

DrugDomain: the evolutionary context of drugs and small molecules bound to domains

Kirill E. Medvedev^{1,*}, R. Dustin Schaeffer¹, Nick V. Grishin^{1,2}

¹Department of Biophysics, University of Texas Southwestern Medical Center, Dallas, TX
75390

²Department of Biochemistry, University of Texas Southwestern Medical Center, Dallas, TX
75390

* Corresponding author

5323 Harry Hines Blvd., Dallas, Texas 75390-9050, Phone: 214-645-5946

E-mail: Kirill.Medvedev@UTSouthwestern.edu

Abstract

Interactions between proteins and small organic compounds play a crucial role in regulating protein functions. These interactions can modulate various aspects of protein behavior, including enzymatic activity, signaling cascades, and structural stability. By binding to specific sites on proteins, small organic compounds can induce conformational changes, alter protein-protein interactions, or directly affect catalytic activity. Therefore, many drugs available on the market today are small molecules (72% of all approved drugs in the last five years). Proteins are composed of one or more domains: evolutionary units that convey function or fitness either singly or in concert with others. Understanding which domain(s) of the target protein binds to a drug can lead to additional opportunities for discovering novel targets. The Evolutionary Classification Of protein Domains (ECOD) classifies domains into an evolutionary hierarchy that focuses on distant homology. Previously, no structure-based protein domain classification existed that included information about both the interaction between small molecules or drugs and the structural domains of a target protein. This data is especially important for multidomain proteins and large complexes. Here, we present the DrugDomain database that reports the interaction between ECOD domains of human target proteins and DrugBank molecules and drugs. The pilot version of DrugDomain describes the interaction of 5,160 DrugBank molecules associated with 2,573 human proteins. It describes domains for all experimentally determined structures of these proteins and incorporates AlphaFold models when such structures are unavailable. The DrugDomain database is available online: <http://prodata.swmed.edu/DrugDomain/>

Keywords: small molecules, drugs, target, domain, protein structure

Introduction

Proteins are vital components of cells, playing essential roles in regulating a myriad of cellular processes. The three-dimensional (3D) structure of a protein provides valuable information about protein interactions and functions. Protein domains are functionally, evolutionarily, and structurally distinct units that ensure evolutionary viability by fulfilling specific functions (Buljan and Bateman 2009; Grishin 2001). The identification and classification of protein domains into a hierarchy of evolutionary relations can lead to a better understanding of protein function by analyzing the known functions of their homologous relatives. Until recently, major structure-based classifications of protein domains have principally focused on the classification of experimentally determined protein structures: SCOP (Andreeva et al. 2020), and CATH (Sillitoe et al. 2021). Our team has developed and maintains the Evolutionary Classification Of protein Domains (ECOD), which primarily groups domains based on homology rather than topology (Cheng et al. 2014; Schaeffer et al. 2019). This feature aids in detecting instances of homology between domains with differing topologies. Another significant aspect of ECOD is its focus on distant homology, resulting in a comprehensive repository of evolutionary connections among categorized domains.

Artificial intelligence provides potent tools for scientific research across various fields, and structural computational biology is no exception. AlphaFold (AF), a recently developed deep learning method, demonstrated the capability to predict protein structure with atomic-level accuracy and has thus become an indispensable tool in structural biology (Jumper et al. 2021). Utilizing AF models, ECOD became one of the first databases to incorporate domain classification both for the entire human proteome (Schaeffer et al. 2023) and the whole proteomes of 48 model organisms (Schaeffer et al. 2024). AlphaFold has significantly expanded

the available tools for computational structural biology related to drug discovery, target prediction, protein-protein and protein-ligand interaction, and prediction of complex structures (Akdel et al. 2022; Medvedev et al. 2023a; Yang et al. 2023).

Many proteins convey their functions through interaction with small organic molecules. These interactions play crucial roles in various biological processes, including enzymatic reactions, signal transduction, and regulation of gene expression. Today the majority of FDA-approved drugs in the market and their generics are small molecules (Makurvet 2021). In the last five years, small molecules accounted for 72% of all FDA-approved drugs (178 out of a total of 247 drugs) (de la Torre and Albericio 2024). Each drug has one or more protein targets with an affinity defined through experimental methods. However, many proteins are multidomain and it is commonly unknown to which domain a drug binds in most cases. Especially in the cases where the target is a large receptor. Some multidomain proteins exhibit multidomain binding sites (Kruger et al. 2012). For example, human prostaglandin D-synthase (PDB: 3EE2) binds Nocodazole (DB08313) through both domains: Thioredoxin-like (ECOD: e3ee2A2) and Repetitive alpha hairpins (ECOD: e3ee2A1) (Weber et al. 2010). Human topoisomerase II beta (PDB: 3QX3) binds anticancer drug Etoposide (DB00773) using two out of four classified domains: helix-turn-helix (ECOD: 3qx3A1) and HAD domain-like (ECOD: e3qx3A11) (Wu et al. 2011). Homologous proteins exhibit highly similar active sites, capable of accommodating similar chemical compounds (Medvedev et al. 2021; Medvedev et al. 2019). Understanding the locations and principles that govern multi-domain drug binding sites may help us to better identify them in homologous proteins and complexes made up of homologous proteins. Chemical compound repositories such as PubChem (Kim et al. 2023) and DrugBank (Wishart et al. 2018) contain information about domains of the target protein derived from Pfam (Mistry et

al. 2021), which is a sequence-based classification. However, no indications are provided about to which domain(s) the compound interacts. Previously several resources were created that described ligand domain mapping using CATH and SCOP (structure-based) domain classification (Bashton et al. 2008; Bashton and Thornton 2010; Chalk et al. 2004) and Pfam database (sequence-based) (Mistry et al. 2021) classification (Kruger et al. 2012). However, these resources are out-of-date and presently unavailable for usage by the scientific community. Currently, no structure-based protein domain classification includes this type of data.

Here, we have developed the DrugDomain database, which reports those ECOD domains of proteins targeted by small molecules and drugs from DrugBank. The current version not only encompasses experimentally defined protein structures but also incorporates AlphaFold models in cases where such structures are unavailable. DrugDomain is available online at: <http://prodata.swmed.edu/DrugDomain/>

DrugDomain database features and statistics

We developed the DrugDomain database (<http://prodata.swmed.edu/DrugDomain/>) with web interfaces that display two types of database hierarchy: protein and molecule-centric. DrugDomain catalogs ECOD domains whose residues are located within 5Å of the DrugBank molecule's atoms. The distribution of DrugBank molecules interacting with ECOD homologous groups and architectures is shown in Figure 1. For the target proteins whose 3D structure was experimentally determined with any DrugBank molecule, the top three ECOD A-groups of the interacting domains include a+b complex topology, a/b three-layered sandwiches, and alpha arrays (Fig. 1A). The majority of small molecules interacting with a+b complex topology domains are associated with protein kinases, which are one of the most druggable protein

domains and is the domain most commonly encoded among genes associated with cancer (Anderson et al. 2023; Medvedev et al. 2023b; Wang et al. 2020). The a/b three-layered sandwiches are mostly represented by Rossmann-like proteins, which were shown to bind the majority of organic molecules superclasses (Medvedev et al. 2021; Medvedev et al. 2019). For the proteins lacking their interacting DrugBank molecule in any experimental structure classified by ECOD, we predicted interacting ECOD domains using the AlphaFill algorithm (Hekkelman et al. 2023) (see below for details). In this case, the three most populated ECOD A-groups containing modeled interacting DrugBank molecules include 1) the alpha bundles with G protein-coupled receptors (GPCRs), one of the most druggable protein domains together with kinases (Wang et al. 2020), followed by 2) the a/b three-layered sandwiches and 3) a+b complex topology (Fig. 1B).



Figure 1. Distribution of DrugBank molecules interacting with ECOD domains of target proteins. (A) Distribution of ECOD domains from experimentally determined 3D protein structures determined with associated DrugBank molecule, stratified by architecture (inside pie) and homologous group (outside donut). (B) Similar distribution for ECOD domains with

predicted small-molecule interactions using AlphaFill tool, AlphaFold models, and experimentally determined 3D structures that do not include DrugBank molecule.

The protein-centric hierarchy begins with the alphabetically sorted list of UniProt accessions (<http://prodata.swmed.edu/DrugDomain/proteins/>). Each accession leads to the dedicated webpage, which lists the UniProt accession, UniProt entry name, gene name, and protein name followed by the list of DrugBank molecules and drugs that target this protein. The list of molecules includes links to DrugDomain data webpages, DrugBank accession, molecule name, InChI Key, and SMILES formula. The molecule-centric hierarchy starts with the list of DrugBank accessions (<http://prodata.swmed.edu/DrugDomain/molecules/>). Each accession leads to the dedicated webpage, which lists the DrugBank accession, molecule name, InChI Key, and SMILES formula followed by the list of human proteins that are targets for this molecule. The list of proteins includes links to DrugDomain data webpages, UniProt accession, and protein names. The last level of our database hierarchy (DrugDomain data webpages) contains DrugBank accession, molecule name, InChI Key, and SMILES formula followed by the table of PDB structures and/or AlphaFold models with target particular protein. The table contains different information depending on the availability of experimentally determined 3D protein structures together with a particular DrugBank molecule. Such structures are available at RCSB Protein Data Bank for 2,149 molecules interacting with 783 target proteins, encompassing 5,702 PDB structures. For these cases, we identify residues whose atoms are located within 5Å of the DrugBank molecule's atoms and map these residues to existing ECOD domains. In such instances DrugDomain data webpage's table includes PDB accession linked to RCSB, downloadable PyMOL script, an indication that this interaction between molecule and protein target was confirmed experimentally and a list of ECOD domains interacting with the molecule

with links to the ECOD database (Fig. 2A). The PyMOL script downloads the PDB structure from RCSB, colors chains by different colors, colors residues interacting with the molecule in magenta and sets their representation as sticks (Fig. 2B). We specified links to DrugDomain data webpages for corresponding domains in the ECOD database (for example: <http://prodata.swmed.edu/ecod/complete/domain/e7du9A1>).

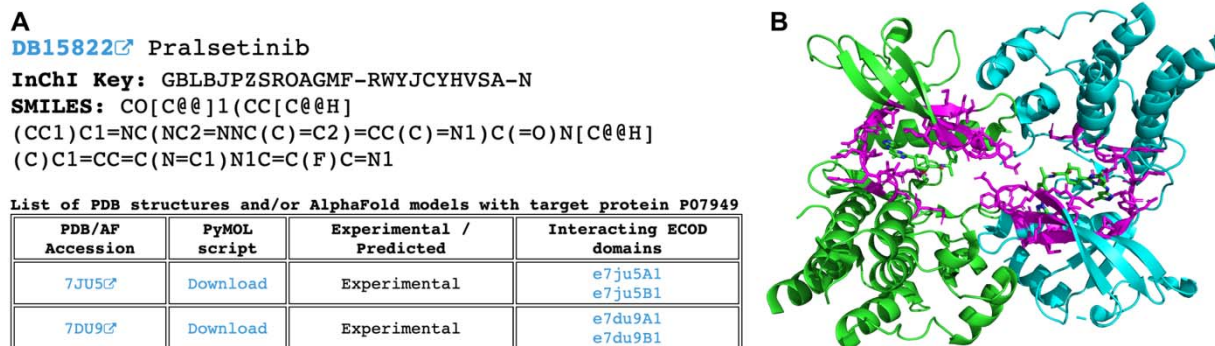


Figure 2. Example of the DrugDomain data webpage for cases with known PDB structure that includes DrugBank molecule and target protein. (A) Basic information about DrugBank molecule and table with experimental PDB structures. (B) Example of PyMOL script result – structure of human Proto-oncogene tyrosine-protein kinase receptor Ret (PDB: 7DU9) in complex with Pralsetinib (DB15822). Interacting residues are shown as sticks and colored in magenta.

The subset of DrugBank molecules that are not present in experimentally defined PDB structures includes 3,480 molecules targeting 2,361 human proteins. Among those proteins, experimentally defined PDB structures exist for 1,776 (75%) but without the DrugBank molecules of interest; only AlphaFold models exist for 573 (24%) proteins and there are no structural data for 12 (<1%) proteins (for example such huge proteins as E3 ubiquitin-protein ligase UBR4 (Q5T4S7), which contains 5,183 amino acids). In these cases, we applied the AlphaFill algorithm

(Hekkelman et al. 2023) that uses sequence and structure similarity to retrieve small molecules and ions from experimentally determined structures to predicted protein models by AlphaFold. Using AlphaFill models we identify residues whose atoms are located within 5Å of the DrugBank molecule's atoms of interest (if present) and map these residues to ECOD domains identified for the whole human proteome using AlphaFold models (Schaeffer et al. 2023; Schaeffer et al. 2024). If the molecule of interest is not present in the AlphaFold model, all other present ligands are considered as such. The DrugDomain data webpage's table in these cases includes AlphaFold accession, downloadable PyMOL script, an indication that this interaction between molecule and protein target was not confirmed experimentally (and was predicted), and a list of ECOD domains interacting with the molecule with links to the AlphaFold-based ECOD database. If experimentally defined PDB structures exist for a particular protein, the table also lists them with domains from the main ECOD database (with experimental structures), which correspond to domains from the AlphaFold model. In these instances, PyMOL script is not provided due to the absence of the molecule of interest in the PDB structure (Fig. 3A). The PyMOL script includes the AlphaFold model, colors it by rainbow, colors residues interacting with the molecule in magenta and sets their representation as sticks (Fig. 3B). Full information about ECOD domains and residues interacting with DrugBank molecules are also available for download in a plain text format.

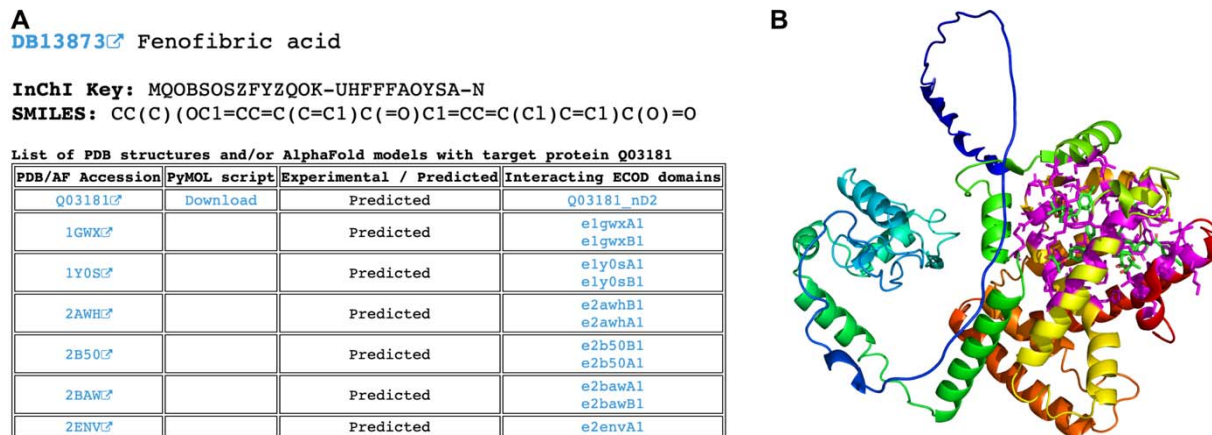


Figure 3. Example of the DrugDomain data webpage for cases without experimental PDB structure that includes DrugBank molecule and target protein. (A) Basic information about DrugBank molecule and table with AlphaFold model and corresponding PDB structures. **(B)** Example of PyMOL script result – AlphaFold model of human Peroxisome proliferator-activated receptor delta (Q03181) in complex with molecules of fenofibric acid (DB13873). Interacting residues are shown as sticks and colored in magenta.

Collection of DrugBank molecules dataset

To obtain all drugs and small molecules that target human proteins we retrieved all DrugBank accessions (Wishart et al. 2018) related to proteins from the reference human proteome (UP000005640) using UniProt KB (UniProt 2019). Overall, using this approach we obtained 6,506 DrugBank molecules (approximately 30% of DrugBank) that are associated with 3,193 human proteins. The DrugDomain database focuses on small molecules, so we have excluded 484 DrugBank entities representing biologics (“biotech” type of molecules in DrugBank). Additionally, 224 molecules were removed from the dataset due to incomplete DrugBank records, for example, birch bark extract (DrugBank accession: DB16536) and KW-6356 (DB17080). Each molecule can have multiple targets, enzymes, transporters, and carriers

associated with it in DrugBank. The purpose of this study is to catalog the interaction between DrugBank small molecules and the structural domain(s) of their target protein. Moreover, as we focus on targets that are human proteins, some molecules do not have targets in our dataset. For example, antibiotics that target bacterial proteins are excluded from the DrugDomain database (Fig. 4). Some molecules include protein groups as targets, which are often represented by various receptors. For example, Amoxapine (DB00543) targets, among others, the GABA receptor, which consists of 16 subunits, each of which is a different protein. In the majority of cases, drug targets include a particular receptor subunit known to bind the drug. However, in some cases, only protein groups are present as drug targets (for example for Mephentermine (DB01365)) and we removed such molecules from our current dataset. Thus, the remaining set, which was used for the development of the DrugDomain database version 1.0, includes 5,160 DrugBank molecules associated with 2,573 human proteins (Fig. 4).

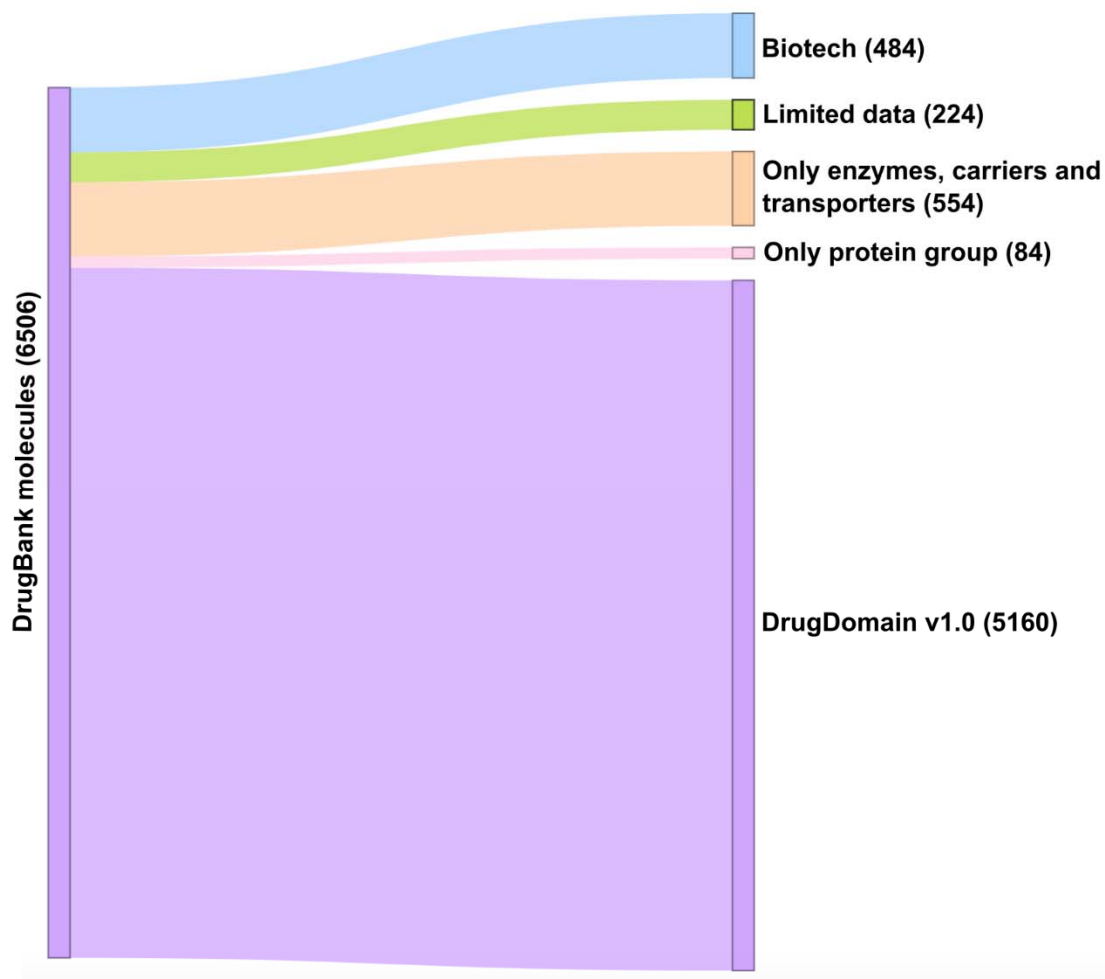


Figure 4. DrugBank molecules statistics for categories excluded and included in the DrugDomain database. The number of molecules for each category is shown in parentheses.

Future perspectives

The DrugDomain database version 1.0 represents the initial step of its development. We plan to significantly expand the range of represented molecules and include all DrugBank entities in our database. Additionally, we are going to incorporate not only the protein targets of molecules but also enzymes, transporters, and carriers associated with these molecules. Considering the conservation rate of residue positions will enhance the accuracy of predicting amino acids that

interact with small molecules. Finally, based on ECOD evolutionary classification and homologous evidence between proteins we anticipate to suggest new potential targets for known drugs. By focusing on evolutionarily conserved domains, we can prioritize targets that are likely to be functionally essential. Homologous proteins exhibit highly similar active sites, capable of accommodating similar chemical compounds. Better understanding along these lines opens up opportunities for discovering novel targets.

Competing interests

The authors declare that there are no competing interests associated with the manuscript.

Funding

The study is supported by grants from the National Institute of General Medical Sciences of the National Institutes of Health GM127390 (to N.V.G.), GM147367 (to R.D.S), the Welch Foundation I-1505 (to N.V.G.), the National Science Foundation DBI 2224128 (to N.V.G.).

CRedit Author Contribution

Kirill E. Medvedev: Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data Curation, Visualization, Writing - Original Draft, Project administration. **R.**

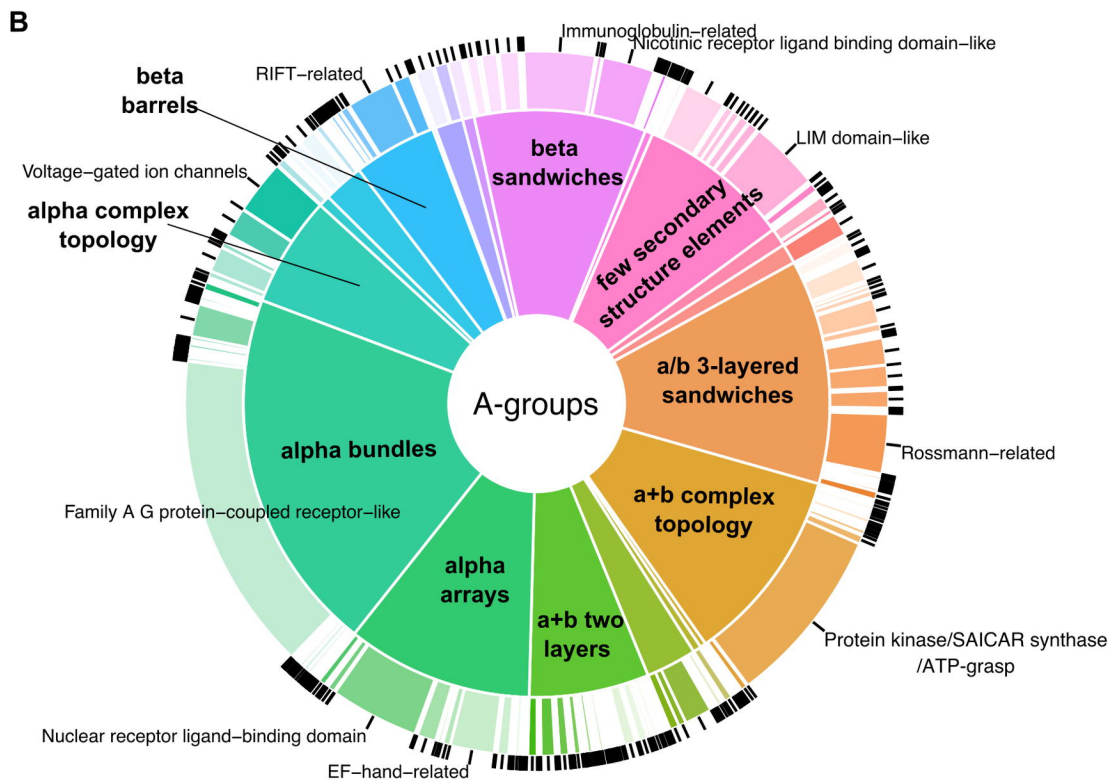
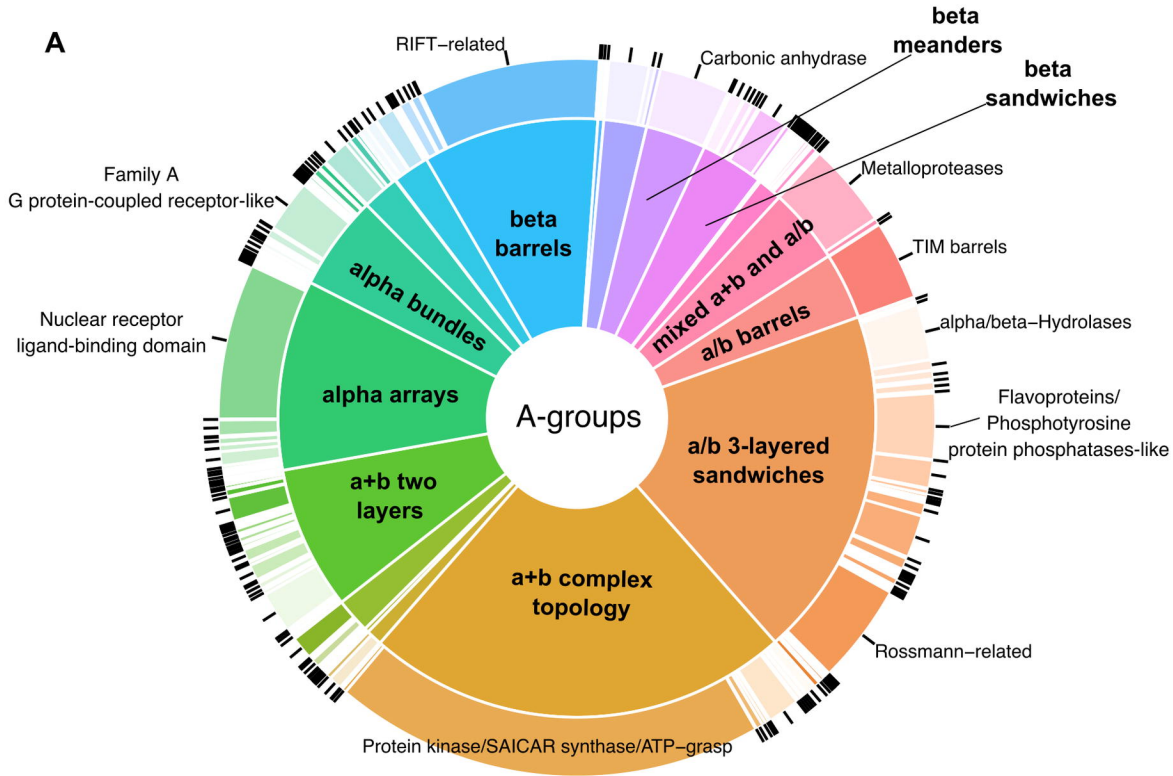
Dustin Schaeffer: Software, Writing - Review & Editing, Funding acquisition. **Nick V.**

Grishin: Conceptualization, Resources, Funding acquisition, Writing - Review & Editing.

References

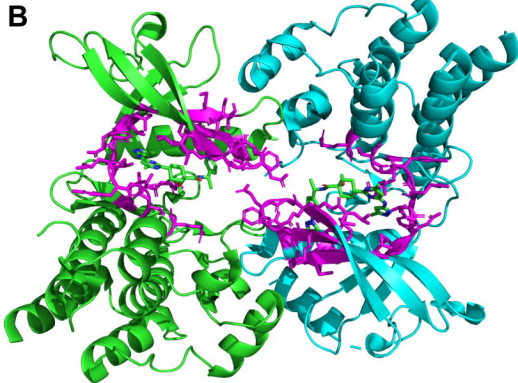
- Akdel M, Pires DEV, Pardo EP, Janes J, Zalevsky AO, Meszaros B et al. (2022) A structural biology community assessment of alphafold2 applications. *Nat Struct Mol Biol.* 29(11):1056-1067.
- Anderson B, Rosston P, Ong HW, Hossain MA, Davis-Gilbert ZW, Drewry DH. (2023) How many kinases are druggable? A review of our current understanding. *Biochem J.* 480(16):1331-1363.
- Andreeva A, Kulesha E, Gough J, Murzin AG. (2020) The scop database in 2020: Expanded classification of representative family and superfamily domains of known protein structures. *Nucleic Acids Res.* 48(D1):D376-D382.
- Bashton M, Nobeli I, Thornton JM. (2008) Procognate: A cognate ligand domain mapping for enzymes. *Nucleic Acids Res.* 36(Database issue):D618-622.
- Bashton M, Thornton JM. (2010) Domain-ligand mapping for enzymes. *J Mol Recognit.* 23(2):194-208.
- Buljan M, Bateman A. (2009) The evolution of protein domain families. *Biochem Soc Trans.* 37(Pt 4):751-755.
- Chalk AJ, Worth CL, Overington JP, Chan AW. (2004) Pdblig: Classification of small molecular protein binding in the protein data bank. *J Med Chem.* 47(15):3807-3816.
- Cheng H, Schaeffer RD, Liao Y, Kinch LN, Pei J, Shi S et al. (2014) Ecod: An evolutionary classification of protein domains. *PLoS Comput Biol.* 10(12):e1003926.
- de la Torre BG, Albericio F. (2024) The pharmaceutical industry in 2023: An analysis of fda drug approvals from the perspective of molecules. *Molecules.* 29(3).
- Grishin NV. (2001) Fold change in evolution of protein structures. *J Struct Biol.* 134(2-3):167-185.
- Hekkelman ML, de Vries I, Joosten RP, Perrakis A. (2023) Alphafill: Enriching alphafold models with ligands and cofactors. *Nat Methods.* 20(2):205-213.
- Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O et al. (2021) Highly accurate protein structure prediction with alphafold. *Nature.* 596(7873):583-589.
- Kim S, Chen J, Cheng T, Gindulyte A, He J, He S et al. (2023) Pubchem 2023 update. *Nucleic Acids Res.* 51(D1):D1373-D1380.
- Kruger FA, Rostom R, Overington JP. (2012) Mapping small molecule binding data to structural domains. *BMC Bioinformatics.* 13 Suppl 17(Suppl 17):S11.
- Makurvet FD. (2021) Biologics vs. Small molecules: Drug costs and patient access. *Medicine in Drug Discovery.* 9:100075.
- Medvedev KE, Kinch LN, Dustin Schaeffer R, Pei J, Grishin NV. (2021) A fifth of the protein world: Rossmann-like proteins as an evolutionarily successful structural unit. *J Mol Biol.* 433(4):166788.
- Medvedev KE, Kinch LN, Schaeffer RD, Grishin NV. (2019) Functional analysis of rossmann-like domains reveals convergent evolution of topology and reaction pathways. *PLoS Comput Biol.* 15(12):e1007569.
- Medvedev KE, Schaeffer RD, Chen KS, Grishin NV. (2023a) Pan-cancer structurome reveals overrepresentation of beta sandwiches and underrepresentation of alpha helical domains. *Sci Rep.* 13(1):11988.
- Medvedev KE, Schaeffer RD, Pei J, Grishin NV. (2023b) Pathogenic mutation hotspots in protein kinase domain structure. *Protein Sci.* 32(9):e4750.
- Mistry J, Chuguransky S, Williams L, Qureshi M, Salazar GA, Sonnhammer ELL et al. (2021) Pfam: The protein families database in 2021. *Nucleic Acids Res.* 49(D1):D412-D419.
- Schaeffer RD, Kinch L, Medvedev KE, Pei J, Cheng H, Grishin N. (2019) Ecod: Identification of distant homology among multidomain and transmembrane domain proteins. *BMC Mol Cell Biol.* 20(1):18.
- Schaeffer RD, Zhang J, Kinch LN, Pei J, Cong Q, Grishin NV. (2023) Classification of domains in predicted structures of the human proteome. *Proc Natl Acad Sci U S A.* 120(12):e2214069120.
- Schaeffer RD, Zhang J, Medvedev KE, Kinch LN, Cong Q, Grishin NV. (2024) Ecod domain classification of 48 whole proteomes from alphafold structure database using dpam2. *PLoS Comput Biol.* 20(2):e1011586.
- Sillitoe I, Bordin N, Dawson N, Waman VP, Ashford P, Scholes HM et al. (2021) Cath: Increased structural coverage of functional space. *Nucleic Acids Res.* 49(D1):D266-D273.

- UniProt C. (2019) Uniprot: A worldwide hub of protein knowledge. *Nucleic Acids Res.* 47(D1):D506-D515.
- Wang J, Yazdani S, Han A, Schapira M. (2020) Structure-based view of the druggable genome. *Drug Discov Today.* 25(3):561-567.
- Weber JE, Oakley AJ, Christ AN, Clark AG, Hayes JD, Hall R et al. (2010) Identification and characterisation of new inhibitors for the human hematopoietic prostaglandin d2 synthase. *Eur J Med Chem.* 45(2):447-454.
- Wishart DS, Feunang YD, Guo AC, Lo EJ, Marcu A, Grant JR et al. (2018) Drugbank 5.0: A major update to the drugbank database for 2018. *Nucleic Acids Res.* 46(D1):D1074-D1082.
- Wu CC, Li TK, Farh L, Lin LY, Lin TS, Yu YJ et al. (2011) Structural basis of type ii topoisomerase inhibition by the anticancer drug etoposide. *Science.* 333(6041):459-462.
- Yang Z, Zeng X, Zhao Y, Chen R. (2023) Alphafold2 and its applications in the fields of biology and medicine. *Signal Transduct Target Ther.* 8(1):115.



A
DB15822 [↗](#) Pralsetinib**InChI Key:** GBLBJPZSROAGMF-RWYJCYHVSA-N**SMILES:** CO[C@@]1(CC[C@@H](CC1)C1=NC(NC2=NNC(C)=C2)=CC(C)=N1)C(=O)N[C@@H](C)C1=CC=C(N=C1)N1C=C(F)C=N1**List of PDB structures and/or AlphaFold models with target protein P07949**

PDB/AF Accession	PyMOL script	Experimental / Predicted	Interacting ECOD domains
7JU5 ↗	Download	Experimental	e7ju5A1 e7ju5B1
7DU9 ↗	Download	Experimental	e7du9A1 e7du9B1



A

DB13873 [↗](#) Fenofibric acid

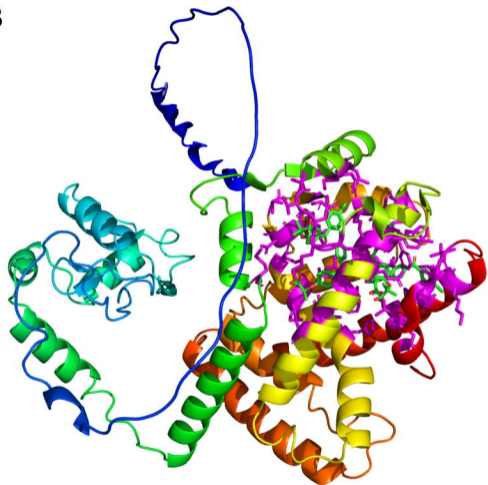
InChI Key: MQOBSOSZFYZQOK-UHFFFAOYSA-N

SMILES: CC(C)(OC1=CC=C(C=C1)C(=O)C1=CC=C(C1)C=C1)C(O)=O

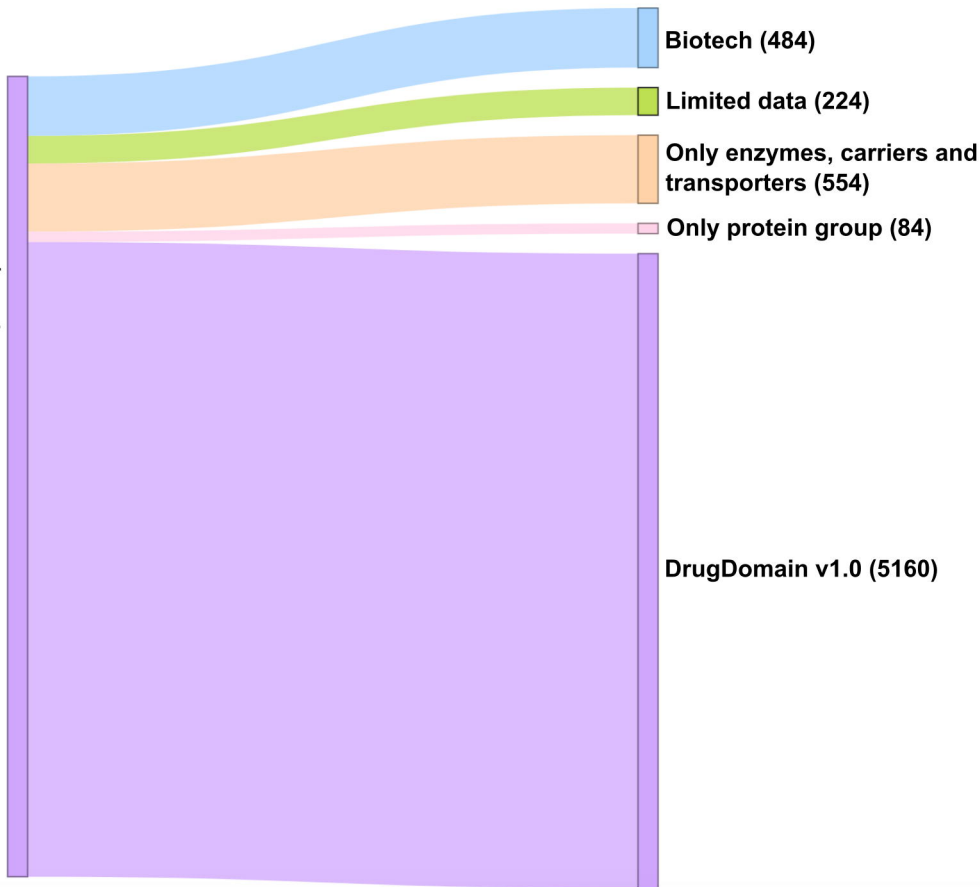
List of PDB structures and/or AlphaFold models with target protein Q03181

PDB/AF Accession	PyMOL script	Experimental / Predicted	Interacting ECOD domains
Q03181 ↗	Download	Predicted	Q03181_nD2
1GWX ↗		Predicted	e1gwxA1 e1gwxB1
1Y0S ↗		Predicted	ely0sA1 ely0sB1
2AWH ↗		Predicted	e2awhB1 e2awhA1
2B50 ↗		Predicted	e2b50B1 e2b50A1
2BAW ↗		Predicted	e2bawA1 e2bawB1
2ENV ↗		Predicted	e2envA1

B



DrugBank molecules (6506)



Biotech (484)

Limited data (224)

Only enzymes, carriers and transporters (554)

Only protein group (84)

DrugDomain v1.0 (5160)