

# DrugPred\_RNA – Structure-based druggability predictions for RNA binding sites

Illimar Hugo Rekand and Ruth Brenk\*

Department of Biomedicine, University of Bergen, Jonas Lies Vei, 5020 Bergen, Norway

\*corresponding author: [ruth.brenk@uib.no](mailto:ruth.brenk@uib.no)

## Abstract

RNA is an emerging target for drug discovery. However, like for proteins, not all RNA binding sites are equally suited to be addressed with conventional drug-like ligands. To this end, we have developed the structure-based druggability predictor DrugPred\_RNA to identify druggable RNA binding sites. Due to the paucity of annotated RNA binding sites, the predictor was trained on protein pockets, albeit using only descriptors that can be calculated for both, RNA and protein binding sites. DrugPred\_RNA performed well in discriminating druggable from less druggable binding sites for the protein set and delivered sensible predictions for selected RNA binding sites. Further, the majority of drug-like ligands contained in a data set of RNA-containing pockets were found in pockets predicted to be druggable, further adding confidence to the performance of DrugPred\_RNA. The method is robust against conformational changes in the binding site and can contribute to direct drug discovery efforts for RNA targets.

## Introduction

The vast majority of targets for approved drugs are proteins.<sup>1,2</sup> However, in recent years it has been increasingly realized that also RNAs constitute promising drug targets as they play a key role in many biological processes, can fold into diverse 3D structures, and specifically recognize small molecules.<sup>3-6</sup> By targeting RNA the functions of currently undruggable protein-mediated pathways and the non-coding transcriptome can be modulated and thus the size of the druggable genome can be increased considerably.<sup>3</sup> A prime example of an RNA drug target is the bacterial ribosome, where protein synthesis is inhibited through binding of small molecules.<sup>7</sup> This is illustrated by linezolid, an FDA-approved antibiotic, which acts by binding to ribosomal RNA (Figure 1).<sup>8</sup> Another active research area are RNA-binding splicing modifiers for the treatment of spinal muscular atrophy with several compounds in clinical trials.<sup>9,10</sup> Riboswitches, which are non-coding RNA structures in the 5'-UTR and regulate gene expression through metabolite binding are new RNA drug targets for antibiotics.<sup>11,12</sup> For example, compounds binding to the flavin mononucleotide (FMN) riboswitch, e. g. ribocil and 5FDQD, have been shown to kill bacteria (Figure 1).<sup>13,14</sup> Riboflavin is known to bind to both the FMN riboswitch and riboflavin kinase. In both binding sites, the ligand is recognized in a similar way forming hydrophobic contacts and hydrogen bonds between the surrounding residues and the pteridine ring system, the dimethylbenzene ring, and the ribose chain. This fact nicely illustrates the capability of RNA to make specific molecular interactions with a wide variety of functional groups and ligand surfaces.<sup>3</sup>

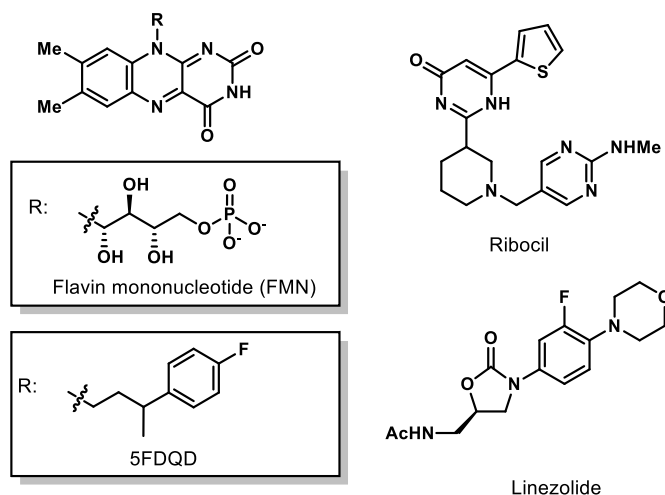


Figure 1 Examples of RNA-binding small molecules

When targeting RNA, the question arises which targets are best suited for drug discovery and where in chemical space to look for potent ligands. Analysis of RNA-binding small molecules has revealed that some RNA ligands have drug-like properties comparable to FDA approved drugs while others lie outside this

space.<sup>4,15</sup> Warner et al. have argued that RNA targets that bind such drug-like molecules and are thus deemed to be “ligandable” hold the greatest promise.<sup>3</sup> Consequently, tools are needed to identify such targets.

Targets are commonly considered to be “ligandable” or “druggable” if they possess binding sites that allow them to bind orally bioavailable drugs with high affinity.<sup>16,17</sup> The terms to name such pockets are hotly debated and several alternative terms such as “bindability”, “tractability” or “chemical tractability” have been proposed.<sup>17</sup> We will use the term “druggability” throughout this manuscript because it is the prevalent term used in literature. Druggability is not an absolute property and for other pockets potent drugs can be developed, albeit larger efforts might be required. According, we will label pockets that are not classified to be druggable as “less-druggable”.

Over the last few years, several methods have been reported that are able to segregate druggable pockets from less-druggable ones based on the 3D structure of the binding site.<sup>17</sup> Typically, these methods use descriptors describing the hydrophobicity, size and shape of the pockets to classify them using machine learning methods. As training and validation sets, protein pockets that have been assigned to either category are used. One of these methods, the DLID (drug-like density) measure,<sup>18</sup> has also been applied to analyse RNA pockets. DLID uses PocketFinder<sup>19</sup> to identify potential binding sites and the descriptors volume, buriedness and hydrophobicity to estimate how likely a pocket is to bind a drug-like molecule. Warner et al. used this approach to illustrate the diversity of selected RNA binding sites.<sup>3</sup> Hewitt et al. conducted a comprehensive analysis of RNA structures in the PDB using the same method and concluded that many RNAs contain pockets that are likely suitable for small molecule binding.<sup>20</sup> However, they did not distinguish between the binding of drug-like ligands and other molecules.

In our group, we have developed DrugPred as a structure-based druggability prediction method.<sup>21,22</sup> DrugPred describes the size and shape of the binding site using a “superligand” as a negative print, which is obtained by merging predicted binding modes of drug molecules that were docked into the pocket using only steric constraints. Descriptors encoding polarity and size of the pocket are subsequently calculated based on the superligand and used to predict the druggability of the binding site. DrugPred was trained and validated on a set of non-redundant druggable and less druggable protein binding sites (NRDLID) which has become a standard in the field. In comparison studies, DrugPred performed at least equally well than other methods and achieved an accuracy of about 90%.<sup>21,23,24</sup>

Here, we adopted DrugPred for druggability predictions of RNA binding sites. One of the original DrugPred descriptors is the hydrophobicity indices of the amino acids lining the binding site.<sup>21,25</sup> As this descriptor cannot be calculated for RNA pockets, we have implemented alternative descriptors and re-trained the

predictor, and thus made a prediction software which is applicable to both protein and RNA binding sites. Compared to the protein field, there is very little data about ligands binding to RNA, and even less data that can be accessed in an efficient way. The NALDB and SMMRNA databases contain affinities of small molecules binding to RNA extracted from the literature.<sup>26,27</sup> However, it is not possible to download the data for further processing. The R-BIND database links binding data also to RNA crystal structures, but for only five of the ligands in this database a complex structure is available in the Protein Data Bank (PDB).<sup>28</sup> Due to the paucity of suitable data, we opted to train and validate our modified DrugPred model, which we termed DrugPred\_RNA, on protein data. Subsequently, DrugPred\_RNA was used for druggability predictions of RNA structures including the ribosome. In the following, we present the construction of DrugPred\_RNA together with its validation on protein and RNA binding sites.

## Methods

Scripts to download crystal structures from the PDB, process them and to calculate ligand and binding site descriptors were written using Python 3.6.8. with the Biopython (1.73) and RDKit (2019.09.1) libraries.<sup>29,30</sup>

### NRDLD set for training and validation

As training and test set, our non-redundant set of druggable and less druggable binding sites (NRDLD) with the most recent modifications was used.<sup>21,22</sup> The binding sites were carved out of the cif-files downloaded from the PDB by keeping all protein residues with an atom within 15 Å of the ligand. The isolated binding sites together with co-factors and metal ions if present were saved in the PDB format and used for descriptor calculation as described below.

### Descriptor calculation

A superligand as a negative print of the binding site was obtained as done previously with minor modifications.<sup>21</sup> In brief, a set of approved drug molecules was docked into the pockets using DOCK 3.6.<sup>31</sup> Since the aim of docking was solely to obtain information about the shape and the volume of the binding sites, all receptor atoms were set to carbon atoms and assigned a partial charge of 0. Subsequently, compounds for which a docking pose were obtained and for which the ratio of van der Waals (VDW) score to number of heavy atoms was  $\leq -1.3$  were merged into a superligand. Only the atoms adhering to all of the following criteria were retained during this process: 1) the atom had to be a non-hydrogen atom 2) at least two atoms coming from different docked compounds had to be closer than 1.2 Å 3) only one of the atoms within 1.2 Å from other atoms was kept in the final superligand. If no docked ligands passed these filters, the ligand contained in the complex structure was used as superligand.

Based on the superligand, binding site and buried superligand atoms were determined. For that purpose, using FreeSASA<sup>32</sup> as implemented in RDKit, the solvent accessible surface area (SASA) of each receptor

and superligand atom in the superligand-bound and superligand-free state was calculated using a 1.0 Å probe radius and ProtOr radii<sup>33</sup>. All receptor atoms for which the SASA differed between superligand-bound and unbound state were assigned as being binding site atoms. Likewise, all superligand atoms for which the SASA changed between the free and complexed state were assigned as buried superligand atoms and all superligand atoms with SASA > 0 Å in the unbound state were assigned as superligand surface atoms.

Using superligand and binding site atoms as input, descriptors describing the size, shape and polarity of the pocket were calculated (Table S1). For shape descriptors that are not based on the surface area or the number of receptor or superligand atoms, the Descriptors3D module of RDKit was used. For calculating polarity descriptors, we considered all carbon, phosphor, and sulphur atoms in addition to nitrogen atoms of the bases that are bound to the ribose to be hydrophobic and all oxygen atoms of amino acids, ribose sugars and phosphate groups in addition to non-aromatic nitrogen atoms of amino acids to be polar. The SASA values of these atoms were calculated with FreeSASA using the same settings as described above. The side chains of histidine and tryptophane residues as well as the RNA bases are known to form hydrogen bonds in the plane of the heterocycles while parallel to this plane they engage in pi-stacking interactions which are more hydrophobic in nature. To account for this ambivalent behaviour, the SASA of endocyclic aromatic nitrogen atoms of bases and amino acid side chains and exocyclic oxygen and nitrogen atoms of the bases was split into a hydrophobic and a polar contribution in the following way. The SASA of these atoms was calculated both in the absence (SASA<sub>total</sub>) and the presence (SASA<sub>pol</sub>) of two blocking carbon atoms which were placed perpendicular to the plane of the aromatic ring with a 1.70 Å distance from the atom of interest. The area SASA<sub>pol</sub> was considered to belong to a polar atom while the difference SASA<sub>total</sub> – SASA<sub>pol</sub> was considered to belong to a hydrophobic atom. Similarly, if more than half of an atom's SASA value considered to be hydrophobic area, the atom was included in the hydrophobic binding site atom count.

Training the predictive model using decision trees

Machine learning was carried using the decision tree algorithm eXtreme Gradient Boosting package (XGBoost)<sup>34</sup> in R<sup>35</sup>. As learning objective, logistic regression for binary classification with output probability was used. Thus, all binding sites obtained a score between 0.0 and 1.0, whereas pockets with a score  $\geq 0.5$  were labelled druggable and pockets with a score < 0.5 were labelled as less druggable. Divergent from the default settings, the following parameters were used for training the model:

- *Max\_depth* = 3 (Maximum depth of trees)
- *Scale\_pos\_weight* = 0.59 (Adjusts for the skewness between training and testing set)

- *Early\_stopping\_rounds* = 20 (Validation metric needs to improve at least once in every 30 rounds to continue training.).

The influence of the descriptors on the model was evaluated with the help of Shapley Additive Explanation (SHAP) values as implemented in the SHAPforxgboost package.<sup>36,37</sup> The same package was also used to make Figure 2. Descriptors included in the final model were chosen by iteratively removing the least impactful descriptors until the predictive performance of the model was negatively affected. To further assess the robustness of the final model (called DrugPred\_RNA), leave-one-out-cross validation was carried out yielding a training and testing error of 0.00342 and 0.127, respectively.

#### Assembly of RNA containing data sets

We selected RNA structures for druggability assessment by querying the PDB for structures containing only RNA and ligands (accessed November 2019). In addition, the PDB was searched for entries containing ligands and the keyword riboswitch to include structures which were excluded in the first query due to the presence of proteins. In total, this yielded 1084 structures. Subsequently, all structures that contained ligands that were detergents, buffer salts or crystallization components were filtered out reducing the data set to 427 unique entries (Table 1, see supplementary material for three letter codes of filtered out ligands). If a crystal structure contained several instances of the same ligand, only the first instance was retained. In addition, all metal ions and water molecules were deleted. This resulted in 465 distinct binding sites spanning 224 unique ligands. A second variant of this set was also prepared. In this variant, metal ions which were not more than 5 Å away from a ligand atom were retained. If a binding site contained several metal ions, several copies of the binding sites each of them containing one of the metal ions were prepared. This variant contained 343 entries. In the following, the first variant is called the metal-free and the second variant the metal-containing set. Further, a data set containing ligand binding sites in ribosome crystal structures was compiled by querying the PDB for structures that contained “ribosome” as keyword. These structures were treated as described above. In addition, the ligands were visually inspected to remove buffer components that had slipped the filter rules. This resulted in 731 binding sites in the metal-free ribosome set and 732 in the metal-containing set.

Table 1: Data sets of RNA and ribosomal binding sites for assessing DrugPred\_RNA.

	RNA data set			Ribosome data set		
	Total	Metal-free entries	Metal- containing entries	Total	Metal- free entries	Metal- containing entries
Unique PDB IDs	427			590		
Binding sites containing small molecule ligands	808	465	343	1463	731	732
Unique ligands	224			247		
Druggable entries	298	172	126	356	215	141

The binding sites were carved out of the original cif-files by keeping all RNA residues with at least one atom within 15 Å of the ligand and potentially metal ions. The isolated binding sites were saved in the PDB format and used for descriptor calculation as described above.

#### Binding site similarity consensus scoring

To investigate the robustness of DrugPred-RNA, binding sites were grouped based on binding site similarity. First, the binding site sequence of each pocket was generated by including all residues that contained at least one binding site atom (identified as described above) in ascending order while for modified nucleic residues the name of the corresponding unmodified residue was used. Subsequently, all binding site sequences were pairwise aligned using BioPython and the global alignment similarity was calculated. If this value was > 85% the pockets were assigned to the same family. As done previously, the consensus of the druggability predictions within each family (C) was calculated using the formula

$$C = \frac{|n_d - n_{ld}|}{N} \times 100\%$$

where  $n_d$  is the number of druggable binding sites within the family,  $n_{ld}$  is the number of less druggable binding sites  $N$  is the total number of family members.<sup>22</sup>

## Results

### Construction of DrugPred\_RNA

Compared to protein data, there is very little data about ligands binding to RNA and a data set of sufficient size composed of druggable or less druggable RNA bindings sites to train a druggability predictor could not be compiled. Therefore, we opted to predict the druggability of RNA binding sites by training a descriptor on protein binding sites and to subsequently apply it to the prediction of RNA pockets. This approach required that only descriptors that can be calculated both for protein and RNA binding sites were used. This was not the case for our previously derived DrugPred model, as it contained the two descriptors “relative occurrence of hydrophobic amino acid” and “hydrophobicity indices of the amino acids”.<sup>21</sup> Thus, a modified DrugPred model, termed DrugPred\_RNA was derived. As training and test set, our non-redundant set of druggable and less druggable binding sites (NRDLD) with the most recent modifications was used.<sup>21,22</sup> For all 110 binding sites in the NRDLD, 23 descriptors describing the size, shape and polarity were calculated (Table S1). Subsequently, the data set was divided into a training and test set as done previously<sup>22</sup> to train and evaluate a predictor. For DrugPred and DrugPred 2.0, partial least squares-discriminant analysis (PLS-DA) was used to model the data. However, using only protein-independent descriptors with PLS-DA resulted in worse predictions (data not shown). Therefore, we retreated to decision tree modelling based on XGBoost.<sup>34</sup> To avoid overfitting, the maximum depth of trees was limited to 2 and the early stopping option was used (Figure S 1). In an iterative process, weak descriptors were removed until the predictive performance of the model was negatively affected. With the final model, termed DrugPred\_RNA, of the 75 binding sites in the training set, 1 druggable pocket was misclassified as less druggable, and of the 35 binding sites in the validation set, 4 were misclassified (2 false positives and 2 false negatives) leading to accuracy, precision and recall values between 0.85 and 1.00 (Table 2 and Figure S2). With DrugPred\_RNA, the error in the test set is slightly better than with DrugPred 2.0 while the error for the training set is slightly worse.

Table 2: Performance of DrugPred\_RNA and DrugPred 2.0 on the training and test set of the NRDL

	Training set [druggable / less druggable]		Test set [druggable / less druggable]	
	DrugPred_RNA	DrugPred 2.0	DrugPred_RNA	DrugPred 2.0
Accuracy	0.99	0.91	0.91	0.94
Precision	1.00 / 0.97	0.92 / 0.89	0.95 / 0.86	0.95 / 0.93
Recall	0.98 / 1.00	0.94 / 0.86	0.91 / 0.92	0.95 / 0.93

The influence of the descriptors on the model output was evaluated with the help of Shapley Additive Explanation (SHAP) values which describe the importance of each descriptor on the model's output taking into account the interaction with other descriptors.<sup>36,38</sup> Each descriptor for each data point (here, a particular binding site) is assigned a positive or negative SHAP value describing the contribution of the descriptor to the model output (here, druggable or less druggable) for that data point. The mean SHAP value formed by all SHAP values for a descriptor for the entire data set indicates the importance of the descriptor for the model (the larger the absolute mean SHAP value, the more important the descriptor, Figure 2A). For DrugPred\_RNA, positive SHAP values imply high druggability probability, while negative SHAP values imply low druggability probability. Further, by plotting the individual SHAP values for a descriptor against the descriptor values, it becomes evident which descriptor values contribute positively or negatively to the model (Figure 2B). The sum of the SHAP values of all descriptors for a single data point indicates the direction of the prediction.

The final DrugPred\_RNA predictor was based on 12 descriptors (Figure 2A, Table S1). According to the

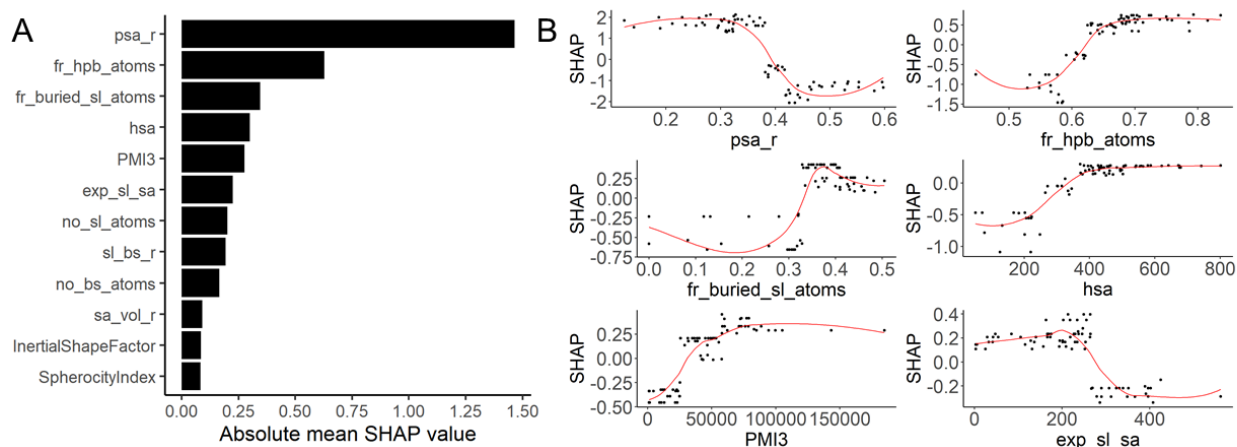


Figure 2: SHAP values for the DrugPred\_RNA model. **A)** Absolute mean SHAP values ranked from highest to lowest impact on the model's output. **B)** Individual SHAP values for each pocket in the training set for the top six descriptors in the model plotted against the descriptor values. Locally estimated scatterplot smoothing (LOESS) curves are overlaid on the descriptor observations (black dots). The midpoint in each curve indicates the cut-off value from where the model's prediction changes the direction. The plots for the remaining descriptors are displayed in Figure S 3.

SHAP values, the two most important descriptors were the relative polar surface area (*psa\_r*, absolute mean SHAP value = 1.46) and the fraction of hydrophobic binding site atoms (*fr\_hpb\_atoms*, absolute mean SHAP value = 0.63), which both describe the polarity of the binding site. As expected, druggable binding sites were less polar than less druggable sites (Figure 2B and Figure S 3). Both the high-ranking descriptor *fr\_buried\_sl\_atoms* (absolute mean SHAP value = 0.34) and the less important descriptor *sa\_vol\_r* (absolute mean SHAP value = 0.09) encode how compact a pocket is with less druggable pockets being more shallow (lower descriptor values for *fr\_buried\_sl\_atoms* and higher values for *sa\_vol\_r*) than druggable ones. Further, two descriptors for the solvent accessibility of the pocket (*exp\_sl\_sa*, absolute mean SHAP value = 0.22 and *sl\_bs\_r*, 0.19) were included in the final model. Here, it was found that druggable binding sites were less solvent accessible than less druggable ones. The descriptor *hsa* was also found to be among the more important ones (absolute mean SHAP value = 0.30). This descriptor describes size of the surface area of hydrophobic binding site atoms and correlates roughly with the size of the pocket. Other descriptors describing the size of the pocket were also included in the model but had less influence on the predictions (*no\_bs\_atoms*, absolute mean SHAP value = 0.17, and *no\_sl\_atoms*, 0.20). In agreement with previous findings, druggable pockets were larger and more hydrophobic than less druggable ones. The descriptors *InertialShapeFactor*, *SphericityIndex* and *PMI3* describing the shape of the superligand as a negative print of the binding site were also included in the final model. Pockets with a superligands with a

larger third moment of inertia (*PMI3*, absolute mean SHAP value = 0.27) and which were less spherical (*SphericityIndex*, absolute mean SHAP value = 0.08, *InertialShapeFactor*, absolute mean SHAP value = 0.08 ) were more likely to be assessed as druggable, albeit the latter two descriptors were determined to be less important.

#### Druggability predictions for RNA containing binding sites

Encouraged by the good performance of DrugPred\_RNA on the NRDLD, we proceeded with druggability predictions for RNA and ribosomal binding sites. No benchmark set for the evaluation of RNA druggability predictions is available in the public domain. Therefore, using the PDB, we compiled two data sets for this purpose, one containing RNA only binding sites and one with ribosome binding sites which in addition to ribosomal RNA could also contain ribosomal proteins. As binding sites, we considered all pockets that contained a ligand that is not a common crystallization buffer component. If a binding site contained metal ions within 5 Å of the ligand, several copies of the binding sites each of them containing one of the metal ions in addition to the metal-free pocket were prepared. In total, the RNA binding site set was composed of 427 unique PDB IDs spanning 808 binding sites (464 metal-free and 343 with metal ions, Table 1). 224 different ligands were found in these pockets. The binding sites were grouped into different families whereas family members were required to have a binding site sequence similarity of > 85%. This resulted in 46 different families in the RNA set (Table S2). The ribosomal binding site set was prepared in a similar fashion resulting in 590 unique PDB IDs with 731 pockets in the metal-free and 732 in the metal-containing subset. 247 different ligands were bound to these pockets and they were grouped into 66 different families.

Subsequently, the druggability of the pockets in all sets was predicted. In the RNA data set, 36% of the binding pockets (metal-containing and metal-free combined) were predicted to be druggable while in the ribosomal data set 24% of the pockets were predicted to be druggable .

To assess the impact of metal ions on the druggability prediction, we compared the predictions of metal-free and metal-containing versions of same parent pocket. In both sets, for the majority of the cases no change in the prediction outcome was found, regardless if there were metal ions present in the binding site or not. Accordingly, metal ions had only a very minor influence on the prediction outcome. In the following, we therefore only present data for pockets which were stripped of metal ions.

#### Assessment of druggability predictions on RNA-containing binding sites

Next, the quality of the predictions of DrugPred\_RNA on RNA-containing binding sites was evaluated. The following aspects were considered for the evaluation 1) the agreement of the predictions with anecdotal examples, 2) the extent to which binding sites which are known to efficiently bind drug-like ligands were

predicted to be druggable, and 3) the robustness of predictions with respect to substitutions and conformational changes in the binding sites.

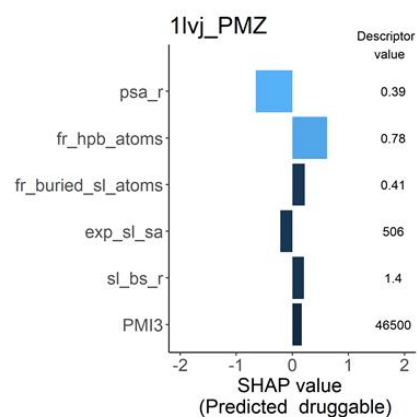
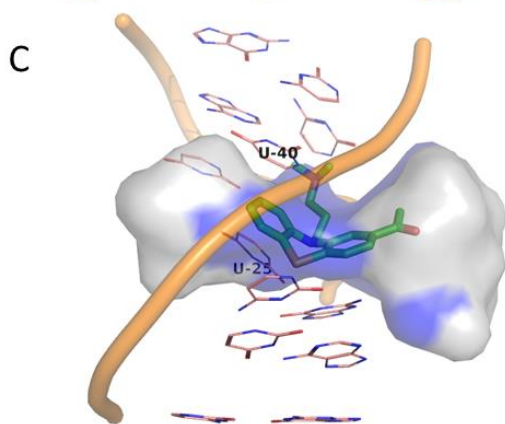
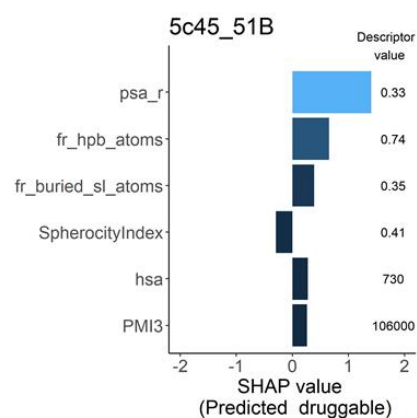
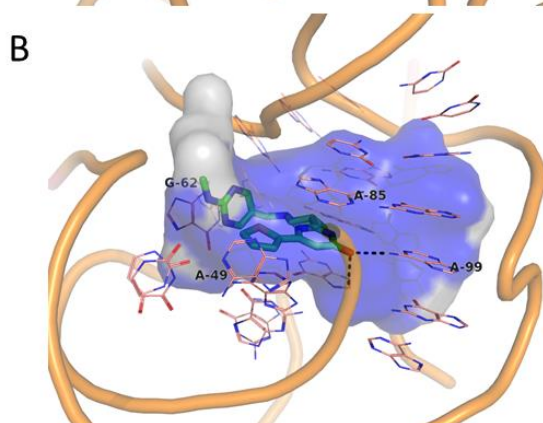
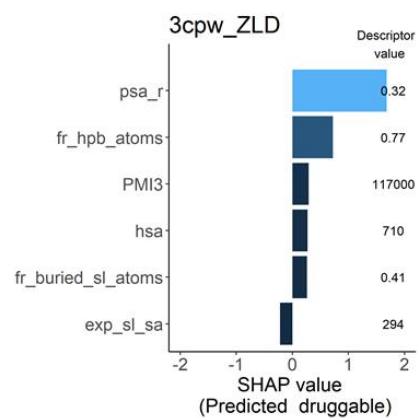
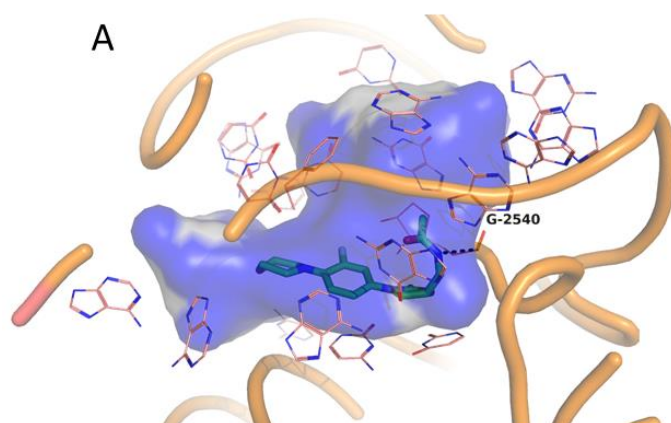
For this assessment, the drug-likeness of the ligands was predicted using the quantitative estimate of drug-likeness (QED) score.<sup>39</sup> This score weighs multiple molecular features (e. g. molecular weight, number of hydrogen bond donors or acceptors, polar surface area, presence of unwanted functionalities) into one single unitless score, which ranges from 0 (undesirable) to 1 (desirable). Albeit this metric does not provide a clear cut-off to distinguish “desirable” from “undesirable” compounds, the authors denoted a mean score of 0.67 for attractive compounds, 0.49 for less attractive and 0.34 for too complex and unattractive compounds. Accordingly, in the following we classified compounds with a QED score  $\geq 0.67$  as drug-like, with a QED score  $\leq 0.49$  as less drug-like and compounds with a score in between as moderate drug-like.

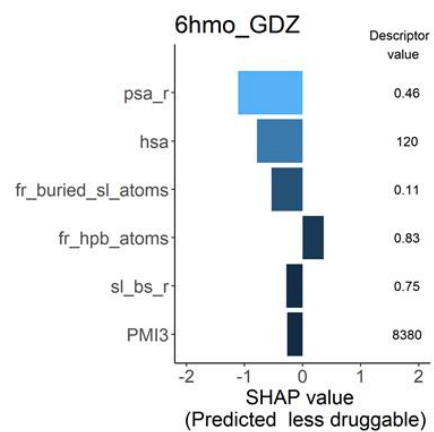
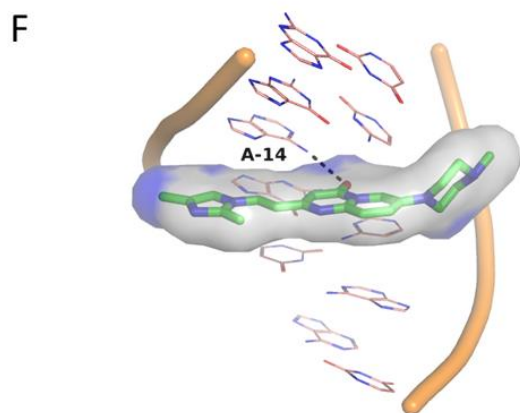
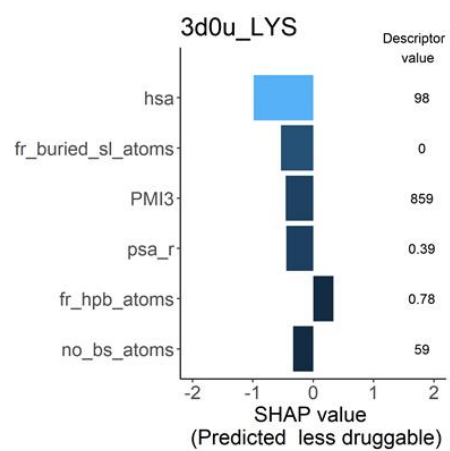
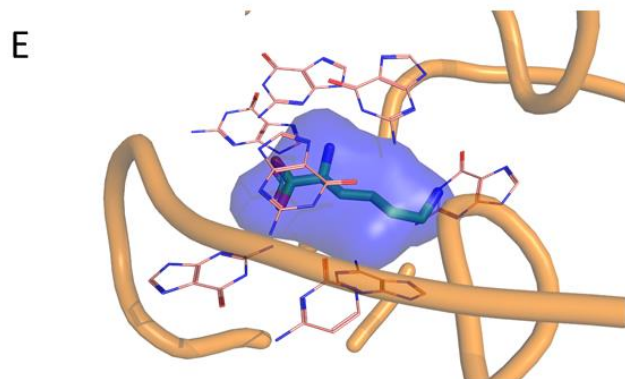
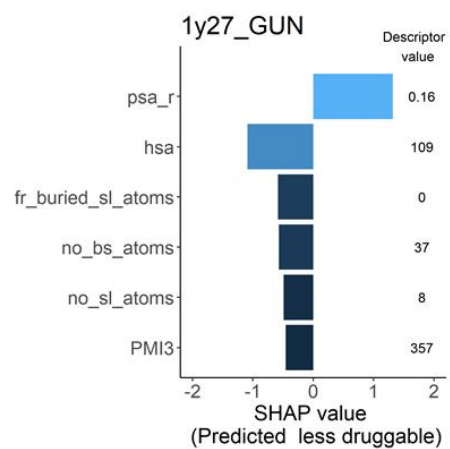
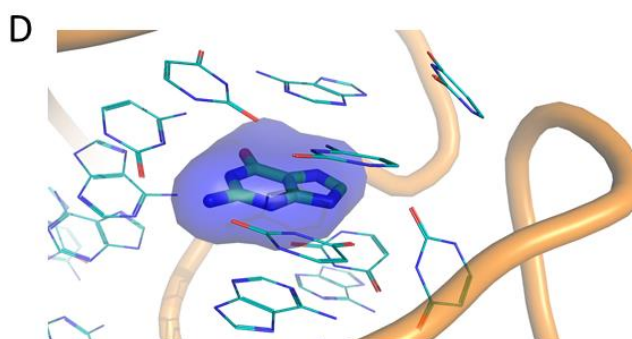
Further, for a binding site to be druggable, it needs to bind drug-like ligands potently.<sup>16,17</sup> Thus, if a potent drug-like ligand is known, one can with certainty say that a binding site is druggable. However, the absence of a potent ligand does not necessarily imply that a binding site is less druggable as there is always the possibility that a ligand can be optimized to increase its binding affinity. To take this into account, we used ligand efficiency (the binding energy normalized by the number of heavy atoms, LE) instead of binding affinity as a measure to judge if a ligand binds potently to its target.<sup>40</sup> Under the assumption that the ligand efficiency stays at its best constant during optimization, we considered ligands to bind potentially to their target if they have a LE of around  $0.30 \text{ kcal}\cdot\text{mol}^{-1}\cdot\text{heavy atom}^{-1}$  which translates to low nanomolar binding affinities of compounds with a molecular weight of maximum 500 Da. If no such ligand is known, we abstained from classifying a binding site on a general basis.

Evaluation of the performance of DrugPred\_RNA based on anecdotal examples

Linezolid is an FDA approved antibiotic targeting the 50S ribosomal subunit (Figure 1).<sup>8</sup> Based on its QED score of 0.89, it is highly drug-like. Its modest affinity of  $20 \mu\text{M}$  translates to a ligand efficiency (LE) of  $0.27 \text{ kcal}\cdot\text{mol}^{-1}\cdot\text{heavy atom}^{-1}$ . The binding site detected by DrugPred\_RNA encompasses linezolid. Linezolid is deeply buried in the pocket and forms mainly hydrophobic contacts in addition to a hydrogen bond to the ribose backbone of G2540 (Figure 3A). DrugPred\_RNA predicted this pocket to be druggable. According to the individual SHAP values of the descriptor values, the druggability was driven by a low relative polar surface area ( $psa_r = 0.29$ ), a large third moment of inertia ( $PMI3 = 11.7 \times 10$ ), and a large size ( $hsa = 748 \text{ \AA}^2$ ). These descriptors outweighed the fraction of hydrophobic atoms ( $fr_{hpb\_atoms} = 0.49$ ), the exposed surface area of the superligand ( $exp\_sl\_sa = 294 \text{ \AA}^2$ ) and the superligand atom/binding site atom ratio ( $sl\_bs_r = 1.2$ ) which were in a range more frequently associated with less druggable pockets. Based on the binding mode of linezolid and the fact that linezolid is a drug-like ligand the prediction that

this binding pocket is druggable appears to be sensible, despite the ligand not binding as potently as expected for a drug.





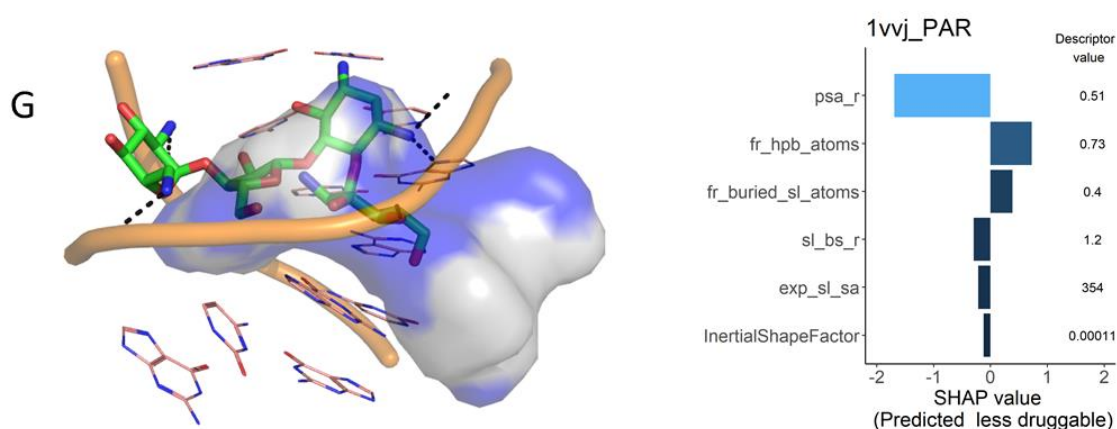


Figure 3: Evaluation of the performance of DrugPred\_RNA based on anecdotal examples. The RNA backbones are shown as orange tubes, nucleobases as thin sticks with carbon atoms coloured pink and ligands as thick sticks with carbon atoms in green. The surface of the superligand created by DrugPred\_RNA as a negative print of the pocket is shown as blobs with the solvent exposed surface area coloured grey and the remaining surface area coloured blue. Hydrogen bonds are indicated as dotted black lines. For each pocket, the individual SHAP values for the six most important descriptors together with the descriptor values are also displayed. The SHAP value plots are labelled with the PDB codes and the three letter codes of the ligands found in each pocket. A) The binding site of linezolid in the 50S ribosomal subunit. B) Ribocil bound to the FMN riboswitch. C) TAR RNA complexed with acetylpromazine. D) Guanine bound to the guanine riboswitch. E) Lysine in the binding site of the lysine riboswitch. F) Splicing site complexed with a splicing site modifier. G) Paromomycin bound to a bacterial ribosome site.

The FMN riboswitch has been validated as a target for the antibiotic compound ribocil, a drug-like small molecule (QED score = 0.71, Figure 1).<sup>13</sup> The affinity for ribocil ( $K_D = 13$  nM) is driven by hydrogen bonding with the base of A99 and the ribose group of A48 as well as stacking interactions with A85, A49 and, G62 (Figure 3B).<sup>41</sup> The binding site was rather deep ( $fr\_buried\_sl\_atoms = 0.35$ ) and characterized by a low relative polar surface area ( $psa\_r = 0.33$ ), a large the fraction of hydrophobic atoms ( $fr\_hpb\_atoms = 0.74$ ), a rather large size of the hydrophobic contact surface area ( $hsa = 730 \text{ \AA}^2$ ) and a large third principal moment of inertia ( $PMI3 = 10.6 \times 10^4$ ). These values drove the site to be predicted as druggable despite its spherocitiy index lying in the less druggable range ( $SpherocitiyIndex = 0.41$ ). This agrees with the site binding drug-like ligands like ribocil with high ligand efficiency ( $LE = 0.41 \text{ kcal}\cdot\text{mol}^{-1}\cdot\text{heavy atom}^{-1}$ ).

A known ligand for the HIV-1 trans activating region (TAR) RNA is the drug acetylpromazine (QED = 0.85, Figure 4).<sup>42</sup> Developed for a different target, the compound binds only with moderate affinity and efficiency to TAR RNA ( $K_D = 270 \mu\text{M}$ ,  $LE = 0.22 \text{ kcal}\cdot\text{mol}^{-1}\cdot\text{heavy atom}^{-1}$ ).<sup>43</sup> In the structure of the complex, the ligand forms stacking interactions with U25 and U40 (Figure 3C). DrugPred\_RNA predicted the ligand binding site to be druggable. As with the examples above, the classification was driven by a large fraction of hydrophobic atoms ( $fr\_hpb\_atoms = 0.78$ ), the depth of the pocket ( $fr\_buried\_sl\_atoms = 0.41$ ), the high proportion of the superligand atoms to binding site atoms ( $sl\_bs\_r = 1.4$ ) and the large third moment of inertia ( $PMI3 = 46.5 \times 10^3$ ). These properties overcame the high solvent accessibility ( $exp\_sl\_sa = 508 \text{ \AA}^2$ ), and the relative high polarity of the binding site ( $psa\_r = 0.39$ ). More potent ligands for HIV TAR RNA are also known albeit structural information about their binding modes is lacking. Examples are a drug-like screening hit (QED = 0.72) and furimidazoline (QED = 0.72) which have affinities of 230 nM and 1  $\mu\text{M}$ , resp. translating to LEs of 0.33 and 0.31  $\text{kcal}\cdot\text{mol}^{-1}\cdot\text{heavy atom}^{-1}$  (Figure 4).<sup>44,45</sup> Assuming that these ligands bind into the same pocket as acetylpromazine, the prediction that this pocket is druggable appears to be reasonable.

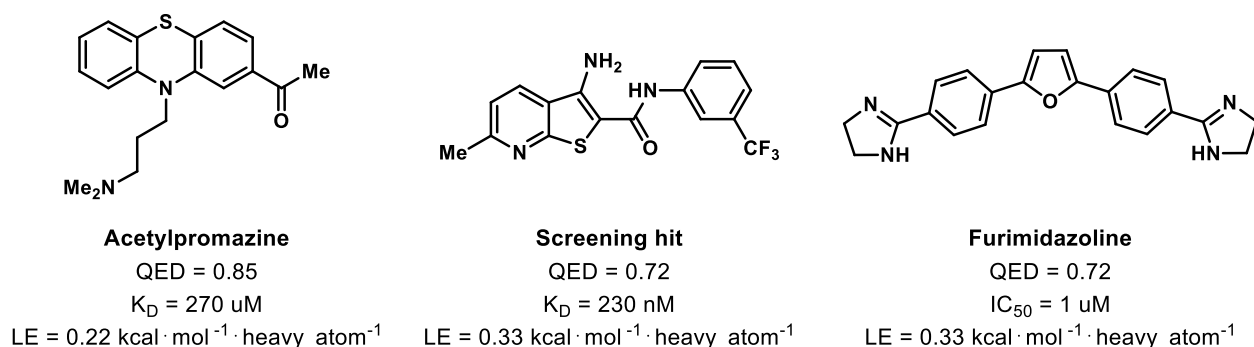


Figure 4: Ligands of HIV-1 TAR RNA

Ligands binding to the guanine and lysine riboswitch have been shown to act as antibiotics.<sup>46,47</sup> In both cases, the pockets are rather small and almost fully enclose the natural ligands (Figure 3D and E). Structure-activity relationships (SAR) are very tight and only small modifications of the ligands are possible without losing binding affinity. DrugPred\_RNA predicted these pockets to be less druggable which agrees with the SAR data. The predictions of the pockets were driven by their low relative polar surface areas ( $psa\_r = 0.16$  and  $0.39$ , resp.), their lack of a sufficiently large hydrophobic surface area ( $hsa = 109 \text{ \AA}^2$  and  $98 \text{ \AA}^2$ , resp.), small third principal moments of inertia ( $PMI3 = 357$  and  $8380$ , resp.), their shallowness ( $fr\_buried\_sl\_atoms = 0.0$  in both cases), and their small size ( $no\_bs\_atoms = 37$  and  $59$ ,  $no\_sl\_atoms = 8$ ,  $13$ , resp.).

Splicing modifiers for the treatment of spinal muscular atrophy are currently in clinical trials.<sup>9,10</sup> In our data set, the ligand SMN-C5 was included (Figure 3F). This ligand is moderate druglike (QED = 0.55) and has

a binding affinity of 28  $\mu\text{M}$  translating to a LE of  $0.22 \text{ kcal}\cdot\text{mol}^{-1}\cdot\text{heavy atom}^{-1}$  for its target RNA. In the NMR structure, the flat ligand is lying in a highly solvent exposed binding site. DrugPred\_RNA predicted this binding site to be less druggable. The prediction was due to the pocket being polar ( $psa_r = 0.46$ ,  $hsa = 120 \text{ \AA}^2$ ), shallow ( $fr\_buried\_sl\_atoms = 0.11$ ), and having an undesirable shape ( $sl\_bs\_r = 0.75$ ,  $PMI3 = 8.38 \times 10^3$ ). The druggability prediction appears to be reasonable based on the properties of the pocket, especially the high solvent exposure, but not on the fact that splicing modifiers are currently in clinical trials. This discrepancy is probably caused by the compounds binding *in vivo* to a ribonucleoprotein-RNA complex with a still unknown structure.<sup>48</sup> Thus, the biological relevant pocket of this type of compounds could not be included in our study.

One class of FDA-approved ribosome binding antibiotics are aminoglycosides. These are relatively polar natural products which break the rule of 5 for drug-like compounds.<sup>49</sup> One example of an aminoglycoside is paromomycin which acts by binding to the 16S ribosomal RNA (Figure 3G). Its low QED score of 0.11 is in agreement with the poor bioavailability of this compound class and the fact that aminoglycosides get into the bacteria by active transport.<sup>50</sup> In the complex of paromomycin bound to the ribosome of *T. thermophilus*, the ligand forms several hydrogen bonds with surrounding binding site residues and water molecules (not shown), with little to hydrophobic interactions. The terminal sugar ring in this ligand is located outside of the superligand created by DrugPred\_RNA, suggesting that this area is a less optimal for ligand binding. The SHAP values suggested that despite the depth of the pocket being in a range beneficial for druggable sites ( $fr\_buried\_sl\_atoms = 0.40$ ), the lack of hydrophobic atoms ( $fr\_hpb\_atoms = 0.45$ ), the large polar surface area ( $psa_r = 0.48$ ), the solvent-exposure ( $exp\_sl\_sa = 354 \text{ \AA}^2$ ) combined with a less ideal shape ( $InertialShapeFactor = 1.10 \times 10^{-4}$ ,  $sl\_bs\_r = 1.2$ , ) contributed to the pocket being predicted as less druggable. This prediction agrees with the nature of the known ligands.

Overall, the results for the selected anecdotal examples looked very promising. The predictions of DrugPred\_RNA generally agreed with what one would await based on the physico-chemical properties and binding efficiencies of known ligands. As expected, pockets predicted to be druggable were generally larger and more hydrophobic while the less druggable sites among the selected examples were more polar and solvent exposed. Encouraged by these results, we went on to investigate the performance of DrugPred\_RNA on a broader level.

#### Assessment of RNA-containing pockets binding drug-like ligands

In the next step, we investigated if the drug-like ligands contained in our RNA test sets bound to pockets predicted to be druggable. In total, the sets contained 331 unique ligands which 22 of them having a QED score  $\geq 0.67$ . Four of these ligands were found in the binding site of the preQ1 riboswitch. Upon closer inspection of these pockets it became evident that some of the bases in these structures were not resolved.

These pockets were therefore not further considered. Out of the remaining ligands, 12 (67%) were found in binding sites assessed by DrugPred\_RNA as druggable (Table S 4) and 6 (23%) in binding sites assessed to be less druggable (Table S 5). As only 16 % of all metal-free binding sites were predicted to be druggable, the drug-like ligands were clearly enriched in druggable binding sites.

For 11 out of the 12 drug-like ligands binding to pockets predicted to be druggable, we could find binding data in the literature. Based on this data, 9 ligands bind efficiently to their target with LEs  $> 0.30 \text{ kcal}\cdot\text{mol}^{-1}\cdot\text{heavy atom}^{-1}$  hinting that these pockets are indeed druggable. The two remaining ligands were linezolid with the 50S ribosomal subunit as target and acetylpromazine binding to HIV-1 TAR RNA. For the reasons discussed above, these pockets also appear to be druggable. Thus, all druggability predictions for the pockets binding the 11 drug-like ligands with known binding data appear to be sensible.

On the other hand, 6 drug-like ligands were found in pockets predicted to be less druggable (Table S 5). For 5 of them we could retrieve affinity data in the literature and all of these bind rather efficiently to their targets ( $\text{LE} \geq 0.29 \text{ kcal}\cdot\text{mol}^{-1}\cdot\text{heavy atom}^{-1}$ ). Three of these ligands are fragments binding to the TPP riboswitch, one ligand binding the influenza A virus promoter region, and one a ligand for the Spinach aptamer. The TPP riboswitch contains a rather large and partially buried binding site (Figure 5). Several other examples of this pockets with slightly different binding site conformations were also contained in the set (Figure 5c). Some of them were predicted to be druggable, discussed in more detail below. That the TPP riboswitch binding site conformation binding efficiently the drug-like fragments were predicted to be less druggable can be considered a false negative prediction. The drug-like ligand of the to influenza A promoter region sits on the surface of the RNA molecule and is almost entirely solvent exposed (Figure 6). It is highly unusual that a ligand with such a binding mode bind that efficiently ( $\text{LE} = 0.29 \text{ kcal}\cdot\text{mol}^{-1}\cdot\text{heavy atom}^{-1}$ ). However, the structure of the complex has been determined by NMR and it is possible that the resolution of the structure is not accurate enough to reveal the actual details of the binding mode.<sup>51</sup> The Spinach aptamer binds a small molecule dye, DFHBI, which forms hydrogen bonding and pi-stacking interactions in the binding site (Figure 7). The top six SHAP-values for this entry showed that while the fraction of hydrophobic atoms ( $fr\_hpb\_atoms = 0.8$ ) and the hydrophobic surface area value ( $hsa = 409 \text{ \AA}^2$ ) were in a range that is favorable for druggable binding sites, the shallow shape of the pocket ( $fr\_buried\_sl\_atoms = 0.17$ ,  $sl\_bs\_r = 0.79$ ), combined with the solvent exposure ( $exp\_sl\_sa = 376 \text{ \AA}^2$ ) and the high relative polar surface area ( $psa\_r = 0.38$ ) drove the site to be predicted as less druggable. Considering the drug-likeness of the ligand together with its efficient binding, this prediction can be considered to be a misclassification by DrugPred\_RNA. Taken together, the druggability predictions for the TPP riboswitch pockets binding the fragments and the Spinach aptamer were likely false negatives, while the prediction for influence A promotor region appeared to be sensible. This could point to DrugPredRNA having a tendency to

misclassify druggable binding sites as less druggable. However, the investigated data set was too small to conclude firmly on this.

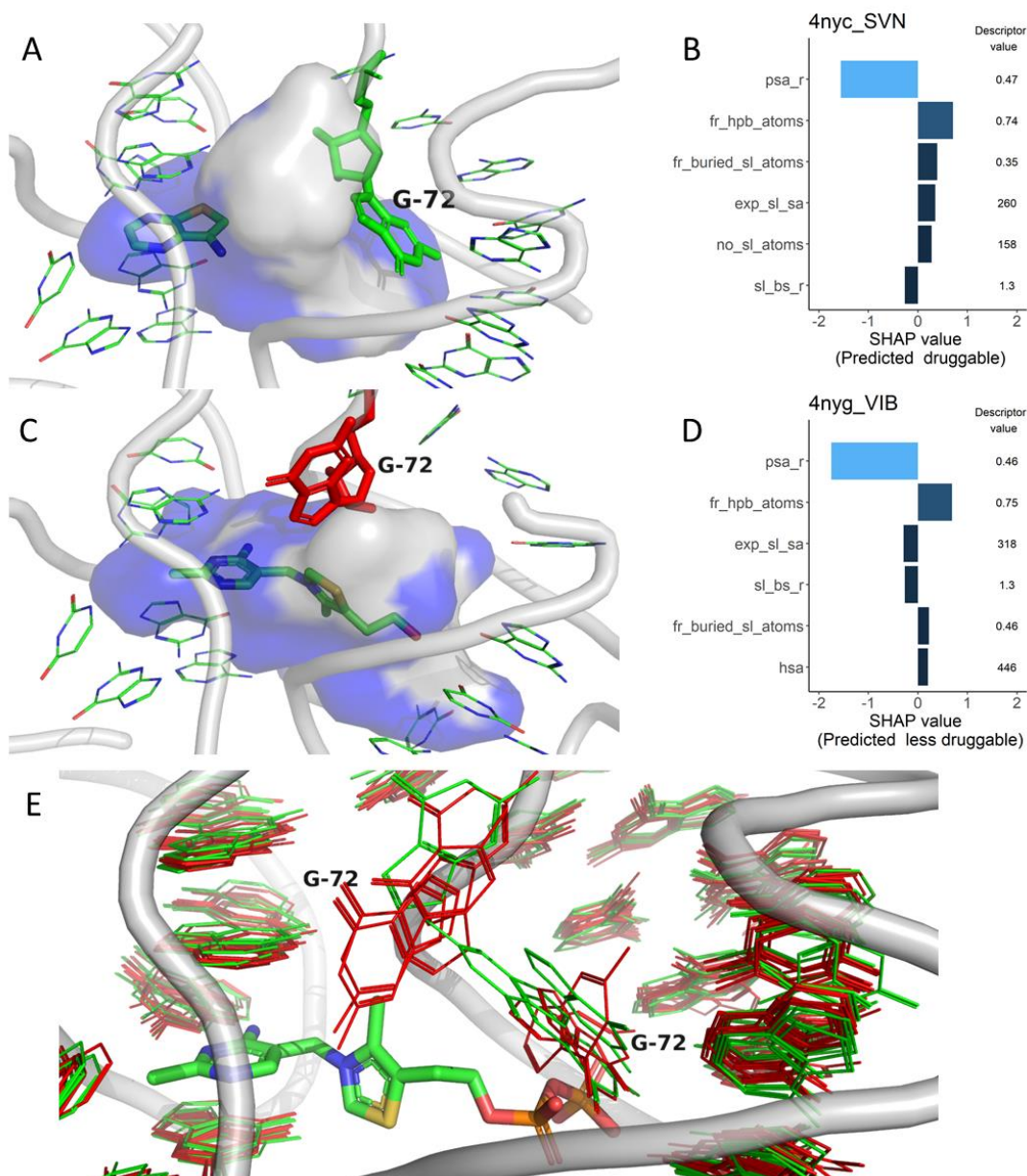


Figure 5: Druggability predictions for TPP riboswitch binding sites. The flexible residue G72 is highlighted in all figures. A) The TPP riboswitch binding site (PDB ID 4nyc) in complex with a fragment screening hit (green sticks) and the superligand generated by DrugPred\_RNA (blob, with solvent-exposed areas colored grey, remaining blue). B,D) The six descriptors with the highest SHAP values for 4nyc\_SVN (B) and 4nyg\_VIB (D). C) The TPP riboswitch binding site (PDB ID 4nyg) in complex with thiamine (VIB). E) Superposition of all *E. coli* TPP riboswitch binding sites. Entries predicted as druggable are colored in green and as less druggable in red. The conformations of the residue G72 influences the prediction outcome.

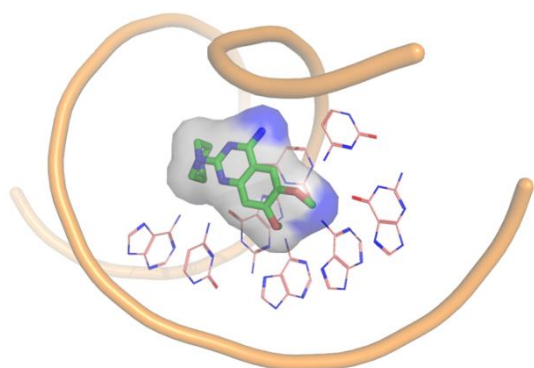


Figure 6: Binder of influenza A promoter region (PDB ID 2lwk). The surface of the superligand created by DrugPred\_RNA as a negative print of the pocket is shown as a blob with the solvent exposed surface area coloured grey and the remaining surface area coloured blue.

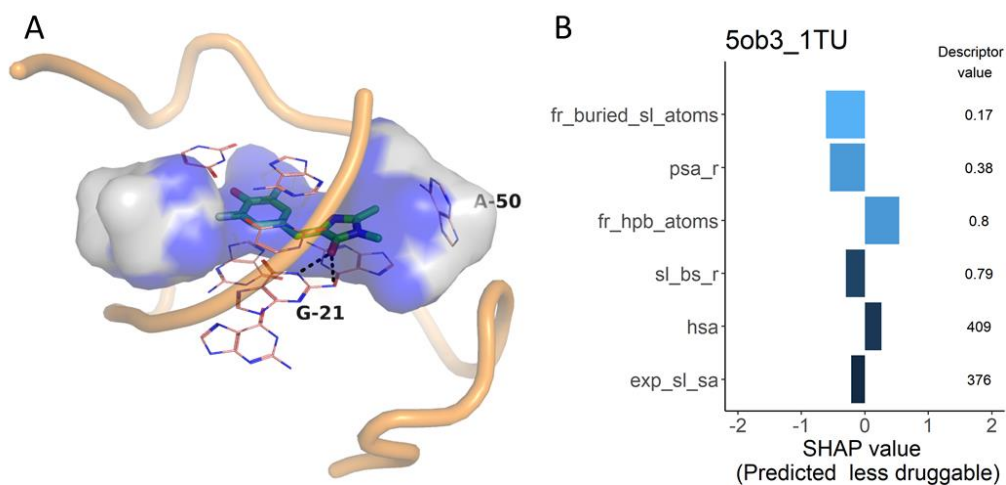


Figure 7: a) The Spinach aptamer (PDB ID 5ob3, orange thick and thin lines) bound to the dye DFHBI (green sticks). The surface of the superligand created by DrugPred\_RNA as a negative print of the pocket is shown as a blob with the solvent exposed surface area coloured grey and the remaining surface area coloured blue. b) Individual SHAP values for the six most important descriptors together with the descriptor values obtained by DrugPred\_RNA.

## Assessment of the robustness of the druggability predictions

Further, we assessed the robustness of the predictions with respect to conformational changes and base composition of the binding sites. For that purpose, all binding sites with a sequence similarity >85% were grouped together resulting in 46 families in the RNA binding site set and 56 in the ribosome set. The families spanned between 2 and 23 members in the RNA dataset (Table S2) and 2 and 56 in the ribosomal set (Table S3). Subsequently, the consensus of the predictions for each family was calculated. The consensus was defined as  $(100 * |\# \text{druggable binding sites} - \# \text{less-druggable binding sites}|) / (\text{total number of predictions})$ .<sup>22</sup> Thus, 100% consensus would be obtained if all pockets in one family were predicted to belong to the same class (druggable or less druggable) and 0% if one half of the pockets was predicted to belong to one class and the other half to the other class. In the RNA set, for 74% of all 46 evaluated families, a consensus of 80% or more was obtained. Out of those, 22 families were predicted to be less druggable with a consensus score of 100% while 7 families were predicted to be druggable with a consensus score  $\geq 80\%$ . In the ribosome set, 75% of the 52 families had a consensus score of  $\geq 80\%$  whereas 21 of these families were predicted to be druggable and the remaining ones to be less druggable. Thus, in most cases using different crystal structures of the same or a related pocket did not change the outcome of the prediction.

The TPP riboswitch family containing pockets from 16 distinct PDB entries, obtained a low consensus score of 12.5% with the majority of the pockets predicted as less druggable (Table S2). Superimposing the pockets, it became evident that there is some plasticity in the binding sites (Figure 5c). One guanine residue (G72 in the *E. coli* TPP riboswitch) can adopt several conformations depending on the bound ligand leading to considerably different superligands (Figure 5a and b). Consequently, the pockets differ in compactness (*fr\_buried\_sl\_atoms*, *sl\_bs\_r*) and solvent exposure (*exp\_sl\_sa*). Some of the binding site conformations were predicted to be druggable and others less druggable (Figure 5C). No distinct pattern emerged of which conformation was favored but it became evident that a position which hindered the size of the superligand was unfavorable.

## Conclusion

RNA is an emerging target for drug discovery.<sup>3-6</sup> However, like for proteins, not all RNA binding sites are equally suited to be addressed with conventional drug-like ligands. To identify pockets that are primed to potentially bind such ligands we have developed the structure-based druggability predictor DrugPred\_RNA. Due to the paucity of annotated RNA binding sites, the predictor was trained on a set of protein pockets, albeit containing only descriptors that can be calculated for both, RNA and protein binding sites. DrugPred\_RNA performed comparable on the protein binding site set as our previous DrugPred 2.0 predictor trained with slightly different descriptors (Table 2) In addition, druggability predictions of

DrugPred\_RNA on selected anecdotal examples appeared to be sensible (Figure 3). Likewise, the majority of the drug-like ligands contained in our RNA binding site sets were found in pockets predicted to be druggable, further adding confidence to the DrugPred\_RNA predictions (Table S 4 and Table S 5). As observed before,<sup>21,52</sup> using different conformations of a binding site could result in opposing druggability predictions (Table S2 and Table S3). However, for the majority of cases consistent predictions were obtained indicating that DrugPred\_RNA is robust towards small changes in binding site conformations. Interestingly, many riboswitches were found among the binding site families that were predicted to be druggable (Table S2). This finding underlines the notation that these promising targets for new antibiotics could be addressed with drug-like ligands.<sup>3,12,20</sup> Further, also in the ribosomal binding site set druggable pockets were contained (Table S3) which can help to direct efforts when targeting the ribosome for the development of drugs to overcome the looming antibiotic crisis.<sup>7,53</sup> Collectively, DrugPred\_RNA is a promising tool for structure-based druggability predictions of RNA binding sites that can be used to decide in which area of chemical space to search for ligands for a given binding site.

## References

- (1) Imming, P.; Sinning, C.; Meyer, A. Drugs, Their Targets and the Nature and Number of Drug Targets. *Nat. Rev. Drug Discov.* **2006**, *5* (10), 821–834.
- (2) Santos, R.; Ursu, O.; Gaulton, A.; Bento, A. P.; Donadi, R. S.; Bologa, C. G.; Karlsson, A.; Al-Lazikani, B.; Hersey, A.; Oprea, T. I.; Overington, J. P. A Comprehensive Map of Molecular Drug Targets. *Nat. Rev. Drug Discov.* **2017**, *16* (1), 19–34. <https://doi.org/10.1038/nrd.2016.230>.
- (3) Warner, K. D.; Hajdin, C. E.; Weeks, K. M. Principles for Targeting RNA with Drug-like Small Molecules. *Nat. Rev. Drug Discov.* **2018**. <https://doi.org/10.1038/nrd.2018.93>.
- (4) Rizvi, N. F.; Smith, G. F. RNA as a Small Molecule Druggable Target. *Bioorg. Med. Chem. Lett.* **2017**. <https://doi.org/10.1016/j.bmcl.2017.10.052>.
- (5) Matsui, M.; Corey, D. R. Non-Coding RNAs as Drug Targets. *Nat. Rev. Drug Discov.* **2017**, *16* (3), 167–179. <https://doi.org/10.1038/nrd.2016.117>.
- (6) Ursu, A.; Childs-Disney, J. L.; Andrews, R. J.; O’Leary, C. A.; Meyer, S. M.; Angelbello, A. J.; Moss, W. N.; Disney, M. D. Design of Small Molecules Targeting RNA Structure from Sequence. *Chem. Soc. Rev.* **2020**, *49* (20), 7252–7270. <https://doi.org/10.1039/D0CS00455C>.
- (7) Wilson, D. N. Ribosome-Targeting Antibiotics and Mechanisms of Bacterial Resistance. *Nat. Rev. Microbiol.* **2013**, *12* (1), 35–48. <https://doi.org/10.1038/nrmicro3155>.
- (8) Ippolito, J. A.; Kanyo, Z. F.; Wang, D.; Franceschi, F. J.; Moore, P. B.; Steitz, T. A.; Duffy, E. M. Crystal Structure of the Oxazolidinone Antibiotic Linezolid Bound to the 50S Ribosomal Subunit. *J. Med. Chem.* **2008**, *51* (12), 3353–3356. <https://doi.org/10.1021/jm800379d>.
- (9) Calder, A. N.; Androphy, E. J.; Hodgetts, K. J. Small Molecules in Development for the Treatment of Spinal Muscular Atrophy. *J. Med. Chem.* **2016**, *59* (22), 10067–10083. <https://doi.org/10.1021/acs.jmedchem.6b00670>.
- (10) Ratni, H.; Karp, G. M.; Weetall, M.; Naryshkin, N. A.; Paushkin, S. V.; Chen, K. S.; McCarthy, K. D.; Qi, H.; Turpoff, A.; Woll, M. G.; Zhang, X.; Zhang, N.; Yang, T.; Dakka, A.; Vazirani, P.; Zhao, X.; Pinard, E.; Green, L.; David-Pierson, P.; Tuerck, D.; Poirier, A.; Muster, W.; Kirchner, S.; Mueller, L.; Gerlach, I.; Metzger, F. Specific Correction of Alternative Survival Motor Neuron 2 Splicing by Small Molecules: Discovery of a Potential Novel Medicine To Treat Spinal Muscular Atrophy. *J. Med. Chem.* **2016**, *59* (13), 6086–6100. <https://doi.org/10.1021/acs.jmedchem.6b00459>.
- (11) Rekand, I.; Brenk, R. Design of Riboswitch Ligands, an Emerging Target Class for Novel Antibiotics. *Future Med Chem* **2017**, *9* (14), 1649–1662.
- (12) Panchal, V.; Brenk, R. Riboswitches as Drug Targets for Antibiotics. *Antibiotics Submitted*.
- (13) Howe, J. A.; Wang, H.; Fischmann, T. O.; Balibar, C. J.; Xiao, L.; Galgoci, A. M.; Malinverni, J. C.; Mayhoo, T.; Villafania, A.; Nahvi, A.; Murgolo, N.; Barbieri, C. M.; Mann, P. A.; Carr, D.; Xia, E.; Zuck, P.; Riley, D.; Painter, R. E.; Walker, S. S.; Sherborne, B.; de Jesus, R.; Pan, W.; Plotkin, M. A.; Wu, J.; Rindgen, D.; Cummings, J.; Garlisi, C. G.; Zhang, R.; Sheth, P. R.; Gill, C. J.; Tang, H.; Roemer, T. Selective Small-Molecule Inhibition of an RNA Structural Element. *Nature* **2015**, *526* (7575), 672–677. <https://doi.org/10.1038/nature15542>.
- (14) Blount, K. F.; Megyola, C.; Plummer, M.; Osterman, D.; O’Connell, T.; Aristoff, P.; Quinn, C.; Chrusciel, R. A.; Poel, T. J.; Schostarez, H. J.; Stewart, C. A.; Walker, D. P.; Wuts, P. G. M.; Breaker, R. R. Novel Riboswitch-Binding Flavon Analog That Protects Mice against *Clostridium Difficile* Infection without Inhibiting Cecal Flora. *Antimicrob. Agents Chemother.* **2015**, *59* (9), 5736–46. <https://doi.org/10.1128/AAC.01282-15>.
- (15) Morgan, B. S.; Forte, J. E.; Culver, R. N.; Zhang, Y.; Hargrove, A. E. Discovery of Key Physicochemical, Structural, and Spatial Properties of RNA-Targeted Bioactive Ligands. *Angew. Chem. Int. Ed Engl.* **2017**, *56* (43), 13498–13502. <https://doi.org/10.1002/anie.201707641>.
- (16) Edfeldt, F. N.; Folmer, R. H.; Breeze, A. L. Fragment Screening to Predict Druggability (Ligandability) and Lead Discovery Success. *Drug Discov. Today* **2011**. [https://doi.org/S1359-6446\(11\)00034-1](https://doi.org/S1359-6446(11)00034-1) [pii] 10.1016/j.drudis.2011.02.002.

- (17) Abi Hussein, H.; Geneix, C.; Petitjean, M.; Borrel, A.; Flatters, D.; Camproux, A.-C. Global Vision of Druggability Issues: Applications and Perspectives. *Drug Discov. Today* **2017**, 22 (2), 404–415. <https://doi.org/10.1016/j.drudis.2016.11.021>.
- (18) Sheridan, R. P.; Maiorov, V. N.; Holloway, M. K.; Cornell, W. D.; Gao, Y. D. Drug-like Density: A Method of Quantifying the “Bindability” of a Protein Target Based on a Very Large Set of Pockets and Drug-like Ligands from the Protein Data Bank. *J. Chem. Inf. Model.* **2010**, 50 (11), 2029–2040. <https://doi.org/10.1021/ci100312t>.
- (19) An, J.; Totrov, M.; Abagyan, R. Pocketome via Comprehensive Identification and Classification of Ligand Binding Envelopes. *Mol. Cell. Proteomics* **2005**, 4 (6), 752–761. <https://doi.org/10.1074/mcp.M400159-MCP200>.
- (20) Hewitt, W. M.; Calabrese, D. R.; Schneekloth, J. S. Evidence for Ligandable Sites in Structured RNA throughout the Protein Data Bank. *Bioorg. Med. Chem.* **2019**, 27 (11), 2253–2260. <https://doi.org/10.1016/j.bmc.2019.04.010>.
- (21) Krasowski, A.; Muthas, D.; Sarkar, A.; Schmitt, S.; Brenk, R. DrugPred: A Structure-Based Approach to Predict Protein Druggability Developed Using an Extensive Nonredundant Data Set. *J. Chem. Inf. Model.* **2011**, 51 (11), 2829–2842.
- (22) Sarkar, A.; Brenk, R. To Hit or Not to Hit, That Is the Question - Genome-Wide Structure-Based Druggability Predictions for *Pseudomonas Aeruginosa* Proteins. *PloS One* **2015**, 10 (9), e0137279. <https://doi.org/10.1371/journal.pone.0137279>.
- (23) Desaphy, J.; Azdimousa, K.; Kellenberger, E.; Rognan, D. Comparison and Druggability Prediction of Protein-Ligand Binding Sites from Pharmacophore-Annotated Cavity Shapes. *J. Chem. Inf. Model.* **2012**, 52 (8), 2287–2299. <https://doi.org/10.1021/ci300184x>.
- (24) Volkamer, A.; Rarey, M. Exploiting Structural Information for Drug-Target Assessment. *Future Med. Chem.* **2014**, 6 (3), 319–331. <https://doi.org/10.4155/fmc.14.3>.
- (25) Kyte, J.; Doolittle, R. F. A Simple Method for Displaying the Hydropathic Character of a Protein. *J. Mol. Biol.* **1982**, 157 (1), 105–132. [https://doi.org/0022-2836\(82\)90515-0](https://doi.org/0022-2836(82)90515-0) [pii].
- (26) Kumar Mishra, S.; Kumar, A. NALDB: Nucleic Acid Ligand Database for Small Molecules Targeting Nucleic Acid. *Database J. Biol. Databases Curation* **2016**, 2016. <https://doi.org/10.1093/database/baw002>.
- (27) Mehta, A.; Sonam, S.; Gouri, I.; Loharch, S.; Sharma, D. K.; Parkesh, R. SMMRNA: A Database of Small Molecule Modulators of RNA. *Nucleic Acids Res.* **2014**, 42 (Database issue), D132–141. <https://doi.org/10.1093/nar/gkt976>.
- (28) Morgan, B. S.; Sanaba, B. G.; Donlic, A.; Karloff, D. B.; Forte, J. E.; Zhang, Y.; Hargrove, A. E. R-BIND: An Interactive Database for Exploring and Developing RNA-Targeted Chemical Probes. *ACS Chem. Biol.* **2019**, 14 (12), 2691–2700. <https://doi.org/10.1021/acscchembio.9b00631>.
- (29) Cock, P. J. A.; Antao, T.; Chang, J. T.; Chapman, B. A.; Cox, C. J.; Dalke, A.; Friedberg, I.; Hamelryck, T.; Kauff, F.; Wilczynski, B.; de Hoon, M. J. L. Biopython: Freely Available Python Tools for Computational Molecular Biology and Bioinformatics. *Bioinformatics* **2009**, 25 (11), 1422–1423. <https://doi.org/10.1093/bioinformatics/btp163>.
- (30) Landrum, G. RDKit: Open-source cheminformatics <http://www.rdkit.org>.
- (31) Lorber, D. M.; Shoichet, B. K. Flexible Ligand Docking Using Conformational Ensembles. *Protein Sci.* **1998**, 7 (4), 938–50.
- (32) Mitternacht, S. FreeSASA: An Open Source C Library for Solvent Accessible Surface Area Calculations. *FI000Research* **2016**, 5. <https://doi.org/10.12688/fi000research.7931.1>.
- (33) Tsai, J.; Taylor, R.; Chothia, C.; Gerstein, M. The Packing Density in Proteins: Standard Radii and Volumes11Edited by J. M. Thornton. *J. Mol. Biol.* **1999**, 290 (1), 253–266. <https://doi.org/10.1006/jmbi.1999.2829>.
- (34) Chen, T.; Guestrin, C. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*; KDD '16; ACM: New York, NY, USA, 2016; pp 785–794. <https://doi.org/10.1145/2939672.2939785>.

- (35) R Core Team. *R: A Language and Environment for Statistical Computing*; R Foundation for Statistical Computing: Vienna, Austria, 2020.
- (36) Lundberg, S. M.; Lee, S.-I. A Unified Approach to Interpreting Model Predictions. In *Advances in Neural Information Processing Systems 30*; Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R., Eds.; Curran Associates, Inc., 2017; pp 4765–4774.
- (37) Liu, Y.; Just, A. *SHAPforxgboost: SHAP Plots for “XGBoost”*; manual; 2019.
- (38) Rodríguez-Pérez, R.; Bajorath, J. Interpretation of Machine Learning Models Using Shapley Values: Application to Compound Potency and Multi-Target Activity Predictions. *J. Comput. Aided Mol. Des.* **2020**, *34* (10), 1013–1026. <https://doi.org/10.1007/s10822-020-00314-0>.
- (39) Bickerton, G. R.; Paolini, G. V.; Besnard, J.; Muresan, S.; Hopkins, A. L. Quantifying the Chemical Beauty of Drugs. *Nat. Chem.* **2012**, *4* (2), 90–98. <https://doi.org/10.1038/nchem.1243>.
- (40) Hopkins, A. L.; Groom, C. R.; Alex, A. Ligand Efficiency: A Useful Metric for Lead Selection. *Drug Discov. Today* **2004**, *9* (10), 430–431.
- (41) Howe, J. A.; Xiao, L.; Fischmann, T. O.; Wang, H.; Tang, H.; Villafania, A.; Zhang, R.; Barbieri, C. M.; Roemer, T. Atomic Resolution Mechanistic Studies of Ribocil: A Highly Selective Unnatural Ligand Mimic of the E. Coli FMN Riboswitch. *RNA Biol.* **2016**, *13* (10), 946–954. <https://doi.org/10.1080/15476286.2016.1216304>.
- (42) Structure of TAR RNA Complexed with a Tat-TAR Interaction Nanomolar Inhibitor That Was Identified by Computational Screening. *Chem. Biol.* **2002**, *9* (6), 707–712. [https://doi.org/10.1016/S1074-5521\(02\)00151-5](https://doi.org/10.1016/S1074-5521(02)00151-5).
- (43) Mayer, M.; James, T. L. NMR-Based Characterization of Phenothiazines as a RNA Binding Scaffold. *J. Am. Chem. Soc.* **2004**, *126* (13), 4453–4460. <https://doi.org/10.1021/ja0398870>.
- (44) Sztuba-Solinska, J.; Shenoy, S. R.; Gareiss, P.; Krumpe, L. R. H.; Le Grice, S. F. J.; O’Keefe, B. R.; Schneekloth, J. S. Identification of Biologically Active, HIV TAR RNA-Binding Small Molecules Using Small Molecule Microarrays. *J. Am. Chem. Soc.* **2014**, *136* (23), 8402–8410. <https://doi.org/10.1021/ja502754f>.
- (45) Gelus, N.; Bailly, C.; Hamy, F.; Klimkait, T.; Wilson, W. D.; Boykin, D. W. Inhibition of HIV-1 Tat-TAR Interaction by Diphenylfuran Derivatives: Effects of the Terminal Basic Side Chains. *Bioorg. Med. Chem.* **1999**, *7* (6), 1089–1096. [https://doi.org/10.1016/S0968-0896\(99\)00041-3](https://doi.org/10.1016/S0968-0896(99)00041-3).
- (46) Blount, K. F.; Wang, J. X.; Lim, J.; Sudarsan, N.; Breaker, R. R. Antibacterial Lysine Analogs That Target Lysine Riboswitches. *Nat. Chem. Biol.* **2007**, *3* (1), 44–49.
- (47) Yan, L.-H.; Le Roux, A.; Boyapelly, K.; Lamontagne, A.-M.; Archambault, M.-A.; Picard-Jean, F.; Lalonde-Seguin, D.; St-Pierre, E.; Najmanovich, R. J.; Fortier, L.-C.; Lafontaine, D.; Marsault, É. Purine Analogs Targeting the Guanine Riboswitch as Potential Antibiotics against *Clostridioides Difficile*. *Eur. J. Med. Chem.* **2018**, *143*, 755–768. <https://doi.org/10.1016/j.ejmech.2017.11.079>.
- (48) Sivaramakrishnan, M.; McCarthy, K. D.; Campagne, S.; Huber, S.; Meier, S.; Augustin, A.; Heckel, T.; Meistermann, H.; Hug, M. N.; Birrer, P.; Moursy, A.; Khawaja, S.; Schmucki, R.; Berntsen, N.; Giroud, N.; Golling, S.; Tzouros, M.; Banfai, B.; Duran-Pacheco, G.; Lamerz, J.; Hsiu Liu, Y.; Luebbers, T.; Ratni, H.; Ebeling, M.; Cléry, A.; Paushkin, S.; Krainer, A. R.; Allain, F. H.-T.; Metzger, F. Binding to SMN2 Pre-mRNA-Protein Complex Elicits Specificity for Small Molecule Splicing Modifiers. *Nat. Commun.* **2017**, *8* (1), 1476. <https://doi.org/10.1038/s41467-017-01559-4>.
- (49) Lipinski, C. A.; Lombardo, F.; Dominy, B. W.; Feeney, P. J. Experimental and Computational Approaches to Estimate Solubility and Permeability in Drug Discovery and Development Settings. *Adv Drug Deliv Rev* **1997**, *23* (1–3), 3–25.
- (50) Krause, K. M.; Serio, A. W.; Kane, T. R.; Connolly, L. E. Aminoglycosides: An Overview. *Cold Spring Harb. Perspect. Med.* **2016**, *6* (6). <https://doi.org/10.1101/cshperspect.a027029>.
- (51) Lee, M.-K.; Bottini, A.; Kim, M.; Bardaro, M. F.; Zhang, Z.; Pellecchia, M.; Choi, B.-S.; Varani, G. A Novel Small-Molecule Binds to the Influenza A Virus RNA Promoter and Inhibits Viral Replication. *Chem. Commun.* **2013**, *50* (3), 368–370. <https://doi.org/10.1039/C3CC46973E>.

- (52) Wehrhan, L.; Hillisch, A.; Mundt, S.; Tersteegen, A.; Meier, K. Druggability Assessment for Selected Serine Proteases in a Pharmaceutical Industry Setting. *ChemMedChem* **2020**, *15* (21), 2010–2018. <https://doi.org/10.1002/cmdc.202000425>.
- (53) Mj, B.; H, V.; A, W.; S, S.; I, S.; C, S.; Rs, B.; Dw, L.; T, L. cryoEM-Guided Development of Antibiotics for Drug-Resistant Bacteria <https://pubmed.ncbi.nlm.nih.gov/30667174/>. <https://doi.org/10.1002/cmdc.201900042>.
- (54) Flinders, J.; DeFina, S. C.; Brackett, D. M.; Baugh, C.; Wilson, C.; Dieckmann, T. Recognition of Planar and Nonplanar Ligands in the Malachite Green–RNA Aptamer Complex. *ChemBioChem* **2004**, *5* (1), 62–72. <https://doi.org/10.1002/cbic.200300701>.
- (55) A Small-Molecule Probe Induces a Conformation in HIV TAR RNA Capable of Binding Drug-Like Fragments. *J. Mol. Biol.* **2011**, *410* (5), 984–996. <https://doi.org/10.1016/j.jmb.2011.03.039>.
- (56) Structural insights into species-specific features of the ribosome from the human pathogen *Mycobacterium tuberculosis* | Nucleic Acids Research | Oxford Academic <https://academic.oup.com/nar/article/45/18/10884/4103595> (accessed Nov 11, 2020).
- (57) Lin, A. H.; Murray, R. W.; Vidmar, T. J.; Marotti, K. R. The Oxazolidinone Eprezolid Binds to the 50S Ribosomal Subunit and Competes with Binding of Chloramphenicol and Lincomycin. *Antimicrob. Agents Chemother.* **1997**, *41* (10), 2127–2131.
- (58) Loubresse, N. G. de; Prokhorova, I.; Holtkamp, W.; Rodnina, M. V.; Yusupova, G.; Yusupov, M. Structural Basis for the Inhibition of the Eukaryotic Ribosome. *Nature* **2014**, *513* (7519), 517–522. <https://doi.org/10.1038/nature13737>.
- (59) Belousoff, M. J.; Venugopal, H.; Wright, A.; Seoner, S.; Stuart, I.; Stubenrauch, C.; Bamert, R. S.; Lupton, D. W.; Lithgow, T. CryoEM-Guided Development of Antibiotics for Drug-Resistant Bacteria. *ChemMedChem* **2019**, *14* (5), 527–531. <https://doi.org/10.1002/cmdc.201900042>.
- (60) Blaha, G.; Gürel, G.; Schroeder, S. J.; Moore, P. B.; Steitz, T. A. Mutations Outside the Anisomycin-Binding Site Can Make Ribosomes Drug-Resistant. *J. Mol. Biol.* **2008**, *379* (3), 505–519. <https://doi.org/10.1016/j.jmb.2008.03.075>.
- (61) Warner, K. D.; Homan, P.; Weeks, K. M.; Smith, A. G.; Abell, C.; Ferré-D'Amaré, A. R. Validating Fragment-Based Drug Discovery for Biological RNAs: Lead Fragments Bind and Remodel the TPP Riboswitch Specifically. *Chem. Biol.* **2014**, *21* (5), 591–595. <https://doi.org/10.1016/j.chembiol.2014.03.007>.
- (62) Warner, K. D.; Chen, M. C.; Song, W.; Strack, R. L.; Thorn, A.; Jaffrey, S. R.; Ferré-D'Amaré, A. R. Structural Basis for Activity of Highly Efficient RNA Mimics of Green Fluorescent Protein. *Nat. Struct. Mol. Biol.* **2014**, *21* (8), 658–663. <https://doi.org/10.1038/nsmb.2865>.
- (63) Edwards, T. E.; Ferré-D'Amaré, A. R. Crystal Structures of the Thi-Box Riboswitch Bound to Thiamine Pyrophosphate Analogs Reveal Adaptive RNA-Small Molecule Recognition. *Structure* **2006**, *14* (9), 1459–1468. <https://doi.org/10.1016/j.str.2006.07.008>.

## Supplementary material

### Table of Content:

#### *Tables*

- *Table S1:* Descriptor names, descriptions/formulae, descriptor type, mean SHAP value
- *Table S2:* Overview of binding site families and consensus score of the RNA dataset as assessed by DrugPred\_RNA
- *Table S3:* Overview of the binding site families and consensus score of the ribosomal dataset as assessed by DrugPred\_RNA
- *Table S4:* Drug-like ligands ( $\text{QED} \geq 0.67$ ) found in binding sites predicted to be druggable
- *Table S5:* Drug-like ligands ( $\text{QED} \geq 0.67$ ) found in binding sites predicted to be less druggable

#### *Figures*

- *Figure S1:* Training and testing set accuracy error during the construction of DrugPred\_RNA
- *Figure S2:* Individual SHAP values for each pocket in the training set for all descriptors in the model plotted against the descriptor values.
- *Figure S3:* Druggability predictions with DrugPred\_RNA for the NRDL training and test set

## Tables

Table S1: Descriptors to describe size, polarity, and shape of ligand binding sites together with mean SHAP values for descriptors included in the final model.

Descriptor name	Description	Descriptor type	Mean absolute SHAP value
<i>csa</i>	Sum of SASA of binding site atoms	size	NA
<i>psa</i>	Sum of SASA of polar binding site atoms	polarity	NA
<i>psa_r</i>	$psa / csa$	polarity	1.46
<i>hsa</i>	Sum of SASA of hydrophobic binding site atoms	polarity/size	0.30
<i>ali_sa_r</i>	Sum of SASA of aliphatic binding site atoms / <i>csa</i>	polarity	NA
<i>exp_sl_sa</i>	Sum of SASA of superligand atoms that are solvent exposed in the superligand-receptor complex	shape	0.224
<i>no_sl_atoms</i>	Number of superligand atoms	size	0.202
<i>no_bs_atoms</i>	Number of binding site atoms	size	0.167
<i>fr_buried_sl_atoms</i>	Number of atoms buried inside the superligand / number of superligand surface atoms	shape	0.345
<i>fr_hpb_atoms</i>	Number of hydrophobic binding site atoms / <i>no_bs_atoms</i>	polarity	0.629
<i>sl_bs_r</i>	$no\_sl\_atoms / no\_bs\_atoms$	shape	0.193
<i>vol</i>	Volume of superligand	size	NA
<i>sa_vol_r</i>	Surface area of superligand / <i>vol</i>	shape	0.0907
<i>PM1</i>	First principal moment of inertia	shape/size	NA
<i>PM2</i>	Second principal moment of inertia	shape/size	NA
<i>PM3</i>	Third principal moment of inertia	shape/size	0.277
<i>NPR1</i>	$PM1 / PM3$	shape	NA
<i>NPR2</i>	$PM2 / PM3$	shape	NA
<i>Asphericity</i>	$0.5 \frac{(PM3 - PM1)^2 + (PM3 - PM1)^2 + (PM2 - PM1)^2}{PM1^2 + PM2^2 + PM3^2}$	shape	NA

<i>Eccentricity</i>	$\frac{\sqrt{PM3^2 - PM1^2}}{PM3^2}$	shape	NA
<i>SpherocityIndex</i>	$\frac{3 \times PM1}{(PM1 + PM2 + PM3)}$	shape	0.0825
<i>RadiusOfGyration</i>	$\sqrt{\frac{2\pi \frac{PM3 \times PM2 \times PM1}{3}}{MW}}$	shape	NA
<i>InertialShapeFactor</i>	$\frac{PM2}{PM1 \times PM3}$	shape	0.0849

Table S2: RNA families based on binding site sequence similarity. For each family, a head with PDB ID and three letter code of the small molecule bound to the pocket are listed. The total number of family members and the consensus score are also given. The druggability column contains the prediction that the majority of the members in each family obtained.

Family	Head	Description	Members	Consensus score	Druggability
1	2gis_SAM	SAM-I RS	23	100.0	druggable
2	1j7t_PAR	Ribosomal binding site	17	41.2	less druggable
3	1y26_ADE	Adenine/Guanine RS	11	100.0	less druggable
4	1o9m_42B	Ribosomal binding site/ HIV-1 Kissing loops	18	88.9	less druggable
5	2yie_FMN	FMN RS	16	100.0	druggable
6	3ds7_GNG	Purine/ deoxyguanosine RS	8	50.0	less druggable
7	4qk8_2BA	c-di-AMP RS	10	45.5	less druggable
8	3irw_C2E	c-di-GMP RS	8	75.0	druggable
9	2ho7_G6P	glmS ribozyme	7	100.0	less druggable
10	2cky_TPP	TPP RS	16	12.5	Less druggable
11	1i9v_NMY	tRNA / Corn aptamer	4	100.0	less druggable
12	3d0u_LYS	Lysine RS	3	100.0	less druggable
13	4ts0_38E	Spinach aptamer	5	60.0	less druggable
14	6qiq_J48	CAG repeats	3	100.0	druggable
15	3e5c_SAM	SAM-III RS	3	100.0	druggable
16	4znp_AMZ	pfl/ZTP/ZMP RS	3	33.3	druggable
17	1byj_GE1	Ribosomal sites	5	100.0	less druggable
18	3owi_GLY	Glycine RS	3	100.0	less druggable
19	1aju_ARG	HIV-2 Trans-Activating Region	3	100.0	less druggable
20	3suh_FFO	THF RS	3	100.0	less druggable
21	2ktz_ISH	HCV IRES	3	33.3	less druggable
22	6e8s_EKJ	Mango aptamer	2	0.0	-
23	6dmc_G4P	ppGpp RS	2	100.0	druggable
24	6c63_EKJ	Mango-II aptamer	2	100.0	druggable
25	5ny8_AGU	Guanidine-III RS	2	100.0	less druggable

26	5eao_CVC	Hammerhead ribozyme	2	100.0	druggable
27	2oe5_AM2	Human ribosomal site	2	100.0	less druggable
28	5ddp_GLN	Glutamine riboswitch	2	100.0	less druggable
29	2o43_ERN	Ribosomal site	2	100.0	druggable
30	4yaz_4BW	cGAMP RS	2	100.0	druggable
31	4qlm_2BA	ydaO RS	2	100.0	druggable
32	1nta_SRY	Streptomycin-binding aptamer	2	100.0	less druggable
33	5bjo_747	Corn aptamer	2	33.3	druggable
34	4l81_SAM	SAM-I/IV RS	2	100.0	druggable
35	4k32_GET	-	2	100.0	less druggable
36	4jf2_PRF	PreQ1-II RS	2	100.0	less druggable
37	5ux3_8OS	RNA hairpin	2	100.0	less druggable
38	3td1_GET	Protozoal cytoplasmic Ribosomal site	2	100.0	less druggable
39	6dlq_PRP	PRPP RS	2	100.0	druggable
40	3sd3_FFO	THF RS	3	100.0	less druggable
41	6e1t_HLV	PreQ1-I RS	2*	0.0	-
42	3nnp_SAH	SAH RS	2	100.0	druggable
43	1f1t_ROS	Malachite green aptamer	2	100.0	druggable
44	3gca_PQ0	PreQ1-I/PreQ0 RS	2*	100.0	less druggable
45	3fu2_PRF	PreQ1-I RS	2*	100.0	less druggable
46	1eht_TEP	Theophylline-binding RNA	2	100.0	druggable

Table S3: Ribosomal binding site families based on binding site sequence similarity. For each family, a head with PDB ID and three letter code of the small molecule bound to the pocket together with its name are listed. The organism associated with the binding site, the total number of family members, and the consensus score are also given. The druggability column contains the prediction that the majority of the members in each family obtained.

Family	Head	Ligand	Organism	Members	Consensus	Druggability
1	1k8a_CAI	Carbomycin	<i>H. morismortui</i>	13	100.0	Druggable
2	1jzx_CLY	Clindamycin	<i>D. radiodurans</i> , <i>T. thermophilus</i>	17	100.0	Druggable
3	1fjg_PAR	Paromomycin	<i>T. Thermophilus</i>	56	41.9	Less druggable
4	1j5a_CTY	Clarithromycin	<i>D. radiodurans</i> , <i>T. thermophilus</i>	13	100	Druggable
5	5jup_GDP	GDP	<i>S. cerevisiae</i>	9	77.8	less druggable
6	1ttt_GNP	GNP	<i>E. coli</i> , <i>T. aquaticus</i>	10	100.0	less druggable
7	4wfa_ZLD	Linezolid	<i>S. aureus</i>	7	100.0	Druggable
8	6ole_MVM	PF846	<i>H. sapiens</i>	6	100.0	Druggable
9	1j7t_PAR	Paromomycin	16S rRNA Synthetic constructs	9	11.1	Druggable
10	1fjg_SRY	Spectinomycin	<i>T. thermophilus</i>	9	100.0	less druggable
11	4u3u_3HE	Cicloheximide	<i>S. cerevisiae</i>	6	100.0	Druggable
12	1ibk_PAR	Paromomycin	<i>T. thermophilus</i>	30	72.4	less druggable
13	6gxm_GCP	GCP	<i>E. coli</i>	4	100.0	less druggable
14	4wpo_GDP	GDP	<i>T. thermophilus</i> , <i>E. coli</i>	6	100.0	less druggable
15	3jap_GCP	GCP	<i>K. lactis</i> , <i>S. cerevisiae</i>	4	100.0	less druggable
16	3id5_SAM	SAM/SAH	<i>S. solfataricus</i>	4	50.0	Less druggable
17	5zq0_SAH	SAH	<i>S. pneumoniae</i>	4	50.0	Druggable
18	4v52_NMY	Neomycin	<i>E. coli</i> , <i>t. thermophilus</i>	4	0.00	less druggable
19	6hiv_GTP	GTP	<i>T. brucei</i>	3	33.3	Druggable
20	4u3m_ANM	Anisomycin	<i>S. cerevisiae</i>	8	100.0	Druggable

21	1m90_SPS	Sparsomycin	<i>H. morismortui</i> , <i>S. cerevisiae</i>	6	66.7	Druggable
22	5jup_SO1	Sordarin	<i>S. cerevisiae</i>	4	50.0	Less druggable
23	6gaw_GSP	GSP	<i>S. scrufa</i> , <i>H. sapiens</i>	3	100.0	less druggable
24	5lzx_GCP	GCP	<i>H. sapiens</i>	3	100.0	less druggable
25	3jct_GTP	GTP	<i>S. cerevisiae</i>	4	100.0	less druggable
26	4wf1_NEG	Negamycin	<i>E. coli</i>	4	100.0	less druggable
27	1kqs_PPU	Puromycin	<i>H. morismortui</i>	3	100.0	Druggable
28	6hiv_UTP	UTP	<i>T. brucei</i>	2	100.0	less druggable
29	4k32_GET	Geneticin	<i>Leishmania</i>	2	100.0	less druggable
30	4ji3_SRY	Streptomycin	<i>T. thermophilus</i>	2	100.0	Less druggable
31	4io9_1F2	Carbomycin A derivative	<i>D. radiodurans</i>	2	100.0	Druggable
32	1hnw_TAC	Tetracycline	<i>T. thermophilus</i>	2	100.0	Less druggable
33	3wru_SJP	Neomycin analogue	-	2	100.0	Less druggable
34	3td1_GET	Geneticin	-	2	100.0	Less druggable
35	5lzb_GNP	GNP	<i>E. coli</i>	2	100.0	less druggable
36	5kcr_6UQ	Avilamycin	<i>E. coli</i>	2	100.0	Less Druggable
37	5jvg_6NO	Avilamycin	<i>D. radiodurans</i>	2	100.0	less druggable
38	5juu_SO1	Sodarin	<i>S. cerevisiae</i>	2	0.0	-
39	3oij_SAH	SAH	<i>S. cerevisiae</i>	2	100.0	less druggable
40	3jcj_GNP	GNP	<i>E. coli</i>	2	0.0	Druggable
41	5aj4_GDP	GDP	<i>S. scrofa</i>	2	0.0	Druggable
42	3jah_ADG	ADP	<i>O. corniculm</i> , <i>H. sapiens</i>	2	100.0	Less drugga
43	3j7a_34G	Enetube	<i>P. falciparum</i>	2	100.0	less druggable
44	3cpw_ZLD	Linezolid	<i>H. morismortui</i>	2	100.0	Druggable
45	4v9o_GCP	GCP	<i>T. thermophilus</i>	2	100.0	less druggable
46	2otj_13T	13-deoxytedanolide	<i>H. morismortui</i>	2	100.0	Druggable
47	4v56_SCM	Spectinomycin	<i>E. coli</i>	2	100.0	less druggable
48	2g5k_AM2	Apramycin	<i>Leishmania</i>	2	100.0	Less druggable

<b>49</b>	4u56_BLS	Blasticidin S	<i>S. cerevisiae</i> , <i>T. Thermophilus</i>	2	100.0	Les druggable
<b>50</b>	6rxt_GTP	GTP	<i>C. thermophilum</i>	2	100.0	Less druggable
<b>51</b>	4u4y_PCY	Pactamycin, amicoumacin	<i>S. cerevisiae</i>	2	100.0	Druggable
<b>52</b>	6sg9_SAH	SAH	<i>T. Brucei</i>	2	100.0	Druggable

Table S 4: Drug-like ligands ( $\text{QED} \geq 0.67$ ) found in binding sites predicted to be druggable

Ligand ID	PDB ID	Receptor name	QED score	$K_D$ [nM]	LE [kcal·mol <sup>-1</sup> ·heavy atom <sup>-1</sup> ]
RNA data set					
MGR	1q8n	Malachite green aptamer	0.76	800 <sup>54</sup>	0.34
6YG	5kx9	FMN RS	0.69	13.4 <sup>41</sup>	0.41
L8H	2l8h	HIV-1 TAR RNA	0.67	NA <sup>#55</sup>	-
PMZ	1lvj	HIV-1 TAR RNA	0.85	27,000 <sup>44</sup>	0.22
Ribosomal data set					
917	5v7q	50S ribosomal subunit	0.94	700 <sup>56</sup>	0.39
ZLD	3cpw	50S ribosomal subunit	0.89	20,000 <sup>57</sup>	0.27
G6M	6ddg	50S ribosomal subunit	0.79	2,600 <sup>53</sup>	0.31
3HE	4u3u	80S ribosome	0.76	140 <sup>58</sup>	0.48
G6V	6ddd	50S ribosomal subunit	0.76	2,600 <sup>59</sup>	0.30
ANM	3cc4	50s ribosomal subunit	0.78	20,000 <sup>60</sup>	0.34
HN8	5on6	80S ribosome	0.71	NA <sup>#</sup>	-
3K8	4u55	80S ribosome	0.71	39	0.32

\* RS = riboswitch, # binding affinity unknown

Table S 5: : Drug-like ligands ( $\text{QED} \geq 0.67$ ) found in binding sites predicted to be less druggable

Ligand ID	PDB code	Receptor name	QED	$K_D$ [nM]	LE [kcal·mol <sup>-1</sup> ·heavy atom <sup>-1</sup> ]
RNA data set					
VIB	4nyg	TPP RS	0.79	1,500 <sup>61</sup>	0.45
2QC	4nyb	TPP RS	0.77	103,000 <sup>61</sup>	0.43
0EC	2lwk	Influenza A	0.86	50,000 <sup>51</sup>	0.29
1TU	5ob3	Spinach aptamer	0.85	530 <sup>62</sup>	0.49
218	2hop	TPP RS*	0.77	6,000 <sup>63</sup>	0.38
Ribosomal data set					
TRP	4v6o	Tryptophan-sensing ribosomal site	0.67	NA#	-

RS = riboswitch, # binding affinity unknown

## Figures

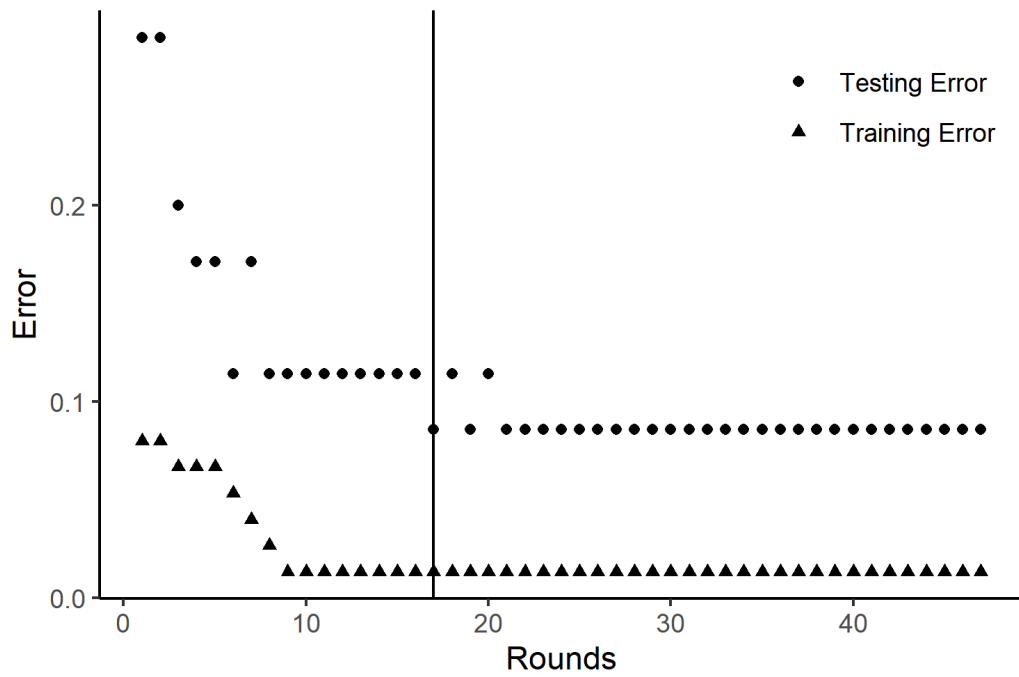


Figure S 1: Training and test set accuracy error vs. rounds of xgboost iteration. The line denotes the earliest iteration (17) where the error on the test set has not improved in the following 20 iterations

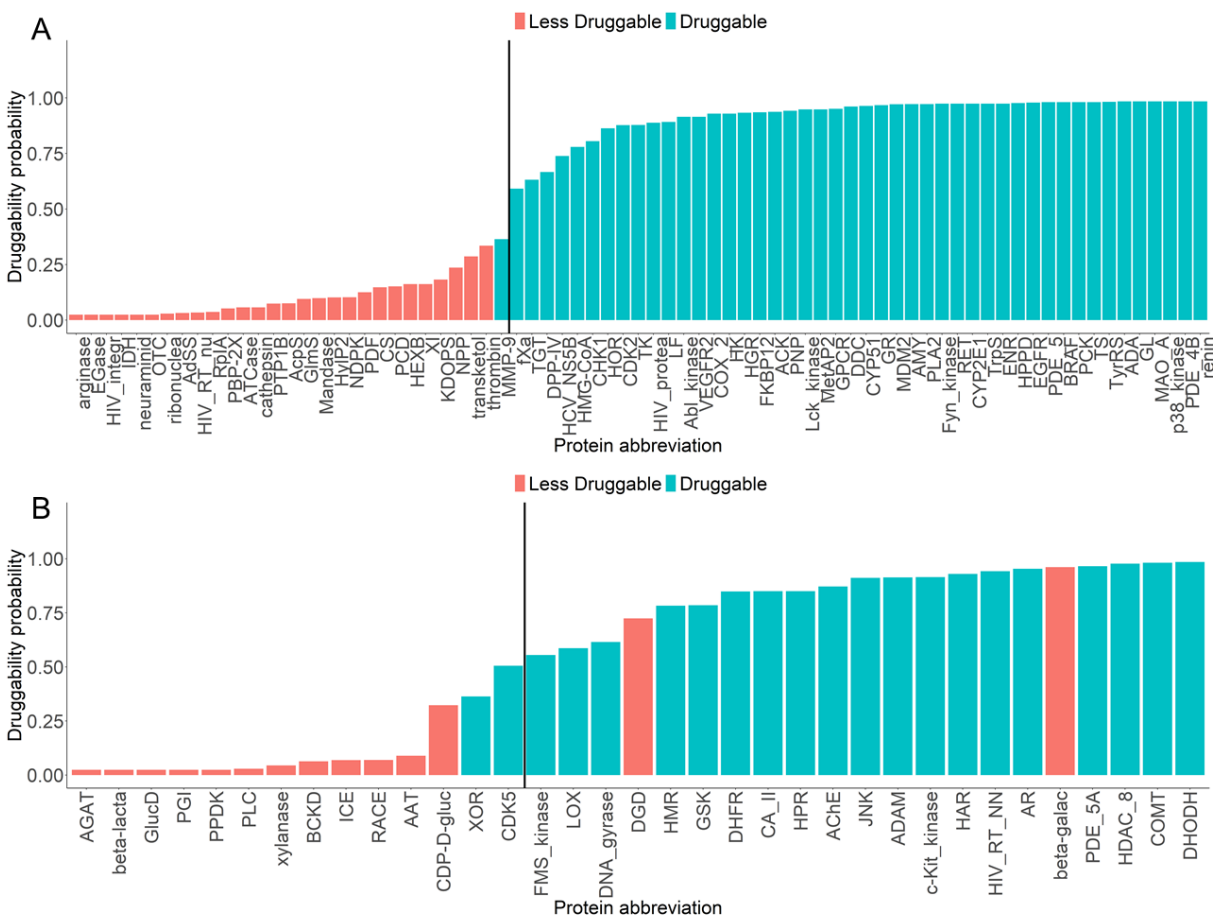


Figure S2: Druggability predictions with DrugPred\_RNA for the NRDLD training (A) and test set (B). Cyan bars represent druggable and red bars less druggable binding sites. All pockets at the right side of the black line are classified as druggable while the pockets on the left side of the line are classified as less druggable. The full names of the proteins are listed in <sup>21</sup>.

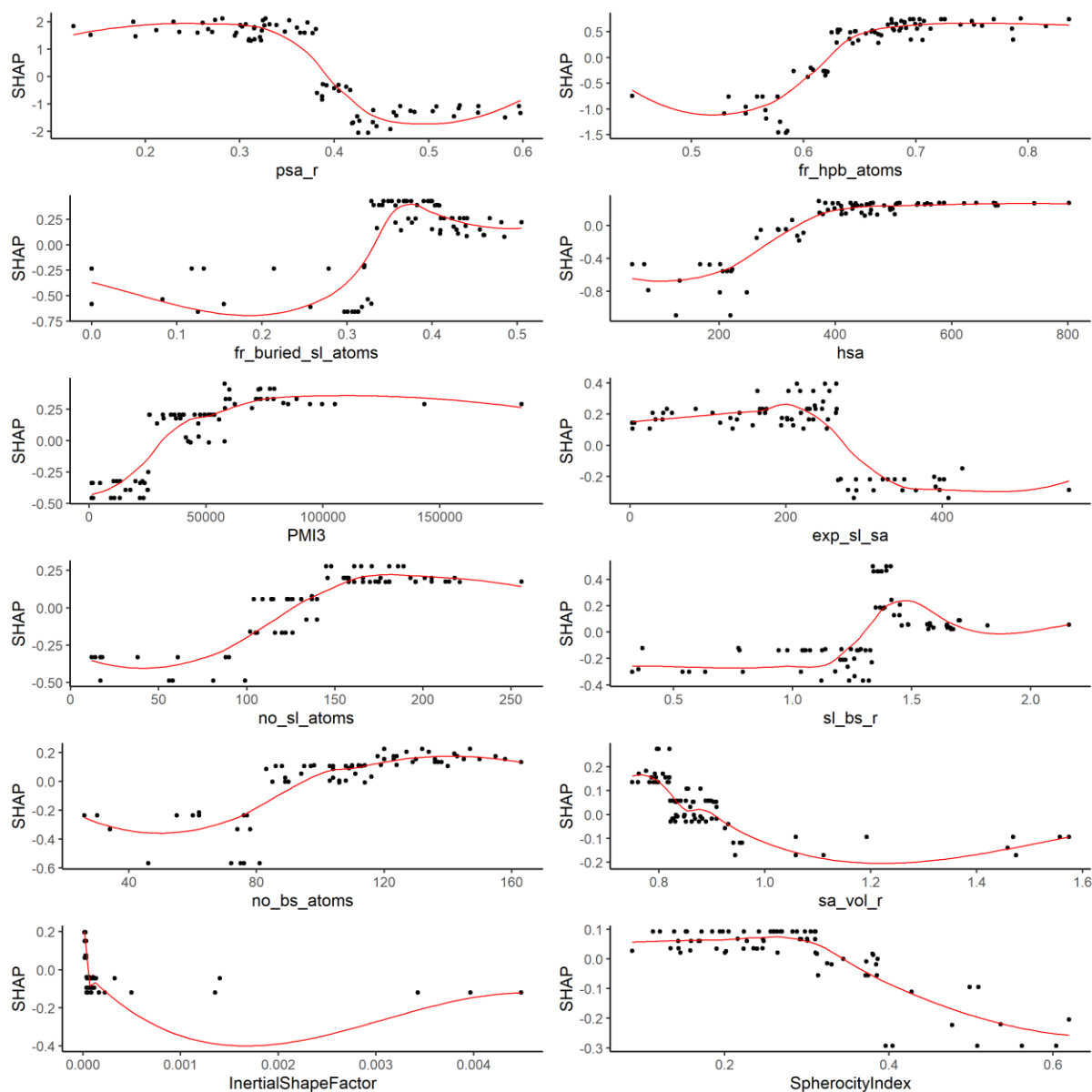


Figure S 3: Individual SHAP values for each pocket in the training set for all descriptors in the final model plotted against the descriptor values. Locally estimated scatterplot smoothing (LOESS) curves are overlaid on the descriptor observations (black dots). The midpoint in each curve indicates the cut-off value from where the model's prediction changes the direction.