


RESEARCH

eModel-BDB: a database of comparative structure models of drug-target interactions from the Binding Database

Misagh Naderi ¹, Rajiv Gandhi Govindaraj ¹ and Michal Brylinski ^{1,2,*}

¹Department of Biological Sciences, Louisiana State University, 202 Life Sciences Bldg, Baton Rouge, LA 70803, USA and ²Center for Computation & Technology, Louisiana State University, 2054 Digital Media Center, Baton Rouge, LA 70803, USA

*Correspondence address. Michal Brylinski, Louisiana State University, 202 Life Sciences Bldg, Baton Rouge, LA 70803, USA; Tel: +(225) 578-2791; Fax: +(225) 578-2597; E-mail: michal@brylinski.org  <http://orcid.org/0000-0002-6204-2869>

Abstract

Background: The structural information on proteins in their ligand-bound conformational state is invaluable for protein function studies and rational drug design. Compared to the number of available sequences, not only is the repertoire of the experimentally determined structures of holo-proteins limited, these structures do not always include pharmacologically relevant compounds at their binding sites. In addition, binding affinity databases provide vast quantities of information on interactions between drug-like molecules and their targets, however, often lacking structural data. On that account, there is a need for computational methods to complement existing repositories by constructing the atomic-level models of drug-protein assemblies that will not be determined experimentally in the near future. **Results:** We created eModel-BDB, a database of 200,005 comparative models of drug-bound proteins based on 1,391,403 interaction data obtained from the Binding Database and the PDB library of 31 January 2017. Complex models in eModel-BDB were generated with a collection of the state-of-the-art techniques, including protein meta-threading, template-based structure modeling, refinement and binding site detection, and ligand similarity-based docking. In addition to a rigorous quality control maintained during dataset generation, a subset of weakly homologous models was selected for the retrospective validation against experimental structural data recently deposited to the Protein Data Bank. Validation results indicate that eModel-BDB contains models that are accurate not only at the global protein structure level but also with respect to the atomic details of bound ligands. **Conclusions:** Freely available eModel-BDB can be used to support structure-based drug discovery and repositioning, drug target identification, and protein structure determination.

Keywords: eModel-BDB; eThread; eFindSite; BindingDB; homology modeling; comparative modeling; binding pocket prediction; similarity-based docking; protein function; drug targets

Background

Structural bioinformatics is becoming an increasingly important component of modern drug discovery. Despite significant advances in experimental methods to acquire protein structures, such as X-ray crystallography, nuclear magnetic resonance, and cryo-electron microscopy, technical limitations and expensive

procedures make it unlikely to have the experimental structures of all known protein sequences in the near future. For example, more than 110 million gene products have been included in the Reference Sequence Database [1] as of June 2018. In contrast, the number of experimentally determined protein structures in the Protein Data Bank (PDB) [2] is 140,824, which re-

Received: 4 October 2017; Revised: 15 March 2018; Accepted: 16 July 2018

© The Author(s) 2018. Published by Oxford University Press. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

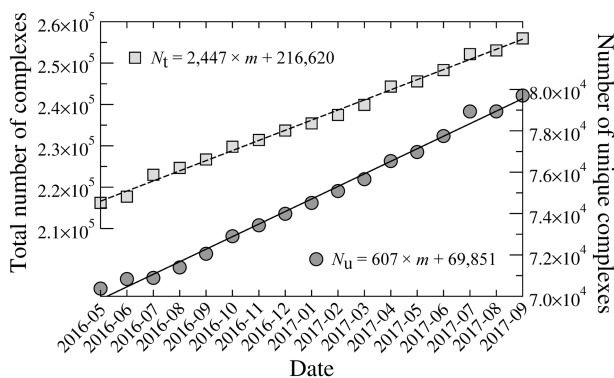


Figure 1: Deposition rate of ligand-bound structures to the Protein Data Bank. The total number of protein chains binding small molecules (light gray squares and a dashed line) is counted at any point in time. The number of unique complex structures is obtained by clustering individual chains at 80% sequence identity (dark gray circles and a solid line). N_t and N_u in the linear regression equations are the total and unique number of ligand-protein complexes, respectively, and m stands for month.

duces to 51,990 structures after removing similar proteins at 95% sequence identity. Genome sequencing currently produces as many as 13 million protein sequences each year, whereas only 8,872 protein structures are solved experimentally at the same time on average. Since this disparity between the number of available sequences and structures will likely continue to grow, high-throughput computational modeling is expected to play a significant role in biomedical sciences by generating three-dimensional models for those proteins whose structures will not be determined in the near future.

In addition to protein sequence and structure repositories, the Binding Database (BindingDB) provides comprehensive information on interactions between small, drug-like molecules and proteins considered to be drug targets collected from affinity measurements [3]. The BindingDB can be used to identify protein targets for small molecules and bioactive compounds for new proteins, as well as to conduct virtual screening with ligand-based methods. As of June 2018, BindingDB contained 1,450,120 binding data; however, only 2,291 ligand-protein crystal structures with BindingDB affinity measurements are available in the PDB. To bridge this gap, we created eModel-BDB, a new database of 200,005 high-quality drug-protein complex models involving 108,363 unique drug-like compounds and 2,791 proteins from the BindingDB. This repository was constructed with a state-of-the-art protocol to generate protein models in their ligand-bound conformational state, employing meta-threading, pocket detection, and protein structure and ligand chemical alignment techniques. eModel-BDB significantly expands the current structural information on known drug-protein complexes.

To fully appreciate the immensity of the structural data included in eModel-BDB, we estimate the time required to solve an equal number of drug-protein assemblies. Figure 1 shows that at the current pace, 2,447 ligand-bound protein structures containing 607 nonredundant complexes are deposited to the PDB each month. Therefore, it would take about 329 months for 200,005 unique complex structures to be determined experimentally. In contrast to other databases comprising protein models in the unbound conformational state generated through traditional structure modeling [4, 5], eModel-BDB includes annotated structure models of drug-protein complexes with known binding affinities. It provides high-quality data to

support structure-based drug discovery as well as repurposing of known drugs based on binding pocket and ligand similarities. In addition, the information provided by eModel-BDB can be utilized to facilitate experimental structure determination by developing protocols to stabilize proteins with ligands. The protocol to construct eModel-BDB described here is based entirely on open-source software to ensure that any researcher is able to produce new holo-protein models as more data become available in the PDB and BindingDB.

Methods

Protein structure modeling

Drug-bound protein complexes in eModel-BDB are generated with a template-based approach. The first phase is to construct structure models for single protein chains 50–999 amino acids in length obtained from BindingDB with eThread [6], which supports both close and remote homology modeling. eThread employs Modeller, a commonly used comparative modeling program (RRID:SCR_008395) [7], to build apo-protein structures based on alignments produced by 3-fold recognition algorithms, HH-suite [8], SparksX [9], and RaptorX [10]. Subsequently, side-chain positions and hydrogen-bonding networks in the initial models are improved with ModRefiner, a program to refine protein structures at the atomic level with a composite physics- and knowledge-based force field [11]. The quality assessment of refined models is carried out with ModelEvaluator [12] in terms of the estimated global distance test score (GDT-score). Of 5,501 BindingDB proteins, 4,906 were assigned an estimated GDT-score of ≥ 0.4 , indicating good-quality models [13, 14].

Ligand-binding site identification

Confident structure models with a GDT-score of ≥ 0.4 are further annotated with binding pockets and residues by eFindSite [15], which also computes a calibrated pocket confidence score. eFindSite detected 2,922 high-, 644 moderate-, and 776 low-confidence pockets in the eThread models of BindingDB targets. At this point, BindingDB drugs can be assigned to the predicted pockets with fingerprint-based virtual screening. Specifically, for a given drug-target pair in the BindingDB, we compute a rank of the drug against pockets detected by eFindSite, where the remaining BindingDB compounds are used as the background library. eFindSite conducts virtual screening with a set of molecular fingerprints and physicochemical properties calculated for ligands extracted from weakly homologous template structures [16]. The top one, two, and three pockets are considered for high-, moderate-, and low-confidence targets, respectively. A drug matches the predicted pocket if it is ranked within the top 20% of the screening library. With this protocol, we matched 108,363 drugs to binding pockets identified in their target proteins.

Similarity-based ligand docking

In the next phase, drug molecules are positioned within the predicted pockets with a two-step similarity-based docking protocol. This procedure exploits a significant structural conservation of ligand binding modes across remote homologs [17]. First, globally similar ligand-bound templates from the PDB, identified by eFindSite to have a similar pocket as the BindingDB protein, are superposed onto the apo-model. Proteins are aligned with Fr-TM-align [18] employing the template modeling score (TM-score) [19] to measure the global structure similarity. Subse-

quently, the BindingDB compound is aligned onto the template-bound ligand in order to place it in the predicted pocket of the apo-model. Here, we use chemical alignments constructed with kcombu [20], which also reports the chemical similarity between the BindingDB compound and the template-bound ligand measured by the Tanimoto coefficient (TC). Since a perfect case corresponds to both a TM-score and a TC of 1.0, we introduce a new metric, the perfect match distance (PMD), combining protein structure and ligand chemical similarity values:

$$PMD = \sqrt{(1 - TM\text{-score})^2 + (1 - TC)^2} \quad (1)$$

PMD is simply the Cartesian distance from the perfect match in the TM-score/TC space. In order to generate only high-quality holo-models, those cases with a PMD of >0.6 are excluded from the modeling process. This PMD cutoff was chosen to ensure that TM-score and TC for the selected templates are always above their individual significance threshold values of 0.4 [19, 20]. Further, for those cases having multiple ligand-bound templates satisfying the PMD criterion of ≤ 0.6 , a template with the shortest PMD is selected to build the holo-model of the BindingDB complex.

Complex structure refinement and assessment

In the final phase, protein models are rebuilt in the presence of the docked BindingDB compounds with Modeller. To make sure that the binding site is remodeled to accommodate the specific ligand, binding residues identified by eFindSite are removed from the initial model while enforcing the presence of secondary structure predicted by PSIPRED [21]. The resulting models are further annotated with the ligand-protein interaction score according to the distance-scaled finite ideal-gas reference (DFIRE) potential [22]. The eModel-BDB database contains atomic-level structure models of 200,005 drug-protein interactions from BindingDB, comprising 2,791 nonredundant proteins and 108,363 drug-like compounds. The list of eModel-BDB complexes is provided in Supplementary File S1.

Analyses

Data quality control

The quality control is pertinent to both protein structure modeling as well as binding site prediction. The quality of protein models is assessed with ModelEvaluator employing various structural features to compute the absolute quantitative score with a support vector regression. This approach assigns the GDT-score to a model by analyzing its secondary structure, relative solvent accessibility, contact map, and β -sheet structure. It has been demonstrated that GDT-scores estimated by ModelEvaluator for template-based models are highly correlated with the actual values with the Pearson correlation coefficient of 0.82 [12]. The first violin in Fig. 2 shows that eModel-BDB contains close and remote homology models with the median target-template sequence identity of 63%. The second violin indicates that the vast majority of these structures are accurate with the median estimated GDT-score for BindingDB proteins of 0.62. Further, as many as 78% of binding sites predicted by eFindSite to match BindingDB ligands have a high confidence of >0.8 . We showed previously that confidence scores of >0.8 assigned by eFindSite correspond to the Matthews correlation coefficient [23] of ≥ 0.6 for predicted binding residues [15]. On that account, we expect the majority of binding sites for BindingDB drugs to be correctly

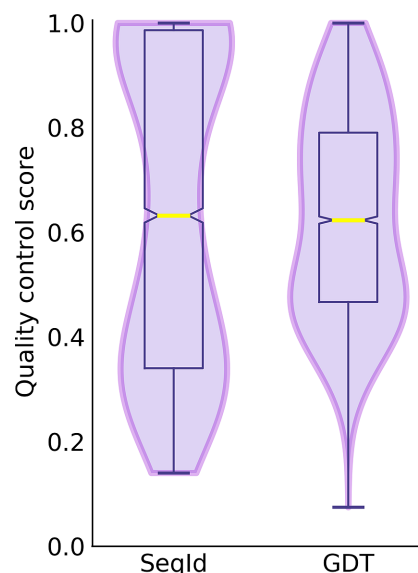


Figure 2: Violin and box plots for model quality control. The distribution of the target-template sequence identity (SeqId) and the global distance test (GDT) score estimated for structure models. Horizontal yellow lines represent median values.

annotated as well. Note that in contrast to other pocket predictors, eFindSite annotations and confidence assignments are, to some extent, independent of the accuracy of protein models.

The quality of complex models is controlled by imposing thresholds on the chemical similarity between BindingDB and PDB ligands as well as the global structure similarity between eThread models and ligand-bound templates from the PDB. Figure 3A shows the distribution of both parameters across eModel-BDB models. Encouragingly, the median TM-score and TC for ligand-bound templates used to build eModel-BDB are as high as 0.81 and 0.67, respectively. Previous studies show that the probability for a protein pair to belong to the same fold is 98% when the TM-score is close to 0.8 [24]. In addition, it was demonstrated that the root-mean-square deviation (RMSD) over ligand non-hydrogen atoms for similarity-based docking conducted with the TC in the range of 0.6–0.8 is typically 2–3 Å [25]. TM-score and TC values are combined into a single assessment score, the PMD, measuring the distance from the perfect match. Therefore, selecting template proteins with a lower TM-score to BindingDB targets requires their ligands to have a high TC and vice versa; selecting PDB ligands with a lower chemical similarity to BindingDB molecules requires a high global structure similarity between proteins. Figure 3B shows that the median PMD for eModel-BDB complex models is 0.46.

Data validation

In addition to the rigorous quality control maintained during dataset generation, eModel-BDB is validated retrospectively against experimental structures recently deposited to the PDB. The structure models of BindingDB interactions have been constructed with the PDB library as of 31 January 2017. Therefore, in order to validate eModel-BDB models, 7,012 experimental structures deposited to the PDB after February 2017 were considered. The validation protocol is made more challenging by including only remote homology models with a template-target sequence identity of $<40\%$. In order to maximize the validation coverage, we use the recently determined structures of eModel-BDB tar-

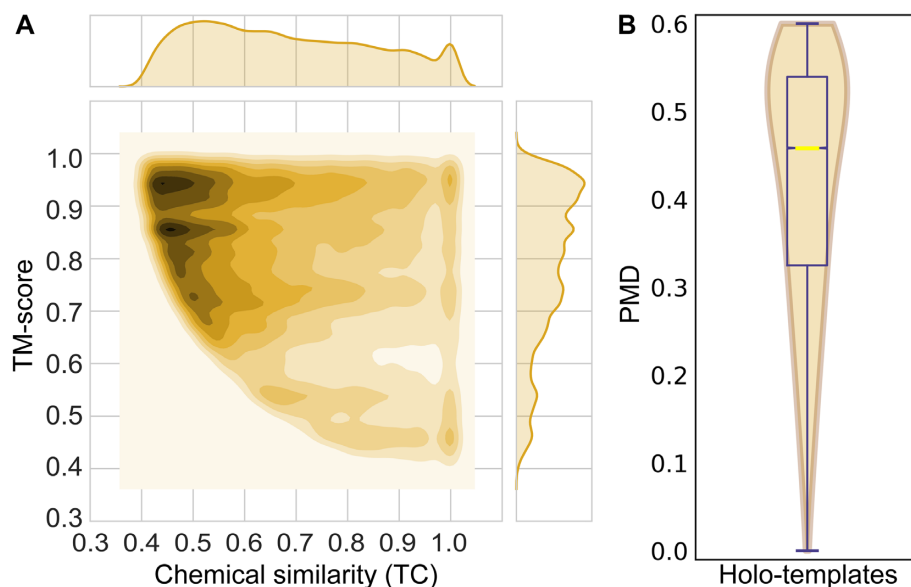


Figure 3: Similarities between target and holo-template proteins. (A) The chemical similarity between BindingDB and PDB ligands measured with the TC is plotted against the global structure similarity of *e*Thread models and ligand-bound templates from the PDB assessed by the TM-score. The two-dimensional contour plot is generated by smoothing the data with the kernel density estimation technique. One-dimensional histograms show the distribution of TC (top) and TM-score (right) values across eModel-BDB models. (B) Violin and box plot for the holo-template PMD combining TC and TM-score. The horizontal yellow line represents the median value.

gets and their homologs with at least 40% sequence identity. After applying these filters, 41 recently solved experimental structures selected from the PDB can be used to validate 161 *e*Thread models and 952 BindingDB reaction set IDs, comprising 39 target proteins, 52 pockets, and 881 compounds. This set is referred to as the validation dataset. The list of validation pairs is given in Supplementary File S2.

Protein structure modeling

The first violin in Fig. 4 shows that the median TM-score of eModel-BDB vs experimental structures is 0.85, with as many as 98.1% of the models having a TM-score of ≥ 0.4 . Clearly, the majority of structures are modeled by *e*Thread with a high accuracy. A representative example of the correctly predicted target structure is dihydrofolate reductase (DHFR) from *Streptococcus pyogenes* build on the crystal structure of DHFR from *Streptococcus pneumoniae* (PDB-ID: 3ix9, chain B, 36% sequence identity to the target) [26]. The *e*Thread model, whose estimated GDT-score is 0.92, was then used to construct a structure model for the BindingDB reactant set ID 00267770 consisting of DHFR complexed with BDBM50329610. This model is validated against the crystal structure of DHFR-UCP1106 from *Staphylococcus aureus* (PDB-ID: 5isp, chain X, 43% sequence identity to the target) released on 28 June 2017 [27]. Figure 5 shows the predicted weakly homologous model of DHFR-BDBM50329610 (purple) superposed on the experimental structure of DHFR-UCP1106 (gold). The eModel-BDB model is indeed highly accurate, with a TM-score of 0.95 and a $C\alpha$ -RMSD of 1.23 Å over 157 aligned residues. In addition, Fig. 6A shows that the estimated GDT-score employed in this study as the confidence measure to control the quality of protein models correlates with the accuracy of final models evaluated with the TM-score. On that account, the estimated GDT-score provides a robust quality assessment measure to control the quality of models in eModel-BDB.

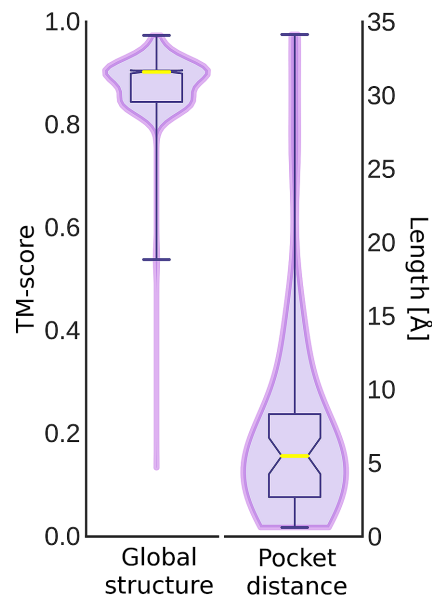


Figure 4: Violin and box plots for the distribution of validation scores. The accuracy is assessed for remote homology complex models in the validation set. The global structure similarity is measured with the TM-score. The pocket distance is measured between the predicted pocket center and the geometric center of the ligand in the experimental structure superposed onto the *e*Thread model. Horizontal yellow lines represent median values.

Binding pocket prediction

The accuracy of pocket prediction is validated by superposing the experimental holo structure onto the eModel-BDB model and then calculating the distance between the geometric center of a bound ligand in the experimental complex and the pocket center predicted by *e*FindSite in the model. The second violin in Fig. 4 shows that the median pocket distance is 5.5 Å, with 59.6%

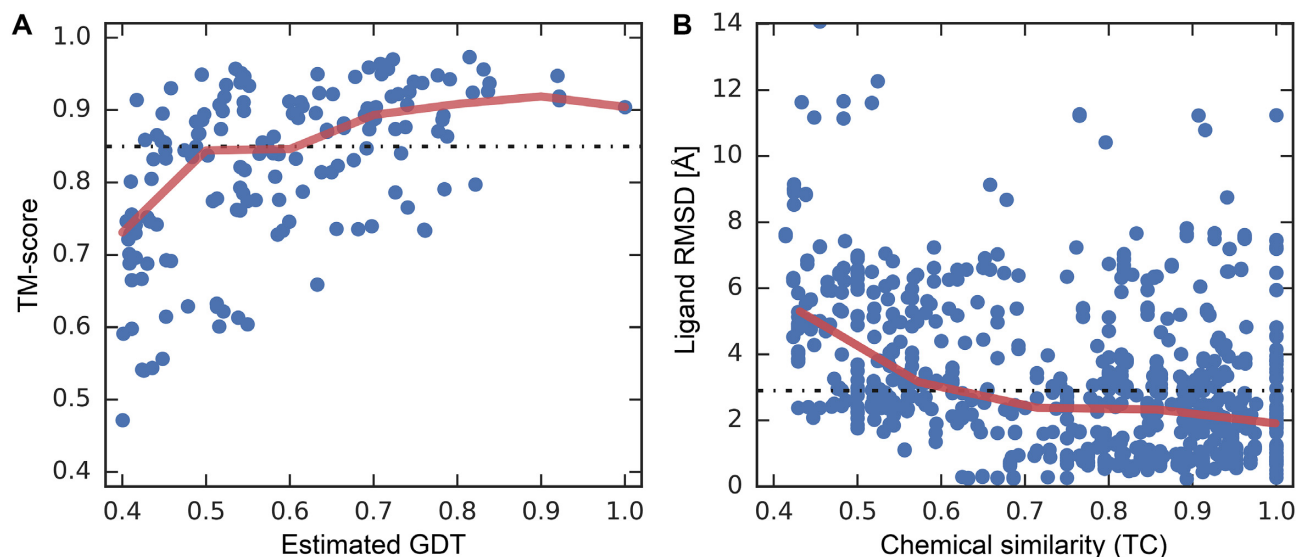


Figure 6: Analysis of structure modeling and ligand docking accuracy. The accuracy is assessed for remote homology complex models in the validation set. **(A)** Accuracy of global structure prediction evaluated by the TM-score with respect to the estimated GDT-score. **(B)** Accuracy of similarity-based docking with respect to the chemical similarity between BindingDB and PDB ligands measured with the TC. The ligand RMSD is calculated over non-hydrogen atoms according to the chemical alignment reported by kcombu. Solid red lines show the average prediction accuracy for binned GDT-score values in A and the chemical similarity in B. Dotted black lines mark the median TM-score in A and RMSD in B across all benchmarking cases.



Figure 5: Representative example of a structure model constructed by eThread. The model of DHFR (purple ribbons) complexed with BDBM50329610 is superposed onto the crystal structure of homologous DHFR from *S. aureus* (gold ribbons) complexed with UCP1106. Ligands bound to target proteins are shown as solid sticks (BDBM50329610 is purple and UCP1106 is gold) with non-carbon atoms colored by atom type (O—red, N—blue).

of pockets predicted within 6 Å; therefore, most eFindSite annotations are accurate. A binding site in the model of vitamin D receptor (VDR) is a representative example of a pocket predicted with eFindSite. This model was constructed for the BindingDB reactant set ID 50662356 based on human retinoic acid receptor RXR-alpha (PDB-ID: 4nqa, chain H, 38% sequence identity to the target) [28]. Although the GDT-score estimated for the VDR model is 0.62, indicating a moderately accurate structure, the top-ranked binding site annotated by eFindSite is assigned a high confidence of 94.2%. Figure 7 shows the VDR model (purple ribbons) superposed onto the crystal structure of vitamin D3 receptor A (gold ribbons) complexed with a synthetic analog of $1\alpha,25$ -dihydroxyvitamin D3 (PDB-ID: 5nky, chain A, 66% sequence identity to the target) released on 24 May 2017 [29]. Not only does the VDR model align well to the experimental struc-

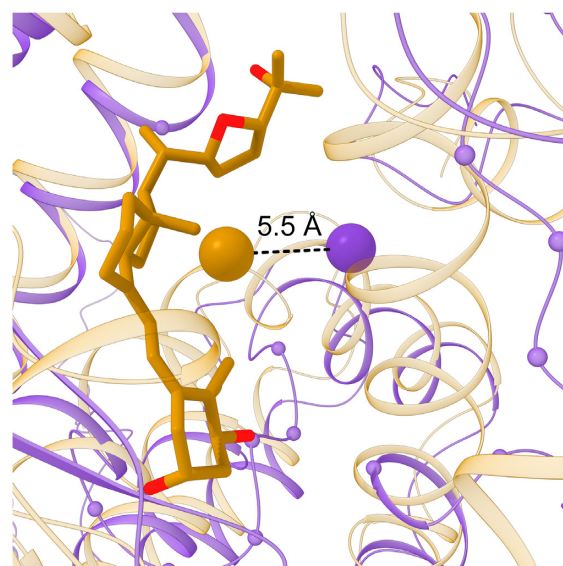


Figure 7: Representative example of a binding site detected by eFindSite. The model of VDR (purple ribbons) is superposed onto the crystal structure of homologous VDR from human (gold ribbons) complexed with a synthetic analog of vitamin D (gold and red sticks). $C\alpha$ atoms of binding residues predicted in the VDR model by eFindSite are shown as small spheres. Large spheres connected by a dashed black line are placed at the location of the predicted pocket center (purple) and the geometric center of vitamin D analog (gold).

ture, with a TM-score of 0.90 and a $C\alpha$ -RMSD of 2.13 Å over 235 residues, but the predicted pocket center (purple sphere) is only 5.5 Å away from the geometric center of vitamin D analog (gold sphere).

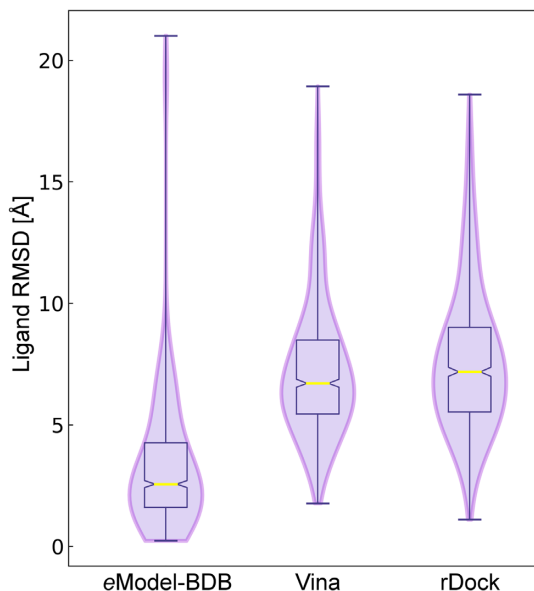


Figure 8: Violin and box plots for the docking accuracy. The accuracy is assessed for remote homology complex models in the validation set. The ligand RMSD is calculated over non-hydrogen atoms according to the chemical alignment reported by kcombu. The performance of similarity-based docking employed to construct eModel-BDB is compared to that of AutoDock Vina and rDock. Horizontal yellow lines represent median values.

Ligand docking

Finally, we calculate the RMSD over non-hydrogen atoms between the BindingDB drug in the eModel-BDB structure and the bound ligand in the superposed experimental complex. Here, we employ a subset of 37 models selected from the validation set whose pocket centers are predicted within 8 Å. The first violin in Fig. 8 shows that the median ligand RMSD is 2.6 Å and it is ≤ 3 Å for 58.1% of BindingDB compounds. The model of the BindingDB reactant set ID 50974033 consisting of tyrosine-protein kinase (TSK) complexed with BDBM50399512 is selected to exemplify the accuracy of complex structures in eModel-BDB. The model of TSK built on the crystal structure of human hemopoietic cell kinase (HCK) (PDB-ID: 1qcf, chain A, 39% sequence identity to the target) [30] by eThread is assigned an estimated GDT-score of 0.61. Subsequently, the complex model of TSK-BDBM50399512 was constructed by similarity-based docking employing the crystal structure of HCK bound to a pyrazolo-pyrimidine inhibitor (PDB-ID: 3vs7, chain B, 37% sequence identity to the target) [31]. This HCK complex was selected as the best ligand-bound template based on the high TM-score of 0.64 and TC of 0.53, yielding the shortest PMD of 0.35. Figure 9 shows the validation of the modeled TSK-BDBM50399512 by the experimental structure of Bruton's tyrosine kinase (Btk) bound to an anti-cancer drug, ibrutinib (PDB-ID: 5p9i, chain A, TM-score: 0.92, 58% sequence identity to the target) released on 24 May 2017 [32]. Kcombu reports a significant chemical alignment between BDBM50399512 and ibrutinib with a TC of 0.68 (Fig. 9A). Upon the superposition of TSK and Btk proteins, the RMSD between BDBM50399512 docked to the model and ibrutinib bound in the experimental structure calculated over the chemical alignment reported by kcombu is 2.62 Å (Fig. 9B). These results verify that the computer-generated TSK-BDBM50399512 model for the BindingDB reactant set ID 50974033 is correct.

The model of the BindingDB reactant set ID 50103430 consisting of cytochrome P450 17A1 (CYP17A1) complexed with BDBM50061174 is selected to exemplify the accuracy of complex structures in eModel-BDB. The model of CYP17A1 built on the crystal structure of human microsomal cytochrome P450 2A6 (PDB-ID: 1z11, chain A, 29% sequence identity to the target) [33] by eThread is assigned an estimated GDT-score of 0.69. Subsequently, the complex model of CYP17A1-BDBM50061174 was constructed by similarity-based docking employing the crystal structure of CYP17A1 bound to abiraterone, a steroidal prostate cancer drug (PDB-ID: 3ruk, chain D) [34]. The CYP17A1-abiraterone complex was selected as the best ligand-bound template based on the high TM-score of 0.84 and TC of 0.89, yielding the shortest PMD of 0.19. Figure 9 shows the validation of the modeled CYP17A1-BDBM50061174 by the experimental structure of CYP17A1-(R)-orterone (PDB-ID: 5irq, chain B, 64% sequence identity to the target) released on 15 March 2017 [35]. Kcombu reports a significant chemical alignment between steroidal BDBM50061174 and nonsteroidal (R)-orterone with a TC of 0.54 (Fig. 9A). Upon the superposition of CYP17A1 proteins, the RMSD between BDBM50061174 docked to the model and (R)-orterone bound in the experimental structure calculated over the chemical alignment reported by kcombu is 2.95 Å (Fig. 9B). These results verify that the computer-generated CYP17A1-BDBM50061174 model for the BindingDB reactant set ID 50103430 is correct.

The similarity-based docking procedure employed to construct ligand-bound structures in eModel-BDB superposes target ligands onto template molecules selected from the PDB according to the chemical alignment reported by kcombu. One may expect that superposing target compounds onto chemically similar template ligands yields more accurate binding poses than those generated from chemically less similar template molecules. Indeed, Fig. 6B shows that the target-template chemical similarity measured with the TC correlates with the docking accuracy evaluated with the RMSD of ligand poses constructed based on target-template alignments. These results are in line with other studies reporting that the average RMSD values for similarity-based docking methods are generally below 2 Å when the target-template similarities are above 0.7 [36]. The performance of similarity-based docking employed to construct eModel-BDB is also compared to that of AutoDock Vina [37] and rDock [38]. In contrast to the median ligand RMSD of 2.6 Å for eModel-BDB complexes, the median RMSD values for BindingDB drugs docked to eFindSite pockets with AutoDock Vina and rDock are 6.7 Å and 7.2 Å, respectively (Fig. 8). We note that similarity-based docking was demonstrated to outperform traditional docking when the target-template similarity is greater than 0.4 [36], which was employed as the TC threshold to construct eModel-BDB complex models. Overall, the quality assessment as well as independently obtained validation results demonstrate that the eModel-BDB database contains high-quality models closely resembling experimentally determined structures, not only at the global structure level but also at the level of binding pockets and bound ligands.

Discussion

eModel-BDB is generated to support rational drug development projects. These data can directly aid structure-based drug discovery pipelines and protein function analysis by providing atomic-level models of a large set of drug-protein interactions with known affinities curated in the BindingDB. An important

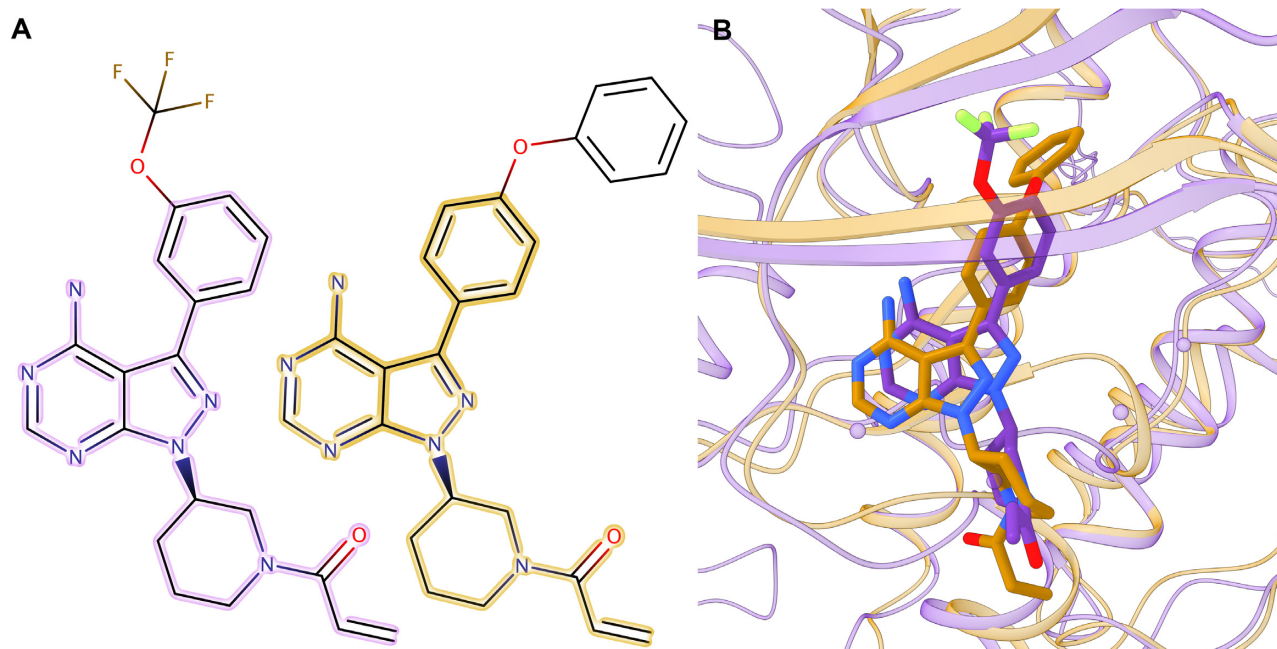


Figure 9: Representative example of a complex structure constructed by similarity-based docking. **(A)** Chemical alignment between BDBM50399512 (left) and ibrutinib (right) reported by kcombu. The 26 equivalent atom pairs that constitute the maximum common substructure are outlined in purple in BDBM50399512 and in gold in ibrutinib. **(B)** The model of TSK (purple ribbons) is superposed onto the crystal structure of HCK (gold ribbons) complexed with ibrutinib (gold sticks). $C\alpha$ atoms of binding residues identified in the TSK model by eFindSite are shown as purple spheres, whereas the target compound, BDBM50399512, docked into the predicted pocket is represented by purple sticks. Non-carbon atoms in BDBM50399512 and ibrutinib are colored by atom type (O—red, N—blue, F—green).

application of eModel-BDB is computational drug repositioning, i.e., finding new indications for existing drugs [39]. Although drug repurposing holds a significant promise to speed up drug development, particularly for diseases considered to be unprofitable, its major bottleneck is the scarce structural information on druggable pockets. On that account, a diverse dataset of small, drug-like molecules bound to high-quality models with accurately annotated pockets provides an invaluable resource for drug repositioning employing sequence order-independent pocket matching algorithms [40–43]. It is noteworthy that computational drug repurposing has suggested new opportunities to combat tuberculosis [44, 45], malaria [46], and rare diseases [47, 48].

Binding sites in eModel-BDB can also be matched to pockets predicted in potential drug targets in order to determine whether these proteins are druggable or not. If a new pocket aligns well with drug-bound pockets in eModel-BDB, then it is likely going to be druggable. That being the case, our data can be utilized at the outset of drug discovery, in the target identification phase. Finally, ligand binding can significantly help stabilize a protein, particularly from the point of view of conformational stability [49]. eModel-BDB can, therefore, inform crystallography efforts by suggesting possible compounds binding to certain protein targets at either the active or allosteric sites in order to increase the chances of successful crystallization.

Availability of supporting data

Structure models in eModel-BDB are named according to the BindingDB reactant set IDs, which can be obtained by searching the BindingDB at <https://www.bindingdb.org>. This procedure is illustrated in Fig. 10. The BindingDB can be searched either by protein and compound names (Fig. 10A) or by the

target sequence (Fig. 10B and 10C). Next, the complex of interest can be selected from the list of hits (Fig. 10D) in order to download the corresponding SDfile of the complex (Fig. 10E). The BindingDB reactant set ID, e.g., 00267770, is stored inside the SDfile (Fig. 10F). The reactant set ID can then be used to find the detailed information on the BindingDB website, e.g., https://www.bindingdb.org/jsp/dbsearch/Summary.ki.jsp?reactant_set_id=00267770 (Fig. 10G) as well as access the structure model in eModel-BDB, e.g., http://brylinski.cct.lsu.edu/pub/eModelBDB.php?reactant_set_id=00267770 (Fig. 10H). A web-based interface to query eModel-BDB is provided at <http://brylinski.org/emodel-bdb-0>. Data are available to download from the GigaScience GigaDB database [50].

Additional files

SupplementaryFileS1.csv

SupplementaryFileS2.csv

Figure 10. Procedure to obtain eModel-BDB complexes via the BindingDB website. The target complex can be identified based on either the protein (red arrows and boxes) or the ligand of interest (blue arrows and boxes). Common actions that a user needs to perform are colored in green. **(A)** Specific ligands and proteins can be searched directly for on the BindingDB website. **(B, C)** Alternatively, target proteins can be found with the blast search. **(D)** A complex of interest can then be selected in order to **(E)** generate and download an SDfile. **(F)** The BindingDB reactant set ID stored inside the SDfile is used to **(G)** view a web page containing detailed information about the target complex as well as **(H)** access the corresponding eModel-BDB structure model named according to the BindingDB reactant set ID.

Abbreviations

BindingDB: Binding Database; Btk: Bruton's tyrosine kinase; DFIRE: distance-scaled finite ideal-gas reference; DHFR: dihydrofolate reductase; GDT-score: global distance test score; HCK: hemopoietic cell kinase; PDB: Protein Data Bank; PMD: perfect match distance; RMSD: root-mean-square deviation; TC: Tanimoto coefficient; TM-score: template modeling score; TSK: tyrosine-protein kinase; VDR: vitamin D receptor.

Competing interests

The authors declare that they have no competing interests.

Funding

Research reported in this publication was supported by the National Institute of General Medical Sciences of the National Institutes of Health (under award R35GM119524).

Author contributions

M.B. prepared protein models, annotated binding pockets, and validated protein structures and pockets. M.N. constructed, refined, and validated drug-bound models. R.G.G. prepared case studies. M.N. and M.B. wrote the article.

Acknowledgements

The authors are grateful to Louisiana State University for providing computing resources.

References

- O'Leary NA, Wright MW, Brister JR, et al. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res* 2016;**44**(D1):D733–45.
- Berman HM, Westbrook J, Feng Z, et al. The Protein Data Bank. *Nucleic Acids Res* 2000;**28**(1):235–42.
- Liu T, Lin Y, Wen X, et al. BindingDB: a web-accessible database of experimentally determined protein-ligand binding affinities. *Nucleic Acids Res* 2007;**35**(Database issue):D198–201.
- Castrignano T, De Meo PD, Cozzetto D, et al. The PMDB Protein Model Database. *Nucleic Acids Res* 2006;**34**(Database issue):D306–9.
- Sanchez R, Pieper U, Mirkovic N, et al. MODBASE, a database of annotated comparative protein structure models. *Nucleic Acids Res* 2000;**28**(1):250–3.
- Brylinski M, Lingam D. eThread: a highly optimized machine learning-based approach to meta-threading and the modeling of protein tertiary structures. *PLoS One* 2012;**7**(11):e50200.
- Sali A, Blundell TL. Comparative protein modelling by satisfaction of spatial restraints. *J Mol Biol* 1993;**234**(3):779–815.
- Remmert M, Biegert A, Hauser A, et al. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nat Methods* 2011;**9**(2):173–5.
- Yang Y, Faraggi E, Zhao H, et al. Improving protein fold recognition and template-based modeling by employing probabilistic-based matching between predicted one-dimensional structural properties of query and corresponding native properties of templates. *Bioinformatics* 2011;**27**(15):2076–82.
- Ma J, Wang S, Zhao F, et al. Protein threading using context-specific alignment potential. *Bioinformatics* 2013;**29**(13):i257–65.
- Xu D, Zhang Y. Improving the physical realism and structural accuracy of protein models by a two-step atomic-level energy minimization. *Biophys J* 2011;**101**(10):2525–34.
- Wang Z, Tegge AN, Cheng J. Evaluating the absolute quality of a single protein model using structural features and support vector machines. *Proteins* 2009;**75**(3):638–47.
- He Z, Alazmi M, Zhang J, et al. Protein structural model selection by combining consensus and single scoring methods. *PLoS One* 2013;**8**(9):e74006.
- Mirabello C, Adelfio A, Pollastri G. Reconstructing protein structures by neural network pairwise interaction fields and iterative decoy set construction. *Biomolecules* 2014;**4**(1):160–80.
- Brylinski M, Feinstein WP. eFindSite: improved prediction of ligand binding sites in protein models using meta-threading, machine learning and auxiliary ligands. *J Comput Aided Mol Des* 2013;**27**(6):551–67.
- Feinstein WP, Brylinski M. eFindSite: enhanced fingerprint-based virtual screening against predicted ligand binding sites in protein models. *Mol Inform* 2014;**33**(2):135–50.
- Brylinski M, Skolnick J. FINDSITE: a threading-based approach to ligand homology modeling. *PLoS Comput Biol* 2009;**5**(6):e1000405.
- Pandit SB, Skolnick J. Fr-TM-align: a new protein structural alignment method based on fragment alignments and the TM-score. *BMC Bioinformatics* 2008;**9**:531.
- Zhang Y, Skolnick J. Scoring function for automated assessment of protein structure template quality. *Proteins* 2004;**57**(4):702–10.
- Kawabata T. Build-up algorithm for atomic correspondence between chemical structures. *J Chem Inf Model* 2011;**51**(8):1775–87.
- Jones DT. Protein secondary structure prediction based on position-specific scoring matrices. *J Mol Biol* 1999;**292**(2):195–202.
- Zhang C, Liu S, Zhu Q, et al. A knowledge-based energy function for protein-ligand, protein-protein, and protein-DNA complexes. *J Med Chem* 2005;**48**(7):2325–35.
- Matthews BW. Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochim Biophys Acta* 1975;**405**(2):442–51.
- Xu J, Zhang Y. How significant is a protein structure similarity with TM-score = 0.5? *Bioinformatics* 2010;**26**(7):889–95.
- Brylinski M. Nonlinear scoring functions for similarity-based ligand docking and binding affinity prediction. *J Chem Inf Model* 2013;**53**(11):3097–112.
- Lee J, Yennawar NH, Gam J, et al. Kinetic and structural characterization of dihydrofolate reductase from *Streptococcus pneumoniae*. *Biochemistry* 2010;**49**(1):195–206.
- Reeve SM, Scocchera E, Ferreira JJ, et al. Charged propargyl-linked antifolates reveal mechanisms of antifolate resistance and inhibit trimethoprim-resistant MRSA strains possessing clinically relevant mutations. *J Med Chem* 2016;**59**(13):6493–500.
- Lou X, Toresson G, Benod C, et al. Structure of the retinoid X receptor alpha-liver X receptor beta (RXRalpha-LXRbeta) heterodimer on DNA. *Nat Struct Mol Biol* 2014;**21**(3):277–81.
- Belorusova AY, Martinez A, Gandara Z, et al. Structure-activity relationship study of vitamin D analogs with oxolane group in their side chain. *Eur J Med Chem* 2017;**134**:86–96.

30. Schindler T, Sicheri F, Pico A, et al. Crystal structure of Hck in complex with a Src family-selective tyrosine kinase inhibitor. *Mol Cell* 1999;3(5):639–48.
31. Saito Y, Yuki H, Kuratani M, et al. A pyrrolo-pyrimidine derivative targets human primary AML stem cells in vivo. *Sci Transl Med* 2013;5(181):181ra52.
32. Bender AT, Gardberg A, Pereira A, et al. Ability of Bruton's tyrosine kinase inhibitors to sequester Y551 and prevent phosphorylation determines potency for inhibition of Fc receptor but not B-cell receptor signaling. *Mol Pharmacol* 2017;91(3):208–19.
33. Yano JK, Hsu MH, Griffin KJ, et al. Structures of human microsomal cytochrome P450 2A6 complexed with coumarin and methoxsalen. *Nat Struct Mol Biol* 2005;12(9):822–3.
34. DeVore NM, Scott EE. Structures of cytochrome P450 17A1 with prostate cancer drugs abiraterone and TOK-001. *Nature* 2012;482(7383):116–9.
35. Petrunak EM, Rogers SA, Aube J, et al. Structural and functional evaluation of clinically relevant inhibitors of steroidogenic cytochrome P450 17A1. *Drug Metab Dispos* 2017;45(6):635–45.
36. Kawabata T, Nakamura H. 3D flexible alignment using 2D maximum common substructure: dependence of prediction accuracy on target-reference chemical similarity. *J Chem Inf Model* 2014;54(7):1850–63.
37. Trott O, Olson AJ. AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *J Comput Chem* 2010;31(2):455–61.
38. Ruiz-Carmona S, Alvarez-Garcia D, Foloppe N, et al. rDock: a fast, versatile and open source program for docking ligands to proteins and nucleic acids. *PLoS Comput Biol* 2014;10(4):e1003571.
39. Haupt VJ, Schroeder M. Old friends in new guise: repositioning of known drugs with structural bioinformatics. *Brief Bioinform* 2011;12(4):312–26.
40. Brylinski M. eMatchSite: sequence order-independent structure alignments of ligand binding pockets in protein models. *PLoS Comput Biol* 2014;10(9):e1003829.
41. Brylinski M. Local alignment of ligand binding sites in proteins for polypharmacology and drug repositioning. *Methods Mol Biol* 2017;1611:109–22.
42. Konc J, Janezic D. ProBiS algorithm for detection of structurally similar protein binding sites by local structural alignment. *Bioinformatics* 2010;26(9):1160–8.
43. Ren J, Xie L, Li WW, et al. SMAP-WS: a parallel web service for structural proteome-wide ligand-binding site comparison. *Nucleic Acids Res* 2010;38(Web Server issue):W441–4.
44. Kinnings SL, Liu N, Buchmeier N, et al. Drug discovery using chemical systems biology: repositioning the safe medicine Comtan to treat multi-drug and extensively drug resistant tuberculosis. *PLoS Comput Biol* 2009;5(7):e1000423.
45. Anand P, Chandra N. Characterizing the pocketome of *Mycobacterium tuberculosis* and application in rationalizing polypharmacological target selection. *Sci Rep* 2014;4:6356.
46. Ramakrishnan G, Chandra N, Srinivasan N. Exploring anti-malarial potential of FDA approved drugs: an in silico approach. *Malar J* 2017;16(1):290.
47. Brylinski M, Naderi M, Govindaraj RG, et al. eRepo-ORP: exploring the opportunity space to combat orphan diseases with existing drugs. *J Mol Biol* 2018;430:2266–73.
48. Govindaraj RG, Naderi M, Singha M, et al. Large-scale computational drug repositioning to find treatments for rare diseases. *NPJ Syst Biol Appl* 2018;4:13.
49. Deller MC, Kong L, Rupp B. Protein stability: a crystallographer's perspective. *Acta Crystallogr F Struct Biol Commun* 2016;72(Pt 2):72–95.
50. Naderi M, Govindaraj RG, Brylinski M. Supporting data for “eModel-BDB: a database of comparative structure models of drug-target interactions from the Binding Database.” *GigaScience Database* 2018 <http://dx.doi.org/10.5524/100396>.