

---

# Towards a more inductive world for drug repurposing approaches

---

Jesus de la Fuente<sup>1,2,†</sup>, Guillermo Serrano<sup>1,2,†</sup>, Uxía Veleiro<sup>1,†</sup>, Mikel Casals<sup>2</sup>  
Laura Vera<sup>1</sup>, Marija Pizurica<sup>3,4</sup>, Antonio Pineda-Lucena<sup>1</sup>, Idoia Ochoa<sup>2,5</sup>  
Silve Vicent<sup>1</sup>, Olivier Gevaert<sup>3,\*</sup> and Mikel Hernaez<sup>1,5,\*</sup>

<sup>1</sup>CIMA University of Navarra, IdiSNA, Pamplona, Spain.

<sup>2</sup>TECNUN, University of Navarra, San Sebastián, Spain.

<sup>3</sup>Stanford Center for Biomedical Informatics Research, Stanford University, California.

<sup>4</sup>Internet technology and Data science Lab (IDLab), Ghent University, Belgium.

<sup>5</sup> Instituto de Ciencia de los Datos e Inteligencia Artificial (DATAI), University of Navarra, Spain.

\* Corresponding authors, † These authors have contributed equally.

jdlfuentec@gmail.com uxveleiro@gmail.com gserranos@unav.es  
ogevaert@stanford.edu mhernaez@unav.es

## Abstract

Drug-target interaction (DTI) prediction is a challenging, albeit essential task in drug repurposing. Learning on graph models have drawn special attention as they can significantly reduce drug repurposing costs and time commitment. However, many current approaches require high-demanding additional information besides DTIs that complicates their evaluation process and usability. Additionally, structural differences in the learning architecture of current models hinder their fair benchmarking. In this work, we first perform an in-depth evaluation of current DTI datasets and prediction models through a robust benchmarking process, and show that DTI prediction methods based on transductive models lack generalization and lead to inflated performance when evaluated as previously done in the literature, hence not being suited for drug repurposing approaches. We then propose a novel biologically-driven strategy for negative edge subsampling and show through *in vitro* validation that newly discovered interactions are indeed true. We envision this work as the underpinning for future fair benchmarking and robust model design. All generated resources and tools are publicly available as a python package.

## 1 Introduction

Drug discovery aims at finding the most effective pharmacological compound that can target a specific disease-causing mechanism while yielding minimal side effects. Traditionally, predicting drug-target interactions (DTIs) has relied on determining physical parameters between both components, such as the dissociation constant or the inhibitory concentration [1, 2]. However, experimental screenings of compounds have a limited success rate and require time, effort, and elevated monetary costs [3], which considerably hinders the process of finding new drugs interacting with the intended targets. High throughput sequencing technologies have unveiled thousands of interesting targets with many potential modulators, making the experimental screening of compounds a daunting challenge.

This new paradigm, fueled by the current availability of large amounts of biological data, has promoted breakthrough deep learning on graphs [4] approaches that have accelerated the first stages of drug discovery pipelines by narrowing down the most promising DTIs [5, 6, 7, 8, 9]. However, there are

currently four critical challenges that prevent proper performance evaluations of newly proposed DTI prediction models: **i)** Current gold standard datasets for evaluating DTI prediction models, such as the well-known Yamanishi dataset [10] are small, outdated and missing many interactions. **ii)** State-of-the-art DTI prediction methodologies require additional information to predict novel DTIs, which is generally not readily available and thus restricts their usage and evaluation. **iii)** Current methods can be divided into inductive or transductive, based on whether they can learn underlying patterns in the data to make predictions on unseen samples (inductive), or whether they directly build a prediction model for the seen ones (transductive). These structural differences make it challenging to fairly compare these methodologies. **iv)** Current techniques for dealing with the existing positive/negative edge imbalance when training DTI prediction models do not incorporate biological information, potentially affecting subsequent experimental validation.

To address these limitations, in this work we perform an in-depth evaluation of current state-of-the-art DTI prediction methodologies, taking into account drug repurposing datasets, transductive and inductive learning and DTIs network splitting and subsampling techniques. We demonstrate that designing DTI prediction methods using transductive-based approaches are not optimal, and recommend utilizing inductive-based ones instead. Specifically, we show that a baseline transductive classifier achieves near optimal performance only due to data leakage. Additionally, we introduce a technique based on Root Mean Square Deviation (RMSD) for subsampling negative edges during the construction of the DTI dataset, and show that it can lead to the discovery of true interactions (validated *in vitro*) otherwise missed. Finally, we provide data and tools to ease the design and the fair evaluation of novel DTI methods as we envision this work as the first step towards a community-driven unified benchmark to assess novel DTI prediction approaches.

## 2 Pearls and pitfalls of current DTI datasets and prediction models

### 2.1 Overview of current DTI datasets

In recent years multiple datasets [11, 12, 13, 14] have been generated for in-silico validation of DTI prediction models. While these datasets typically contain a set of targets and their interacting drugs, they strongly differ in the origin of the data, as well as the topology and size of the network (Appx. Table A1, Appx. Note 1). For example, the current gold standard datasets were defined by Yamanishi et al. in 2008, which consist of four precise, albeit small datasets (less than 100 edges) divided by protein families: enzymes (E), ion channels (IC), G-protein-coupled receptors (GPCR) and nuclear receptors (NR) [10]. This contrasts with recent datasets derived from DrugBank, such as DrugBank-DTI [11] and BIOSNAP [12] which contain more than 15000 edges. In addition, data from drug-target binding affinity experiments has also been used for DTI prediction tasks (e.g., DAVIS [13] and BindingDB [14] datasets), with the caveat that data must be previously binarized at an arbitrary threshold of affinity.

**Chemically-driven evaluation of current DTI datasets.** To enable an accurate DTI prediction and avoid introducing a bias towards certain chemical drug categories, datasets should encompass drugs with high chemical diversity and high promiscuity (i.e., high capability to interact with multiple protein families) [15, 16]. When assessing the aforementioned datasets for these properties, we found that drugs within datasets are indeed chemically diverse (i.e., their pair-wise Tanimoto distance followed a 0-skewed distribution, Appx. Fig. A1, A2), and promiscuous (Appx. Fig. A3, A4). Further, datasets should comprise diverse protein families to enable generalization capabilities of DTI models. However, the included protein families are highly variant across datasets, with some containing a wide range (e.g., DrugBank) and others being family-specific (e.g., Yamanishi Enzymes, Appx. Fig. A5, A6, A7 and A8). This analysis revealed that while the latest DTI datasets such as DrugBank are suitable for training DTI prediction models, the still-considered gold standard such as Yamanishi should be used with caution since it may introduce bias towards certain protein families.

### 2.2 Overview of current DTI prediction models

Depending on their learning process on the DTI graph, current DTI models can be classified into two groups: inductive and transductive. Inductive graph learning involves using a set of labeled nodes/edges to learn the underlying data structure, aiming to make predictions on unseen samples leveraging the knowledge acquired during training in the form of weights. Transductive graph learning, on the other hand, does not build a predictive model from seen nodes/edges, as there are

no weights that can be used to predict a set of unseen samples. Instead, it uses every sample in the dataset to generate the desired prediction. The following models have been recently shown to achieve state-of-the-art performance [17, 18]: DTINet [19], DDR [7], DTiGEMS+ [20] and DTi2Vec [21] which fall under the transductive category, and NeoDTI [22], Moltrans [23], Hyper-Attention-DTI [9] and EEG-DTI [8] which fall into the inductive category (Appx. Table A2, Appx. Note 2).

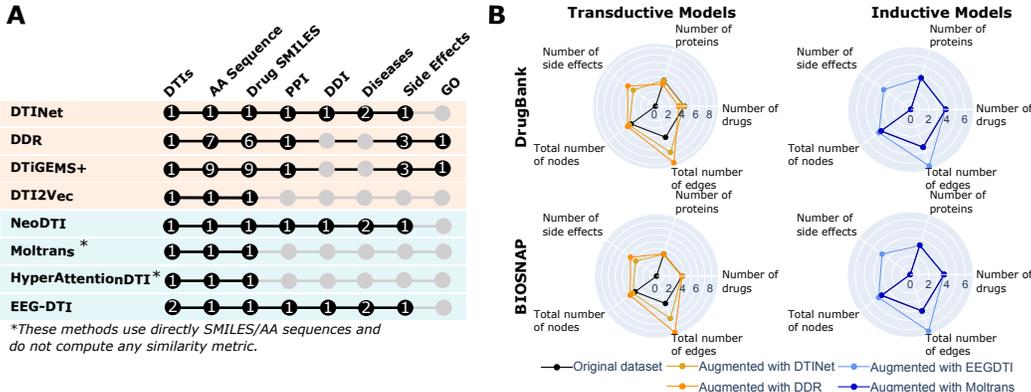


Figure 1: **Analysis of required resources and network augmentation.** **A.** Number of different side information matrices used by each model. **B.** Radar plots depicting the original and modified number of nodes and edges (log10) for DrugBank and BIOSNAP, when used as input for different models.

**Data requirements for current DTI prediction methods hinder their applicability.** DTI prediction methods typically augment the above-mentioned DTI datasets to include additional information beyond DTIs, such as protein-protein interactions or side-effect-drug associations. For example, methods like DTINet, DDR or DTiGEMS+ [19, 7, 20] require collecting information from several complex data sources, such as Side Effects from SIDER [24], or Diseases from CTD [25] which hinders their usability (Fig. 1-A, Appx. Note 1). Generating these heterogeneous networks requires accessing information that may not be always readily available due to the inconsistency of identifiers across databases. Further, the absence of drug-target pairs in any required additional matrix precludes some models from including such pairs in the final graph. Indeed, the original number of proteins and drugs when using DrugBank, BindingDB, and NR datasets are considerably shrunk when used in high-demanding side-information models, losing up to 82% of the drug nodes and 72% of the protein nodes for DDR in DrugBank (Fig. 1-B). In approaches that require less demanding side information, such as Moltrans, the dataset size is maintained.

**A publicly available resource of augmented DTI datasets.** To enable robust benchmarking across DTI prediction models with different augmented datasets, we built an augmented version of the most-used DTI datasets, including the gold standard. We computed all complementary matrices with the latest data releases required by every evaluated DTI prediction model (Fig. 1-A, Supp. Note 1).

Table 1: **AUC (mean and std) for the evaluated DTI prediction models.** AUC shown for the eight state-of-the-art models (first four are transductive, second four are inductive) for every dataset.

Method	DrugBank	BIOSNAP	BindingDB	DAVIS	Yamanishi E	Yamanishi IC	Yamanishi GPCR	Yamanishi NR
DTINet	0.815 ± 0.000	0.856 ± 0.001	0.853 ± 0.005	0.813 ± 0.008	0.904 ± 0.004	0.737 ± 0.010	0.792 ± 0.003	0.803 ± 0.022
DDR	OOT	OOT	OOT	0.934 ± 0.020	0.975 ± 0.008	0.986 ± 0.006	0.957 ± 0.023	0.917 ± 0.049
DTi-GEMS	OOT	OOT	0.922 ± 0.034	0.843 ± 0.160	0.968 ± 0.230	0.973 ± 0.080	0.777 ± 0.107	0.935 ± 0.120
DTi2Vec	OOT	OOT	1.000 ± 0.004	0.862 ± 0.001	0.999 ± 0.009	0.996 ± 0.007	0.992 ± 0.004	0.968 ± 0.007
NeoDTI	OOT	OOT	OOT	OOT	OOT	0.974 ± 0.006	0.924 ± 0.021	0.825 ± 0.094
Moltrans	0.798 ± 0.008	0.792 ± 0.008	0.896 ± 0.009	0.693 ± 0.013	0.883 ± 0.011	0.875 ± 0.012	0.737 ± 0.037	0.544 ± 0.073
HyperAttentionDTI	0.862 ± 0.003	0.862 ± 0.003	0.954 ± 0.003	0.738 ± 0.003	0.953 ± 0.002	0.949 ± 0.006	0.804 ± 0.003	0.459 ± 0.001
EEG-DTI	0.889 ± 0.005	0.902 ± 0.010	0.893 ± 0.021	0.741 ± 0.034	0.951 ± 0.009	0.940 ± 0.017	0.870 ± 0.041	0.765 ± 0.084

OOT and OOR mean *out of time* ( $\geq 7$  days) and *out of RAM* ( $\geq 250$  GBs), respectively.

**Evaluation of current DTI prediction models.** Using these augmented datasets, we then evaluated the above-mentioned methods following the originally proposed evaluation benchmark (Appx. Table A2, default splitting column). Transductive models yielded significantly better AUC and AUPRC results (Table 1, Appx. Table A3, A4). However, except for DTINet, they did not converge on the two largest networks (DrugBank and BIOSNAP), potentially due to their large DTI network size, which gets further augmented with the needed additional matrices. On the other hand, inductive models such as Moltrans and HyperAttentionDTI obtained low AUCs in the smallest network, NR, indicating that the size of the network may be hampering the model learning capabilities.

### 3 Graph ablation studies are crucial for fair and robust benchmarking

Since transductive methodologies can present data leakage during feature generation [26], we hypothesized that the significant AUC discrepancies among inductive and transductive methods shown in the previous section could be a consequence of data leakage. This would indicate that the high AUC values achieved by these methods are not representative of their true ability to predict interactions as the model evaluation setting artificially raises the performance. As such, transductive methodologies should be carefully applied when designing DTI prediction approaches as information used in generating node embeddings is shared with the test set. Therefore, it becomes remarkably challenging to establish fair benchmarking guidelines as one must ensure that each model leverages the data in its intended way while ensuring fairness in cross-model comparisons.

**Previously proposed ablation studies.** To start addressing these challenges, recent reviews have proposed various scenarios regarding how drugs and proteins should be distributed among train-test splits [27, 7]:  $S_p$ , where drugs and proteins are shared within train-test splits,  $S_d$ , where only proteins are shared, and  $S_t$ , where only drugs are shared. Despite these efforts to provide a homogeneous benchmarking practice, current DTI prediction models are typically evaluated considerably differently by their authors which makes it challenging to compare the performance across DTI prediction models (Appx. Table A2).

**Ablation studies on state-of-the-art DTI prediction models.** We next performed an evaluation of the above-mentioned DTI prediction models using the different split scenarios for all the generated augmented DTI datasets (see previous section). As hypothesized, the overall results showed that the transductive models suffered a higher loss of AUC compared to their default split while the inductive models achieved a similar accuracy (Fig. 2). Interestingly, the transductive approach DTi2Vec, which exhibited exceptionally high performance during the default run by attaining a 0.999 AUC score on the Yamanishi’s Enzyme dataset, yielded a performance decline to 0.61 when utilizing the  $S_p$  split. The same behavior was observed when comparing the AUC scores within the models;  $S_d$  and  $S_t$  accuracies tended to be lower than those for  $S_p$  (Fig. 2). Part of this variability appears to depend on the evaluated model, as demonstrated by the performance of different models on the DAVIS dataset. Specifically, the AUC scores were consistently higher for both  $S_p$  and  $S_t$  splits and lower for the  $S_d$  split across most models. However, significant performance discrepancies were observed for the EEG-DTI and DTi2Vec models.

Furthermore, smaller datasets (e.g., Yamanishi’s NR) yielded a significant decrease of accuracy. This phenomenon was especially notable in DTi2Vec and DTiGEMS+, where removing edges hindered training on  $S_p$  and aggravated both training and testing in  $S_d$  and  $S_t$ , as only very few edges remain in the network. A similar trend was observed in inductive methods like Moltrans and HyperAttentionDTI. While these methods can still be trained, their AUC score falls below 0.5, suggesting an inability to perform the task effectively, thereby emphasizing the need for larger networks.

**Summary.** The performed ablation studies based on different train-test types of splits offer a more complete and realistic evaluation of DTI prediction approaches and suggest a lack of generalization capabilities of current transductive methodologies.

### 4 Train-to-Test data leakage in transductive DTI models prevents them from generalizing and yields inflated performance

We noticed that the best-performing transductive models, DTiGEMS+ and DTi2Vec, shared the use of *node2vec* (N2V) to generate the node embeddings for the DTI network [28]. As the node embeddings in N2V are built by local neighborhood visits within the network, it requires to be rerun

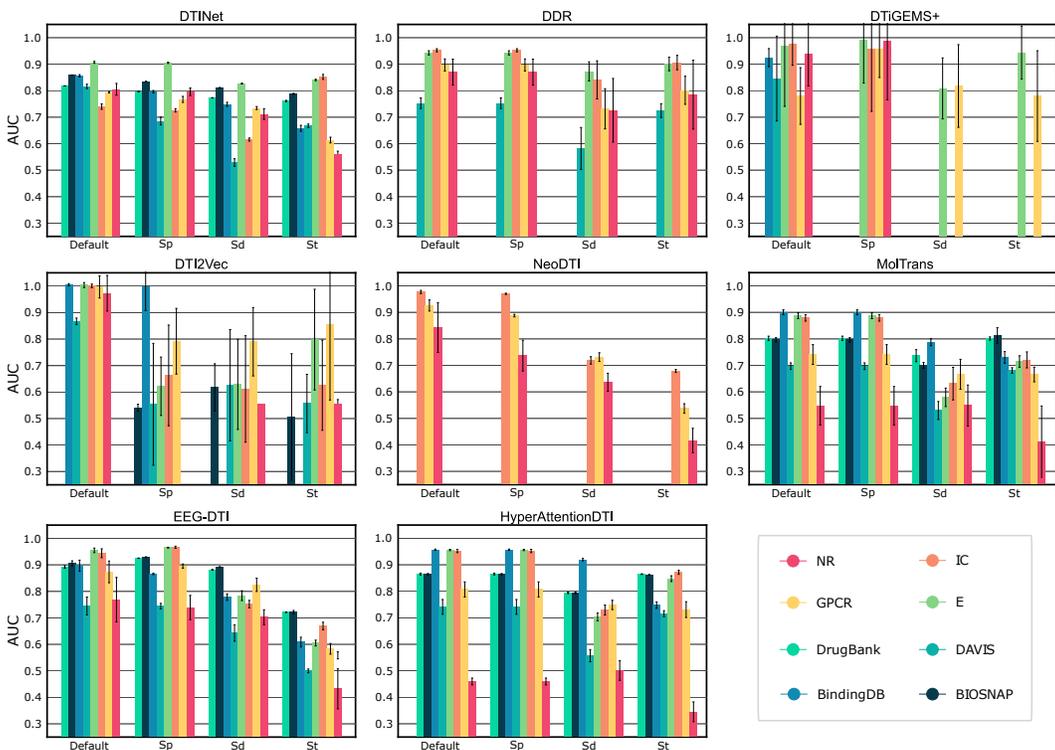


Figure 2: **DTI prediction models benchmarking.** AUC results for each dataset and model (see Appx. Tables A5, A6, A7 for AUPRC). Results correspond to an average of 5-Times 10-Fold Cross Validation. DDR, DTiGEMS+, MolTrans and HyperAttentionDTI models use  $S_p$  as the default split.

whenever a new sample is included in the dataset. Thus, when used in DTI prediction methodologies, if the DTI network embedding using N2V occurs prior to the network splitting, it can promote data leakage issues when performing traditional train/test folds evaluation.

**Design of a baseline transductive model.** To delve into this potential data leakage, we designed a baseline model (Fig. 3-A) based on N2V followed by a shallow neural network (SNN). We then performed a grid-search over multiple model parameters (Appx. Table A8) following a train/val/test evaluation setup, across all assessed datasets. We report the test AUC scores, and found that for all datasets there is an optimal embedding dimension for which the variance of the AUC scores is minimized, while maintaining a high AUC score (Appx. Fig. A9). This variance is mostly influenced by the size of the datasets, as it decreases for larger networks.

**Assessment of generalization capabilities of transductive DTI prediction models.** We trained and tested the baseline model on different DTI networks, generating an AUC matrix (Fig. 3-B). The matrix’s diagonal represents the test AUC when both the training and testing data come from the same dataset, aligning with the benchmarking process for the DTI prediction models (Table 1). These results raised concerns about their reliability, as they consistently outperformed other evaluated methods across all datasets without leveraging any additional biological information.

Furthermore, this near-perfect performance drastically compares with the poor performance of the upper and lower triangles (where the train and test were constructed using different datasets, Appx. Note 3). This behavior aligns with what we observed for transductive models in the previous section and reaffirms that their inflated performance (Table 1) may be due to data leakage, as information from the test fold is present on the node’s embeddings used in the train fold. Also, this analysis reveals one major drawback of the N2V approach for building DTI prediction models: the embedding process generally yields considerably different node embeddings for every network, hindering its capability to translate to unseen data. This also complicates generating embeddings on training and test folds separately to prevent data leakage, as the change in the graph topology produced by the splitting plus the transductive nature of N2V will heavily influence the generated embeddings.

**Summary.** The conclusions drawn from our baseline model can be transferable to transductive models such as the path-category-based technique in DDR, or the network diffusion algorithm (random walk with restart) in DTINet, which reduces the confidence of their obtained results. These findings also promote the adoption of inductive models for DTI prediction tasks, as the predictive models they construct during training enable testing on unseen graphs, mitigating the risk of data leakage and rendering them more suitable for a production environment.

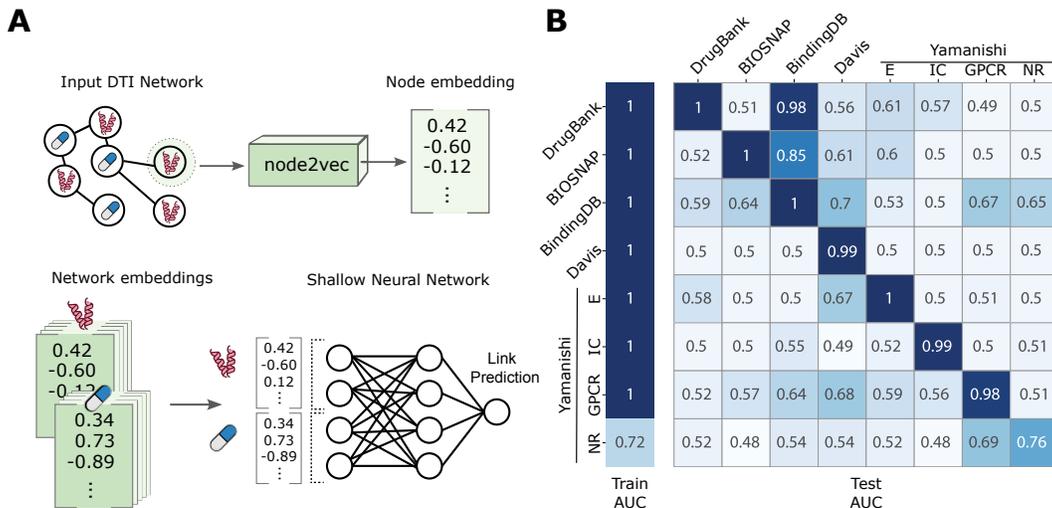


Figure 3: **Evaluation of the designed *node2vec*-based DTI prediction model.** **A.** Baseline classifier schematic. N2V embeddings are generated solely from the DTI network and drug-target pairs are fed into a SNN classifier. **B.** AUC matrix built by training on each dataset (left) and testing on another (right) (see Appx. Note 3).

## 5 Protein structure-based metric leads to improved accuracy over currently used random subsampling

When training DTI prediction models, the choice of positive and negative DTIs is still a challenging task. The sparse nature of current DTI networks, when used for classification tasks, yields very unbalanced datasets. The true edges are few and the negative edges, which are defined by all other possible connections, are orders of magnitude greater in number (see sparsity ratio in Appx. Table A1). Random subsampling is the preferred method to balance negative and positive edges. However, this can hamper the prediction task, as it is likely that the model is not trained on hard-to-classify negative samples. To address this issue, we propose a novel way to subsample negative DTIs that relies on the target’s structural information to find hard-to-classify negative DTIs (Fig. 4-A).

**Identifying informative negative edges via Root Mean Square Deviation (RMSD) between backbone alpha carbons of two proteins.** Since evolution preserves protein structure more than the sequence itself [29], we consider those drug-target pairs with potential structural interaction to be plausible (uncovered) edges, measured using the RMSD between backbone C-alpha of two proteins. Hence, this metric ranks, for each positive DTI, all the negative pairs containing the same drug according to the structural similarity, which enabled us to identify hard-to-classify samples (negative pairs with low RMSD) and select high-quality negatives (negative pairs with RMSD within a defined window) (see Appx. Note 3 for further details on protein structure and RMSD calculations). Thus, the proposed subsampling scheme consists of two differentiated steps: the proposed ranking of edges via the RMSD metric, and a selection of the negative edges for training.

**Assessment of the proposed negative sampling criteria.** Since the analysis of subsampling methods becomes particularly relevant when dealing with larger DTI datasets, we tested the proposed sampling methodology on the largest datasets: BIOSNAP and BindingDB. For the evaluation methods, we discarded approaches based on N2V because of the aforementioned potential lack of generalization, as well as slow or hard-to-evaluate methodologies. From the remaining models, we chose MolTrans

and HyperAttentionDTI, because of their inductive nature, their easiness of use, and their good performance in our previous analysis (Fig. 4-B, Appx. Fig. A10).

The RMSD criteria helped both models to generalize and obtain more robust results, for almost every selected window (see Fig. 4 and Appx. Note 3) across datasets and methodologies. Further, the AUC tends to decrease as we increase the RMSD, i.e., relax the similarity criteria. This is an expected behavior as including more easy-to-classify negative pairs into the folds can potentially decrease the total AUC score when tested on hard-to-classify DTI pairs. In summary, these results emphasize the importance of considering biologically driven criteria for future DTI prediction models' design and show the potential of the proposed negative-edge selection process for increasing the chances of uncovering novel DTIs.

**All previous analysis have been integrated into the python package *GraphGuest*.** It enables 1) easy ablation studies of input DTI graphs, 2) computation of the RMSD-based score for negative edge selection, and 3) seamless access to the generated augmented DTI datasets (Appx. Note 4).

## 6 RMSD-based subsampling enhances models capabilities for identifying novel interactions

The utilization of RMSD-based selection for negative edges leads to an improved AUC, potentially facilitating the discovery of novel DTI interactions. To further investigate this hypothesis, we examined the excluded negative DTIs (2.5 to 5 Å) within the largest selected network, BIOSNAP, using Moltrons trained on the highest yield AUC window (5 to 6 Å). From the later, we specifically chose a set of DTIs that showed a high confident prediction when using RMSD, across five independent runs. From those, we selected EGFR and GSK3 $\beta$  targets, as their inactivation elicited an antiproliferative effect, hence facilitating *in vitro* validation. Finally, we compared the RMSD-based probabilities with the random subsampling ones, finding that the proposed metric consistently reported higher and more reliable positive predictions (Fig. 4-B).

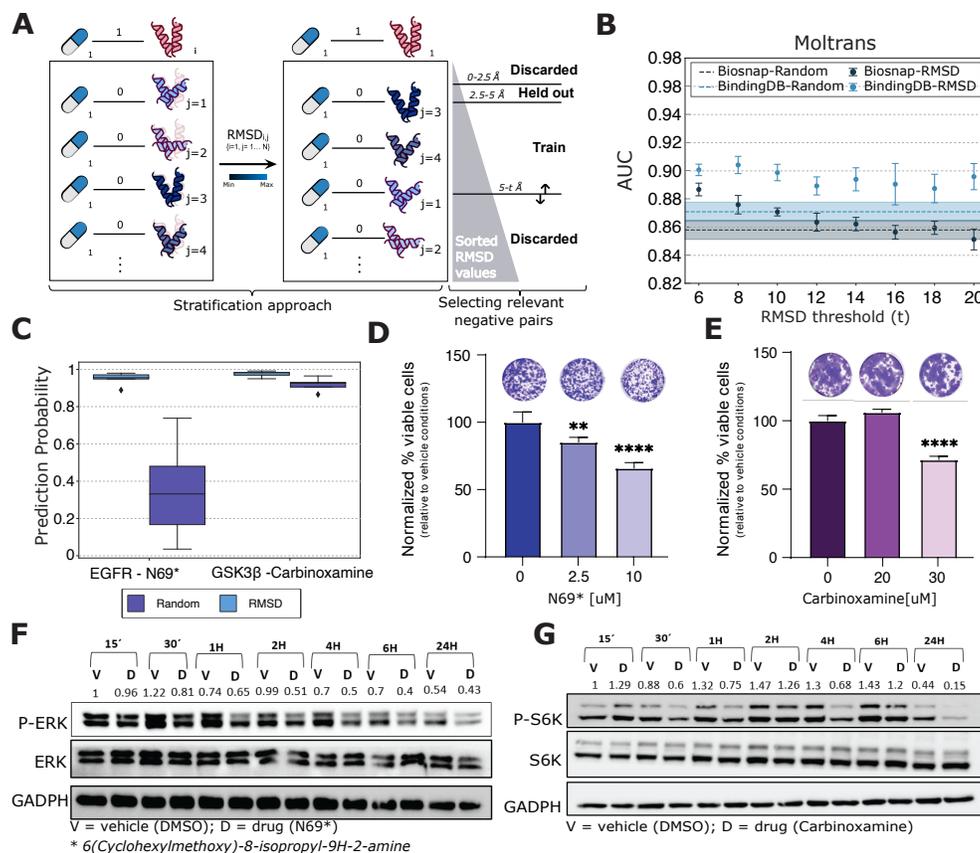
### 6.1 In vitro validation of newly uncovered drug-target interactions.

We conducted *in vitro* validation experiments to further assess their accuracy. The selected DTIs for subsequent validation were EGFR with 2-Amino-6-cyclohexylmethoxy-8-isopropyl-9H-purine (N69) and GSK3 $\beta$  with carbinoxamine due to the readily availability of both reactives at our lab.

**In vitro validation of EGFR-N69 interaction.** We used a cell line with active basal EGFR levels. We first assessed cell viability in the presence of the predicted drug. Cell viability decreased upon N69 treatment, particularly on the highest concentration of the compound, suggesting a potential interaction (Fig. 4-C). To understand if this effect was specific to the inactivation of the EGFR pathway, we examined the activation level of downstream effectors within the target pathway (Appx. Fig. A11). Specifically, we examined the activation state of the RAF-MEK-ERK pathway in HPAFII cells using antibodies specific to phospho-ERK1/2. We found that while the protein levels of ERK, as well as those of a protein loading control (GAPDH) remained unchanged in the presence of the drug, the phosphorylated active version (P-ERK) decreased as early as 30 min post-treatment. Thus, the reduction in P-ERK may be a consequence of the observed decrease in cell viability (Fig. 4-E) [30].

**In vitro validation of GSK3 $\beta$ -Carbinoxamine interaction.** Similarly, we tested the GSK3 $\beta$ -Carbinoxamine interaction, revealing a decrease in cell viability at a 30  $\mu$ M concentration (Fig. 4-D). When looking into the downstream proteins reported to mediate GSK3 $\beta$  signaling (Appx. Fig. A12), we found that the levels of phosphorylated S6 kinase, a protein involved in protein synthesis when activated through phosphorylation, were reduced 30 minutes after Carbinoxamine treatment while its non-phosphorylated version and the loading protein control remained unchanged (Fig. 4-F) [31].

**Summary.** These observations suggest that the antiproliferative effect caused by N69 and Carbinoxamine may indeed be related to inactivation of protein effectors downstream of the predicted targets EGFR and GSK3 $\beta$ . Thus, the RMSD selection method may increase the likelihood of discovering positive DTIs compared to random subsampling.



**Figure 4: Biological criteria as an alternative to random subsampling.** **A.** Alternative criteria proposed for negative subsampling based on RMSD protein structure comparison. **B.** AUCs of the test split for the RMSD-based (threshold from 6 to 20 Å) and random subsampling techniques when using Moltrons model, on Biosnap and BindingDB datasets. **C.** Prediction probabilities (across five independent runs) of Moltrons for interactions EGFR - N69, and GSKB - Carbinoxamine, when using random-based and RMSD-based subsampling. **D-E.** Percentage of cell viability and representative images of crystal violet-stained cells for **(D)** KRASG12D pancreatic cell line (HPAFII) after three days of treatment with N69 (0-10 μM) and **(E)** KRASG12C LUAD cell (H1792) after five days of treatment with Carbinoxamine (0-30 μM). **F-G.** Protein expression at different time points after N69 10 μM **(F)** or Carbinoxamine 30 μM **(G)** drugs, or DMSO vehicle treatment. Relative P-Erk  $\frac{1}{2}$  **(F)** and P-S6 Kinase **(G)** densitometry quantification is shown. 25 μg of protein were loaded per sample. GAPDH is shown as loading control.

## Discussion

Previous in-silico drug repurposing methodologies have often required high-demanding additional information, exhibited significant disparities in their evaluation framework, and employed structurally distinct learning architectures, which has resulted in a lack of a standardized benchmarking approach to determine the most suitable model. In this study, we started by assessing currently used datasets in DTI prediction problems, and generated a valuable resource of augmented DTI datasets that will enable accessible and robust future benchmarking of DTI prediction models. Using this newly generated data resource, we benchmarked diverse drug repurposing models by first using the traditional approach and then following graph-aware train-test splitting techniques. The latter revealed that methods employing transductive feature generation exhibited over-performance. This motivated further assessment of transductive approaches, which allowed uncovering data leakage issues that could be avoided if using inductive approaches.

To improve the predictive capabilities of inductive DTI models we proposed a subsampling method based on structural differences across proteins. This revealed improved accuracy when compared with traditional random subsampling, increasing the reliability of uncovering novel DTIs. Importantly, we then performed *in vitro* validation suggesting a direct interaction between drugs and protein targets leading to potential pathway inactivation, as revealed by variations in the activation levels of canonical downstream effectors of the targeted proteins.

**Conclusion.** This study emphasizes the significance of larger and diverse DTI databases, accessible drug repurposing models, data-leakage-free evaluation, and biologically driven subsampling techniques. It also presents the *GraphGuest* python package that will ease the design of drug repurposing approaches.

## Appendix

### Appendix Notes

#### 1 Datasets

Along the work, multiple DTI networks have been evaluated. Here we briefly describe them all:

- **DrugBank** [11]. DTIs collected from DrugBank Database Release 5.1.9. It has undergone significant upgrades since its first release in 2006.
- **BIOSNAP** [12]. Dataset created by Stanford Biomedical Network Dataset Collection. It contains proteins targeted by drugs on the U.S. market from DrugBank release 5.0.0 using MINER [32].
- **BindingDB** [14]. Database that consists of measured binding affinities, focusing on protein interactions with small molecules. The binarization of the dataset was done by considering interactions as positive if their K<sub>d</sub> was lower than 30 units. Data was downloaded from Therapeutics Data Commons (TDC) [33].
- **DAVIS** [13]. Dataset of kinase inhibitors to kinases covering more than 80% of the human catalytic protein kinome. The binarization of the dataset was done by considering as positive those interactions with a K<sub>d</sub> lower than 30 units. Data downloaded from Therapeutics Data Commons (TDC) [33].
- **Yamanishi et al.** [10]. It is composed of four subsets of different protein families: enzymes (E), ion channels (IC), G-protein-coupled receptors (GPCR) and nuclear receptors (NR). The Yamanishi dataset has been considered the gold standard dataset for DTI prediction and has been used in several published models [34, 8, 35]. DTIs in this dataset come from KEGG BRITE [36], BRENDA [37], SuperTarget [38] and DrugBank. Compounds with molecular weights lower than 100 are excluded from the dataset. In the Enzyme group all the ligands are inhibitors or activators, and co-factors are not included.

Also, complementary datasets were used for building the augmented networks:

- **CTD** [25]. Comparative Toxicogenomics Database for disease-drug and disease-protein associations.
- **DrugBank** [11]. DrugBank Database can be used for extracting other information such as drug-drug interaction.
- **FDA Adverse Event Reporting System (FAERS)** [39]. The FAERS is a database that contains adverse event reports, medication error reports and product quality complaints resulting in adverse events that were submitted to FDA.
- **HPRD** Human Protein Reference Database [40] for human protein-protein interactions.
- **SIDER** Side Effect Resource Database [24] aggregates information from side effects

Further, other databases have been used to change between identifier types, e.g., KEGG Drug ID to PubChem ID, such as STITCH [41], bioMART [42], ChemBL [43].

## 2 Related work

In what follows we briefly describe the selected state-of-the-art DTI models, where the first four are transductive and the second four are inductive. Note that referring to a model as transductive or inductive concerns the feature (embedding) generation process and not the prediction task itself.

- **DTINet** [19]. DTINet considers a heterogeneous graph with four node types (drugs, proteins, side effects and diseases) and six edge types (DTIs, protein-protein interaction, drug-drug interaction, drug-disease association, protein-disease association, drug-side-effect association, plus similarity edges between drugs and proteins). After compact feature learning (based on a random walk with restart) on each drug and protein network, it calculates the best projection from one space onto another using a matrix completion method, and then infers interactions according to the proximity criterion. The matrices generated are known as “Luo Dataset”.
- **DDR** [7]. DDR uses a heterogeneous graph built from known DTIs, multiple drug-drug similarities, and several protein-protein similarities. Firstly, DDR performs a pre-processing step where a subset of similarities is selected in a heuristic process to obtain an optimized combination of similarities. Then, DDR applies a non-linear similarity fusion method to combine different similarities. Finally, from these combined similarities, a path-category-based feature extraction method is applied, and these features are fed into a random forest model.
- **DTiGEMS+** [20]. The information of the interaction within drugs and proteins coming from diverse matrices is selected and integrated to create a heterogeneous graph alongside the DTI information. Simultaneously, a second graph is created by applying *node2vec* to the DTI graph, obtaining the features for each node and augmenting the interactions based on the similarity of the calculated features. Multiple paths are extracted from both graphs and feed to a supervised ML classifier after a feature selection process.
- **DTI2Vec** [21]. DTI2Vec stems from the previous and more complex model DTiGEMS+, trying to improve the precision of the predictions while reducing the amount of side information needed. This method only uses the similarity matrices within drugs and proteins to increase the number of connections on the DTI network. The nodes of this augmented network are used as input to *node2vec*, and the resulting embeddings are combined to create a feature vector and feed a classifier.
- **NeoDTI** [22]. NeoDTI aims to automatically learn a network topology-preserving node-level embedding to facilitate DTI prediction. First, neighborhood information aggregation and node embedding update processes ensure that each node within the heterogeneous network generates a new feature representation by integrating its neighborhood information with its own features. Then, they enforce the node embeddings to preserve the network topology, aiming to reconstruct the original individual networks. Finally, from these embeddings they extract the node features and use them for the DTI prediction.
- **MolTrans** [23]. MolTrans uses unlabeled data to decompose drugs and proteins into high-quality substructures. Then it creates an augmented embedding for each using a transformer and a map of interactions, allowing it to predict which substructures contribute most to the overall interaction.
- **HyperAttentionDTI** [9]. HyperAttentionDTI embeds each character of the different sequences into vectors. Then the model makes use of an attention mechanism and convolutional neural networks (CNNs) to make DTI predictions. It models the complex non-covalent intermolecular interactions between atoms and amino acids using the attention mechanism.
- **EEG-DTI** [8]. EEG-DTI considers a heterogeneous graph using the same type of dataset as DTINet. It first generates low-dimensional embeddings for drugs and proteins with three graph convolutional networks (GCN) layers and concatenates them separately. Then, it calculates their inner product to get a protein-drug score.

Our **baseline classifier** (denoted as **N2V+NN**) is based on *node2vec* to embed the DTI network so that it solely relies on the topology of the network. From the generated embeddings, positive edges and a random subsampling of negative edges are used to train and validate a 2-layer neural network  $\Psi$ . Being  $X \in \mathbb{R}^{K \times 2d}$  the batched input matrix, and  $\mathbf{W}_1 \in \mathbb{R}^{2d \times n}$  and  $\mathbf{W}_2 \in \mathbb{R}^{n \times 1}$  the associated

weight matrices, our model  $\Psi$  will generate the output  $h \in \mathbb{R}^{K \times 1}$  as:

$$h = \sigma(\mathbf{W}_2 \cdot f(\mathbf{W}_1 \cdot X)),$$

where  $d$  is the selected *node2vec* embedding dimension for each node,  $K$  is the number of samples per batch,  $f$  is a ReLU activation function,  $\sigma$  is a sigmoid activation function and  $n$  is the number of neurons of the first layer. In order to solve the DTI classification problem, we use a loss that combines the sigmoid of the output layer and the binary cross entropy loss in a single function. This combination takes advantage of the log-sum-exp trick for numerical stability [44]. For each sample  $x_k$  in a given batch ( $k \in [1, K]$ ), the loss is given by:

$$l_k = -w_k [y_k \cdot \log h_k + (1 - y_k) \cdot \log (1 - h_k)],$$

where  $w_k$  is a manual rescaling weight,  $y_k$  is the associated label for sample  $x_k$ , and  $h_k$  is the model output for sample  $x_k$ . The final loss  $L$  is then computed as the average of  $(l_1, \dots, l_K)$ . We performed a train/validation/test (0.75, 0.15, 0.1) splitting prior performing a hyperparameter tuning, varying several architectures, loss functions, epochs and batch sizes to select the model with the highest test AUROC for every evaluated dataset (Appendix Table A8).

### 3 Evaluation setup

#### A fair evaluation scheme: graph embedding splitting approach

The following evaluation scheme consisting of constructing three different train-test splits ( $S_p$ ,  $S_d$ , and  $S_t$ ) was used:

- $S_p$  Related to pairs. Any protein or drug may appear both in the train and test set, but interactions cannot be duplicated in the two sets.
- $S_d$  Related to drug nodes. Drug nodes are not duplicated in the train and test set, i.e., a node evaluated during training does not appear in the test set.
- $S_t$  Related to targets. Protein nodes are not duplicated in the train and test set, each protein seen during training does not appear in the test set.

If the model to be compared uses three splits (train/val/test), the criterion is applied the same way as if they were just two splits (train/test), but applying an extra split to the train fold, yielding train/val/test folds. Hence, train and validation will be evaluated together when verifying  $S_p$ ,  $S_d$  and  $S_t$  splits.

Note that most assessed models have not been evaluated on these splitting criteria, but perform a *traditional* split. This consists of a random splitting of the DTI network without constraining the DTI distribution, which may lead to repetition of drug or protein nodes across folds. As the  $S_p$ ,  $S_d$  and  $S_t$  splits impose certain constraints not assumed by the authors and may result in lower performance than what was initially reported, we also provide, for each model, the results following the originally proposed evaluation benchmark.

Furthermore, the  $S_c$  split, related to a couple of different DTI networks that do not have common drugs nor proteins [9, 27], involves training a model initially on one dataset and then testing the trained model on another dataset. This split can assist in assessing the methods' generalization capabilities, potentially revealing data leakage concerns. However, the limited reproducibility of most methods have complicated the application of this evaluation scheme to the evaluated ones. Nonetheless, we validated our hypothesis regarding *node2vec*-based methods by applying this split to our baseline DTI classifier (see Appendix Note 2).

#### Building train and test splits for N2V evaluation

In assessing the generalization capability of *node2vec*-based drug repurposing models across multiple DTI networks, we evaluated the designed baseline model using train/test splits. First, node embeddings for each network were constructed individually using *node2vec*. Next, for each network, a balanced dataset was created by selecting all positive pairs and randomly pairing them with negatives in a 1:1 ratio. Finally, the baseline model was trained on embeddings from one dataset and tested on a different one, yielding both train and test AUROC and AUPRC values. When the same network is used for both training and testing (as shown in Fig. 3-B matrix's diagonal), the dataset was constructed as previously described, with a 70/30 train-test split.

## Considering a biological driven criteria for negative subsampling

Here now we describe the process of DTI stratification and hard-to-classify pairs selection. First, for each known DTI interaction (labeled as 1), we compute the RMSD between the selected protein and every other protein available in the data set. Then, in order to generate a balanced dataset, for each positive DTI, we select a protein to form a negative interaction, based on the computed RMSD between the known target and every other protein in the network. The selection is made by sorting the proteins' RMSD, and we will select or discard them based on three different windows. The first window ranges from 0 to 2.5 Å, where proteins in this interval are discarded, as in this range we may include small structures or very simple proteins that align non-specifically to others. Proteins lying on the second window, from 2.5 to 5 Å, are held out for validation, as they are very similar to the actual target but are labelled as 0, so they can generate false positives, potentially hinder the model's training. The third window, ranged from 5 to  $t$  Å ( $t \in [6, 20]$ ), will be the selected one as the train split, as is populated with proteins that are closely enough to the target to be a difficult train event, but different enough to assure the potentially true negativity of the data.

## Tanimoto Similarity

The pairwise drug similarity, calculated with Tanimoto metric was calculated in RDKit [45] creating fingerprints in default configuration using *RDKFingerprint* (with 2048 bits) function.

## Protein Structures and RMSD Calculation

Protein structures were obtained from the PDB Database [46] and Alpha Fold [47, 48], considering X-Ray structures with resolution lower than 2 Å and a per-residue confidence score higher than 70 on average, respectively. The RMSD was calculated using an adapted script from PyMOL [49], considering superimposition mode and all objects aligned using the alpha carbons (C-alpha) of the backbone of the two proteins and the default configuration of 5 cycles. See Appendix Figs. A1, A2, A5, A6, for distribution of pairwise RMSD in all datasets.

## Hardware

All simulations were performed on a workstation with 64 cores Intel xeon gold 6130 2.1Ghz and 754Gb of RAM. A Quadro RTX 4000 GPU was also used, with driver version 460.67 and cuda 11.2. version.

## 4 *In vitro* validation

### Cell lines

Human mut KRAS (H1792) LUAD cell line and human mut KRAS (HPAF II) PDAC cell line were used. All these cell lines were obtained from American Type Culture Collection (ATCC) and authenticated by the Genomics Unit at CIMA using Short Tandem Repeat profiling (AmpFLSTR Identifiler Plus PCR Amplification Kit). Human cells were grown according to ATCC specifications.

### Reagents

Carbinoxamine maleate (PHR2802) was purchased from Merck and 6-(Cyclohexylmethoxy)-8-isopropyl-9H-purin-2-amine was synthesized and obtained from Wuxi.

### Western blotting

Western blot methodology was performed as previously published [50]. For these experiments cells were treated with DMSO (vehicle, control conditions) or drug. In this last case, we used a final concentration of 30uM for Carbinoxamine and 10uM for 6-(Cyclohexylmethoxy)-8-isopropyl-9H-purin-2-amine. Antibodies used: GAPDH (1:5,000, ab9484, Abcam), ERK1/2 (1:1,000, #9102, Cell Signaling Technology), p-ERK1/2 (1:1,000, #9101, Cell Signaling Technology), p70S6K (1:1,000, #2708, Cell Signaling Technology), p-p70S6K (1:1,000, #9205, Cell Signaling Technology), EGFR (1:1,000, #2232, Cell Signaling Technology), p-EGFR (1:1,000, #2236, Cell Signalling Technology),

GSK3 $\beta$  (1:1,000, ab31826, Abcam), p-GSK3 $\beta$  (1:1,000, #9336, Cell Signalling Technology) and p-4E-BP1 (1:1,000, #9451, Cell Signaling Technology).

### **Drug studies in vitro**

To determine the number of viable cells in proliferation and the potential cytotoxicity of drugs in cell lines, cells were seeded in triplicate into 96-well plates (range: 500 to 1,800 cells per well depending on the cell line). The next day, cells were cultured in the absence or presence of rising concentrations of single drugs (Carbinoxamine 0-30 $\mu$ M; 6-(Cyclohexylmethoxy)-8-isopropyl-9H-purin-2-amine 0-10 $\mu$ M) during 3 or 5 days. At these time points, remaining cells were fixed with 4% formaldehyde (Panreac) for 15 minutes at RT, stained with crystal violet solution (Sigma-Aldrich) (1% crystal violet in H<sub>2</sub>O) for 15 minutes and photographed using a digital scanner (EPSON Perfection v850 Pro). Relative growth was quantified by measuring absorbance at 570 nm in a spectrophotometer (SPECTROstar Nano – BMG Labtech) after extracting crystal violet from the stained cells using 20% of acetic acid (Sigma).

### **Protein and Drug Annotation**

Proteins were annotated using Molecular Function Keywords from Uniprot [51] and drugs with Classyfire [52]. Annotated heatmaps were generated to check whether proteins cluster per molecular function and drugs by chemical classification.

### **Code availability**

Dockers for all evaluated models are available for in DockerHub. The repository containing all the developed tools and code, along with the *GraphGuest* Python Package are available at <https://github.com/ubioinformat/GraphEmb>.

### **Acknowledgments**

We would like to thanks to Oier Azurmendi for his support along the development of the work.

### **Funding**

This work was supported by the following grants: DoD of the US - CDMR Programs [W81XWH-20-1-0262], Ramon y Cajal contracts [MCIN/AEI RYC2021-033127-I] [RYC2019-028578-I], DeepCTC [MCIN/AEI TED2021-131300B-I00], Gipuzkoa Fellows [2022-FELL-000003-01], the Spanish MCIN (PID2021-126718OA-I00), Fulbright Predoctoral Research Program [PS00342367], and FEDER/MCIN - AEI (PID2020-116344-RB-100/MCIN/AEI/10.13039/501100011033).

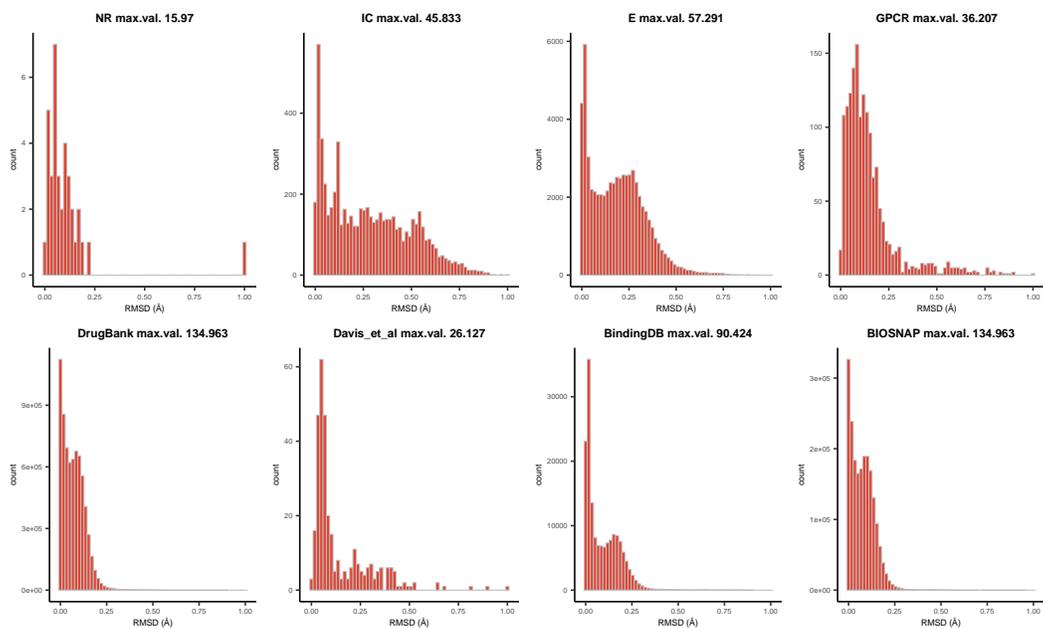
### **Competing interests**

The authors declared no competing interests.

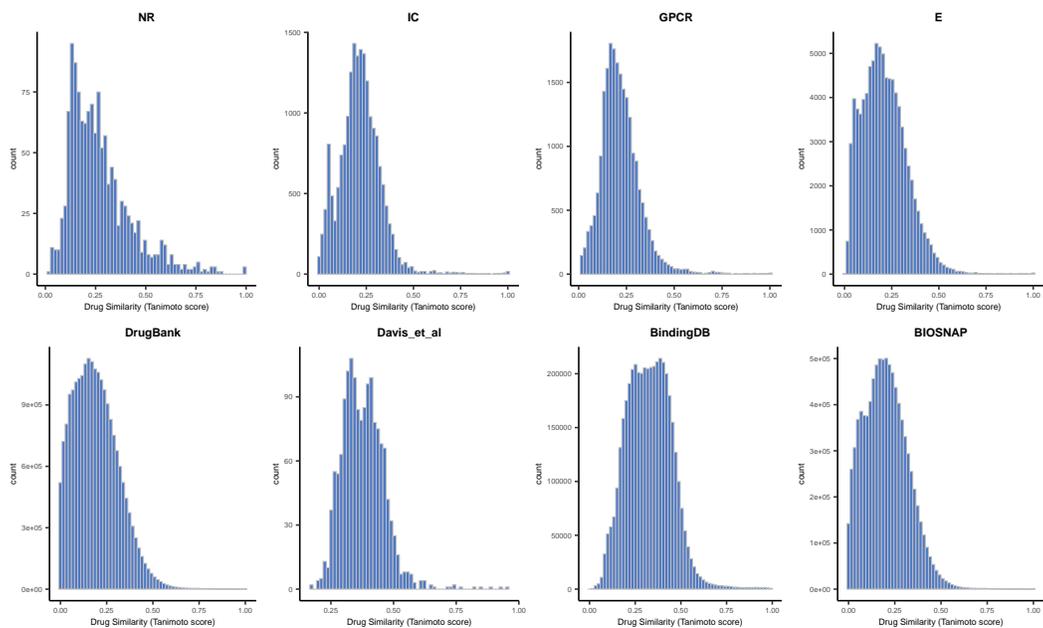
### **Author contributions**

Conceptualization: J.F., G.S., U.V., I.O., O.G., and M.H.; methodology: J.F., G.S., U.V., I.O., and M.H.; software: J.F., G.S., U.V. and M.C.; formal analysis: J.F., G.S., U.V. and M.C.; investigation: J.F., G.S., U.V.; validation: J.F., G.S., U.V., M.C and S.V.; supervision: I.O., S.V., O.G., and M.H.; writing-original draft: J.F., G.S., U.V., I.O., S.V. O.G., and M.H.; visualization: J.F., G.S., U.V.; writing-review & editing: all authors.

## Appendix Figures

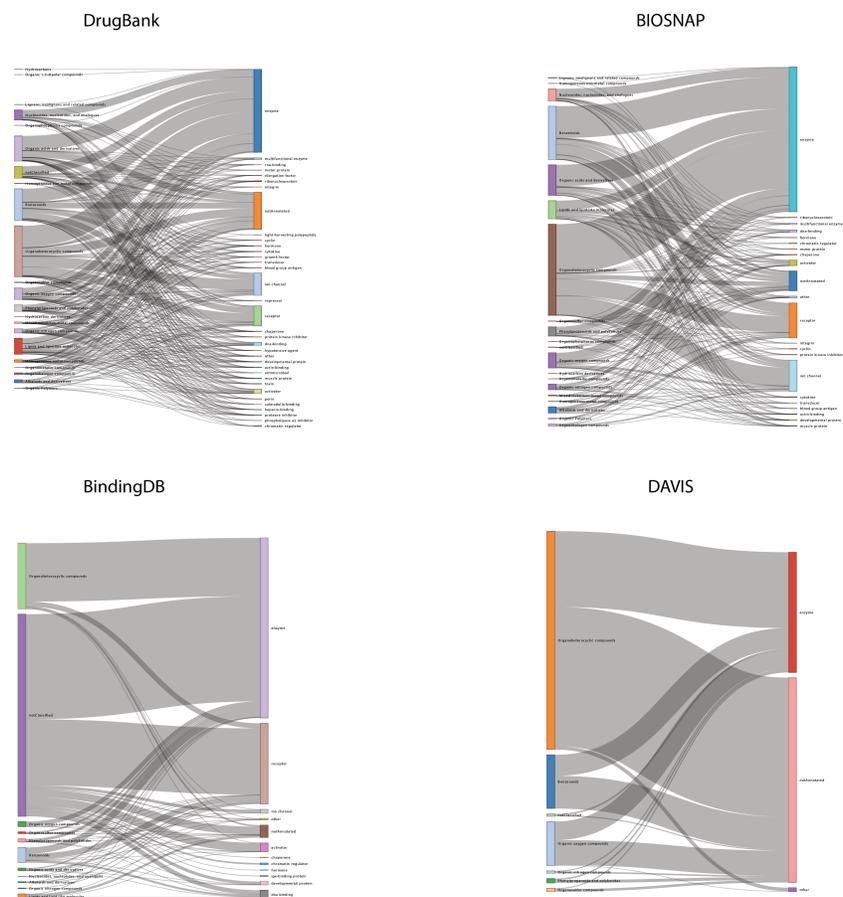


Appendix Figure A1: Histograms of the distribution of the pairwise Root Mean Square Deviation (RMSD) of all proteins for each dataset, regardless of whether they come from PDB or AlphaFold.

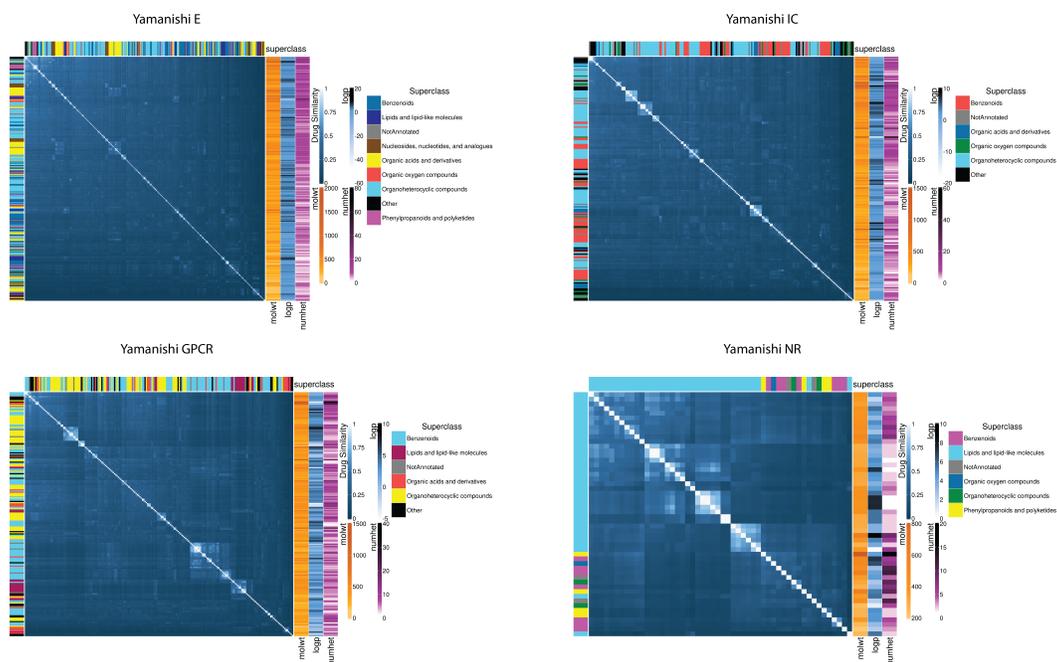


Appendix Figure A2: Histograms of the distribution of Tanimoto score calculated pairwise over all drugs for each dataset. The distribution of chemical similitude shift to the left indicates that drugs are chemically diverse

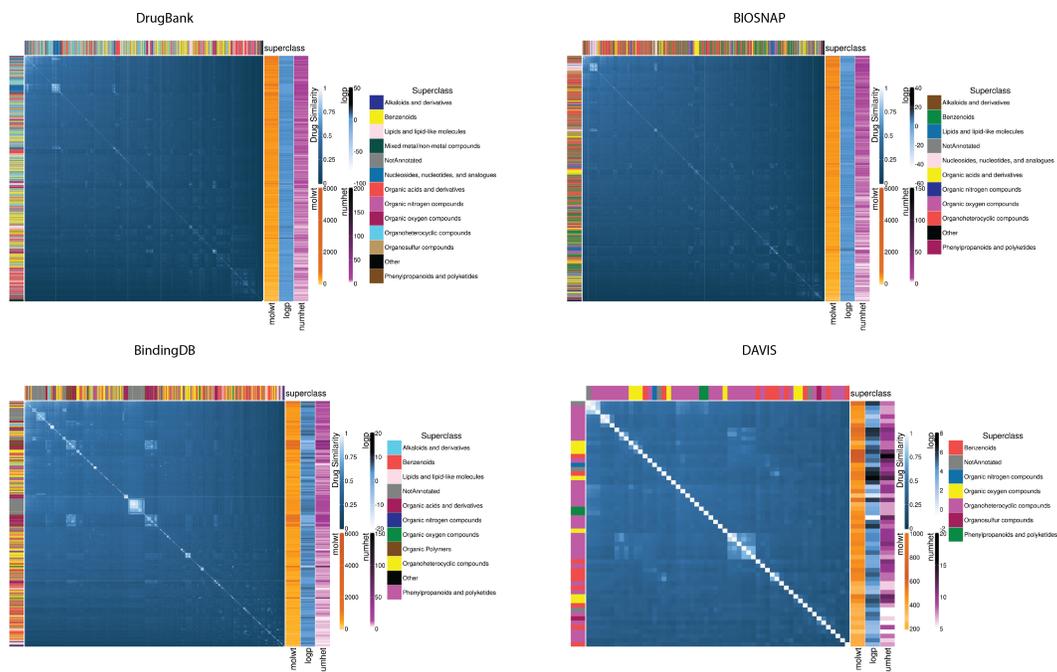




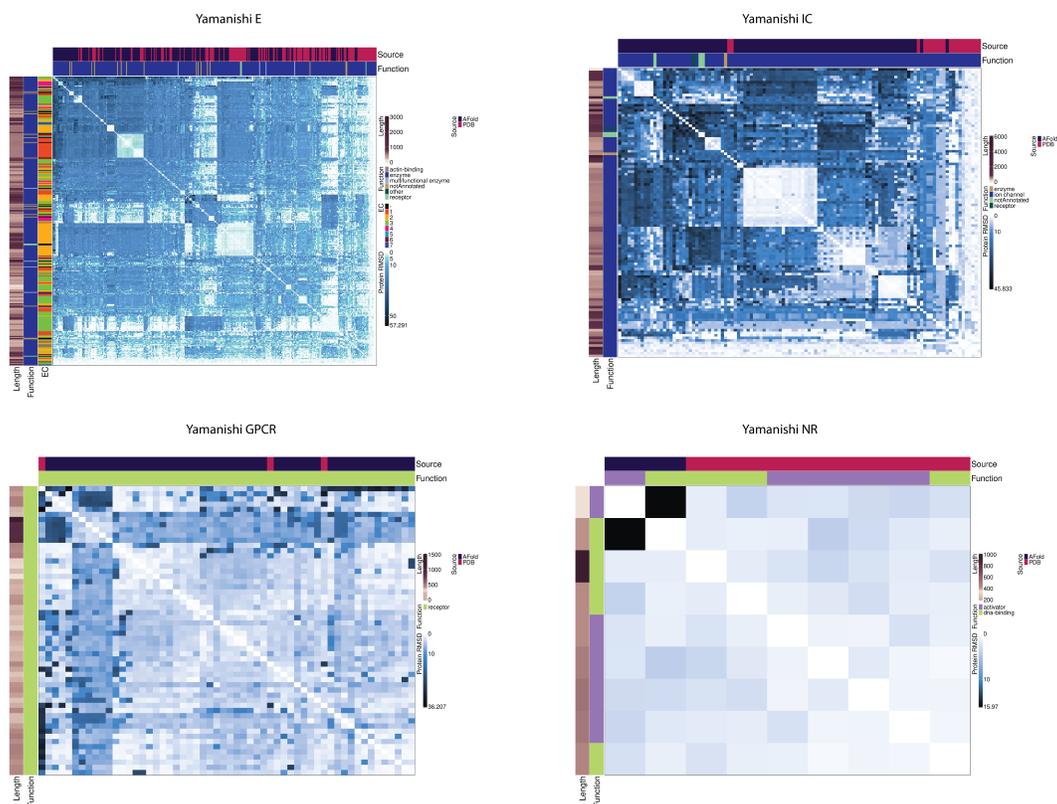
Appendix Figure A4: Sankey plots of DrugBank, BIOSNAP, BindingDB and DAVIS datasets. These Sankey plots connect the drug chemical category with the protein family, remarking that each connection is an existing DTI in the dataset.



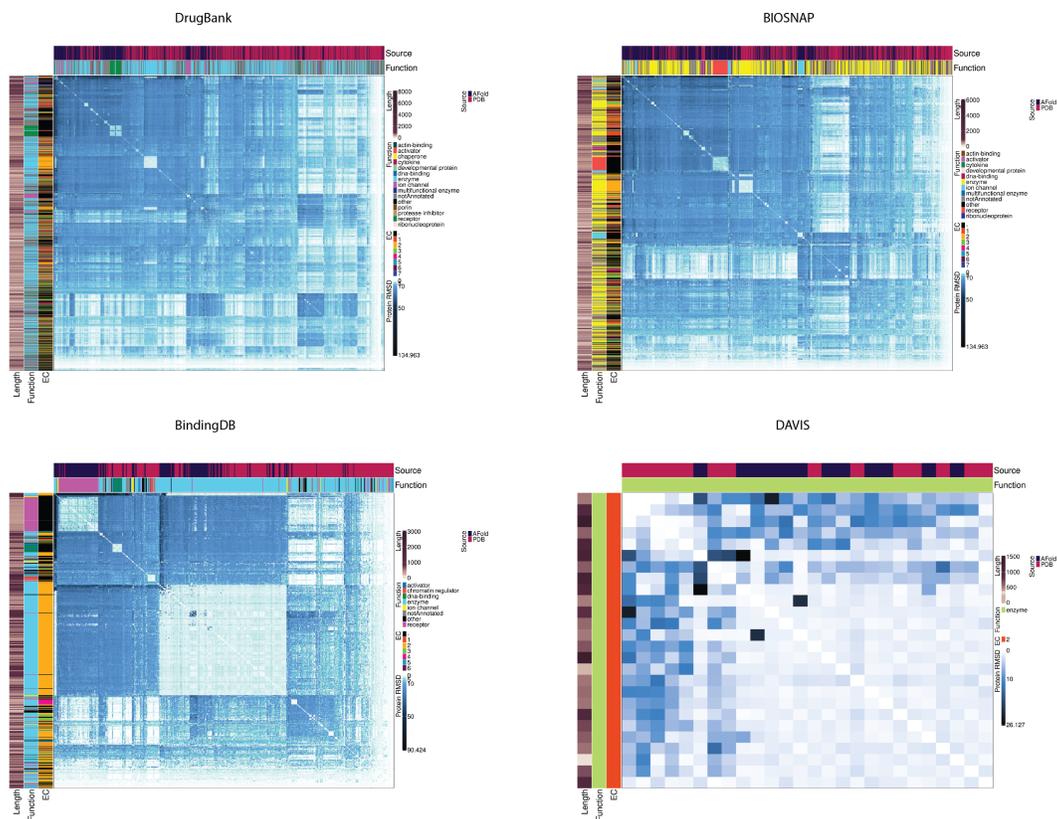
Appendix Figure A5: Heatmaps of Tanimoto score for drugs of Yamanishi datasets calculated with euclidean distances. It is appreciated how some drugs form small clusters. The annotation represents the chemical classification of the drug (by superclass in Classyfire), and three different molecular descriptors the molecular weight (molwt), the logP, and the number of heteroatoms (numhet).



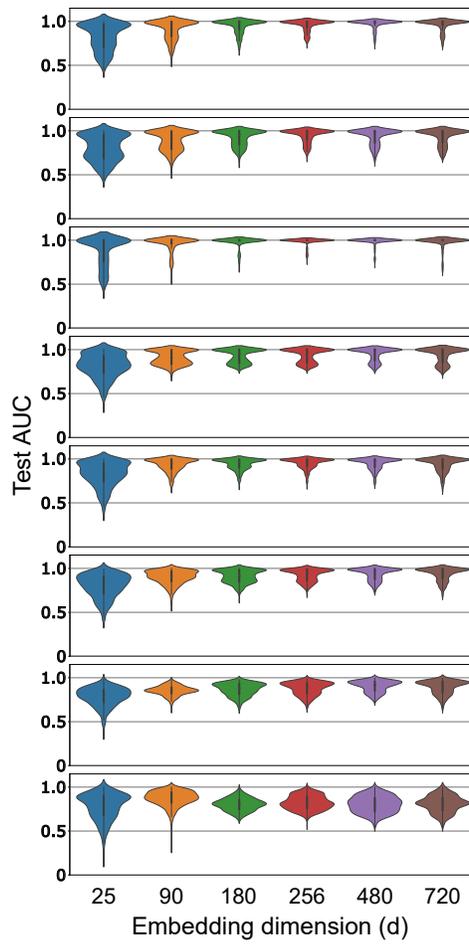
Appendix Figure A6: Heatmaps of Tanimoto score for drugs of DrugBank, BIOSNAP, BindingDB and Davis datasets calculated with euclidean distances. It is appreciated how some drugs form small clusters. The annotation represents the chemical classification of the drug (by superclass in Classyfire), and three different molecular descriptors the molecular weight (molwt), the logP, and the number of heteroatoms (numhet).



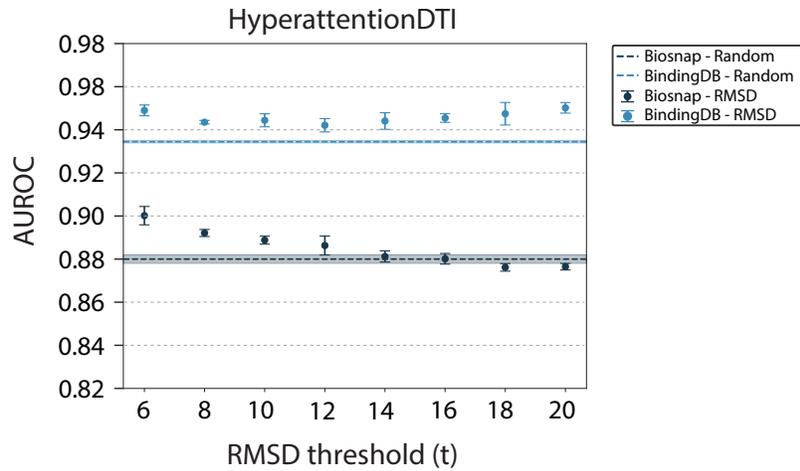
Appendix Figure A7: Heatmaps of RMSD score for proteins of Yamanishi datasets calculated with Euclidean distances. The annotation represents the source of the protein, i.e., from where we downloaded the structure (PDB/AlphaFold), the molecular function of the protein, and the enzyme classification for the E dataset. Further, we added the sequence length in the left annotation in brown. In the heatmaps, proteins form several clusters. In E, this occurs especially for EC-1 (oxidoreductases) and EC-2 (transferases), further, smaller clusters of EC-3 (hydrolases) and EC-4 (lyases) also appear. Further, while in E and GPCR we found mixed the protein source, it does not happen for IC and NR.



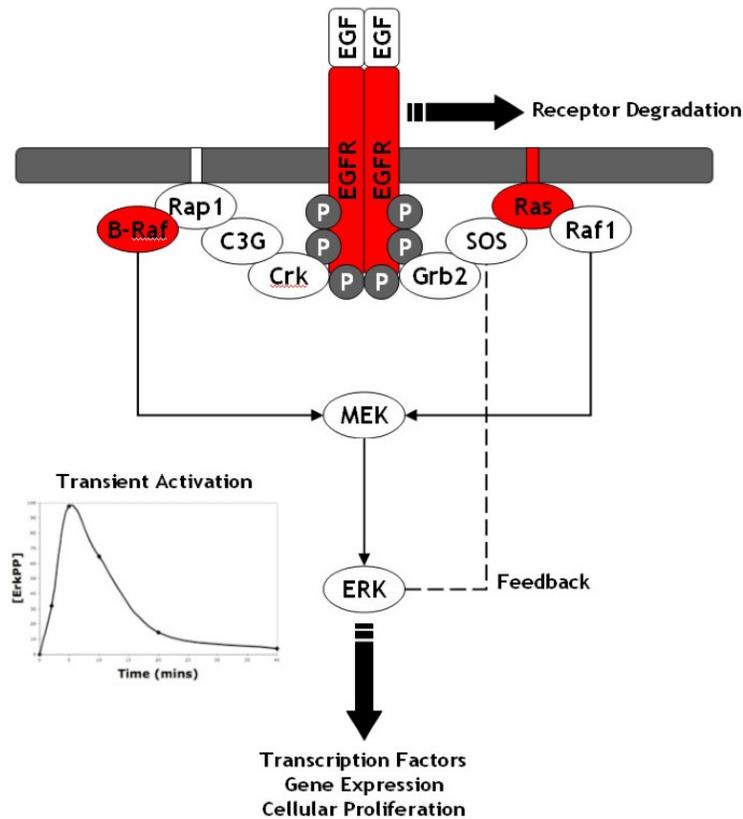
Appendix Figure A8: Heatmaps of RMSD score for proteins of DrugBank, BIOSNAP, BindingDB and DAVIS datasets calculated with Euclidean distances. The annotation represents the source of the protein, i.e., from where we downloaded the structure (PDB/AlphaFold), the molecular function of the protein, and the enzyme classification for the E dataset. Further, we added the sequence length in the left annotation in brown. For the four datasets, we found a mixed distribution over clusterings of proteins independently of the source. Further, we shown how proteins clustering by protein family, which can be specially appreciated in BindingDB for receptors and enzymes, with brenda classification EC-2.



Appendix Figure A9: Test AUC distribution across evaluated datasets, grouped by node2vec embedding dimension



Appendix Figure A10: AUROCs of the test split for the RMSD-based (threshold from 6 to 20 Å) and random subsampling techniques when using HyperAttentionDTI model, on Biosnap and BindingDB datasets.



Appendix Figure A11: EGF-ERK-pathway: This schematic illustrates the EGF-activated ERK pathway, starting with EGF binding to EGFR and concluding with ERK activation. Activated ERK has multiple targets in the cytoplasm and nucleus, including numerous transcription factors, therefore directly affecting gene expression and influence cellular growth. Source: <https://doi.org/10.1186/1752-0509-3-100>



## Appendix Tables

Appendix Table A1: **State-of-the-art DTI networks statistics**. The density of a graph represents the proportion of edges present in the graph to the total number of edges that could possibly exist in the graph. The number of connected components represents the count of isolated subgraphs within the network.

	DrugBank	BIOSNAP	BindingDB	DAVIS	Yamanishi E	Yamanishi IC	Yamanishi GPCR	Yamanishi NR
Number of drugs	8042	5017	3085	65	445	210	223	54
Number of proteins	5141	2324	719	314	664	204	95	26
Total number of nodes	13183	7341	3804	379	1109	414	318	80
Total number of edges	27861	15138	5938	1048	2926	1476	635	90
Density (%)	0.03	0.06	0.08	1.46	0.48	1.73	1.26	2.85
# of connected components	490	205	232	1	44	3	19	10

Appendix Table A2: **Comparison across DTI prediction methods**. Evaluation considers various criteria, such as the language used, the splitting or the utilized validation technique. Upper and lower groups illustrate transductive and inductive methods, respectively.

Method	Language	Complete Code	Issue-free Code	GitHub Medal	Default Splitting	Data Availability	Time Consumption	Validation (type)
DTiNet	Matlab	✓	✓	×	Unconstrained 10 folds CV	×	Very Fast	Experimental, In-sillico
DDR	Python	×	✓	🏅	Sp, Sd, St	×	Slow	Bibliographic
DTi-GEMS	Python	×	×	×	Unconstrained 10 times TVT	×	Fast	In-sillico
DTi2Vec	Python	×	×	×	Unconstrained 10 times TVT	×	Slow	Bibliographic
-----								
NeoDTI	Python	×	×	×	Unconstrained 10 times TVT	✓	Slow	Bibliographic
Moltrans	Python	×	✓	×	Unconstrained 5 times TVT	×	Fast	×
HyperAttentionDTI	Python	×	×	×	Sp, Sd, St, Sh	×	Fast	×
EEG-DTI	Python	×	×	×	Unconstrained 10 folds CV	✓	Fast	×

TVT: train-validation-test, CV: Cross-Validation

Appendix Table A3: Default Splits AUPRC benchmarking for the evaluated DTI prediction models.

Method/Dataset	DrugBank	BIOSNAP	BindingDB	Davis et al	E	IC	GPCR	NR
DDR	OOT	OOT	OOT	0.679 ± 0.049	0.911 ± 0.010	0.934 ± 0.013	0.818 ± 0.036	0.785 ± 0.040
DTi2Vec	OOD	OOD	0.920 ± 0.005	0.390 ± 0.060	0.970 ± 0.012	0.960 ± 0.009	0.830 ± 0.058	0.750 ± 0.072
DTi-GEMS	OOD	OOD	0.323 ± 0.451	0.386 ± 0.210	0.820 ± 0.278	0.873 ± 0.092	0.731 ± 0.150	0.626 ± 0.170
DTiNet	0.857 ± 0.001	0.886 ± 0.001	0.873 ± 0.006	0.812 ± 0.010	0.918 ± 0.004	0.764 ± 0.010	0.803 ± 0.006	0.730 ± 0.022
-----								
NeoDTI	OOT	OOT	OOT	OOT	OOT	0.917 ± 0.005	0.752 ± 0.045	0.450 ± 0.145
Moltrans	0.662 ± 0.005	0.645 ± 0.004	0.803 ± 0.005	0.530 ± 0.039	0.800 ± 0.004	0.772 ± 0.007	0.530 ± 0.053	0.421 ± 0.018
HyperAttentionDTI	0.776 ± 0.046	0.772 ± 0.063	0.910 ± 0.051	0.576 ± 0.019	0.922 ± 0.062	0.917 ± 0.082	0.648 ± 0.057	0.326 ± 0.073
EEG-DTI	0.880 ± 0.007	0.905 ± 0.010	0.890 ± 0.030	0.700 ± 0.036	0.955 ± 0.009	0.948 ± 0.015	0.882 ± 0.047	0.801 ± 0.061

OOT and OOR mean *out of time* and *out of RAM*, respectively.

Appendix Table A4: Analyzed methodologies time consumption.

Model/Dataset	DrugBank-DTI	BioSNAP	BindingDB	Davis et al	E	IC	GPCR	NR
DDR	OOT	OOT	OOT	1h 33m 45s	5d 18h	3h 36m 21s	1h 16m 10s	9m 17s
DTI2Vec	OOR	OOR	11h 50m	3h 17m 42s	3d 17h 18m	2h 36m	1h 50m 49s	1m 39s
DTI-GEMS	OOR	OOR	5h 25min	25m	22m	2min	1m 13s	20s
DTINet	39m 57s	21m 8s	3m 27s	1m 4s	6m 5s	2m 11s	1m 11s	8s
NeoDTI	OOT	OOT	OOT	OOT	OOT	1d 19h	1d 19h	1d 15h
Moltrans	6h 39m 18s	3h 36m 46s	1h 13m 47s	17m 49s	49m 14s	24m 13s	10m 42s	1m 34s
HyperAttentionDTI	10h 10m 27s	5h 32m 50s	2h 31m 37s	31m 27s	1h 18m 29s	45m 59s	18m 17s	2m 21s
EEG-DTI	3d 13h 6m 45s	3h 26m 57s	50m 21s	54m 39s	4h 30m 22s	10m 16s	328m 30s	7m 47s

OOT and OOR mean *out of time* and *out of ram*, respectively.

Appendix Table A5: Mean and standard deviation of the AUPRC values for the  $S_p$  split, for each evaluated model and dataset.

Model/Dataset	Yamanishi-NR	Yamanishi-IC	Yamanishi-GPCR	Yamanishi-E	DrugBank	DAVIS	BindingDB	BioSNAP
DDR	0.785 ± 0.040	0.934 ± 0.013	0.818 ± 0.035	0.910 ± 0.010	OOT	0.679 ± 0.049	OOT	OOT
DTI2Vec	NEET	0.620 ± 0.130	0.760 ± 0.128	0.580 ± 0.171	OOR	0.58 ± 0.098	0.999 ± 0.112	0.57 ± 0.132
DTI-GEMS	0.977 ± 0.249	0.953 ± 0.254	0.953 ± 0.157	0.996 ± 0.128	OOR	OOR	OOR	OOR
DTINet	0.710 ± 0.0133	0.767 ± 0.003	0.767 ± 0.009	0.914 ± 0.002	0.830 ± 0.001	0.689 ± 0.014	0.800 ± 0.002	0.857 ± 0.001
NeoDTI	0.470 ± 0.190	0.893 ± 0.018	0.748 ± 0.032	OOT	OOT	OOT	OOT	OOT
Moltrans	0.4206 ± 0.014	0.771 ± 0.007	0.530 ± 0.053	0.800 ± 0.004	0.620 ± 0.005	0.529 ± 0.039	0.803 ± 0.005	0.645 ± 0.004
HyperAttentionDTI	0.326 ± 0.073	0.917 ± 0.082	0.648 ± 0.057	0.922 ± 0.062	0.776 ± 0.046	0.576 ± 0.019	0.910 ± 0.051	0.772 ± 0.063
EEG-DTI	0.732 ± 0.044	0.962 ± 0.007	0.887 ± 0.006	0.962 ± 0.002	0.925 ± 0.001	0.747 ± 0.015	0.859 ± 0.003	0.928 ± 0.002

NEET, OOR and OOT stands for *Not Enough Edges to Train*, *Out of RAM* and *Out of Time*, respectively.

Appendix Table A6: Mean and standard deviation of the AUPRC values for the  $S_d$  split, for each evaluated model and dataset.

Model/Dataset	Yamanishi-NR	Yamanishi-IC	Yamanishi-GPCR	Yamanishi-E	DrugBank	DAVIS	BindingDB	BioSNAP
DDR	0.642 ± 0.220	0.660 ± 0.095	0.608 ± 0.090	0.699 ± 0.650	OOT	0.462 ± 0.014	OOT	OOT
DTI2Vec	0.73 ± 0.179	0.620 ± 0.012	0.785 ± 0.210	0.59 ± 0.132	OOR	0.630 ± 0.100	OOR	0.640 ± 0.140
DTI-GEMS	NEET	NEET	0.796 ± 0.280	0.980 ± 0.120	OOR	OOR	OOR	OOR
DTINet	0.635 ± 0.011	0.664 ± 0.005	0.741 ± 0.011	0.844 ± 0.004	0.813 ± 0.001	0.556 ± 0.012	0.762 ± 0.008	0.841 ± 0.001
NeoDTI	0.118 ± 0.029	0.213 ± 0.060	0.271 ± 0.024	OOT	OOT	OOT	OOT	OOT
Moltrans	0.421 ± 0.038	0.437 ± 0.056	0.474 ± 0.071	0.411 ± 0.043	0.574 ± 0.025	0.333 ± 0.038	0.653 ± 0.044	0.530 ± 0.010
HyperAttentionDTI	0.343 ± 0.020	0.638 ± 0.016	0.591 ± 0.036	0.611 ± 0.012	0.684 ± 0.005	0.395 ± 0.014	0.859 ± 0.007	0.680 ± 0.004
EEG-DTI	0.747 ± 0.020	0.764 ± 0.017	0.832 ± 0.024	0.781 ± 0.020	0.892 ± 0.001	0.660 ± 0.033	0.779 ± 0.016	0.902 ± 0.002

NEET, OOR and OOT stands for *Not Enough Edges to Train*, *Out of RAM* and *Out of Time*, respectively.

Appendix Table A7: Mean and standard deviation of the AUPRC values for the  $S_t$  split, for each evaluated model and dataset.

Model/Dataset	Yamanishi-NR	Yamanishi-IC	Yamanishi-GPCR	Yamanishi-E	DrugBank	DAVIS	BindingDB	BioSNAP
DDR	0.622 ± 0.300	0.786 ± 0.031	0.658 ± 0.103	795 ± 0.027	OOT	0.618 ± 0.44	OOT	OOT
DTI2Vec	0.740 ± 0.117	0.580 ± 0.017	0.660 ± 0.170	0.690 ± 0.127	OOR	0.540 ± 0.130	OOR	0.440 ± 0.320
DTI-GEMS	NEET	NEET	0.731 ± 0.297	0.976 ± 0.190	OOR	OOR	OOR	OOR
DTINet	0.493 ± 0.017	0.856 ± 0.008	0.623 ± 0.006	0.860 ± 0.002	0.796 ± 0.001	0.649 ± 0.008	0.660 ± 0.010	0.815 ± 0.002
NeoDTI	0.111 ± 0.051	0.362 ± 0.12	0.163 ± 0.032	OOT	OOT	OOT	OOT	OOT
Moltrans	0.266 ± 0.069	0.551 ± 0.064	0.488 ± 0.036	0.604 ± 0.038	0.660 ± 0.020	0.477 ± 0.018	0.561 ± 0.300	0.681 ± 0.049
HyperAttentionDTI	0.294 ± 0.038	0.827 ± 0.012	0.565 ± 0.023	0.795 ± 0.022	0.775 ± 0.002	0.548 ± 0.023	0.617 ± 0.014	0.771 ± 0.002
EEG-DTI	0.501 ± 0.048	0.690 ± 0.017	0.592 ± 0.011	0.604 ± 0.015	0.743 ± 0.001	0.512 ± 0.002	0.606 ± 0.012	0.721 ± 0.008

NEET, OOR and OOT stands for *Not Enough Edges to Train*, *Out of RAM* and *Out of Time*, respectively.

Appendix Table A8: Grid-search parameters used for the baseline (N2V + NN) model. For the grid-search, they were used 6 different embedding space dimensions, 4 type of architectures, 4 different number of epochs, *BinaryCrossEntropy* or *FocalDistance* as the loss function, and 1/16 and 1/64 of the dataset size as the batch size.

Parameter	Value
Embedding dimension	d25, d90, d180, d256, d480, d720
Architecture	Type 1, Type 2, Type 3, Type 4
Epochs	2, 5, 10, 50
Loss Function	BCE, Focal
Batch Size	1/16, 1/64

## References

- [1] Simon G. Patching. Surface plasmon resonance spectroscopy for characterisation of membrane protein–ligand interactions and its potential for drug discovery. *Biochimica et Biophysica Acta (BBA) - Biomembranes*, 1838(1, Part A):43–55, 2014. Structural and biophysical characterisation of membrane protein–ligand binding.
- [2] Suzanne B. Shuker, Philip J. Hajduk, Robert P. Meadows, and Stephen W. Fesik. Discovering high-affinity ligands for proteins: Sar by nmr. *Science*, 274(5292):1531–1534, 1996.
- [3] Joseph A. DiMasi, Henry G. Grabowski, and Ronald W. Hansen. Innovation in the pharmaceutical industry: New estimates of r&d costs. *Journal of Health Economics*, 47:20–33, 2016.
- [4] Michael M Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.
- [5] Mingyue Zheng, Xian Liu, Yuan Xu, Honglin Li, Cheng Luo, and Hualiang Jiang. Computational methods for drug design and discovery: focus on china. *Trends Pharmacol Sci*, 34(10):549–559, Oct 2013.
- [6] Antonio Lavecchia and Carmen Cerchia. In silico methods to address polypharmacology: current status, applications and future perspectives. *Drug Discov Today*, 21(2):288–298, Feb 2016.
- [7] Rawan S Olayan, Haitham Ashoor, and Vladimir B Bajic. DDR: efficient computational method to predict drug–target interactions using graph mining and machine learning approaches. *Bioinformatics*, 34(7):1164–1173, 11 2017.
- [8] Jiajie Peng, Yuxian Wang, Jiaojiao Guan, Jingyi Li, Ruijiang Han, Jianye Hao, Zhongyu Wei, and Xuequn Shang. An end-to-end heterogeneous graph representation learning-based framework for drug–target interaction prediction. *Briefings in Bioinformatics*, 22(5), 2021.
- [9] Qichang Zhao, Haochen Zhao, Kai Zheng, and Jianxin Wang. HyperAttentionDTI: improving drug–protein interaction prediction by sequence-based deep learning with attention mechanism. *Bioinformatics*, 38(3):655–662, 10 2021.
- [10] Yoshihiro Yamanishi, Michihiro Araki, Alex Gutteridge, Wataru Honda, and Minoru Kanehisa. Prediction of drug–target interaction networks from the integration of chemical and genomic spaces. *Bioinformatics*, 24(13):i232–i240, 2008.
- [11] David S Wishart, Craig Knox, An Chi Guo, Savita Shrivastava, Murtaza Hassanali, Paul Stothard, Zhan Chang, and Jennifer Woolsey. Drugbank: a comprehensive resource for in silico drug discovery and exploration. *Nucleic acids research*, 34(suppl\_1):D668–D672, 2006.
- [12] Sagar Maheshwari Marinka Zitnik, Rok Sosič and Jure Leskovec. BioSNAP Datasets: Stanford biomedical network dataset collection. <http://snap.stanford.edu/biodata>, 2018.
- [13] Mindy I Davis, Jeremy P Hunt, Sanna Herrgard, Pietro Ciceri, Lisa M Wodicka, Gabriel Pallares, Michael Hocker, Daniel K Treiber, and Patrick P Zarrinkar. Comprehensive analysis of kinase inhibitor selectivity. *Nature biotechnology*, 29(11):1046–1051, 2011.
- [14] Tiqing Liu, Yuhmei Lin, Xin Wen, Robert N Jorissen, and Michael K Gilson. Bindingdb: a web-accessible database of experimentally determined protein–ligand binding affinities. *Nucleic acids research*, 35(suppl\_1):D198–D201, 2007.
- [15] Timothy B. Dunn, Gustavo M. Seabra, Taewon David Kim, K. Eurídice Juárez-Mercado, Chenglong Li, José L. Medina-Franco, and Ramón Alain Miranda-Quintana. Diversity and chemical library networks of large data sets. *Journal of Chemical Information and Modeling*, 62(9):2186–2201, 2022.
- [16] Simon K. Mencher and Long G. Wang. Promiscuous drugs compared to selective drugs (promiscuity can be a virtue). *BMC Clinical Pharmacology*, 5(1):3, 2005.

- [17] Ming Wen, Zhimin Zhang, Shaoyu Niu, Haozhi Sha, Ruihan Yang, Yonghuan Yun, and Hongmei Lu. Deep-learning-based drug–target interaction prediction. *Journal of proteome research*, 16(4):1401–1409, 2017.
- [18] Hakime Öztürk, Arzucan Özgür, and Elif Ozkirimli. Deepdta: deep drug–target binding affinity prediction. *Bioinformatics*, 34(17):i821–i829, 2018.
- [19] Yunan Luo, Xinbin Zhao, Jingtian Zhou, Jinglin Yang, Yanqing Zhang, Wenhua Kuang, Jian Peng, Ligong Chen, and Jianyang Zeng. A network integration approach for drug–target interaction prediction and computational drug repositioning from heterogeneous information. *Nat Commun*, 8(1):573, 2017.
- [20] Maha A. Thafar, Rawan S. Olayan, Haitham Ashoor, Somayah Albaradei, Vladimir B. Bajic, Xin Gao, Takashi Gojobori, and Magbubah Essack. Dtigems+: drug–target interaction prediction using graph embedding, graph mining, and similarity-based techniques. *Journal of Cheminformatics*, 12(1):44, 2020.
- [21] Maha A. Thafar, Rawan S. Olayan, Somayah Albaradei, Vladimir B. Bajic, Takashi Gojobori, Magbubah Essack, and Xin Gao. Dti2vec: Drug–target interaction prediction using network embedding and ensemble learning. *Journal of Cheminformatics*, 13(1):71, 2021.
- [22] Fangping Wan, Lixiang Hong, An Xiao, Tao Jiang, and Jianyang Zeng. NeoDTI: neural integration of neighbor information from a heterogeneous network for discovering new drug–target interactions. *Bioinformatics*, 35(1):104–111, 07 2018.
- [23] Kexin Huang, Cao Xiao, Lucas M Glass, and Jimeng Sun. MolTrans: Molecular Interaction Transformer for drug–target interaction prediction. *Bioinformatics*, 37(6):830–836, 10 2020.
- [24] Michael Kuhn, Ivica Letunic, Lars Juhl Jensen, and Peer Bork. The sider database of drugs and side effects. *Nucleic Acids Res*, 44, 2016.
- [25] Allan Peter Davis, Cynthia J Grondin, Robin J Johnson, Daniela Sciaky, Jolene Wieggers, Thomas C Wieggers, and Carolyn J Mattingly. Comparative toxicogenomics database (ctd): update 2021. *Nucleic Acids Research*, 49(D1):D1138–D1143, 2020.
- [26] Yungki Park and Edward M Marcotte. Flaws in evaluation schemes for pair-input computational predictions. *Nature methods*, 9(12):1134–1136, 2012.
- [27] Tapio Pahikkala, Antti Airola, Sami Pietilä, Sushil Shakyawar, Agnieszka Sz wajda, Jing Tang, and Tero Aittokallio. Toward more realistic drug–target interaction predictions. *Briefings in Bioinformatics*, 16(2):325–337, 04 2014.
- [28] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864, 2016.
- [29] Jessica Siltberg-Liberles, Johan A Grahnén, and David A Liberles. The evolution of protein structures and structural ensembles under functional constraint. *Genes*, 2(4):748–762, 2011.
- [30] Richard J Orton, Michiel E Adriaens, Amelie Gormand, Oliver E Sturm, Walter Kolch, and David R Gilbert. Computational modelling of cancerous mutations in the egfr/erk signalling pathway. *BMC systems biology*, 3(1):1–17, 2009.
- [31] Hui H Zhang, Alex I Lipovsky, Christian C Dibble, Mustafa Sahin, and Brendan D Manning. S6k1 regulates gsk3 under conditions of mtor-dependent feedback inhibition of akt. *Molecular cell*, 24(2):185–197, 2006.
- [32] Stanford-SNAP-Group. Miner: Gigascale multimodal biological network. *GitHub Repository*, 2017.
- [33] Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor W. Coley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Therapeutics data commons: Machine learning datasets and tasks for drug discovery and development, 2021.

- [34] Nansu Zong, Rachael Sze Nga Wong, Yue Yu, Andrew Wen, Ming Huang, and Ning Li. Drug–target prediction utilizing heterogeneous bio-linked network embeddings. *Briefings in Bioinformatics*, 22(1):568–580, 12 2019.
- [35] Junjun Zhang and Minzhu Xie. Graph regularized non-negative matrix factorization with prior knowledge consistency constraint for drug–target interactions prediction. *BMC Bioinformatics*, 23(1):564, 2022.
- [36] Minoru Kanehisa, Susumu Goto, Masahiro Hattori, Kiyoko F Aoki-Kinoshita, Masumi Itoh, Shuichi Kawashima, Toshiaki Katayama, Michihiro Araki, and Mika Hirakawa. From genomics to chemical genomics: new developments in kegg. *Nucleic Acids Res*, 34(Database issue):D354–7, Jan 2006.
- [37] Ida Schomburg, Antje Chang, Christian Ebeling, Marion Gremse, Christian Heldt, Gregor Huhn, and Dietmar Schomburg. Brenda, the enzyme database: updates and major new developments. *Nucleic Acids Res*, 32(Database issue):D431–3, Jan 2004.
- [38] Stefan Günther, Michael Kuhn, Mathias Dunkel, Monica Campillos, Christian Senger, Evangelia Petsalaki, Jessica Ahmed, Eduardo Garcia Urdiales, Andreas Gewiess, Lars Juhl Jensen, Reinhard Schneider, Roman Skoblo, Robert B Russell, Philip E Bourne, Peer Bork, and Robert Preissner. Supertarget and matador: resources for exploring drug-target relationships. *Nucleic Acids Res*, 36(Database issue):D919–22, Jan 2008.
- [39] US Food and Drug Administration. Questions and answers on fda’s adverse event reporting system (faers). *Washington: US Department of Health and Human Services*, 2018.
- [40] T S Keshava Prasad, Renu Goel, Kumaran Kandasamy, Shivakumar Keerthikumar, Sameer Kumar, Suresh Mathivanan, Deepthi Telikicherla, Rajesh Raju, Beema Shafreen, Abhilash Venugopal, Lavanya Balakrishnan, Arivusudar Marimuthu, Sutopa Banerjee, Devi S Somanathan, Aimy Sebastian, Sandhya Rani, Somak Ray, C J Harrys Kishore, Sashi Kanth, Mukhtar Ahmed, Manoj K Kashyap, Riaz Mohmood, Y L Ramachandra, V Krishna, B Abdul Rahiman, Sujatha Mohan, Prathibha Ranganathan, Subhashri Ramabadran, Raghothama Chaerkady, and Akhilesh Pandey. Human protein reference database–2009 update. *Nucleic Acids Res*, 37, 2009.
- [41] Michael Kuhn, Christian von Mering, Monica Campillos, Lars Juhl Jensen, and Peer Bork. Stitch: interaction networks of chemicals and proteins. *Nucleic acids research*, 36(suppl\_1):D684–D688, 2007.
- [42] Damian Smedley, Syed Haider, Benoit Ballester, Richard Holland, Darin London, Gudmundur Thorisson, and Arek Kasprzyk. Biomart–biological queries made easy. *BMC genomics*, 10(1):1–12, 2009.
- [43] Anna Gaulton, Louisa J Bellis, A Patricia Bento, Jon Chambers, Mark Davies, Anne Hersey, Yvonne Light, Shaun McGlinchey, David Michalovich, Bissan Al-Lazikani, et al. ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic acids research*, 40(D1):D1100–D1107, 2012.
- [44] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. *PyTorch: An Imperative Style, High-Performance Deep Learning Library*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- [45] Greg Landrum, Paolo Tosco, Brian Kelley, Gedeck Sriniker, and Gedeck. Rdkit: Open-source cheminformatics. version 2022.09.1. 2022.
- [46] Helen M. Berman, John Westbrook, Zukang Feng, Gary Gilliland, T. N. Bhat, Helge Weissig, Ilya N. Shindyalov, and Philip E. Bourne. The Protein Data Bank. *Nucleic Acids Research*, 28(1):235–242, 01 2000.
- [47] Mihaly Varadi, Stephen Anyango, Mandar Deshpande, Sreenath Nair, Cindy Natassia, Galabina Yordanova, David Yuan, Oana Stroe, Gemma Wood, Agata Laydon, Augustin Žídek, Tim Green, Kathryn Tunyasuvunakool, Stig Petersen, John Jumper, Ellen Clancy, Richard Green,

- Ankur Vora, Mira Lutfi, Michael Figurnov, Andrew Cowie, Nicole Hobbs, Pushmeet Kohli, Gerard Kleywegt, Ewan Birney, Demis Hassabis, and Sameer Velankar. AlphaFold Protein Structure Database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Research*, 50(D1):D439–D444, 11 2021.
- [48] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstern, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [49] Schrödinger, LLC. The PyMOL molecular graphics system, version 1.8. November 2015.
- [50] Adrian Vallejo, Naiara Perurena, Elisabet Guruceaga, Pawel K. Mazur, Susana Martinez-Canarias, Carolina Zandueta, Karnele Valencia, Andrea Arricibita, Dana Gwinn, Leanne C. Sayles, Chen-Hua Chuang, Laura Guembe, Peter Bailey, David K. Chang, Andrew Biankin, Mariano Ponz-Sarvisé, Jesper B. Andersen, Purvesh Khatri, Aline Bozec, E. Alejandro Sweet-Cordero, Julien Sage, Fernando Lecanda, and Silve Vicent. An integrative approach unveils fosl1 as an oncogene vulnerability in kras-driven lung and pancreatic cancer. *Nature Communications*, 8(1):14294, 2017.
- [51] The UniProt Consortium. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids Research*, 51(D1):D523–D531, 11 2022.
- [52] Yannick Djoumbou Feunang, Roman Eisner, Craig Knox, Leonid Chepelev, Janna Hastings, Gareth Owen, Eoin Fahy, Christoph Steinbeck, Shankar Subramanian, Evan Bolton, Russell Greiner, and David S. Wishart. Classyfire: automated chemical classification with a comprehensive, computable taxonomy. *Journal of Cheminformatics*, 8(1):61, 2016.