FANTASTIC ALLOSTERIC BINDING SITES AND WHY DEEP LEARNING CANNOT FIND THEM

Dhvani S. Vora*

Department of Biochemistry and Molecular Biology Mayo Clinic Rochester, MN 55905, USA {vora.dhvani}@mayo.edu

Shashank Yadav*

Department of Biomedical Engineering University of Arizona Tucson, AZ 85721, USA {shashank}@arizona.edu

Abstract

The discovery of druggable and structurally distinct allosteric sites across various protein classes has introduced new avenues for small molecules to modulate protein activity and, hence, cellular functions. Ligands that target allosteric sites may provide advantages like enhanced selectivity and often exhibit the possibility of targeting existing drug-resistant mutations. However, recent deep learning approaches show limited effectiveness in predicting allosteric sites, as demonstrated in the present study. We compare the performance of two deep learning methods, PUResNetV2.0 and VNEGNN, with Fpocket, a traditional geometry-based method and P2Rank, a geometry and machine learning ensemble approach.

1 INTRODUCTION

Proteins contain different types of functional sites that play critical roles in their biological activitiestwo key types being orthosteric and allosteric sites. Orthosteric sites are the primary binding sites where substrates or inhibitors directly interact with the protein, often leading to a functional or catalytic response. Allosteric sites, in contrast, are located distally from the orthosteric site and regulate protein function by binding to effectors that induce conformational changes (Christopoulos, 2002; Lu et al., 2019). Identifying these sites is essential for drug discovery. Multiple in silico approaches to predict ligand-binding sites (LBS) on three-dimensional protein structures have emerged in the past years, employing various approaches from geometry-based techniques like usage of grids, spheres, or tessellations; energy-based methods, conservation-based methods, template-based as well as ensemble methods (Xia et al., 2024; Zhao et al., 2020). While existing computational methods predict multiple LBS, predicting allosteric sites is still faced with significant challenges. In this study, we compare the performance of four popular methods for LBS prediction- Fpocket (Le Guilloux et al., 2009), P2Rank (Krivák & Hoksza, 2018), VNEGNN (Sestak et al., 2024) and PUResNetV2.0 (Jeevan et al., 2024). Fpocket is a geometry-based method, P2Rank a geometry and machine learning-based method, while VNEGNN and PUResNetV2.0 are deep learning (DL) methods. While all methods perform comparably on orthosteric sites of our dataset, the DL methods perform worse on allosteric sites. Our study provides compelling evidence to question if DL approaches truly outperform traditional geometry-based methods, especially on allosteric LBS.

2 Methods

All human proteins from the Allosteric Database (ASD) with binding sites of known allosteric function were compiled (He et al., 2024). An additional 22 proteins with experimentally reported orthoand allosteric sites were added to the dataset after a manual search of recent PDB entries, which were not included in the ASD. The protein structures with complexed orthosteric and allosteric ligands were downloaded from RCSB PDB (Berman et al., 2000), stripped of water molecules and ions. The ligands were then removed to generate a "clean" PDB structure used for pocket prediction. The Fpocket, P2Rank and VNEGNN methods were installed and run according to instructions on the GitHub repositories. PUResNetV2.0 predictions were obtained from the webserver. The predicted sites were manually verified with the ground truth. The results for each protein with all

^{*}Equal contribution

binding sites and predictions are saved as PyMOL session files (DeLano et al., 2002) accessible at https://doi.org/10.5281/zenodo.14977798.

3 RESULTS

Table 1: Model Recall Comparison

A manually curated dataset of 60 proteins with 152 LBS was compiled to evaluate the different prediction methods. Out of 152, 58 sites were orthosteric, while 94 were allosteric in nature. The predictions were counted as true positives if the predicted pocket completely or partially overlaps the ligand in the PDB structure. An example of the ligand-bound protein structure overlayed with prediction visualizations is

-	Ligand Binding Site Type							
Model	All	Orthosteri	c Allosteric					
PUResNet	0.460	0.845	0.223					
VNEGNN	0.328	0.707	0.096					
Fpocket	0.954	0.983	0.936					
P2Rank	0.789	1.000	0.659					

shown in Figure 1, and additional examples are provided in Appendix A1. A comparison of the models' predictions is summarized in Table 1 based on recall, considering it is impossible to obtain true negatives and false negatives without experimental support. A detailed site-wise comparison is included in Appendix A2. The geometry-based Fpocket ranks highest in predicting ligandable cavities on the protein across all types of pockets, as well as for predicting allosteric sites. However, P2Rank proves to be the superior method in predicting orthosteric sites. The more recent deep learning based PUResNetV2.0 and VNEGNN are not state-of-the-art for LBS prediction.



Figure 1: The tyrosine protein kinase ABL1, represented in green, complexed with asciminib (at the allosteric or A-site) and nilotinib (orthosteric or O-site), with the reference PDB structure 5MO4, overlayed with predictions of (a) PUResNetV2.0, (b) VNEGNN, (c) Fpocket and (d) P2Rank predicted pockets. The ligands are shown in blue, and the predicted pockets are shown in violet.

4 DISCUSSION

Despite the advanced capabilities of DL, which is quoted for its ability to model complex patterns, the results from this study suggest that DL techniques struggle to accurately predict allosteric LBS. PUResNetV2.0 may often predict allosteric sites that are present in close proximity to the orthosteric site, while distant sites are missed. The poor performance of VNEGNN and PUResNetV2.0 may be attributed to the absence of extensive, high-quality training data. In contrast, methods that rely on surface topology and shape analysis show more robust performance possibly due to their reliance on inherent structural features. DL methods typically require large amounts of annotated examples to generalize effectively, and the relatively small number of experimentally validated allosteric sites available for training further exacerbates this issue. Allosteric regulation often involves conformational changes that are difficult for models to capture, as these changes may not be directly observable from static protein structures alone. Without sufficient diversity in training examples, DL models may fail to learn the subtle patterns necessary for accurate allosteric LBS prediction. The high structural and functional group complexity inherent in the allosteric mechanisms may also be inadequately represented, which leads to DL methods being unable to fully capture such intricate biological processes. Geometry-based methods excel by detecting physical pocket characteristics and conformational flexibility without being overly dependent on training data. DL methods may improve on allosteric prediction by integrating molecular dynamics features, expanding allosteric datasets and explicitly modeling structural flexibility. The findings from this study highlight the need for investigation into the limitations of deep learning approaches and emphasize that, in certain contexts, traditional methods may still offer competitive performance in allosteric site prediction.

URM STATEMENT

The authors acknowledge that at least one key author of this work meets the URM criteria of the Tiny Papers Track at the "I Can't Believe It's Not Better" Workshop @ ICLR 2025.

REFERENCES

- Helen M Berman, John Westbrook, Zukang Feng, Gary Gilliland, Talapady N Bhat, Helge Weissig, Ilya N Shindyalov, and Philip E Bourne. The protein data bank. *Nucleic acids research*, 28(1): 235–242, 2000.
- Arthur Christopoulos. Allosteric binding sites on cell-surface receptors: novel targets for drug discovery. *Nature reviews Drug discovery*, 1(3):198–210, 2002.
- Warren L DeLano et al. Pymol: An open-source molecular graphics tool. *CCP4 Newsl. Protein Crystallogr*, 40(1):82–92, 2002.
- Jixiao He, Xinyi Liu, Chunhao Zhu, Jinyin Zha, Qian Li, Mingzhu Zhao, Jiacheng Wei, Mingyu Li, Chengwei Wu, Junyuan Wang, et al. Asd2023: towards the integrating landscapes of allosteric knowledgebase. *Nucleic Acids Research*, 52(D1):D376–D383, 2024.
- Kandel Jeevan, Shrestha Palistha, Hilal Tayara, and Kil T Chong. Puresnetv2. 0: a deep learning model leveraging sparse representation for improved ligand binding site prediction. *Journal of Cheminformatics*, 16(1):1–16, 2024.
- Radoslav Krivák and David Hoksza. P2rank: machine learning based tool for rapid and accurate prediction of ligand binding sites from protein structure. *Journal of cheminformatics*, 10:1–12, 2018.
- Vincent Le Guilloux, Peter Schmidtke, and Pierre Tuffery. Fpocket: an open source platform for ligand pocket detection. *BMC bioinformatics*, 10:1–11, 2009.
- Shaoyong Lu, Qiancheng Shen, and Jian Zhang. Allosteric methods and their applications: facilitating the discovery of allosteric drugs and the investigation of allosteric mechanisms. *Accounts* of Chemical Research, 52(2):492–500, 2019.
- Florian Sestak, Lisa Schneckenreiter, Johannes Brandstetter, Sepp Hochreiter, Andreas Mayr, and Günter Klambauer. Vn-egnn: E (3)-equivariant graph neural networks with virtual nodes enhance protein binding site identification. *arXiv preprint arXiv:2404.07194*, 2024.
- Ying Xia, Xiaoyong Pan, and Hong-Bin Shen. A comprehensive survey on protein-ligand binding site prediction. *Current Opinion in Structural Biology*, 86:102793, 2024.
- Jingtian Zhao, Yang Cao, and Le Zhang. Exploring the computational methods for protein-ligand binding site prediction. *Computational and structural biotechnology journal*, 18:417–426, 2020.

A APPENDIX

A.1 EXTENDED FIGURES



Figure 2: The beta2 adrenoceptor, represented in green, complexed with allosteric (A-site) and orthosteric (O-site) ligands shown in blue, with the reference PDB structure 5X7D, overlayed with predictions of (a) PUResNetV2.0, (b) VNEGNN, (c) Fpocket and (d) P2Rank predicted pockets. The predicted pockets are shown in violet.



Figure 3: The caspase-1, represented in green, complexed with allosteric (A-site) and orthosteric (O-site) ligands shown in blue, with the reference PDB structure 2FQQ, overlayed with predictions of (a) PUResNetV2.0, (b) VNEGNN, (c) Fpocket and (d) P2Rank predicted pockets. The predicted pockets are shown in violet.



Figure 4: The mu-type opioid receptor-G protein complex, represented in green, complexed with allosteric (A-site) and orthosteric (O-site) ligands shown in blue, with the reference PDB structure 8K9L, overlayed with predictions of (a) PUResNetV2.0, (b) VNEGNN, (c) Fpocket and (d) P2Rank predicted pockets. The predicted pockets are shown in violet.



Figure 5: The G-protein-coupled receptor 40 (GPR40), represented in green, complexed with allosteric (A-site) and orthosteric (O-site) ligands shown in blue, with the reference PDB structure 5KW2, overlayed with predictions of (a) PUResNetV2.0 (no pocket predicted), (b) VNEGNN, (c) Fpocket and (d) P2Rank predicted pockets. The predicted pockets are shown in violet.

A.2 DETAILED COMPARISON OF PREDICTIONS ON EACH LIGAND BINDING SITE

Table 2:	Summary	of pi	oteins,	ligand	types,	and	binding	infor-
mation.	Abbreviatio	ons:	PUR—	PUResN	letV2.0,	VN	N—VNE	EGNN,
FPK—FP	OCKET, P2	R—P2	2RANK	κ, Ο—Ο	rthosteri	ic, A-	-Alloste	eric, 1-
Correct, 0	-Incorrect.							

#	Gene ID	PDB Base	Multi- mer	LBS Index	LBS Type	Ligand Bound	PUR	VNN	FPK	P2R
1	CDK2	8VQ4	No	1	0	JWS	1	1	1	1
	-CCNE1			2	A	A1AC5	0	0	1	1
2	GBA1	8P3E	Yes	1	0	WYC	0	1	1	1
				2	0	WYC	0	0	0	1
				3	0	PN8	1	1	1	1
				4	0	PN8	1	1	1	1
				5	A	WSI	0	0	0	0
3	PTPN2	9C56	No	1	0	527	1	1	1	1
				2	A	FRJ	0	0	1	0
4	MALT1	8V4X	Yes	1	A	A1A	1	0	1	1
				2	A	A1A	1	0	1	1
				3	A	A1A	1	0	1	1
				4	A	A1A	1	0	1	1
5	MOR	8K9L	No	1	0	7V7	1	1	1	1
				2	A	VV9	0	0	1	1
6	MAT2A	8XAM	Yes	1	0	SAM	1	1	1	1
				2	0	SAM	0	1	1	1
				3	A	XRH	0	0	1	1
				4	A	VUO	0	0	1	0
7	KCNC1	8QUC	Yes	1	0	PCF	0	0	1	1
				2	A	WY9	0	0	1	1
8	ABL1	5MO4	No	1	0	STI	1	1	1	1
				2	A	AY7	0	0	1	1
9	CNR1	8IKH	No	1	0	FMN	1	1	1	1
				2	0	8D0	1	0	1	1

#	Gene ID	PDB Base	Multi- mer	LBS Index	LBS Type	Ligand Bound	PUR	VNN	FPK	P2R
				3	A	Q2L	0	0	1	0
10	ADRB2	5X7D	No	1	0	CAU	1	1	1	1
				2	A	8VS	0	0	1	1
				3	A	M3J	0	0	1	0
				4	A	KBY	0	0	1	0
11	AKT1	4GV1	No	1	0	UCB	0	1	1	1
				2	A	0XZ	1	0	1	1
12	AMD1	3EPA	No	1	0	MAO	0	1	1	1
				2	A	PUT	0	0	1	1
13	AD	2PIQ	No	1	0	TES	1	1	1	1
				2	A	17W	0	0	1	0
				3	A	RB1	0	0	1	0
14	ALB	2BX8	No	1	A	MYR	0	0	1	1
				2	A	MYR	0	0	1	1
				3	A	MYR		0	1	1
				4	A	MYR	0	0		
				5	A	RWF				
				0	A	AZQ				
				0	A			0		
15		5012	No	0	A			1	1	
15	AUKKA	SDNS	INO			5DN				
16	DDAE	11133/11	N.		A			1		
10	BKAF	IUWH	INO			ACP	1	1		
17	CASD1	2500	N		A		1	1		
17	CASPI	2FQQ	INO							
10	CASDZ	10111	Vee		A	FIU NVN	0	1		
18	CASP/	15HJ	res		A					
				$\begin{vmatrix} 2 \\ 3 \end{vmatrix}$		NYN	0	1	1	1
						SF1	0	0	1	1
10	CPS1	5DOU	Ves	1	0		1	1	1	1
17	CISI	5000	103	$\frac{1}{2}$	0		1	0	1	1
				3	A	NLG	0	0	1	1
				4	A	374	0	Ő	1	1
20	CRF1R	4K5Y	Yes	1	A	105	1	0	1	1
20			105	2	A	105	1	Ő	1	1
21	CTSK	5.194	No	1	0	INA	1	1	1	1
				2	Ā	6HM	0	0	0	0
22	DHPS	6PGR	Yes	1	A	8XY	1	0	1	1
				2	A	8XY	1	0	1	1
23	EP300	6PF1	No	1	0	COA	1	1	1	1
_		-		2	A	OJ7	1	0	1	1
24	FDPS	5DGM	No	1	0	ZOL	1	1	1	1
				2	Ā	7AM	1	0	1	0
				3	A	7AM	0	0	1	0
25	FFAR1	5KW2	No	1	0	MK6	0	1	1	1
-				2	A	6XQ	0	0	1	1
26	FLT3	4RT7	No	1	0	P30	1	1	1	1
	-			2	Α	CXS	0	0	0	0
27	GLS	5JYO	Yes	1	0	GLU	1	0	1	1
	-			2	A	ZBS	0	1	1	1
				3	A	63J	0	1	1	1
28	GMDS	5IN4	Yes	1	0	GDP	1	0	1	1
				2	0	NAP	1	1	1	1

ш	C		N/14			T		TANT	EDIZ	
#	Gene	PDB Base	Multi- mer	LBS Index	LBS Type	Ligand	PUK	VININ	FPK	P2R
20	CD) (5	5000	NT	3	A	6CK	0	0	1	
29	GRM5	SCGC	No		0	MES	0	0		
20	UD C1	and		2	A			1	1	
30	H3CI	61Q4	No		A	XIS	0	0		0
21	UDD	40.01	N	2	A	DUX	0	0	1	
31	нвв	4ROL	Yes		0	HEM		0		
20		5706	N	2	A	CND		1	1	
32	нказ	JZCO	INO			UNP		1	1	
22			Vac		A			1	1	
33	ШПІ	40MA	168			1NAF 50D		1	1	1
24		61.57	Vac			NAL	1	1	1	1
54	IDHS		168				1	1	1	1
35	KDM4A	5DGW	No			02V	1	0	1	1
26	KDM4A		No	1			1	1	1	1
30	KIF11	ΟΠΚΛ	ies			ADP 4A2		1	1	
						GCE	1	0	1	1
37	KIE18A	50CU	No	1			1	0	1	1
57	KIITOA	5000	INU	$\frac{1}{2}$		9V5		0	1	0
38	TUBB	50CU	No	1		GDP	1	0	1	1
50	TODD	5000	110	$\frac{1}{2}$			0	0	1	0
30		6NCE	No	1			1	1	1	1
57	ALOAS		110	2	A	AE7	0	0	1	1
40	MVH7	400	Ves	1	0	CRO	0	0	1	1
-10	141111/		103	2	A	20W	0	0	1	1
41	NRAS	8VM2	Yes	1	0	GTP	1	1	1	1
71	1111115	0 1112	103	2	A	EZZ	0	0	1	0
42	P2RX3	9IK1	Yes	1	0	ATP	1	0	1	1
	1	,	100	2	Ă	128	0	Ő	1	1
				3	Α	A1L	1	0	1	1
				4	A	PG4	0	0	1	1
43	PAK1	4ZJI	Yes	1	0	40Q	1	0	1	1
				2	A	59U	0	0	0	1
44	PANK3	3MK6	Yes	1	0	AN2	1	0	1	1
				2	A	PKZ	0	0	1	1
45	PARP1	6BHV	Yes	1	0	DQV	1	0	1	1
				2	A	09L	0	0	1	0
46	PCSK9	6U3X	No	1	A	63	1	0	1	1
47	MAPK14	8X3M	No	1	0	LBE	1	1	1	1
				2	A	B8Z	0	0	1	1
48	PGYB	5IKP	No	1	0	PLP	1	1	1	1
				2	A	AMP	0	0	1	1
49	PIK3CA	9C15	No	1	0	GNP	1	0	1	1
				2	A	70S	0	0	1	0
				3	A	A1A	0	0	1	1
				4	A	71K		0		
-	DOCT	2015		5	A	FBI	0	0		
50	PGC1A	3GN8	No		U O	MOF				
				$\frac{2}{2}$	A		0	0		
E 1		5710	N-		A	ILA IV1		1	1	
51	PIPNI	5119	INO					1		
50	DTDMF	61100	No				1	1	1	
52	LILIND	υπδΚ	INO	1	U	UI	1	1	1	1

#	Gene ID	PDB Base	Multi- mer	LBS Index	LBS Type	Ligand Bound	PUR	VNN	FPK	P2R
				2	A	FWB	0	0	0	1
53	PTPN11	5XZR	No	1	0	JZG	1	1	1	1
				2	A	50D	1	0	1	0
				3	A	DZV	0	0	1	0
54	DUSP3	8TK3	No	1	0	PO4	1	1	1	1
				2	A	I2X	0	0	1	0
				3	A	I2X	0	0	1	0
55	SIRT1	5BTR	Yes	1	0	4I5	1	1	1	1
				2	A	STL	0	0	1	1
				3	A	8QF	0	0	1	0
56	SIRT6	5MFP	Yes	1	0	AR6	1	1	1	1
				2	A	8L9	0	0	1	1
57	SMO	5L7I	Yes	1	A	VIS	1	1	1	1
58	THRB	2PIN	Yes	1	0	4HY	1	1	1	1
				2	A	LEG	0	0	1	1
59	TTR	5EZP	Yes	1	0	IPJ	1	0	1	1
				2	A	AJU	0	0	1	1
				3	A	04B	0	0	1	0
60	SMYD3	6YUH	No	1	0	SAM	1	1	1	1
				2	A	POW	0	0	1	0