel} as a criterion; **D** ROC100 curve (cumulative TP against cumulative FP until 100 FP are reached); **E** precision curve for the top-1000 predictions of each method across the LIGYSIS dataset, $Precision_{1K}$. Error bars represent 95% CI of the recall (A–C) and precision (E), which is 100× proportion. Numbers at the right of the panels indicate groups or blocks of methods that perform similarly for each metric. Asterisks (*) indicate outlier methods, or methods that perform very differently than the rest

LIGYSIS dataset are sorted and cumulative precision is plotted for the top-1000 scoring predictions. Methods group into two clear blocks. Newer (machine learning-based) methods VN-EGNN, GrASP, IF-SitePred, P2Rank_{CONS}, DeepPocket_{SEG}, fpocket_{PRANK}, DeepPocket_{RESC} and PURESNet present a $Precision_{1K}$ of 80–95%. Earlier (geometry/energy-based) methods Ligsite⁺, fpocket, PocketFinder⁺ and Surfnet⁺ present lower $Precision_{1K}$ of 40–50%. DeepPocket_{RESC} and fpocket_{PRANK} take fpocket (geometry-based) predictions as a starting point and achieve much higher # TP_{100FP} (+500) as well as $Precision_{1K}$ (+30%). This is further evidence that performance can be boosted with

a solid scoring scheme and agrees with previous studies [37, 38, 40, 87].

Table 4 summarises these results and shows the mean relative residue overlap (RRO) and relative volume overlap (RVO) which measure how well predicted sites align with observed ones in shape. VN-EGNN and IF-SitePred present the smallest RRO and RVO, but it is important to note that these methods do not report pocket residues and so residues were taken within 6 Å of their centroid, or pocket spheres. PocketFinder⁺, Ligsite⁺, and Surfnet⁺ present the highest RRO and RVO. This is a consequence of the large size of their predicted cavities, that rather than overlap with the observed site, fully contain and

Table 4 Pocket level evaluation summary table

Method	% Recall _{top-N}	% Recall _{top-N+2}	% Recall _{max}	% Precision _{1K}	# TP _{100 FP}	% RRO	% RVO
(d) VN-EGNN	27.5 (#11)	40.9 (#12)	49.3 (#10)	92.5⁺ (#1)	1301⁺ (#1)	32.8⁻ (#12)	27.6⁻ (#11)
(d) IF-SitePred	19.8⁻ (#12)	25.7⁻ (#13)	52.1 (#6)	91.0 (#2)	961 (#3)	46.5 (#11)	40.4 (#9)
(d) GrASP	48.0 (#2)	49.9 (#5)	50.0 (#8)	92.5⁺ (#1)	1017 (#2)	54.5 (#7)	59.8 (#6)
(d) PURESNet	40.6 (#6)	41.1 (#11)	41.1⁻ (#12)	81.6 (#6)	534 (#8)	61.0 (#4)	63.9 (#4)
(d) DeepPocket _{SEG}	35.4 (#10)	43.8 (#10)	56.5 (#5)	82.6 (#4)	670 (#5)	57.5 (#5)	60.3 (#5)
(d) DeepPocket _{RESC}	46.6 (#4)	58.1 (#2)	89.3 (#2)	81.7 (#5)	637 (#6)	53.1 (#9)	38.2 (#10)
(d) P2Rank _{CONS}	48.8⁺ (#1)	53.9 (#3)	57.0 (#4)	90.7 (#3)	932 (#4)	56.4 (#6)	43.8 (#8)
(d) P2Rank	46.7 (#3)	51.9 (#4)	57.0 (#3)	79.2 (#7)	586 (#7)	54.4 (#8)	58.2 (#7)
(d) fpocket _{PRANK}	48.8⁺ (#1)	60.4⁺ (#1)	91.3⁺ (#1)	81.7 (#5)	526 (#9)	52.6 (#10)	38.2 (#10)
(d) fpocket	38.8 (#8)	46.5 (#8)	91.3⁺ (#1)	47.3 (#9)	94 (#11)	52.6 (#10)	38.2 (#10)
(d) PocketFinder ⁺	39.2 (#7)	47.8 (#7)	50.5 (#7)	42.0 (#10)	64 (#12)	72.3 (#2)	75.9 (#2)
(d) Ligsite ⁺	41.3 (#5)	48.4 (#6)	49.7 (#9)	52.3 (#8)	115 (#10)	77.6⁺ (#1)	77.0⁺ (#1)
(d) Surfnet ⁺	37.7 (#9)	45.8 (#9)	48.9 (#11)	39.5⁻ (#11)	61⁻ (#13)	71.7 (#3)	72.0 (#3)

These metrics correspond to the default modes of the thirteen methods covered in this work, indicated by (d) preceding the methods' names. % Recall for each method considering top-*N*, *N*+2 and *all* predictions (max) without taking rank into consideration, i.e., maximum recall. Precision of the method for the top-1000 scored predictions. Number of TP reached for the first 100 FP (# TP_{100FP}). Mean % relative residue overlap (RRO) for those sites correctly predicted and % relative volume overlap (RVO) only for correctly predicted sites that have a volume, i.e., are pockets or cavities, and not exposed sites, which don't have a volume. These last two metrics represent the overlap in residues and volume relative to the observed site. See "Methods" section for definitions of RRO and RVO. Within each cell, the numbers following a dash (#) indicate the rank of each method according to the metric in the column. Bold font indicates the best ("⁺") and worst ("⁻") performing methods for each metric

are much larger than it. This might not be convenient in the context of pocket finding for drug discovery where more clearly defined drug-like sites might be of interest. GrASP, PURESNet and DeepPocket_{SEG} present high values of RRO ≈ 60% and RVO ≈ 60% whilst presenting a size distribution more like LIGYSIS (Fig. 5B) and provide the best representation of the observed sites in terms of shape similarity.

Residue level evaluation

Ligand binding site prediction tools can also be evaluated at the residue level. F1 Score and Matthews correlation coefficient (MCC) are utilised to do so. For each protein chain, F1 and MCC are calculated, distributions graphed and means reported (Table 5). Binary labels are employed to calculate these scores: 1 if the residue is found in a pocket, and 0 otherwise and compared to the ground truth, i.e., whether a residue binds a ligand in the LIGYSIS set. For VN-EGNN, IF-SitePred, PocketFinder⁺, Ligsite⁺, and Surfnet⁺, which do not report pocket

residues (Table 1), pocket residues were obtained by taking those residues within 6 Å of the pocket centroid, cloud points or grid points, respectively. DeepPocket_{RESC} is not considered for this analysis as the predictions are re-scored and re-ranked fpocket predictions.

Figure 7A, B illustrate the distributions of the F1 score and MCC for each method on the 2775 protein chains of the final LIGYSIS set. Both metrics agree that PURESNet (F1=0.41, MCC=0.39), GrASP (F1=0.39, MCC=0.33) and P2Rank_{CONS} (F1=0.36, MCC=0.30) are the top-3 performing methods in this task of binary classification into pocket (1) and non-pocket residues (0). fpocket presents the lowest F1=0.23 and MCC=0.12 since it predicts many unobserved pockets (residues) that will count as FP here.

IF-SitePred does not report a residue ligandability score beyond a binary label (0, 1). Nevertheless, in this work, a score was computed by utilising the scores returned by the 40 prediction models of IF-SitePred. These scores range 0–1 and can be averaged as probabilities (Eq. 13).

(See figure on next page.)

Fig. 7 Ligand binding site prediction benchmark at the residue level. DeepPocket_{RESC} predictions are not included in the F1 and MCC analyses as these are re-scored and re-ranked fpocket predictions and the results would be the same as fpocket's. **A** F1 score distributions; **B** MCC distributions. In both panels, each data point in the distribution corresponds to the score obtained from all residues in a protein chain; **C** mean ROC curve for methods that report a residue score. Dashed line represents the baseline, 1 FP for each TP, i.e., diagonal and AUC=0.5; **D** mean PR curve. Dashed line represents the baseline, i.e., proportion of observed binding residues=0.1; Mean ROC and PR curves are calculated by averaging the curves of the 2775 LIGYSIS protein chains. **E** Distribution of residue ligandability scores for IF-SitePred, GrASP and P2Ranks. PocketFinder⁺, Ligsite⁺ and Surfnet⁺ are not included as their scores do not range 0–1, and a small number of scores can reach values > 25. (d) indicates that these predictions originate from the original methods run with default parameters

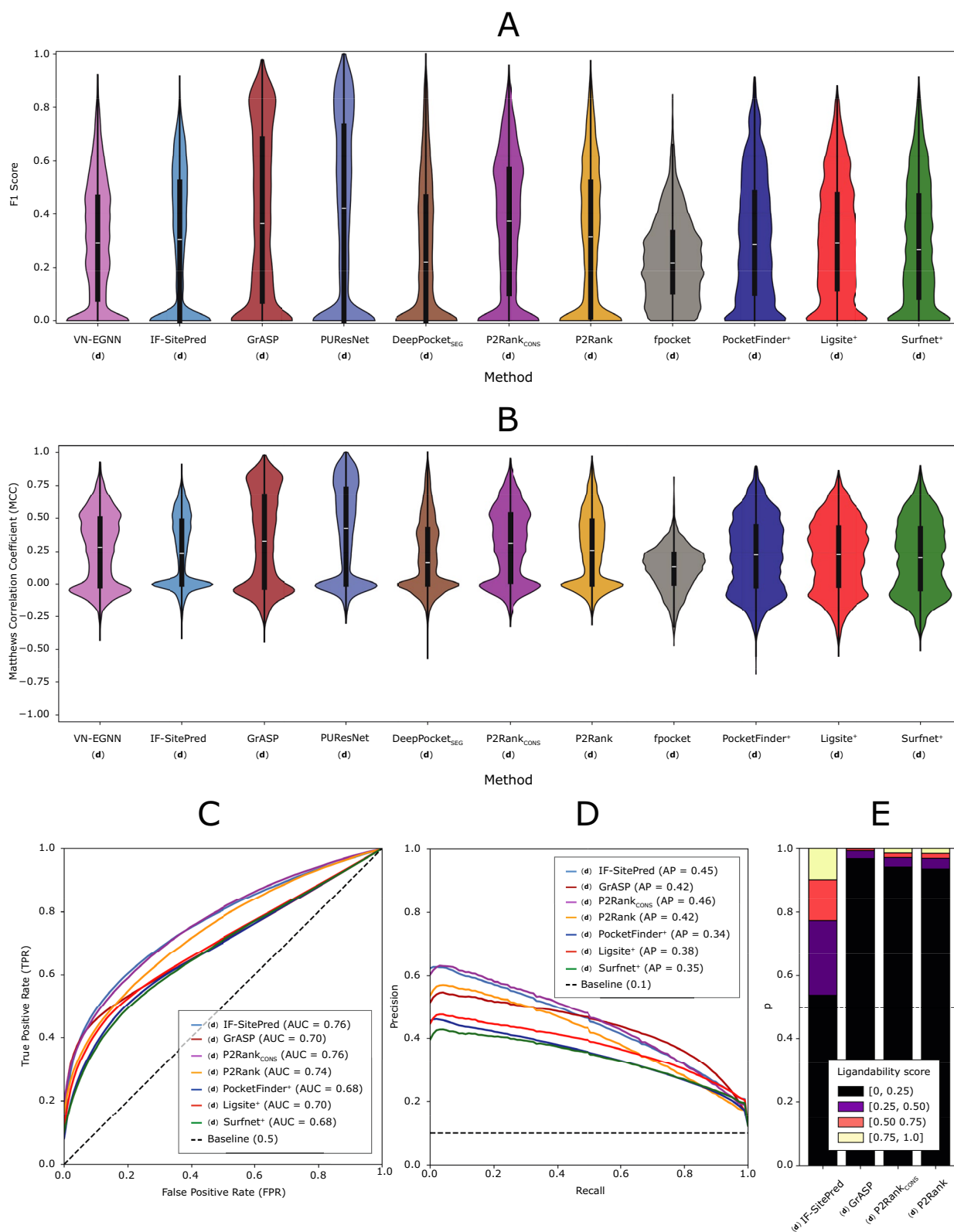


Fig. 7 (See legend on previous page.)

Table 5 Residue level evaluation summary table

Method	F1	MCC	AUC	AP
(d) VN-EGNN	0.29 (#8)	0.26 (#4)	–	–
(d) IF-SitePred	0.29 (#9)	0.24 (#6)	0.76 (#2)	0.45 (#2)
(d) GrASP	0.39 (#2)	0.34 (#2)	0.70 (#4)	0.42 (#3)
(d) PURESNet	0.41⁺ (#1)	0.39⁺ (#1)	–	–
(d) DeepPocket _{SEG}	0.27 (#10)	0.21 (#9)	–	–
(d) P2Rank _{CONS}	0.36 (#3)	0.30 (#3)	0.76⁺ (#1)	0.46⁺ (#1)
(d) P2Rank	0.31 (#4)	0.26 (#5)	0.74 (#3)	0.42 (#3)
(d) fpocket	0.23⁻ (#11)	0.12⁻ (#11)	–	–
(d) PocketFinder ⁺	0.31 (#5)	0.22 (#7)	0.68 (#6)	0.34⁻ (#6)
(d) Ligsite ⁺	0.31 (#6)	0.21 (#8)	0.70 (#5)	0.38 (#4)
(d) Surfnet ⁺	0.29 (#7)	0.20 (#10)	0.68⁻ (#7)	0.35 (#5)

These results come from default predictions, indicated by (d). DeepPocket_{RESC} is not considered in this analysis as their predictions are fpocket's but re-scored and re-ranked. Ligand binding site prediction benchmark at the residue level was calculated from 2775 protein chains in the LIGYSIS dataset. Mean F1 score, mean Matthews correlation coefficient (MCC), mean ROC area under the curve (AUC) and mean precision recall (PR) curve average precision (AP). Within each cell, the numbers following a dash (#) indicate the rank of each method according to the metric in the column. Bold font indicates the best ("⁺") and worst ("⁻") performing methods for each metric. Pocket binary labels (0, 1) were employed for the calculation of F1 and MCC and obtained from predicted pockets. Residue ligandability scores were employed to calculate ROC/AUC and PR/AP. Reported AUC and AP are means resulting from the average across the 2775 LIGYSIS protein chains. This was not possible for VN-EGNN, PURESNet, DeepPocket_{SEG} and fpocket as these methods do not provide such scores, indicated by a dash (–)

This will now be referenced as IF-SitePred ligandability score. For IF-SitePred, GrASP, P2Rank_{CONS}, P2Rank, PocketFinder⁺, Ligsite⁺, and Surfnet⁺, which report a residue level score (beyond a binary label), ROC and PR curves were plotted (Fig. 7C) and mean area under the curve (AUC) and mean average precision (AP) reported. This was not possible for VN-EGNN, PURESNet, DeepPocket_{SEG}, DeepPocket_{RESC} and fpocket as they do not report residue ligandability scores. See Supplementary Figure 10 for details on the variation in ROC and AUC across LIGYSIS protein chains. P2Rank_{CONS} and IF-SitePred, with the ligandability score calculated in this work (Eq. 13), present the highest mean AUC = 0.76, closely followed by P2Rank (AUC = 0.74). Surfnet⁺ presents the lowest mean AUC = 0.68. Figure 7D shows the mean PR curves, which agree with ROC AUC and highlight P2Rank_{CONS} as the method with the highest average precision = 0.46, followed by IF-SitePred (with Eq. 13

scoring) (AP = 0.45) and PocketFinder⁺ the lowest with (AP = 0.34). See Supplementary Figure 11 for details on the variation in PR and AP across LIGYSIS protein chains. Figure 7E shows IF-SitePred presenting a different residue ligandability score distribution to GrASP, P2Rank_{CONS}, and P2Rank. The IF-SitePred ligandability score, resulting from averaging the scores from the 40 IF-SitePred models, is the most "generous" with ≈ 20% of the residues presenting a score > 0.5, in contrast with GrASP which residue scoring is very strict, $P(LS \geq 0.5) = 0.75\%$, and P2Ranks (≈ 3%). This difference, combined with the ROC and PR curves further supports the use of the IF-SitePred ligandability score proposed in this work to define the predicted binding sites in this method. It also suggests that GrASP might benefit from a less strict residue level scoring scheme. PocketFinder⁺, Ligsite⁺, and Surfnet⁺ were not included in this analysis as their scores do not range 0–1 and very high scores (> 25) can be obtained.

Improving current methods by redundancy removal and pocket re-scoring

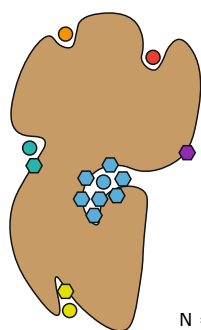
We define pocket prediction redundancy as the prediction of pockets with centroids very close in space ($d \leq 5 \text{ \AA}$) or with overlapping residues ($JI \geq 0.75$). This indicates multiple predictions of the same potential ligand binding site. Most ligand site prediction tools predict not only the location of the pocket by means of a centroid or pocket residues, but also a pocket confidence, and an associated rank among all the predicted pockets. Ligand site predictors tend to be evaluated by considering the top- N , or top- $N+2$ ranking pockets, where N is the number of observed sites for a given protein. Figure 8A, B illustrates how the redundant prediction of pockets can lead to an overestimation of precision and underestimation of recall (top- $N+2$). Figure 8C showcases human creatine kinase S-type, mitochondrial (PDB: 4Z9M) as an example of this phenomenon, where VN-EGNN and IF-SitePred redundantly predict the same pocket 7 and 33 times, whereas PURESNet returns a single prediction. All three methods correctly predict the site, just the difference is in the number of returned predictions.

Figure 5D shows that prediction redundancy is an issue particularly for VN-EGNN, IF-SitePred, and to a lesser extent DeepPocket_{SEG}. To assess the effect that

(See figure on next page.)

Fig. 8 The issue of redundancy in ligand binding site prediction. **A** A set of predictions where 6/10 predictions are redundant, resulting in a low recall (1/5) and inflated precision (7/7) considering the top- $N+2$ predictions; **B** when redundancy is removed, only four predictions remain and recall increases to 3/5 (60%) and precision decreases to 3/4 (75%); **C** predictions by VN-EGNN, IF-SitePred and PURESNet, on chain D of PDB: 4Z9M [88], where ADP binds. For this ADP binding site, VN-EGNN reports 7 predictions, IF-SitePred 33, and PURESNet a single prediction. These three methods correctly predict this site, however, VN-EGNN and IF-SitePred report redundant pocket predictions, which centroids are very close in space ($\leq 5 \text{ \AA}$) and residues overlap considerably (≥ 0.75)

A



Redundant set of predictions

Pocket	Score	Rank	Considered
●	0.99	1	✓
●	0.98	2	✓
●	0.97	3	✓
●	0.96	4	✓
●	0.95	5	✓
●	0.94	6	✓
●	0.93	7	✓
●	0.84	8	✗
●	0.82	9	✗
●	0.77	10	✗

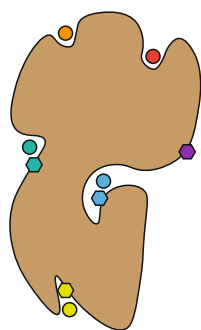
N = 5 observed pockets

Considering top-N+2 = 7 pockets

$$\text{Recall} = \frac{\text{Observed}}{\text{Predicted}} = \frac{5}{7} = 0.71$$

$$\text{Precision} = \frac{\text{Correct Predictions}}{\text{Total Predictions}} = \frac{7}{7} = 1$$

B



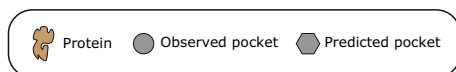
Non-redundant set of predictions

Pocket	Score	Rank	Considered
●	0.99	1	✓
●	0.84	2	✓
●	0.82	3	✓
●	0.77	4	✓

Considering top-N+2 = 7 pockets

$$\text{Recall} = \frac{\text{Observed}}{\text{Predicted}} = \frac{5}{7} = 0.71$$

$$\text{Precision} = \frac{\text{Correct Predictions}}{\text{Total Predictions}} = \frac{4}{5} = 0.8$$



C

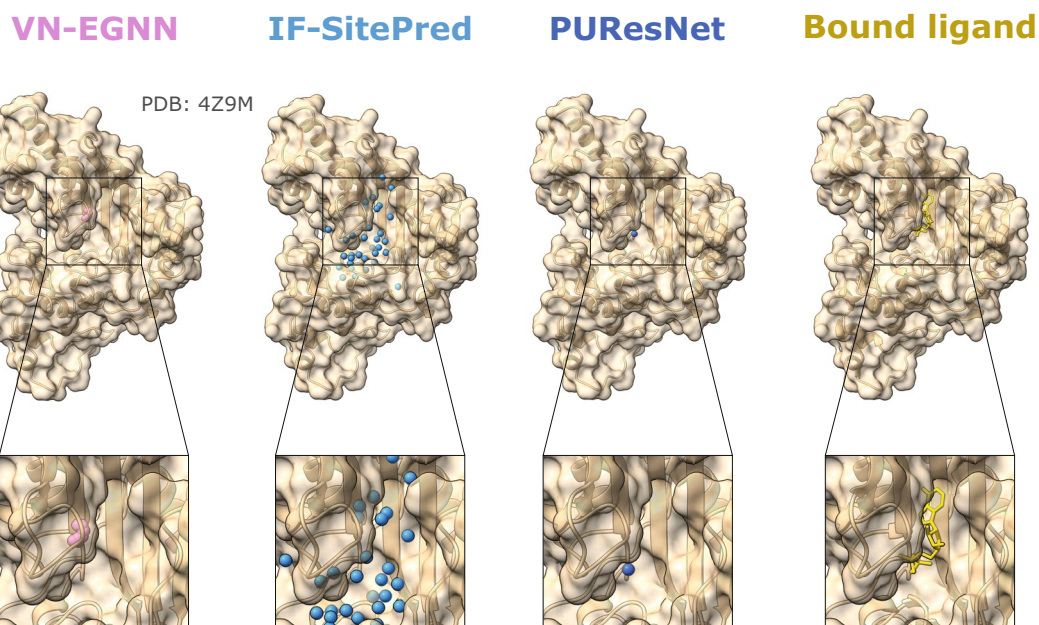


Fig. 8 (See legend on previous page.)

redundancy has on the performance of these methods, non-redundant subsets of predictions were obtained and labelled with the subscript “_{NR}”. Redundancy filtering was carried out for each method and adjacent pockets ($d \leq 5 \text{ \AA}$) or sharing more than $\frac{3}{4}$ of their residues ($JI \geq 0.75$) were dropped, keeping always the higher scoring pocket. Redundancy (%) was calculated as the proportion of redundant pockets relative to the original total number of pockets. VN-EGNN presents the highest percentage of redundant pockets with 9066/13,582 (67%) redundant pockets, followed by IF-SitePred with 22,232/44,948 (49%), and DeepPocket_{SEG} with 6744/21,718 (31%). For other methods, redundancy was minimal (<1%).

New ligand binding site prediction methods tend to score their pockets to rank them based on their likelihood of being a real binding site. A robust pocket scoring scheme is helpful to rank predicted pockets in a meaningful order, so pockets at the top are more likely to be real pockets. In other words, a pocket scoring scheme that captures the nature of ligand sites is key to maximising the benchmark performance and usability of these tools. IF-SitePred opts for a simple scheme, using the number of clustered cloud points as the pocket score. Other methods use more sophisticated approaches, such as taking the output of a machine learning model (VN-EGNN, DeepPocket), summing the model scores of atoms, residues or SAS points (GrASP, P2Rank_{CONS}, P2Rank) or using multiple pocket-level features (fpocket). Finally, methods like PURESNet, PocketFinder⁺, Ligsite⁺ or Surfnet⁺ do not report a pocket score.

As detailed in “**Methods**”, re-scored and re-ranked variants of predictions by IF-SitePred were generated. Additionally, pockets were scored and ranked for PURESNet, PocketFinder⁺, Ligsite⁺ and Surfnet⁺. Three different approaches were investigated: (1) sorting by number of pocket amino acids “_{AA}” (larger pockets ranking higher); (2) running PRANK to score pockets [37] “_{PRANK}”; (3) recalculating a score with the sum of square of cloud points (IF-SitePred_{RESC}) or grid points “_{SS}” for PocketFinder⁺, Ligsite⁺ or Surfnet⁺, in a similar manner as done by Kriv k et al. [37], and later Smith et al. [54]. Refer to Supplementary Note 2 and Supplementary Figure 3 for the effect that redundancy and different pocket scoring strategies have on pocket ranking.

In this section, we explore the effect that redundancy removal and pocket scoring/re-scoring have on % Recall (top- $N+2$), % Precision_{1K}, and # TP_{100FP}. Figure 9 illustrates how removing redundant predictions and re-ranking those remaining resulted in a significant ($p \leq 0.05$) increase in recall of +5.2% for VN-EGNN (Fig. 9A), +13.4% for IF-SitePred (Fig. 9B) and +5.6% for DeepPocket_{SEG} (Fig. 9C). Re-scoring IF-SitePred

predictions using Eq. 13 resulted in a further increase of +2.1%. After removing redundancy, the number of unique pockets by each method reduces drastically for VN-EGNN (−67%), IF-SitePred (−49%) and DeepPocket_{SEG} (−31%) (Fig. 9D–F), thus limiting the recall of these methods. Re-scoring and re-ranking of PURESNet (0.0%), PocketFinder⁺ (1.1%), Ligsite⁺ (0.6%), and Surfnet⁺ (1.6%) predictions did not have a significant effect ($p > 0.05$). For more details, see Supplementary Note 3, and Supplementary Figures 4, 5.

Supplementary Note 4 and Supplementary Figure 6 illustrate the powerful effect that removing redundancy, and a more robust pocket scoring scheme have when ranking predictions across multiple proteins, resulting in increases of up to +600 TP by 100 FP. Moreover, whilst re-scoring fpocket predictions also results in an increase in # TP, DeepPocket_{RESC}, DeepPocket_{SEG-NR} and fpocket_{PRANK}, all starting from fpocket candidates, achieve no more than ≈ 600 TP, indicating that the maximum possible increase in recall achieved by pocket re-scoring is limited by the starting predictions, which are the same for these three methods. Finally, despite benefiting from a simple scoring of their pockets by taking the number of amino acids, PocketFinder⁺, Ligsite⁺ and Surfnet⁺ present the lowest # TP_{100FP}. This is potentially because their purpose is the detection of cavities and not the prediction of ligand sites. Despite being similar tasks, they are not the same. This is reflected in the number of cavities they report and their size, which explain the low number of # TP_{100FP}. Moreover, these methods are not machine learning-based and were not trained to learn from and recapitulate observed ligand binding sites.

Redundancy removal for VN-EGNN, IF-SitePred and DeepPocket_{SEG} does not have a significant effect ($p \geq 0.05$) on the precision of these methods. However, (re-)scoring predictions of PURESNet, fpocket, PocketFinder⁺, Ligsite⁺, and Surfnet⁺ increased precision on the top-1000 predictions by +11.7%, +34.4%, 23.3%, 16.5% and 29.1%, respectively, further proving the importance of the employed pocket scoring scheme (Supplementary Note 5 and Supplementary Figure 7). Supplementary Figure 8 and Supplementary Table 2 summarise the results of this analysis.

Discussion

In this work we have carried out the most complete comparative analysis of ligand binding site prediction methods to date, spanning three decades of methods development. Firstly, predictions from the thirteen methods as well as observed sites from our reference dataset, LIGYSIS, were compared in terms of number of proteins methods predict on, number of predicted sites per protein, their size, distance and overlap between the sites.

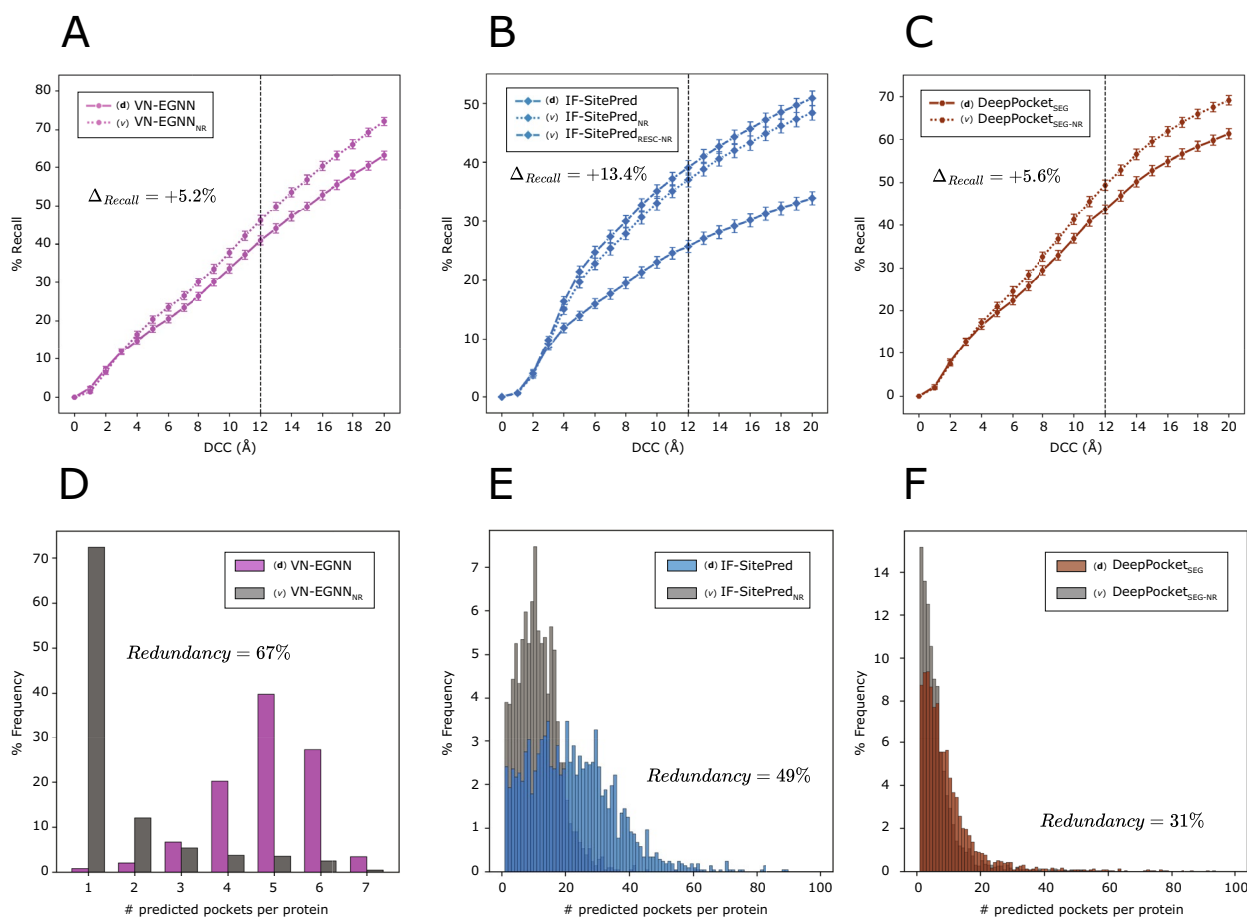


Fig. 9 Improvement on the current methods. Solid lines represent top- $N+2$ recall, i.e., proportion of observed ligand binding sites that are correctly identified by each default version for VN-EGNN (A), IF-SitePred (B) and DeepPocket_{SEG} (C). Dashed/dotted lines represent the non-redundant and re-scored (for IF-SitePred) variants of predictions by these methods generated in this work. Removing redundancy increases significantly the recall of VN-EGNN, IF-SitePred and DeepPocket_{SEG} by > 5%. Vertical black dashed lines indicate the threshold used as reference in this work: DCC = 12 Å. Comparison of number of predicted pockets per protein between standard and non-redundant predictions for VN-EGNN (D), IF-SitePred (E) and DeepPocket_{SEG} (F)

This analysis provides insight into how the different methods work and hints at potential limitations or room for improvement, e.g., the prediction of a fixed number of sites per protein, or considerable proximity and overlap between the predictions. Secondly, predictions from thirteen canonical ligand site prediction methods, as well as fifteen method variants first introduced in this work are objectively evaluated using the LIGYSIS set. This evaluation considers predictions at the residue level by F1 score, MCC, ROC/AUC and PR/AP, as well as at the pocket level by Recall considering top- N , $N+2$ and *all* predictions, Precision_{1K}, # TP_{100 FP}, RRO, RVO and redundancy. This is the first independent ligand site prediction benchmark for over 10 years, since Chen et al. [67] and the largest to date both in terms of reference dataset size (2775 chains), number of methods compared (13 original + 15 variants) and metrics employed (> 10).

We have shown how redundancy in prediction, i.e., predicting multiple times the same observed site, can underestimate the recall, and overestimate the precision of the methods, therefore providing a misleading assessment of the methods' performance. Redundancy removal and subsequent pocket re-ranking can yield a significant increase in recall. A robust pocket scoring scheme can have a strong impact in performance, both in the recall and precision of the methods and emphasis should be put into this area. Even if a single site is predicted per protein, a pocket score can be highly useful when ranking pockets across different proteins, e.g., when having a list of potential drug targets and deciding which protein might be best to target therapeutically.

Recall (% of sites that are correctly predicted) is more informative than precision (% of predictions that are correct), particularly, recall considering top- $N+2$ ranked

predictions. In most cases, not all the existing binding pockets are observed with a ligand bound. In other words, the reference data are incomplete, with 33–50% of existing sites yet to be observed with ligands bound in a structure, as conjectured by Krivák and Hoksza [40]. Considering only the top- N predicted pockets assumes that there are exactly N real pockets for a given protein, which might not be the case. A method could predict a *real* pocket that is yet to be observed and rank it before other predicted and observed pockets. By considering the top- $N+2$ pockets, we are controlling to some extent for this noise in the reference data and obtaining a more accurate representation of the performance of a method. In a context of discovery where the true ligand binding sites of the target are unknown, it is more useful to have multiple predictions that might or might not correspond to real sites (lower precision), rather than a single or few predictions that are very precise but are missing other likely sites (lower recall). Most methods do well in predicting the most obvious (orthosteric) site. This site, however, might not be available for therapeutic targeting and it is convenient to have other predicted sites that could modulate function acting as allosteric sites. Precision, though a metric that provides valuable insight, and covered in this work, must always be contextualised with recall. Our results show that the most precise methods do not correspond to higher recalling methods. A method predicting the most obvious site, that could be identified by eye, might be 90% precise, but present a lower recall, e.g., 30%. This said, methods predicting fewer pockets with higher precision might prove more advantageous when users aim to study a particular region of interest in a protein, a few high-priority sites are needed for experimental validation, or false predictions are costly in a downstream analysis.

Some methods define “success rate” as the precision of the top-1 or top-3 scoring predictions, which is not a very representative performance assessment metric. For this reason, we encourage method developers not only to share the code of their approach, but also of the benchmarking analysis. Furthermore, the definition of *success rate* must be standardised as recall, as some methods use recall, whereas others use precision, both under the same term of success rate. This can be confusing when comparing the results from different analyses. Moreover, due to the inherent noise in the reference data, i.e., not all existing pockets are observed, recall considering top- $N+2$ is more informative than taking top-1, top-3 or top- N predictions. In any case, success rate must be clearly defined, so readers can fully understand the implications of the metric employed on a benchmark.

It is clear from our results that a DCC threshold of 4 Å is too conservative, and a more flexible DCC threshold

of 10–12 Å should be used for comparable performance with $DCA=4$ Å. According to our results and our reference set, predictions with DCC 4–12 Å overlap or are adjacent to observed sites and should be considered as correct predictions. The reason for this is the inherent noise in the ground truth, observed data or reference dataset, i.e., a ligand binding to a cavity might not be representative of all ligands that could bind to it. For most proteins, not all existing ligand sites are characterised, and as different ligands can bind to the same region, it is unrealistic to use such a small DCC threshold. Our results show several examples of predictions of an observed cavity where $DCC > 4$ Å.

$fpocket_{PRANK}$ (60%) and $DeepPocket_{RESC}$ (58%) present the highest top- $N+2$ recall of the methods reviewed in this work. $P2Rank_{CONS}$ and $P2Rank$ follow closely with 54% and 52% recall, then GrASP (50%), $DeepPocket_{SEG-NR}$, $Ligsite^+_{AA}$ and $PocketFinder^+_{AA}$ with 49%, $Surfnet^+_{AA}$ (47%), $VN-EGNN_{NR}$ (46%), $PUResNet_{PRANK}$ (41%) and $IF-SitePred_{RESC-NR}$ (39%). $fpocket$ is the method that predicts the most pockets per protein, reaching a maximum recall between 80 and 90% (considering all pockets regardless of the rank). $P2Rank_{CONS}$ comes second with a maximum recall of 50–60%. The rest of the methods range 40–55%. This indicates that whilst there are still some pockets unpredicted by $fpocket$ (10–20%), the maximum recall of this method is 20–30% higher than any other method. However, considering top- $N+2$ pockets, $fpocket$ only recalls 47% of the observed pockets. $fpocket_{PRANK}$ and $DeepPocket_{RESC}$ gain >10% in recall by simply re-scoring those predictions. This difference also applies to $\#TP_{100FP}$, where PRANK re-scoring results in +432 TP, and $Precision_{1K}$, with an increase of +34.4%. This highlights the paramount importance of a robust scoring scheme, which captures well the nature of binding sites and places those with a higher probability of being real binding sites at the top of the ranking. Newer methods like $VN-EGNN$, $IF-SitePred$, GrASP and $PUResNet$ are the most precise methods, however because of redundancy in predictions ($VN-EGNN$, $IF-SitePred$), or low number of predicted pockets per protein ($VN-EGNN$, GrASP and $PUResNet$) are limited in their recall. Their high precision indicates that their models learn and capture well the nature of ligand binding sites and so they represent a great venue to pursue in the field of ligand binding site prediction. Whilst removing redundancy post-prediction has a significant improvement in performance ($VN-EGNN_{NR}$ and $IF-SitePred_{NR}$), approaching this issue before prediction would be more beneficial. For $VN-EGNN$, which predicts a maximum of 8 sites, ensuring these 8 (or more) predictions are non-redundant is more desirable than removing redundant predictions

ending up with 1/8 predictions. The same applies to IF-SitePred, where non-overlapping starting predictions are more convenient than dealing with redundancy post-prediction.

The usefulness of residue-level metrics as F1 score or MCC is limited, as methods that precisely and correctly predict the clearest sites, such as PURESNet (high precision, low recall) perform better on these metrics, whilst methods that predict more pockets (lower precision, higher recall) such as fpocket and re-scored versions, will obtain worse results. Pocket-level metrics, particularly recall, are more representative of the ability to predict ligand binding sites.

Whilst datasets like PDBbind, binding MOAD or the brand new PLINDER [89] (will) prove extremely useful to train, validate and test deep learning models tackling problems such as rigid body docking [90], flexible pocket docking [91], or pocket-conditioned ligand generation [92], they might not be ideal as a test set for ligand binding site prediction. LIGYSIS analyses all unique, biologically relevant protein–ligand interfaces, including ions, across the biological assembly from multiple experimentally determined structures of a given protein. It then clusters these ligands based on their interactions with the protein, resulting in the observed binding sites. Beyond considering biological assemblies and unique protein–ligand interfaces, the greatest innovation in LIGYSIS is leveraging the extensive structural data on the PDBe-KB to aggregate ligand-binding interactions across different structures of the same protein, thus better capturing the ligand-binding capabilities than just taking a single structure of a protein–ligand complex. In doing so, LIGYSIS represents the most complete and non-redundant protein–ligand complex dataset to date. The benchmark is performed on LIGYSIS, which includes ion binding sites. When these are removed, all methods, except fpocket, experience an increase in (top- $N+2$) recall of 5–10% but the overall ranking of methods does not change (see Supplementary Figure 9). Due to its integrative approach, features, diversity and size (covering >30% of PDB and >20% of BioLiP) LIGYSIS is the most inclusive and representative dataset of protein–ligand interactions.

Aggregating protein–ligand interactions across structures of the same protein is likely to be beneficial not only for testing, but also when training these methods. Most current methods train on datasets where a protein is represented by a single structure interacting with a single ligand. For example, in 100% of sc-PDB and 50% of entries for binding MOAD training sets. Methods consider as ligand binding, and therefore TP, those residues within a certain distance of the ligand and TN all other residues. In doing so, residues of the same protein that bind ligands on other structures, but not the one present

on these sets, will be incorrectly labelled as “non-ligand-binding” (FN). This mislabelling of residues could lead to a lower prediction performance. This issue is to a certain extent approached by P2Rank and GrASP, which enriched their training datasets by including ligands from other chains, or homologous structures. This noise in the training dataset might be more prevalent for DeepPocket, PURESNet and VN-EGNN, which seem to rely fully on 1:1 interactions. The usage of LIGYSIS, or any other data set that aggregates ligand interactions across structures, might alleviate this issue and hints at potential room for improvement in the field of ligand binding site prediction.

Conclusions

The conclusions resulting from our analysis are as follows:

- Ligand binding site prediction methods differ significantly in the number of predicted sites, their size, proximity and overlap, which offers insight into how the methods work.
- Redundancy in ligand binding site prediction leads to an underestimate of recall and an overestimate of precision. The removal of such redundancy and subsequent re-ranking of the remaining pockets results in a drastic increase in recall.
- A robust pocket scoring scheme is crucial for the correct ranking and prioritisation of predicted sites in downstream analysis, e.g., docking, simulation. Additionally, it has a significant positive effect on both precision and recall.
- Recall is a more informative measure of the performance of a ligand site prediction tool, rather than precision and so it must be reported. Precision, though a useful metric, should always be contextualised with recall.
- All authors of ligand site prediction tools should use top- $N+2$ recall as “success rate” for consistency. Benchmarking code should also be shared by the authors for the sake of reproducibility.
- Pocket-level metrics (recall, precision) are a more adequate representation of the capabilities of ligand site prediction methods than residue-level metrics (F1, MCC).
- A DCC threshold of 4 Å is too conservative, and to obtain comparable results between DCA and DCC recall, a threshold of DCC of 10–12 Å should be employed.
- Re-scoring of fpocket predictions, as fpocket_{PRANK} or DeepPocket_{RESC} present the highest (top- $N+2$) recall (60%) among the methods reviewed in this analysis.

- Methods that systematically predict a low number of pockets, e.g., VN-EGNN, GrASP or PURESNet, are very precise (>90%), however their recall is low, and might not be as useful in a discovery context.
- IF-SitePred benefits significantly from pocket re-scoring, and suggests that protein embeddings, which aren't directly dependent of structure, represent great promise in the field of ligand site prediction.
- The use of duplicated protein–ligand interfaces in asymmetric units results in an overestimate of both precision and recall when benchmarking ligand site predictors. Only unique protein–ligand interfaces in biological units should be considered for a more accurate benchmark of the performance of these methods.
- LIGYSIS aggregates non-redundant biologically relevant protein–ligand interactions across multiple structures for a protein and sets a new test set standard for the benchmark of ligand binding site prediction tools.
- This work objectively evaluates the performance of 13 canonical ligand binding site prediction methods and 15 non-redundant and scoring variants using 10 different metrics. This analysis represents the largest benchmark of ligand binding site prediction to date.

Methods

Derivation of the LIGYSIS dataset

There are 20,423 human reviewed proteins in UniProt [93]. Only 7640 (37.4%) of these proteins present experimentally determined three-dimensional structures deposited in the Protein Data Bank (PDB) [94–96]. 5455 (71.4%) proteins present at least one ligand-binding structure. Only biologically relevant ligands, as defined

by BioLiP [66] were considered, yielding a set of 3513 proteins including 4037 structural segments, as defined on the PDBE-KB. These segments represent a UniProt sequence region with one or more structurally overlapping chains [97]. A protein might be represented by multiple segments which usually correspond to different domains. Protein chains mapping to a given UniProt ID were obtained from the PDBE aggregated API endpoint: <https://www.ebi.ac.uk/pdbe/graph-api/uniprot/superposition/> [98]. Transformation matrices to superpose protein chains for a structural segment were downloaded from the PDBE FTP site at <http://ftp.ebi.ac.uk/pub/databases/pdbe-kb/superposition> [99] and used to structurally align a total of 64,498 protein chains across 33,715 structures for the 4037 segments (see Supplementary Figures 12 and 13 for superposition examples and variation across chains for ligand binding residues).

Preferred biological assemblies, as defined by PISA [68], were downloaded from PDBE via ProIntVar [100]. Protein–ligand contacts were determined with *pdbe-arpeggio* [101]. Figure 10 illustrates the ligand site definition approach used to obtain our reference dataset: LIGYSIS. For a pair of ligands, L_A, L_B , fingerprints A, B are defined as sets containing the UniProt residue numbers of the amino acids interacting with each ligand. PDB residues were cross-referenced to UniProt using the SIFTS mapping present in the mmCIF files, located under the `_atom_site.pdbx_sifts_xref_db` fields [102, 103]. Relative intersection, I_{rel} (Eq. 1) is a similarity metric that quantifies how similar these fingerprints are [5]. Subtracting I_{rel} from 1 gives a distance, D (Eq. 2), which takes the value of 0 when A and B share all the binding residues and 1 when they share none. For a given protein segment, interacting with M biologically meaningful ligands across N chains, ligand fingerprints are clustered

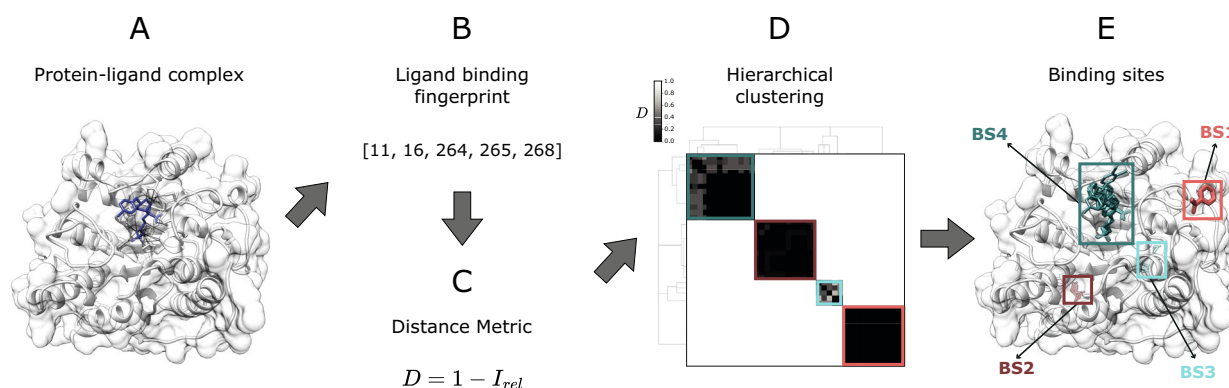


Fig. 10 LIGYSIS ligand binding site definition algorithm. For a given protein–ligand interaction complex (A), a ligand binding fingerprint (B) is obtained as the set of unique UniProt sequence residue numbers interacting with the ligand as defined by *pdbe-arpeggio*; C fingerprints from different ligands binding to the same protein can be compared and a distance calculated; D this distance can be used to perform hierarchical clustering, which groups the different ligands in distinct clusters binding to the same region of the protein or binding site (E)

using average linkage with SciPy [104] and ligand sites obtained by cutting the tree at $D=0.5$. This resulted in 8244 ligand binding sites, i.e., sets of UniProt residues. This method for the definition of ligand binding sites is therefore independent of structure superposition and superposition is merely used for visualisation purposes.

$$I_{rel} = \frac{A \cap B}{A \cap B_{max}} \quad (1)$$

$$D = 1 - I_{rel} \quad (2)$$

Comparison of datasets

Training and test datasets were downloaded for all machine learning based methods reviewed in this work. Datasets were compared to our reference set, LIGYSIS, in terms of number of sites per protein, ligand-interacting chains, chain lengths, size of the sites (number of amino acids), ligand composition, size and diversity. Ligand diversity was quantified by Shannon's Entropy [105] (Eq. 3) where p_i represents the proportion of each ligand i of the R ligands observed in the dataset. Ligand data was extracted from the Chemical Component Dictionary (CCD) [106]. An overlap (%) was calculated for each dataset as the proportion of LIGYSIS binding sites that were covered by at least one ligand in a dataset. A simplistic approach was adopted by calculating the intersection of ligand IDs between LIGYSIS and each dataset. Ligand IDs were defined as a string of PDB ID + "_" + ligand ID, e.g., "6GXT_GTP" corresponds to the guanosine-5'-triphosphate of the PDB entry with ID: 6GX7 [107]. LIGYSIS considers all BioLiP relevant ligands across structures for a protein and dismisses ligands that are a byproduct of crystallisation. Most predictors also tend to discard such ligands for training purposes. For this reason, despite its simplicity, the approach here employed captures well the overlap between training, test sets and LIGYSIS.

$$H' = - \sum_{i=1}^R p_i \ln(p_i) \quad (3)$$

Training datasets

VN-EGNN trains on a subset [44] of the sc-PDB (v2017) [108–112] (sc-PDB_{SUB}). sc-PDB is a comprehensive database of pharmacological ligand–protein complexes. The database is comprised by proteins in complex with buried, biologically relevant synthetic or natural ligands deposited in the PDB. sc-PDB contains unique non-repeating protein–ligand pairs, which means that

only one ligand is considered per PDB structure entry. Smith et al. [54] enriched this dataset with 9000 extra ligands resulting in a version of sc-PDB which we call sc-PDB_{RICH}, which GrASP trained on. Unfortunately, this dataset is not publicly accessible and therefore not considered in our analysis. DeepPocket used the full sc-PDB set to train on, sc-PDB_{FULL}. IF-SitePred uses a sequence identity-filtered version of the non-redundant subset of the binding mother of all databases (MOAD) [113–116], which considers only protein family leaders. The binding MOAD, here referred to as bMOAD_{SUB}, is a large collection of crystal structures with clearly identified biologically relevant ligands with binding data extracted from the literature. Finally, PRANK and P2Rank used the CHEN11 dataset to train, which aimed to cover all SCOP [117–119] families of ligand binding proteins in a non-redundant manner [67]. P2Rank utilised the JOINED dataset for validation. CHEN11 not only considers the ligands in each structure but is enriched with ligands binding to homologous structures. JOINED is a combined dataset formed by other smaller datasets: ASTEX [120], UB48 [30], DT198 [121] and MP210 [34], which represent diverse collections of protein–ligand complexes, including bound/unbound states, drug–target complexes and other ligand site predictor benchmark sets.

Test datasets

The majority of ligand binding site predictors published since 2018 have been using two datasets that were first presented by Krivák et al. [40]: COACH420 and HOLO4K, or subsets of them. COACH420 is comprised by a set of 420 single-chain structures binding a mix of drug-like molecules and naturally occurring ligands which is disjoint with the CHEN11 and JOINED datasets. COACH420 is a modified version of the original COACH test set [26, 27]. HOLO4K is a larger set, $N \approx 4000$, based on the list by Schmidtke et al. [122], which includes a mix of single- and multi-chain complexes, also disjoint with P2Rank training (CHEN11) and validation (JOINED) datasets. PRANK employed the small datasets comprising the JOINED set for testing. VN-EGNN, DeepPocket and GrASP use the Mlig and Mlig+ subsets of the COACH and HOLO4K datasets, which comprise strictly biologically relevant ligands as defined by the binding MOAD. IF-SitePred is tested on the HOLO4K-AlphaFold2 Paired (HAP) and HAP-small sets. HAP is a subset of the HOLO4K dataset which presents high quality models in the AlphaFold database [123]. HAP-small is a smaller subset of HAP that only contains proteins with sequence identity lower than 25% to proteins in the P2Rank training set. VN-EGNN uses the refined version of PDBbind (v2020), referred here as PDBbind_{REF} [124–129], as a third test set. Like binding

MOAD, the PDBbind database provides a comprehensive collection of experimentally measured binding affinity data for macromolecular complexes. Specifically, the refined set includes those protein–ligand complexes for which binding data was obtained with the literature and met certain experimental quality thresholds. Lastly, SC6K is a dataset presented by Aggarwal et al. [48] containing 6000 protein–ligand pairs from PDB entries submitted from 01/01/2018 to 28/02/2020.

Protein chain alignment

For each protein chain, atomic coordinates are translated to be centred at the origin, $O=(0, 0, 0)$, and rotated using a rotation matrix, R . The two principal components of the coordinate space pc_1 and pc_2 are obtained using principal component analysis (PCA) [130]. A third component, pc_{\perp} , (Eq. 4) is obtained with the cross-product of the other two, to ensure orthogonality. A rotation matrix P is constructed from these vectors (Eq. 5). By placing the main component pc_1 on the second row of P , we ensure the Y axis will be the major axis, representing the *height* of the protein chain. The second largest axis will be the X axis, representing the *width* of the protein, and lastly the *depth* will be represented by the smaller magnitude of the Z axis. The final rotation matrix R (Eq. 6) is obtained by multiplying P by the negative identity matrix NI (Eq. 7). This was done to maintain the left-handedness of the protein chains whilst ensuring a consistent alignment on the major axes.

$$pc_{\perp} = pc_1 \times pc_2 \quad (4)$$

$$P = \begin{bmatrix} pc_2 \\ pc_1 \\ pc_{\perp} \end{bmatrix} \quad (5)$$

$$R = P \cdot NI \quad (6)$$

$$NI = -1 \cdot I_3 = -1 \cdot \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{bmatrix} \quad (7)$$

Protein chain characterisation

For a protein chain with N amino acid residues, the centre of mass, CM , is calculated by averaging the coordinates,

r_i , of all atoms (Eq. 8), and from it, the radius of gyration, R_g , is derived (Eq. 9) [131]. As the protein chains are already aligned on the axis and centred on $O=(0, 0, 0)$ the dimensions of the protein chain can be obtained as the magnitude of the PCA components or *eigenvectors*, i.e., the *eigenvalues*. The dimensions represent width, height, and depth for the X , Y and Z axes, respectively.

$$CM = \frac{1}{n} \sum_{i=1}^n r_i \rightarrow CM = O = (0, 0, 0) \quad (8)$$

$$R_g = \sqrt{\frac{1}{n} \sum_{i=1}^n (r_i - CM)^2} \rightarrow R_g \\ = \sqrt{\frac{1}{n} \sum_{i=1}^n (r_i - O)^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n r_i^2} \quad (9)$$

Protein chain volumes were calculated using Protein-Volume [132]. A sphere enclosing the protein and centred on the protein centre of mass was obtained. The radius of this sphere is the maximum Euclidean distance between the protein atoms and the CM (Eq. 10). The volume of the sphere is calculated using Eq. 11. Proteins were classified into four different groups based on their shape and size. Protein chains with ≤ 100 amino acids were classified as “tiny”. Regarding the shape, protein chains were classified into “elongated” if their protein to sphere volume ratio ≤ 0.08 (VR) (Eq. 12), i.e., the protein volume contains no more than 8% of the sphere volume. This threshold was derived empirically by the visual inspection of all 3448 protein chains on the LIGYSIS set. Otherwise, proteins were considered globular (Fig. 11). In this manner, protein chains were classified into *globular* ($N=2104$; 61%), *elongated* ($N=670$; 19%), *elongated tiny* ($N=341$; 10%) and *globular tiny* ($N=333$; 10%).

$$R = \max \|r_i - CM\| \quad (10)$$

$$Volume_{Sphere} = \frac{4}{3} \pi R^3 \quad (11)$$

$$VR = Volume_{Protein} / Volume_{Sphere} \quad (12)$$

(See figure on next page.)

Fig. 11 Protein chain shape and size classification approach. **A** The volume of the sphere enclosing the protein chain as well as the protein chain volumes are calculated, and their ratio obtained (VR). Globular proteins present more spherical shapes and therefore occupy a higher portion of the sphere volume, resulting in higher volume ratios. Non-globular, elongated or fibrous proteins on the other hand do not and present lower VRs. After extensive visual examination, a threshold was established at $VR=0.08$, and so proteins classified in these two groups. Proteins were classified as “tiny” if their chain was ≤ 100 amino acids; **B** eight examples of each protein chain group to illustrate the outcome of the approach

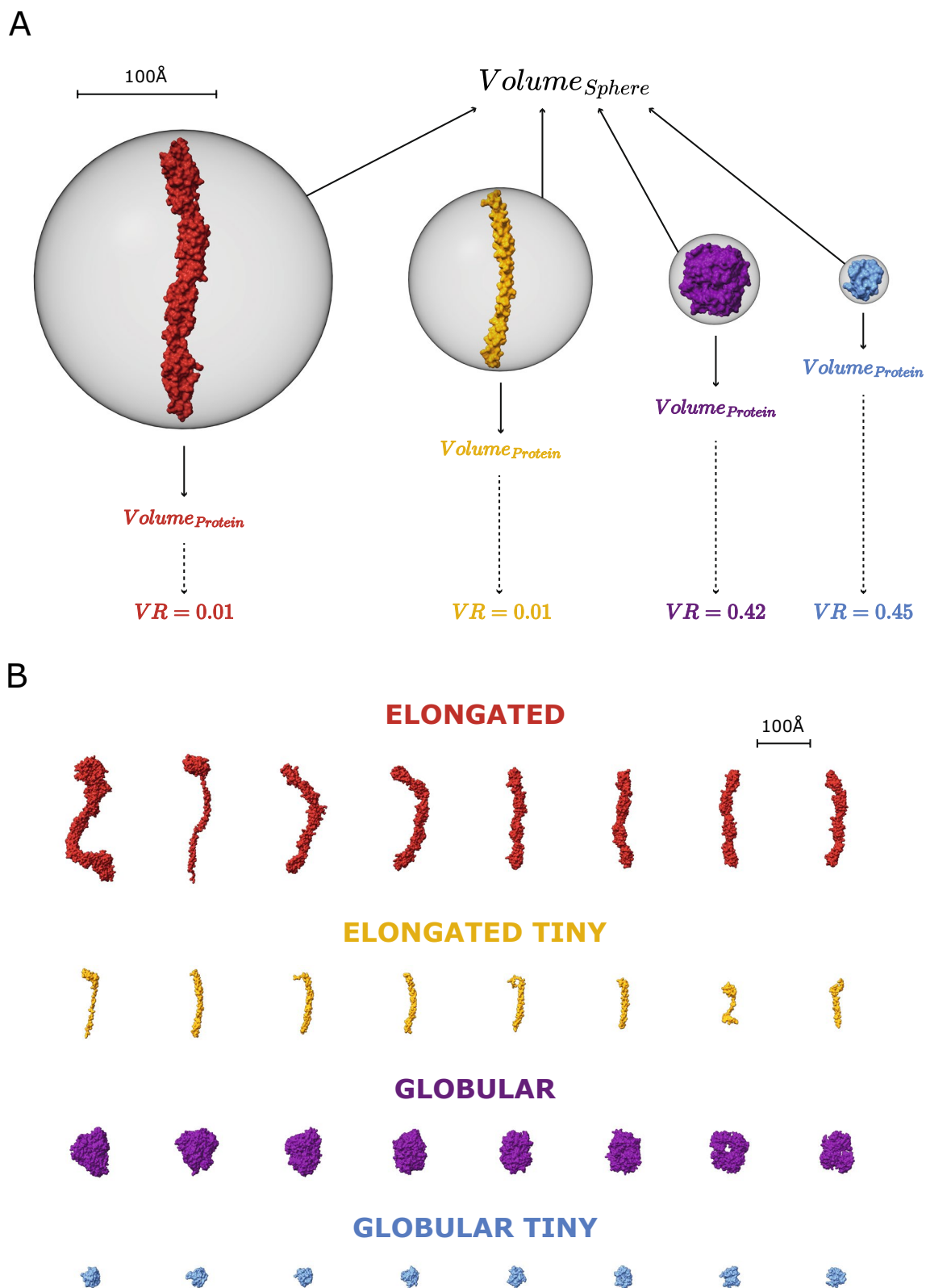


Fig. 11 (See legend on previous page.)

Ligand binding site prediction

For each segment in the LIGYSIS dataset, the representative chain as defined in the PDBE-KB was selected. Structures were cleaned using the *clean_pdb.py* script [133]. Eleven different ligand binding site prediction tools, with their default settings, were employed to predict on the 3448 representative chains: VN-EGNN [56], IF-SitePred [55], GrASP [54], PURESNet [44], DeepPocket [48], P2Rank [37, 40], PRANK [37], fpocket [7, 134], PocketFinder⁺ [15], Ligsite⁺ [8], and Surfnet⁺ [9]. Conservation scores were obtained from PrankWeb <https://prankweb.cz/> and used for further prediction. This variant of P2Rank employing amino acid conservation is referred to as P2Rank_{CONS} [63, 135]. When running DeepPocket, the *-r* threshold was removed and so all fpocket candidates were passed to the CNN-based segmentation module for pocket shape estimation. fpocket predictions re-scored by DeepPocket are referred to as DeepPocket_{RESC}, whereas pockets extracted by the segmentation module of DeepPocket are referred to as DeepPocket_{SEG}. fpocket predictions were also re-scored with PRANK [37] (fpocket_{PRANK}), as introduced in previous studies [37, 38, 40, 87]. Re-implementations by Capra et al. [33] were used for PocketFinder, Ligsite and Surfnet, indicated by the “+” superscript. VN-EGNN, IF-SitePred, PocketFinder⁺, Ligsite⁺ and Surfnet⁺ do not provide a list of residues for each pocket, but a list of centroids and their scores for the first two, and a list of grid points for each predicted pocket for the last three. For VN-EGNN, residues within 6 Å of the virtual nodes were considered pocket residues. For 429 predicted pockets (≈3%) no residues were found within this threshold. For IF-SitePred, residues within 6 Å of the clustered cloud points that resulted on a predicted pocket centroid were considered as pocket residues. Pocket residues were obtained in a similar manner for PocketFinder⁺, Ligsite⁺ and Surfnet⁺, by taking those residues within 6 Å of pocket grid points. In total, thirteen methods are considered on this analysis: VN-EGNN, IF-SitePred, GrASP, PURESNet, DeepPocket_{RESC}, DeepPocket_{SEG}, P2Rank_{CONS}, P2Rank, fpocket_{PRANK}, fpocket, PocketFinder⁺, Ligsite⁺ and Surfnet⁺.

Seven of the considered methods provide residue *ligandability* scores. P2Rank and P2Rank_{CONS} report calibrated probabilities of residues being ligand-binding. Similarly, GrASP predicts the likelihood for any given heavy atom to be part of a binding site. A residue-level score was obtained for GrASP predictions by taking the maximum score of the residue atoms. For IF-SitePred, a residue ligandability score *LS* can be obtained by averaging the 40 independently predicted probabilities of a residue being ligand-binding (Eq. 13). Though calculated in a different way, these three scores range 0–1 and represent

the likelihood of a residue binding a ligand and can therefore be compared. PocketFinder⁺, Ligsite⁺, and Surfnet⁺ also provide residue scores which maximum value can be > 1.

$$LS = \frac{1}{40} \sum_{i=1}^{40} p_i \quad (13)$$

VN-EGNN, PURESNet, DeepPocket_{RESC}, DeepPocket_{SEG}, fpocket_{PRANK} and fpocket do not report residue-level scores. However, binary labels represent whether a residue is part of a pocket (1) or not (0), in the same manner as all other methods.

Throughout this work the terms “site” and “pocket” are used indistinctly. Across all figures, tables and legends, methods are sorted in chronological order.

Binding site characterisation

Radius of gyration, R_g , was calculated for pockets as it was done for whole protein chains (Eq. 9). Distance between pockets was calculated as the Euclidean distance between their centroids and overlap between pocket residues with the Jaccard Index (JI), or intersection over union (IOU) (Eq. 14) [136, 137]. POVME 2.0 was employed for pocket volume calculation [138–140]. A single inclusion region was used for each pocket. This region is defined by the smallest rectangular prism containing all pocket atoms. The prism is centred on the pocket centroid and its dimensions are determined by the distance between the two farthest atomic coordinates on each axis. No exclusion regions were used. Points outside the convex hull were deleted. A contiguous-points region was defined as a 5 Å-radius sphere on the pocket centroid (Fig. 12).

$$JI(A, B) = \frac{|A \cap B|}{|A \cup B|} \quad (14)$$

Prediction evaluation

LIGYSIS binding sites consist of sets of UniProt residue numbers to which ligands bind across the multiple structures of a protein. The thirteen ligand binding site predictors benchmarked in this work predict only on the representative chains for each protein. These representative structures are defined in the PDBE-KB based on three criteria: data quality, sequence coverage and resolution [98]. Despite this, representative chains might still be missing some residues present in other structures. To compare LIGYSIS binding sites to predicted sites on the representative chains, UniProt sequence mappings are needed for each residue in the LIGYSIS-defined sites. For this reason, LIGYSIS entries with ligand-binding residues missing UniProt residue mappings on the protein's

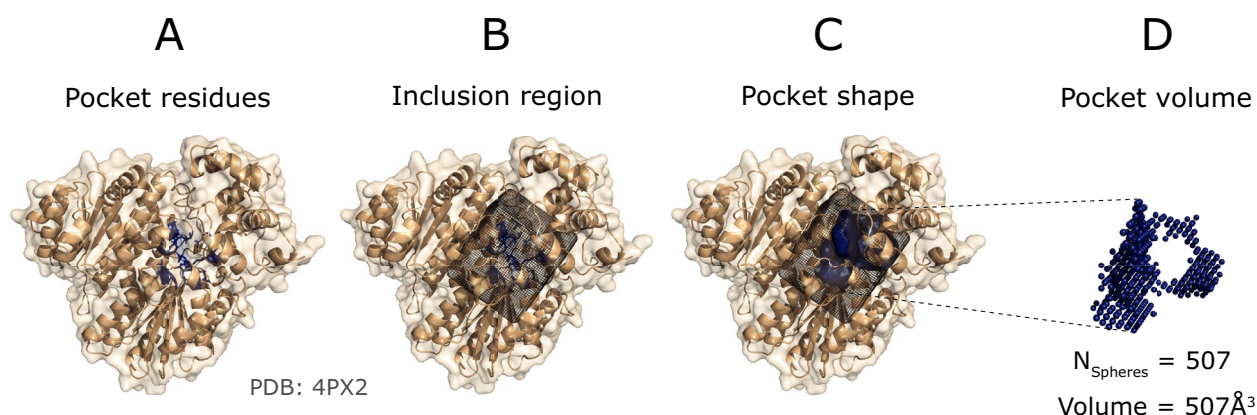


Fig. 12 Pocket volume calculation algorithm. **A** PURESNet predicted pocket for PDB: 4PX2 [141]. Pocket residues are coloured in blue and have their side chains displayed; **B** an inclusion region is determined using the coordinates of the pocket residue atoms; **C** POVME 2.0 calculates the shape of the pocket defined by the residues and contained within the inclusion region; **D** The pocket shape is defined by a series of unit-volume (1 \AA^3) spheres. The volume of the pocket is calculated as the addition of the sphere volumes or the number of spheres within the pocket. Structure visualisation with PyMOL v2.5.2 [61]

representative chain were removed from our set, resulting in a set of 3048 human proteins, including 3448 segments. After predicting on these 3448 LIGYSIS chains, only chains where all residues across predicted sites presented UniProt residue mappings were kept. This resulted in a final set of 2775 protein chains which was employed to assess the performance of the methods. These filters are necessary to compare LIGYSIS-defined sites to predicted sites as UniProt residue overlap between defined and predicted sites is one of the criteria to evaluate the methods.

The performance of ligand binding site prediction methods can be evaluated at two different levels: *residue* level, and *pocket* level. Prediction at the *residue* level involves the discrimination of those residues that are likely to interact with a ligand, whereas the aim of *pocket*-level prediction is to define distinct regions on a protein, i.e., pockets where a ligand is likely to bind. This region can either be defined by a (pocket) centroid, a group of cloud/grid points, a set of pocket residues, or a combination of these. Some methods are *residue*-centric, and first predict at the residue-level, use a threshold to select high-probability ligand-binding residues, and then cluster them into pockets. Residue-centric methods include IF-SitePred, or GrASP. Other methods (*pocket*-centric) directly predict the location or shape of the pocket, without the need of predicting at the residue level first. Some of these methods can use their pocket-level prediction to report residue *ligandability* scores, e.g., P2Rank_{CONS}, P2Rank, PocketFinder⁺, Ligsite⁺ or Surfnet⁺, and others, such as VN-EGNN, PURESNet, DeepPocket, or fpocket do not report residue ligandability scores.

Prior to evaluation, non-redundant sets of predictions were generated for VN-EGNN, IF-SitePred and DeepPocket_{SEG}, as these methods generate a considerable proportion of *redundant* predictions. A predicted pocket i is considered redundant if there exists a pocket $j \neq i$ so that the distance between their centroids $D_{ij} \leq 5 \text{ \AA}$ or their residue overlap $J_{ij} > 0.75$, i.e., they share at least $\frac{3}{4}$ of their residues. Refer to Supplementary Figure 14 for the closest predicted sites for each method.

PURESNet, PocketFinder⁺, Ligsite⁺ and Surfnet⁺ do not score, nor explicitly rank their pockets, and so they were taken in the order given by their pocket ID. This means that when sorting across the dataset, the order of all pockets with the same rank is arbitrary. To obtain a score for these pockets, multiple strategies were employed. Firstly, a pocket score was obtained as the number of pocket amino acids, resulting in variants PURESNet_{AA}, PocketFinder⁺_{AA}, Ligsite⁺_{AA} and Surfnet⁺_{AA}. Secondly, PRANK pocket scoring was employed, resulting in variants PURESNet_{PRANK}, PocketFinder⁺_{PRANK}, Ligsite⁺_{PRANK} and Surfnet⁺_{PRANK}. IF-SitePred uses a simple pocket scoring scheme, which assigns to each centroid the number of clustered cloud points it results from. In this work, newly defined IF-SitePred pocket scores were calculated as the sum of squares (SS) of the ligandability scores (LS_i), calculated with (Eq. 13), of the K residues on a site (Eq. 15) resulting in IF-SitePred_{RESC}. For PocketFinder⁺, Ligsite⁺, and Surfnet⁺ the same was done but instead of residue scores, grid point scores (GS_i) were used (Eq. 16). This resulted in further variants PocketFinder⁺_{SS}, Ligsite⁺_{SS}, and Surfnet⁺_{SS}. This is the same approach introduced by Krivák, et al. [37] and later adopted by Smith et al. [54].

$$SS_{IF-SitePred} = \sum_{i=1}^K LS_i^2 \quad (15)$$

$$SS_{PocketFinder^+, Ligsite^+, Surfnet^+} = \sum_{i=1}^K GS_i^2 \quad (16)$$

Residue level predictions

GrASP, P2Rank_{CONS}, P2Rank, PocketFinder⁺, Ligsite⁺ and Surfnet⁺ all offer residue ligandability scores. Additionally, a ligandability score was derived for IF-SitePred using Eq. 13. Prediction at the residue level is a binary classification problem: binding (1) or non-binding (0). Given a ligandability threshold t_{LS} , a residue with a ligandability score LS_i is classified as “positive” if $LS_i > t_{LS}$. Conversely, the residue is classed as “negative” if $LS_i \leq t_{LS}$. Further stratification results from comparing the predictions to the LIGYSIS reference dataset.

- True positive (TP): residue classified as *positive* that binds a ligand according to the reference.
- False positive (FP): residue classified as *positive* that does not bind a ligand in the reference.
- True Negative (TN): residue classified as *negative* that does not bind a ligand.
- False Negative (FN): residue classified as *negative* but is known to bind a ligand according to the reference.

With these four classes, true positive rate (TPR) (Eq. 17), false positive rate (FPR) (Eq. 18), precision (Eq. 19) and recall (Eq. 20) can be calculated and receiver operating characteristic (ROC), precision-recall (PR) curves plotted. ROC and PR curves were obtained for each of the LIGYSIS protein chains. Using these curves, mean ROC and PR curves, representative of the variation across proteins for these metrics were obtained by taking the mean TPR and FPR (ROC) and mean precision and recall (PR) at each score interval. Mean Area under the curve (AUC) and mean Average Precision (AP) were calculated by averaging the areas and precisions across curves. Baselines for these are 0.5 and the proportion of true binding residues (0.1), respectively. ROC and AUC can't be calculated for VN-EGNN, PURESnet, DeepPocket, fpocket_{PRANK} and fpocket as these methods do not provide residue ligandability scores.

$$TPR = \frac{TP}{TP + FN} \quad (17)$$

$$FPR = \frac{FP}{FP + TN} \quad (18)$$

$$precision = \frac{TP}{TP + FP} \quad (19)$$

$$recall = \frac{TP}{TP + FN} \quad (20)$$

Pocket binary labels (0: no pocket residue; 1: pocket residue) can also determine TP, FP, TN and FN for each residue in a protein chain P_i . VN-EGNN, IF-SitePred, PocketFinder⁺, Ligsite⁺ and Surfnet⁺ do not report pocket residues. For these methods, residues within 6 Å of the pocket centroid, cloud points and grid points (3×), respectively, were labelled as pocket residues (1). All other residues were labelled as non-binding (0). Across all residues in P_i , an F1 score FI_i is computed, which combines precision and recall into a unified metric, capturing the accuracy and completeness of predictions at the residue level (Eq. 21). The Matthews Correlation Coefficient (MCC) [142] (Eq. 22) was also calculated. The mean F1 score and MCC across LIGYSIS proteins is reported for each method.

$$F1_i = \frac{2 \times Precision_i \times Recall_i}{Precision_i + Recall_i} \quad (21)$$

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (22)$$

Pocket level predictions

Ligand binding site prediction at the pocket level is a multi-instance prediction problem. There are no *negatives* predicted, only *positives*. A positive is a predicted pocket, which will be true (TP) or false (FP) depending on whether it is observed in the reference data. False negatives are those pockets observed in the reference data that are not predicted. They are the pockets the method fails to predict, and therefore, are not scored. A true negative would be a “non-pocket” that is not predicted. This can't be quantified easily and even if it was, it would not be scored by the method, as it is not predicted. For this reason, in this context, neither TPR, nor FPR can be calculated. Consequently, ROC/AUC can't be utilised to assess ligand binding site prediction at the pocket level. False negatives are known, but not scored, and therefore PR/AUC is not an option either. What can be calculated is the recall given a certain criterion. In this case, because of the nature of the LIGYSIS dataset, where defined sites result from the clustering of multiple ligands, the distance between the predicted pocket centroid to the observed binding site (DCC) was chosen as a criterion.

For each observed binding site in our reference dataset, the “best” prediction for each method is chosen. This

is defined as the prediction with the minimum Euclidean distance to the observed pocket centroid or DCC. Once the observed-predicted pairs were obtained, only those with $DCC \leq 12 \text{ \AA}$ were considered as correct predictions. A threshold of 12 \AA was chosen as 4 \AA is too strict a threshold when using DCC (Supplementary Note 6 and Supplementary Figures 15–19). A threshold of 4 \AA works well for the distance to closest ligand atom (DCA) but does not for DCC. The top- N and $N+2$ ranking predictions were considered to calculate success rate, or recall (Eq. 23), and maximum recall was calculated by considering *all* predictions, regardless of their score or rank. N represents the number of observed sites for a given protein.

$$\text{Success rate(recall)} = \frac{\text{observed sites with predicted site } DCC \leq 12 \text{ \AA}}{\text{observed sites}} \quad (23)$$

Additionally, instead of conventional ROC, ROC100 [143, 144] can be used to measure the predictive performance of the methods. To do this, for each method, *all* predictions across dataset proteins were ranked based on the pocket score and cumulative true positives were plotted against cumulative false positives until 100 false positives were reached. In a similar way, a precision curve can be calculated by taking the top- N , in this case $N=1000$, predictions, which measures how precision changes as more predictions with lower scores are considered. This is indicative of how informative pocket scores are.

Precision and recall are key measures for evaluating the performance of ligand binding site prediction methods. However, these indicators are calculated and interpreted slightly differently depending on the context a prediction is analysed, i.e., pocket vs residue level, as well as the metric employed, e.g., F1 score, MCC, ROC or PR curves. At the residue level, the prediction is a binary classification task, where each residue is classified as binding (1) or non-binding (0). Here, precision reflects the proportion of residues predicted as binding that are true, i.e., observed in the reference data. Recall measures the proportion of true binding residues that are correctly identified. For the calculation of F1 and MCC, a residue is labelled “positive” or “negative” depending on whether it is part of a predicted pocket. However, for ROC and PR curves, the positive and negative labels are derived

based on a ligandability threshold, t_{LS} . Prediction at the pocket level represents a multi-instance prediction task. Precision indicates the proportion of predicted pockets that are observed in the reference data whilst recall represents the proportion of true binding pockets that are correctly predicted. It is important to keep this in mind to correctly interpret precision and recall across different contexts.

To measure the similarity in shape and residue membership between the predicted and observed pockets, relative residue overlap (RRO) and relative volume overlap (RVO) were employed. For an observed-predicted pocket pair, RRO represents the proportion of observed ligand-binding residues (R_o) that are covered by the pre-

dicted pocket residues (R_p) (Eq. 24). The POVME output was utilised for the calculation of RVO (Fig. 13). POVME defines the volume of a pocket as a series of equidistantly spaced spheres of unit volume. As predictions by the different methods were on the same coordinate reference, these pocket volume spheres were already aligned, and the volume overlap was calculated simply as the proportion of spheres in the observed pocket (V_o) that overlap with the predicted pocket spheres (V_p) (Eq. 25).

$$RRO = \frac{|R_p \cap R_o|}{R_o} \quad (24)$$

$$RVO = \frac{|V_p \cap V_o|}{V_o} \quad (25)$$

Statistics and reproducibility

VN-EGNN was installed and run locally from <https://github.com/ml-jku/vnegnn>. Likewise, for IF-SitePred: <https://github.com/annacarbery/binding-sites>. GrASP was obtained from <https://github.com/tiwarylab/GrASP> and predictions generated using their Google Colab Notebook. PURESNet predictions were obtained through the PURESNet v2.0 web server: <https://nslbio.jbnu.ac.kr/tools/jmol>. DeepPocket was installed and executed

(See figure on next page.)

Fig. 13 Relative volume overlap (RVO) calculation. **A** Example of two very accurate predictions by PURESNet and P2Rank on PDB: 4PX2 [141]. Pocket volumes are calculated with POVME 2.0 [139] and represented by coloured surfaces. These volumes result from the addition of unit-volume spheres on a grid. To obtain the RVO, the intersection of these spheres between predicted and observed site is divided by the number of observed pocket spheres. Both predictions cover the entirety of the observed pocket volume; **B** GrASP and VN-EGNN predictions of a site on PDB: 2ZOX [145]. The volumes of these predicted sites overlap less with the observed site: RVO=0.67 for GrASP and RVO=0.11 for VN-EGNN

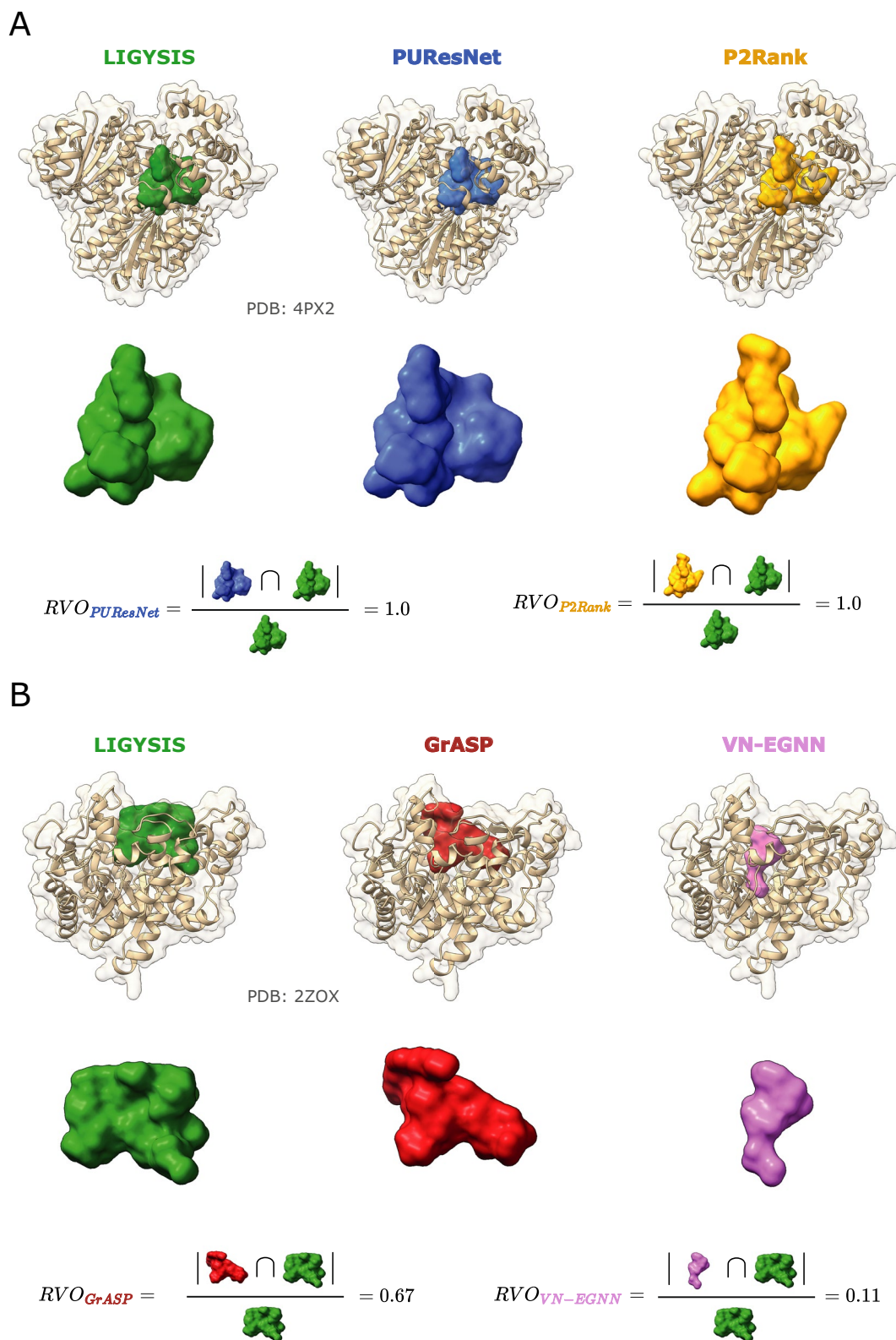


Fig. 13 (See legend on previous page.)

locally: <https://github.com/devalab/DeepPocket>. P2Rank v2.4.2 was used to run all predictions as well as PRANK scoring: <https://github.com/rdk/p2rank>. fpocket v4.0 was installed via Conda: <https://anaconda.org/conda-forge/fpocket>. For PocketFinder, Ligsite and Surfnet, the ConCavity v0.1 “+” re-implementations were employed: <https://compbio.cs.princeton.edu/concavity/>.

Other recent methods including RefinePocket [53], EquiPocket [52], GLPocket [51], SiteRadar [50], NodeCoder [49], RecurPocket [47], PointSite [46], DeepSurf [45], Kalasanty [43], BiteNet [42] GRaSP [41], or DeepSite [39] were not included in this analysis due to technical reasons. Peer-reviewed open-source methods with publicly accessible code, clear installation instructions, well defined dependencies, accessible command line interfaces and trained machine learning models were prioritised in this work. This set of thirteen methods, counting _{RESC} and _{SEG} variants of DeepPocket is representative of the state-of-the-art within the field.

ChimeraX v1.7.1 [146] was used for structural visualisation in all figures unless otherwise stated, in which case PyMOL v2.5.2 was employed [61].

Abbreviations

0XR	Ethyl caffeate
017	Darunavir
3D	Three-dimensional
A32	BMSC-0010
AA	Amino acid
ADE	Adenine
ADP	Adenosine-5'-diphosphate
AP	Average precision
ATP	Adenosine-5'-triphosphate
AUC	Area under the curve
BGC	Glucose
C5P	Cytidine-5'-monophosphate
CI	Confidence interval
CLA	Chlorophyll A
CLR	Cholesterol
CM	Centre of mass
CNN	Convolutional neural network
DBSCAN	Density-based spatial clustering of applications with noise
DCA	Distance to closest ligand atom
DCC	Distance centroid to centroid
DGD	Digalactosyl diacyl glycerol
DNN	Deep neural network
DRN	Deep residual network
EGNN	Equivariant graph neural network
FAD	Flavin-adenine dinucleotide
FMN	Flavin mononucleotide
FN	False negative
FP	False positive
FUC	Fucose
G2P	Phosphomethylphosphonic acid guanylate ester
GAI	Guanidine
GAL	Galactose
GAT	Graph attention network
GFB	GDP-L-fucose
GNN	Graph neural network
GSH	Glutathione
HAP	HOLO4K-AlphaFold2 paired
HEM	Haem
IOU	Intersection over union

IQR	Interquartile range
I _{rel}	Relative intersection
JI	Jaccard Index
LightGBM	Light gradient boosting machine
LJ	Lennard–Jones
MAN	Mannose
MCC	Matthews correlation coefficient
MCD	Minimum centroid distance
MRO	Maximum residue overlap
NAD	Nicotinamide-adenine-dinucleotide
NAP	Nicotinamide-adenine-dinucleotide phosphate
NI	No-ions
N-mer	Protein peptide of N amino acids
NR	Non-redundant
PCA	Principal component analysis
PID	Peridinin
PLP	Vitamin B6 phosphate
PR	Precision recall curve
R _g	Radius of gyration
ROC	Receiver operating characteristic
RRO	Relative residue overlap
RVO	Relative volume overlap
RSA	Relative solvent accessibility
SD	Standard deviation
SAH	S-Adenosyl-L-homocysteine
SAS	Solvent accessible surface
SFG	Sinefungin
SS	Sum of squares
TN	True negative
TP	True positive
VN	Virtual node
VR	Volume ratio
XYP	Xylose

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13321-024-00923-z>.

Supplementary Material 1.

Acknowledgements

We thank Florian Sestak, Rishal Aggarwal, Jeevan Kandel, Drs. Zachary Smith, Radoslav Krivák, and Profs. Günter Klambauer, Pratyush Tiwary and Deva Priyakumar, for their insightful discussions and help on how to run and understand their methods. We also appreciate the authors of the other methods reviewed in this paper for their contributions to the field. We extend our gratitude to Drs. James Abbot, Michele Tinti, Marek Gierliński, Stuart MacGowan, James Procter and James Webster for their suggestions. We also thank the IT service of the University of Dundee for their support of the HPC infrastructure this study was carried out on.

Author contributions

G.J.B. and J.S.U., conceived, designed, and developed the research. J.S.U. analysed the data. J.S.U. developed the software. J.S.U. and G.J.B. wrote, reviewed and edited the manuscript. G.J.B. secured funding and supervised.

Funding

This work was supported by grants to G.J.B. from UKRI-Biotechnology and Biological Sciences [BB/J019364/1; BB/R014752/1] and Wellcome Trust [101651/Z/13/Z; 218259/Z/19/Z]. J.S.U. was supported by a BBSRC EASTBIO Ph.D. Studentship [BB/J01446X/1].

Data availability

The main results tables and files necessary to replicate the analysis described in this paper can be found here: <https://doi.org/10.5281/zenodo.13121414> [147].

Code availability

Software developed to carry out this analysis is found in our GitHub repository: <https://github.com/bartongroup/LB5-comparison> [148].

Declarations**Competing interests**

The authors declare no competing interests.

Received: 2 August 2024 Accepted: 28 October 2024

Published online: 11 November 2024

References

- Congreve M et al (2003) A 'rule of three' for fragment-based lead discovery? *Drug Discov Today* 8(19):876–877
- Rees DC et al (2004) Fragment-based lead discovery. *Nat Rev Drug Discov* 3(8):660–672
- Murray CW, Rees DC (2009) The rise of fragment-based drug discovery. *Nat Chem* 1(3):187–192
- Schiebel J et al (2016) Six biophysical screening methods miss a large proportion of crystallographically discovered fragment hits: a case study. *ACS Chem Biol* 11(6):1693–1701
- Utgés JS et al (2024) Classification of likely functional class for ligand binding sites identified from fragment screening. *Commun Biol* 7(1):320
- Volkamer A et al (2010) Analyzing the topology of active sites: on the prediction of pockets and subpockets. *J Chem Inf Model* 50(11):2041–2052
- Le Guilloux V, Schmidtke P, Tuffery P (2009) Fpocket: an open source platform for ligand pocket detection. *BMC Bioinform* 10:168
- Hendlich M, Rippmann F, Barnickel G (1997) LIGSITE: automatic and efficient detection of potential small molecule-binding sites in proteins. *J Mol Graph Model* 15(6):359–363, 389
- Laskowski RA (1995) SURFNET: a program for visualizing molecular surfaces, cavities, and intermolecular interactions. *J Mol Graph* 13(5):323–330, 307–308.
- Levitt DG, Banaszak LJ (1992) POCKET: a computer graphics method for identifying and displaying protein cavities and their surrounding amino acids. *J Mol Graph* 10(4):229–234
- Kleywegt GJ, Jones TA (1994) Detection, delineation, measurement and display of cavities in macromolecular structures. *Acta Crystallogr D Biol Crystallogr* 50(Pt 2):178–185
- Liang J, Edelsbrunner H, Woodward C (1998) Anatomy of protein pockets and cavities: measurement of binding site geometry and implications for ligand design. *Protein Sci* 7(9):1884–1897
- Brady GP Jr, Stouten PF (2000) Fast prediction and visualization of protein binding pockets with PASS. *J Comput Aided Mol Des* 14(4):383–401
- Weisel M, Proschak E, Schneider G (2007) PocketPicker: analysis of ligand binding-sites with shape descriptors. *Chem Cent J* 1:7
- An J, Totrov M, Abagyan R (2005) Pocketome via comprehensive identification and classification of ligand binding envelopes. *Mol Cell Proteom* 4(6):752–761
- Goodford PJ (1985) A computational procedure for determining energetically favorable binding sites on biologically important macromolecules. *J Med Chem* 28(7):849–857
- An J, Totrov M, Abagyan R (2004) Comprehensive identification of "druggable" protein ligand binding sites. *Genome Inform* 15(2):31–41
- Laurie AT, Jackson RM (2005) Q-SiteFinder: an energy-based method for the prediction of protein–ligand binding sites. *Bioinformatics* 21(9):1908–1916
- Gherzi D, Sanchez R (2009) EasyMIFS and SiteHound: a toolkit for the identification of ligand-binding sites in protein structures. *Bioinformatics* 25(23):3185–3186
- Ngan CH et al (2012) FTSite: high accuracy detection of ligand binding sites on unbound protein structures. *Bioinformatics* 28(2):286–287
- Armon A, Graur D, Ben-Tal N (2001) ConSurf: an algorithmic tool for the identification of functional regions in proteins by surface mapping of phylogenetic information. *J Mol Biol* 307(1):447–463
- Pupko T et al (2002) Rate4Site: an algorithmic tool for the identification of functional regions in proteins by surface mapping of evolutionary determinants within their homologues. *Bioinformatics* 18(Suppl 1):S71–S77
- Xie ZR, Hwang MJ (2012) Ligand-binding site prediction using ligand-interacting and binding site-enriched protein triangles. *Bioinformatics* 28(12):1579–1585
- Zvelebil MJ et al (1987) Prediction of protein secondary structure and active sites using the alignment of homologous sequences. *J Mol Biol* 195(4):957–961
- Wass MN, Kelley LA, Sternberg MJ (2010) 3DLigandSite: predicting ligand-binding sites using similar structures. *Nucleic Acids Res* 38(Web server issue):W469–W473
- Roy A, Yang J, Zhang Y (2012) COFACTOR: an accurate comparative algorithm for structure-based protein function annotation. *Nucleic Acids Res* 40(Web server issue):W471–W477
- Yang J, Roy A, Zhang Y (2013) Protein–ligand binding site recognition using complementary binding-specific substructure comparison and sequence profile alignment. *Bioinformatics* 29(20):2588–2595
- Lee HS, Im W (2013) Ligand binding site detection by local structure alignment and its performance complementarity. *J Chem Inf Model* 53(9):2462–2470
- Brylinski M, Feinstein WP (2013) eFindSite: improved prediction of ligand binding sites in protein models using meta-threading, machine learning and auxiliary ligands. *J Comput Aided Mol Des* 27(6):551–567
- Huang B, Schroeder M (2006) LIGSITEcsc: predicting ligand binding sites using the Connolly surface and degree of conservation. *BMC Struct Biol* 6:19
- Glaser F et al (2006) A method for localizing ligand binding pockets in protein structures. *Proteins* 62(2):479–488
- Halgren TA (2009) Identifying and characterizing binding sites and assessing druggability. *J Chem Inf Model* 49(2):377–389
- Capra JA et al (2009) Predicting protein ligand binding sites by combining evolutionary sequence conservation and 3D structure. *PLoS Comput Biol* 5(12):e1000585
- Huang B (2009) MetaPocket: a meta approach to improve protein ligand binding site prediction. *OMICS* 13(4):325–330
- Bray T et al (2009) SitesIdentify: a protein functional site prediction tool. *BMC Bioinform* 10:379
- Brylinski M, Skolnick J (2009) FINDSITE: a threading-based approach to ligand homology modeling. *PLoS Comput Biol* 5(6):e1000405
- Krivák R, Hoksza D (2015) Improving protein–ligand binding site prediction accuracy by classification of inner pocket points using local features. *J Cheminform* 7:12
- Krivák R, Hoksza D (2015) P2RANK: knowledge-based ligand binding site prediction using aggregated local features. *Algorithms for computational biology*. Springer International Publishing, Cham
- Jimenez J et al (2017) DeepSite: protein-binding site predictor using 3D-convolutional neural networks. *Bioinformatics* 33(19):3036–3042
- Krivák R, Hoksza D (2018) P2Rank: machine learning based tool for rapid and accurate prediction of ligand binding sites from protein structure. *J Cheminform* 10(1):39
- Santana CA et al (2020) GRaSP: a graph-based residue neighborhood strategy to predict binding sites. *Bioinformatics* 36(Suppl_2):i726–i734
- Kozlovskii I, Popov P (2020) Spatiotemporal identification of druggable binding sites using deep learning. *Commun Biol* 3(1):618
- Stepniewska-Dziubinska MM, Zielenkiewicz P, Siedlecki P (2020) Improving detection of protein–ligand binding sites with 3D segmentation. *Sci Rep* 10(1):5035
- Kandel J, Tayara H, Chong KT (2021) PURESNet: prediction of protein–ligand binding sites using deep residual neural network. *J Cheminform* 13(1):65
- Mylonas SK, Axenopoulos A, Daras P (2021) DeepSurf: a surface-based deep learning approach for the prediction of ligand binding sites on proteins. *Bioinformatics* 37(12):1681–1690
- Yan X et al (2022) PointSite: a point cloud segmentation tool for identification of protein ligand binding atoms. *J Chem Inf Model* 62(11):2835–2845

47. Li P et al (2022) RecurPocket: recurrent Lmser network with gating mechanism for protein binding site detection. In: 2022 IEEE international conference on bioinformatics and biomedicine (BIBM)
48. Aggarwal R et al (2022) DeepPocket: ligand binding site detection and segmentation using 3D convolutional neural networks. *J Chem Inf Model* 62(21):5069–5079
49. Abdollahi N et al (2023) NodeCoder: a graph-based machine learning platform to predict active sites of modeled protein structures. arXiv [q-bio.QM]
50. Evteev SA, Ereshchenko AV, Ivanenkov YA (2023) SiteRadar: utilizing graph machine learning for precise mapping of protein–ligand-binding sites. *J Chem Inf Model* 63(4):1124–1132
51. Li P et al (2023) GLPocket: a multi-scale representation learning approach for protein binding site prediction. In: Elkind E (ed) Proceedings of the thirty-second international joint conference on artificial intelligence, IJCAI-23. International Joint Conferences on Artificial Intelligence Organization. pp 4821–4828
52. Zhang Y et al (2023) EquiPocket: an E(3)-equivariant geometric graph neural network for ligand binding site prediction. arXiv [q-bio.BM]
53. Liu Y et al (2023) RefinePocket: an attention-enhanced and mask-guided deep learning approach for protein binding site prediction. *IEEE/ACM Trans Comput Biol Bioinf* 20(5):3314–3321
54. Smith Z et al (2024) Graph attention site prediction (GrASP): identifying druggable binding sites using graph neural networks with attention. *J Chem Inf Model* 64(7):2637–2644
55. Carbery A et al (2024) Learnt representations of proteins can be used for accurate prediction of small molecule binding sites on experimentally determined and predicted protein structures. *J Cheminform* 16(1):32
56. Sestak F et al (2024) VN-EGNN: E(3)-equivariant graph neural networks with virtual nodes enhance protein binding site identification. arXiv [cs.LG]
57. Jeevan K et al (2024) PURESNetV2.0: a deep learning model leveraging sparse representation for improved ligand binding site prediction. *J Cheminform* 16(1):66
58. Rives A et al (2021) Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proc Natl Acad Sci USA* 118(15):e2016239118
59. Hsu C et al (2022) Learning inverse folding from millions of predicted structures. In: Kamalika C et al (ed) Proceedings of the 39th international conference on machine learning. PMLR: Proceedings of Machine Learning Research. pp 8946–8970
60. Ke G et al (2017) LightGBM: a highly efficient gradient boosting decision tree. In: Guyon I et al (eds) Advances in neural information processing systems. Curran Associates, Inc., New York
61. Schrödinger LLC (2015) The PyMOL molecular graphics system, version 1.8
62. Ester M et al (1996) A density-based algorithm for discovering clusters in large spatial databases with noise. In: Proceedings of the second international conference on knowledge discovery and data mining. AAAI Press, Portland. pp 226–231
63. Jendele L et al (2019) PrankWeb: a web server for ligand binding site prediction and visualization. *Nucleic Acids Res* 47(W1):W345–W349
64. Capra JA, Singh M (2007) Predicting functionally important residues from sequence conservation. *Bioinformatics* 23(15):1875–1882
65. Jones JE, Chapman S (1924) On the determination of molecular fields.—I. From the variation of the viscosity of a gas with temperature. *Proc R Soc Lond Ser A Contain Pap Math Phys Character* 106(738):441–462
66. Yang J, Roy A, Zhang Y (2013) BioLiP: a semi-manually curated database for biologically relevant ligand–protein interactions. *Nucleic Acids Res* 41(Database issue):D1096–D1103
67. Chen K et al (2011) A critical comparative assessment of predictions of protein-binding sites for biologically relevant organic compounds. *Structure* 19(5):613–621
68. Krissinel E, Henrick K (2007) Inference of macromolecular assemblies from crystalline state. *J Mol Biol* 372(3):774–797
69. Armstrong DR et al (2020) PDBE: improved findability of macromolecular structure data in the PDB. *Nucleic Acids Res* 48(D1):D335–D343
70. Utgés JS (2024) bartongroup/FRAGSYS: second release. Zenodo. <https://doi.org/10.5281/zenodo.10606595>
71. Xu Q, Dunbrack RL Jr (2019) Principles and characteristics of biological assemblies in experimentally determined protein structures. *Curr Opin Struct Biol* 55:34–49
72. Pickens JC et al (2002) Anchor-based design of improved cholera toxin and *E. coli* heat-labile enterotoxin receptor binding antagonists that display multiple binding modes. *Chem Biol* 9(2):215–224
73. Hofmann E et al (1996) Structural basis of light harvesting by carotenoids: peridinin–chlorophyll–protein from *Amphidinium carterae*. *Science* 272(5269):1788–1791
74. Williams LK et al (2012) Order and disorder: differential structural impacts of myricetin and ethyl caffeate on human amylase, an anti-diabetic target. *J Med Chem* 55(22):10177–10186
75. Wilson EB (1927) Probable inference, the law of succession, and statistical inference. *J Am Stat Assoc* 22(158):209–212
76. Chen CI et al (2012) Structure of human POFUT2: insights into thrombospondin type 1 repeat fold and O-fucosylation. *EMBO J* 31(14):3183–3197
77. Liu RJ et al (2017) Structural basis for substrate binding and catalytic mechanism of a human RNA:m5C methyltransferase NSun6. *Nucleic Acids Res* 45(11):6684–6697
78. Ramirez-Rios S et al (2023) VASH1-SVBP and VASH2-SVBP generate different deetyrosination profiles on microtubules. *J Cell Biol* 222(2):e202205096
79. Lim SM et al (2016) The structure and catalytic mechanism of human sphingomyelin phosphodiesterase like 3a—an acid sphingomyelinase homologue with a novel nucleotide hydrolase activity. *FEBS J* 283(6):1107–1123
80. Finer-Moore J et al (2015) Crystal structure of the human tRNA m(1)A58 methyltransferase-tRNA(3)(Lys) complex: refolding of substrate tRNA allows access to the methylation target. *J Mol Biol* 427(24):3862–3876
81. Lee HJ, Cho HJ, Kang BS (2017) Crystal structure of human pyridoxal 5'-phosphate phosphatase (Chronophin) mutant - C221S. <https://doi.org/10.2210/pdb5gyn/pdb>
82. Blasiak LC et al (2006) Crystal structure of the non-haem iron halogenase SyrB2 in syringomycin biosynthesis. *Nature* 440(7082):368–371
83. Mascarenhas R et al (2023) Architecture of the human G-protein-methylmalonyl-CoA mutase nanoassembly for B(12) delivery and repair. *Nat Commun* 14(1):4332
84. Milani M et al (2011) FAD-binding site and NADP reactivity in human renalase: a new enzyme involved in blood pressure regulation. *J Mol Biol* 411(2):463–473
85. Liu L et al (2022) Autophosphorylation transforms DNA-PK from protecting to processing DNA ends. *Mol Cell* 82(1):177–189.e4
86. Yang H et al (2017) Mechanisms of mTORC1 activation by RHEB and inhibition by PRAS40. *Nature* 552(7685):368–373
87. Comajuncosa-Creus A et al (2024) Comprehensive detection and characterization of human druggable pockets through novel binding site descriptors. bioRxiv. <https://doi.org/10.1101/2024.03.14.584971>
88. Rabe WM, Tempel W, Nedyalkova L, Landry R, Arrowsmith CH, Edwards AM, Bountra C, Bochkarev A, Park H, Structural Genomics Consortium (SGC) (2015) Crystal structure of human sarcomeric mitochondrial creatine kinase. <https://doi.org/10.2210/pdb4z9m/pdb>
89. Durairaj J et al (2024) PLINDER: the protein–ligand interactions dataset and evaluation resource. bioRxiv. <https://doi.org/10.1101/2024.07.17.603955>
90. Stärk H et al (2022) EquiBind: geometric deep learning for drug binding structure prediction. arXiv [q-bio.BM]
91. Qiao Z et al (2024) State-specific protein–ligand complex structure prediction with a multiscale deep generative model. *Nat Mach Intell* 6(2):195–208
92. Schneuing A et al (2023) Structure-based drug design with equivariant diffusion models. arXiv [q-bio.BM]
93. UniProt C (2019) UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res* 47(D1):D506–D515
94. Berman H, Henrick K, Nakamura H (2003) Announcing the worldwide Protein Data Bank. *Nat Struct Biol* 10(12):980
95. WWpdb consortium (2019) Protein Data Bank: the single global archive for 3D macromolecular structure data. *Nucleic Acids Res* 47(D1):D520–D528

96. Varadi M et al (2022) PDBe and PDBe-KB: providing high-quality, up-to-date and integrated resources of macromolecular structures to support basic and applied research and education. *Protein Sci* 31(10):e4439
97. Ellaway JJ et al (2024) Identifying protein conformational states in the Protein Data Bank: toward unlocking the potential of integrative dynamics studies. *Struct Dyn* 11(3):034701
98. PDBe-KB consortium (2020) PDBe-KB: a community-driven resource for structural and functional annotations. *Nucleic Acids Res* 48(D1):D344–D353
99. PDBe-KB consortium (2022) PDBe-KB: collaboratively defining the biological context of structural data. *Nucleic Acids Res* 50(D1):D534–D542
100. MacGowan SA et al (2020) The Dundee resource for sequence analysis and structure prediction. *Protein Sci* 29(1):277–297
101. Jubb HC et al (2017) Arpeggio: a web server for calculating and visualising interatomic interactions in protein structures. *J Mol Biol* 429(3):365–371
102. Velankar S et al (2013) SIFTS: structure integration with function, taxonomy and sequences resource. *Nucleic Acids Res* 41(Database issue):D483–D489
103. Dana JM et al (2019) SIFTS: updated structure integration with function, taxonomy and sequences resource allows 40-fold increase in coverage of structure-based annotations for proteins. *Nucleic Acids Res* 47(D1):D482–D489
104. Virtanen P et al (2020) SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods* 17(3):261–272
105. Shannon CE (1948) A mathematical theory of communication. *Bell Syst Tech J* 27(3):379–423
106. Westbrook JD et al (2015) The chemical component dictionary: complete descriptions of constituent molecules in experimentally determined 3D macromolecules in the Protein Data Bank. *Bioinformatics* 31(8):1274–1278
107. Campanacci V et al (2019) Insight into microtubule nucleation from tubulin-capping proteins. *Proc Natl Acad Sci USA* 116(20):9859–9864
108. Paul N et al (2004) Recovering the true targets of specific ligands by virtual screening of the protein data bank. *Proteins* 54(4):671–680
109. Kellenberger E et al (2006) sc-PDB: an annotated database of drugable binding sites from the Protein Data Bank. *J Chem Inf Model* 46(2):717–727
110. Meslamani J, Rognan D, Kellenberger E (2011) sc-PDB: a database for identifying variations and multiplicity of 'druggable' binding sites in proteins. *Bioinformatics* 27(9):1324–1326
111. Desaphy J et al (2015) sc-PDB: a 3D-database of ligandable binding sites—10 years on. *Nucleic Acids Res* 43(Database issue):D399–D404
112. Mottin M et al (2017) Molecular dynamics simulations of Zika virus NS3 helicase: insights into RNA binding site activity. *Biochem Biophys Res Commun* 492(4):643–651
113. Hu L et al (2005) Binding MOAD (mother of all databases). *Proteins* 60(3):333–340
114. Benson ML et al (2008) Binding MOAD, a high-quality protein–ligand database. *Nucleic Acids Res* 36(Database issue):D674–D678
115. Ahmed A et al (2015) Recent improvements to binding MOAD: a resource for protein–ligand binding affinities and structures. *Nucleic Acids Res* 43(Database issue):D465–D469
116. Smith RD et al (2019) Updates to binding MOAD (mother of all databases): polypharmacology tools and their utility in drug repurposing. *J Mol Biol* 431(13):2423–2433
117. Hubbard TJ et al (1997) SCOP: a structural classification of proteins database. *Nucleic Acids Res* 25(1):236–239
118. Hubbard TJ et al (1998) SCOP, structural classification of proteins database: applications to evaluation of the effectiveness of sequence alignment methods and statistics of protein structural data. *Acta Crystallogr D Biol Crystallogr* 54(Pt 6 Pt 1):1147–1154
119. Lo Conte L et al (2000) SCOP: a structural classification of proteins database. *Nucleic Acids Res* 28(1):257–259
120. Hartshorn MJ et al (2007) Diverse, high-quality test set for the validation of protein–ligand docking performance. *J Med Chem* 50(4):726–741
121. Zhang Z et al (2011) Identification of cavities on protein surface using multiple computational approaches for drug binding site prediction. *Bioinformatics* 27(15):2083–2088
122. Schmidtke P et al (2010) Large-scale comparison of four binding site detection algorithms. *J Chem Inf Model* 50(12):2191–2200
123. Varadi M et al (2022) AlphaFold protein structure database: massively expanding the structural coverage of protein–sequence space with high-accuracy models. *Nucleic Acids Res* 50(D1):D439–D444
124. Wang R et al (2004) The PDBbind database: collection of binding affinities for protein–ligand complexes with known three-dimensional structures. *J Med Chem* 47(12):2977–2980
125. Wang R et al (2005) The PDBbind database: methodologies and updates. *J Med Chem* 48(12):4111–4119
126. Cheng T et al (2009) Comparative assessment of scoring functions on a diverse test set. *J Chem Inf Model* 49(4):1079–1093
127. Li Y et al (2014) Comparative assessment of scoring functions on an updated benchmark: 1. Compilation of the test set. *J Chem Inf Model* 54(6):1700–1716
128. Liu Z et al (2015) PDB-wide collection of binding data: current status of the PDBbind database. *Bioinformatics* 31(3):405–412
129. Liu Z et al (2017) Forging the basis for developing protein–ligand interaction scoring functions. *Acc Chem Res* 50(2):302–309
130. Hotelling H (1933) Analysis of a complex of statistical variables into principal components. *J Educ Psychol* 24(6):417–441
131. Fixman M (1962) Radius of gyration of polymer chains. *J Chem Phys* 36(2):306–310
132. Chen CR, Makhatadze GI (2015) ProteinVolume: calculating molecular van der Waals and void volumes in proteins. *BMC Bioinform* 16(1):101
133. Jubb H (2019) pdbtools
134. Schmidtke P et al (2010) fpocket: online tools for protein ensemble pocket detection and tracking. *Nucleic Acids Res* 38(Web Server issue):W582–W589
135. Jakubec D et al (2022) PrankWeb 3: accelerated ligand-binding site predictions for experimental and modelled protein structures. *Nucleic Acids Res* 50(W1):W593–W597
136. Jaccard P (1901) Distribution de la florine alpine dans la Bassin de Dranses et dans quelques regions voisines. *Bull Soc Vaud Sci Nat* 37:241–272
137. Jaccard P (1912) The distribution of the flora in the alpine zone.1. *New Phytol* 11(2):37–50
138. Durrant JD, de Oliveira CA, McCammon JA (2011) POVME: an algorithm for measuring binding-pocket volumes. *J Mol Graph Model* 29(5):773–776
139. Durrant JD et al (2014) POVME 2.0: an enhanced tool for determining pocket shape and volume characteristics. *J Chem Theory Comput* 10(11):5047–5056
140. Wagner JR et al (2017) POVME 3.0: software for mapping binding pocket flexibility. *J Chem Theory Comput* 13(9):4584–4592
141. Jordan SR and Chmait S (2015) Human GKRP bound to AMG2882 and Sorbitol-6-Phosphate. <https://doi.org/10.2210/pdb4px2/pdb>
142. Matthews BW (1975) Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochim Biophys Acta* 405(2):442–451
143. Webber C, Barton GJ (2003) Increased coverage obtained by combination of methods for protein sequence database searching. *Bioinformatics* 19(11):1397–1403
144. Scott MS, Barton GJ (2007) Probabilistic prediction and ranking of human protein–protein interactions. *BMC Bioinform* 8:239
145. Noguchi J et al (2008) Crystal structure of the covalent intermediate of human cytosolic beta-glucosidase. *Biochem Biophys Res Commun* 374(3):549–552
146. Pettersen EF et al (2021) UCSF ChimeraX: structure visualization for researchers, educators, and developers. *Protein Sci* 30(1):70–82
147. Utgés JS (2024) LBS-comparison results. Zenodo. <https://doi.org/10.5281/zenodo.13121414>
148. Utgés JS (2024) bartongroup/LBS-comparison: publication release. Zenodo. <https://doi.org/10.5281/zenodo.14033286>

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.