

Learning the Drug-Target Interaction Lexicon

Rohit Singh^{a,1}, Samuel Sledzieski^{a,1}, Lenore Cowen^{b,2}, and Bonnie Berger^{a,c,2}

^aComputer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA 02139; ^bDepartment of Computer Science, Tufts University, Medford, MA 02155; ^cDepartment of Mathematics, Massachusetts Institute of Technology, Cambridge, MA 02139

This manuscript was compiled on December 6, 2022

1 **Sequence-based prediction of drug-target interactions has the
2 potential to accelerate drug discovery by complementing exper-
3 imental screens. Such computational prediction needs to be
4 generalizable and scalable while remaining sensitive to subtle
5 variations in the inputs. However, current computational tech-
6 niques fail to simultaneously meet these goals, often sacrific-
7 ing performance on one to achieve the others. We develop
8 a deep learning model, ConPlex, successfully leveraging the
9 advances in pre-trained protein language models (“PLex”) and
10 employing a novel protein-anchored contrastive co-embedding
11 (“Con”) to outperform state-of-the-art approaches. ConPlex
12 achieves high accuracy, broad adaptivity to unseen data, and
13 specificity against decoy compounds. It makes predictions
14 of binding based on the distance between learned representa-
15 tions, enabling predictions at the scale of massive compound
16 libraries and the human proteome. Furthermore, ConPlex is
17 interpretable, which enables us to visualize the drug-target lex-
18 icon and use embeddings to characterize the function of hu-
19 man cell-surface proteins. We anticipate ConPlex will facil-
20 itate novel drug discovery by making highly sensitive and inter-
21 pretable in-silico drug screening feasible at genome scale. Con-
22 PLex is available open-source at <https://github.com/samsledje/>
23 ConPlex.**

Drug discovery | Protein language models | Contrastive learning

1 In the drug discovery pipeline, a key rate-limiting step is the
2 experimental screening of potential drug molecules against a
3 protein target of interest. Thus, fast and accurate computational
4 prediction of drug-target interactions (DTIs) could be extremely
5 valuable, accelerating the drug discovery process. One important
6 class of computational DTI methods, molecular docking, uses
7 3D structural representations of both the drug and target. While
8 the recent availability of high-throughput accurate 3D protein
9 structure prediction models (1–3) means that these methods can
10 be employed starting only from a protein’s amino acid sequence,
11 the computational expense of docking (4) and other structure-
12 based approaches (e.g., rational design (5), active site modeling
13 (6), template modeling (7, 8)) unfortunately remains prohibitive
14 for large-scale DTI screening. An alternative class of DTI predic-
15 tion methods use 3D structure only implicitly, making rapid DTI
16 predictions when the inputs consist only of a molecular descrip-
17 tion of the drug (such as the SMILES string (9)) and the amino
18 acid sequence of the protein target. This class of *sequence-based*
19 DTI approaches enables scalable DTI prediction, but there have
20 been barriers to matching the levels of accuracy obtained by
21 structure-based approaches.

22 In this paper we introduce ConPlex, a new rapid purely

sequence-based DTI prediction method that leverages protein
language models (PLMs), and show it can produce state of the
art performance on the DTI prediction task at scale. The advance
provided by ConPlex comes from two main ideas that together
overcome some of the limitations of previous approaches. While
many methods have been proposed for the sequence-based setting
of the DTI problem (10) (e.g., using secure multi-party computation
(11), convolutional neural networks (12) or transformers
(13)), their protein and drug representations are constructed
solely from DTI ground truth data. The high level of diversity
among the DTI inputs, combined with the limited availability of
DTI training data, limit the accuracy of these methods and their
generalizability beyond their training domain. Furthermore, the
methods that do generalize often do so by sacrificing fine-grained
specificity, i.e., are unable to distinguish true-positive binding
compounds from false positives with similar physico-chemical
properties (“decoys”).

In contrast, the “PLex” (Pre-trained Lexographic) part of
ConPlex helps alleviate the problem of limited DTI training
data. As we showed in our preliminary work (14), one way to get
around the limited size of DTI data sets that has hampered the
quality of the representations learnt by previous methods is to
transfer learned proteins representations from pre-trained protein

Significance Statement

In time and money, one of the most expensive steps of
the drug discovery pipeline is the experimental screen-
ing of small molecules to see which will bind to a protein
target of interest. Therefore, accurate high-throughput
computational prediction of drug-target interactions would
unlock significant value, guiding and prioritizing promis-
ing candidates for experimental screening. We introduce
ConPlex, a machine learning method for predicting drug-
target binding which achieves state-of-the-art accuracy on
many types of targets by using a pre-trained protein lan-
guage model. The approach co-locates the proteins and
the potential drug molecules in a shared feature space
while learning to contrast true drugs from similar non-
binding “decoy” molecules. ConPlex is extremely fast,
which allows it to rapidly shortlist candidates for deeper
investigation.

Please provide details of author contributions here.

The authors have no competing interests to declare.

¹RS and SS contributed equally to this work.

²Correspondence should be addressed to: cowen@cs.tufts.edu, bab@mit.edu

language models to the DTI prediction task. These representations encode structural insights and benefit from being trained on the much larger corpus of single protein sequences (14). Starting with the PLM models, our second insight directly addresses the fine-grained specificity problem in our architecture by using the “Con” (Contrastive learning) part: a novel, protein-anchored contrastive co-embedding that co-locates the proteins and the drugs into a shared latent-space. We show that this co-embedding enforces separation between true interacting partners and decoys to achieve both broad generalization and high specificity.

Putting these two ideas together gives us ConPlex, a novel representation learning approach that enables both broad generalization and high specificity. We show that ConPlex enables more accurate prediction of DTIs than competing methods while avoiding many of the pitfalls suffered by currently available approaches. Thus, our work constitutes a concrete demonstration of the power of a well-designed transfer learning approach that adapts foundation models for a specific task (15, 16). In particular, we found that the performance of existing sequence-based DTI prediction methods could be sensitive to variation in drug-vs-protein coverage in the data set, whereas ConPlex performs well in multiple coverage regimes. Indeed, ConPlex performs especially well in the zero-shot prediction setting where no information is available about a given protein or drug at training time.

ConPlex can also be adapted beyond the binary cases to make predictions about binding affinity. Furthermore, the shared representation also offers advantages beyond prediction accuracy. The co-embedding of both proteins and drugs in the same space offers interpretability, and we show that distances in this space meaningfully reflect protein domain structure and binding function: we leverage ConPlex representations to functionally characterize cell-surface proteins from the Surfaceome database (17), a set of 2,886 proteins localized to the external plasma membrane that participate in signaling and are likely able to be easily targeted by ligands.

ConPlex is extremely fast: as a proof-of-concept, we make predictions for the human proteome against all drugs in ChEMBL (18) ($\approx 2 \times 10^{10}$ pairs) in just under 24 hours. Thus, ConPlex has the potential to be applied for tasks which would require prohibitive amounts of computation for purely structure based approaches or less efficient sequence-based methods, such as genome-scale side-effect screens, identifying drug re-purposing candidates via massive compound libraries searches, or *in silico* deep mutational scans to predict variant effects on binding with currently approved or potential new therapeutics. We note that most DTI methods require significant computation on pairs of drugs and targets (i.e., have quadratic time-complexity). Because ConPlex predictions rely only on the distance in the shared space, predictions can be made highly efficiently once embeddings (which have linear time-complexity) are computed.

Distinguishing between low- and high-coverage DTI prediction. We benchmark performance of ConPlex and competing methods in two different regimes, which we term low-coverage and high-coverage DTI prediction (Figure 1c). We show that ConPlex outperforms its competitors in both set-

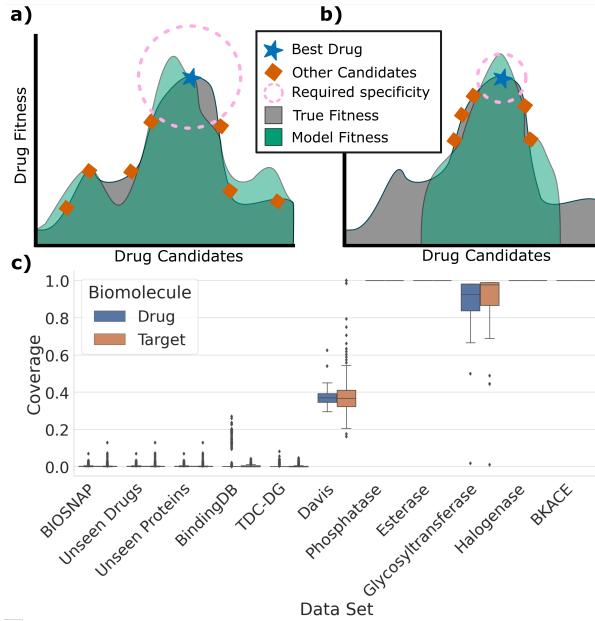


Fig. 1. Drug-target interaction benchmarks display highly variable levels of coverage. Coverage is defined as the proportion of drugs or targets for which a data point (positive or negative) exists in that data set. High- vs. low-coverage benchmarks tend to reward different types of model performance. (a) In this cartoon example of a low coverage data set, drug candidates cover the full diversity of the space, and no two drugs are highly similar. A successful model does not need to be highly specific, but must accurately model the entire fitness landscape to generalize to all candidates. (b) For low-coverage data sets, drugs tend to be much more similar. A successful model does not need to generalize nearly as widely, but must achieve high specificity to differentiate between similar drugs. (c) In a review of existing popular DTI benchmark data sets, we find widely varying coverage, from data sets with nearly zero coverage (each drug/target is represented only a few times) to nearly full coverage (all drug-by-target pairs are known in the data).

tings, but note that separating the two regimes helps clarify an often-seen issue in the field: methods whose performance varies substantially across different proposed DTI benchmarks. Several prior attempts have been made to standardize DTI benchmarking and develop a consistent framework for model evaluation (19, 20). However, much of this work has overlooked a key aspect of benchmarking that we find to significantly affect model performance— differing per-biomolecule data coverage. We define coverage as the average proportion of drugs or targets for which a data point exists in that data set, whether that is a positive or negative interaction (Methods). Depending on the per-biomolecule data coverage of the benchmark data set, we claim that these benchmarks are looking at very different problems. In particular, low-coverage data sets (Figure 1a) tend to measure the broad strokes of the DTI landscape, containing a highly diverse set of drugs and targets. Such data sets can present a modeling challenge due to the diverse nature of targets covered, but allow for a broad assessment of compatibility between classes of compounds and proteins. High-coverage data sets (Figure 1b) represent the opposite trade-off: they contain limited diversity in drug or target type, but report a dense set of potential pairwise interactions. Thus, they capture the fine-grained details of a spe-

124 specific sub-class of drug-target binding and enable distinguishing
125 between similar biomolecules in a particular context.

126 The two coverage regimes correspond to different usage cases.
127 The low-coverage regime is relevant when applying DTI models
128 for large-scale scans to predict interactions for a potential target
129 against a large compound library (e.g. for drug re-purposing as
130 in Dönertaş et al. (21) and Morselli et al. (22)), or for scanning
131 a candidate drug against an entire proteome to identify potential
132 adverse and off-target effects (as in Huang et al. (23, 24)). Data
133 at this scale is often low coverage, with only a small number
134 of known interactions for each unique biomolecule. Thus, it is
135 important that DTI models used for these tasks are broadly ap-
136 plicable and can accurately generalize to many different families
137 of proteins and drugs. However, this generalization often comes
138 at the cost of specificity, resulting in models that are unable to
139 distinguish between highly similar drugs or proteins.

140 The high-coverage regime is relevant when optimizing a par-
141 ticular interaction. Here, models can be trained to be highly
142 specific to a protein family or class of drugs, so much so that
143 a per-drug or per-target model is trained to capture the precise
144 binding dynamics of that biomolecule (25). While such models
145 can be effective for lead optimization, they require high coverage
146 on the biomolecule of interest to make accurate predictions; this
147 may not always be available. Additionally, such models lack
148 the capacity to generalize beyond the training domain and thus
149 cannot be used for genome- or drug bank-scale prediction.

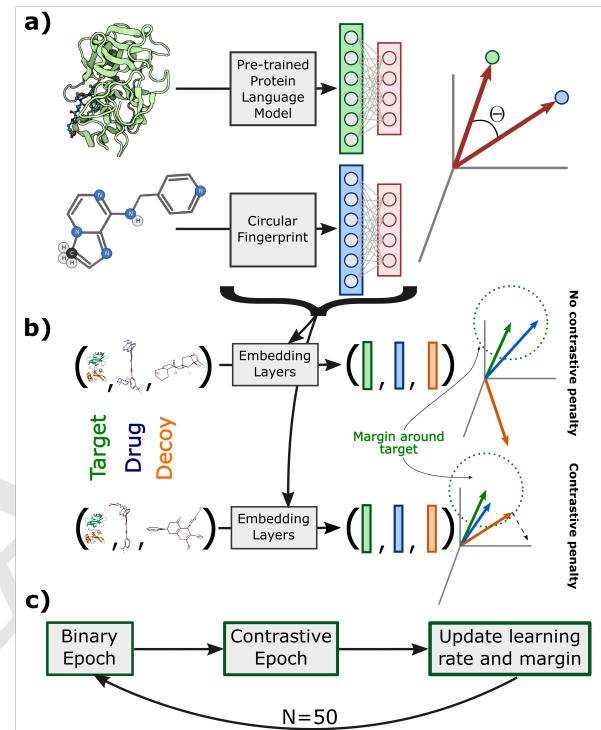
150 The PLM approach of ConPLex enables strong performance
151 in both regimes. In the low coverage regime, the strength is
152 coming mostly from the “PLex” part, where it can leverage the
153 effective generalization of language models to achieve state-
154 of-the-art performance. On high-coverage data sets, the “Con”
155 part also becomes important, since it becomes feasible to train
156 drug- or target-specific models with high accuracy, and such
157 models often outperform more generic models. We find that
158 while single-task models do perform well given available data,
159 ConPLex is able to achieve extremely high specificity in low-
160 diversity, high-coverage scenarios, while remaining broadly ap-
161 plicable to protein targets with limited data. Thus, ConPLex is
162 applicable for both large-scale compound or target screens and
163 fine-grained, highly specific binding prediction. We discuss the
164 issue of matching the right model to the problem domain with
165 respect to coverage further in the **Discussion**.

166 Results

167 **Model Overview.** To achieve both generalizability and speci-
168 ficity, ConPLex leverages advances in both protein language
169 modeling and metric learning. We start with pre-trained represen-
170 tations and learn a non-linear projection of these representations
171 to a shared space (\mathbb{R}^{d_h}). We guide the learning by alternating
172 between two objectives over multiple iterations: a coarse-grained
173 objective of accurately classifying DTIs, and a fine-grained ob-
174 jective of distinguishing decoys from drugs. The coarse-grained
175 objective is evaluated over a low-coverage data set, which trains
176 the model to distinguish between broad classes of drug and tar-
177 get, and makes initial predictions in the right “neighborhood”
178 of the DTI space. The fine-grained objective is evaluated over

179 a high-coverage data set, which fine-tunes the model to distin-
180 guish between true and false positive interactions in the same
181 “neighborhood” and achieve high specificity within a class.

182 To featurize the inputs, here we use the Morgan fingerprint
183 (26) for small molecules, and embeddings from a pre-trained
184 ProtBert model (27) for proteins. We investigate other choices
185 for features, including several other foundation PLMs in Supple-
186 mentary S2. We note that our framework is flexible to different
187 methods of featurization, and make recommendations on the
188 selection of informative representations in the **Discussion**.



189 **Fig. 2. Outline of the ConPLex model architecture and training framework.**
190 ConPLex is trained in two phases, to optimize both generalizability and speci-
191 ficity. (a) Protein features are generated using a pre-trained protein language
192 model (here ProtBert (27)), and drug features are generated using the Morgan
193 fingerprint (26). These features are transformed into a shared latent space
194 by a learned non-linear projection. The prediction of interaction is based on
195 the cosine distance in this space, and the parameters of the transformation
196 are updated using the binary cross-entropy on a low-coverage data set. (b) In
197 the contrastive phase, triplets of a target, drug, and decoy are transformed in
198 the same way into the shared space. Here, the transformation is treated as
199 a metric learning problem. Parameters are updated using the triplet distance
200 loss on a high-coverage data set to minimize the target-drug distance while
201 maximizing the target-decoy distance. No additional penalty is applied if the
202 target-decoy distance is greater than the target-drug distance plus some margin.
203 (c) ConPLex is trained in alternating epochs of the binary and contrastive phase
204 to simultaneously optimize both objectives. After each round, learning rates
205 and the contrastive margin are updated according to an annealing scheme.

206 **ConPLex achieves state-of-the-art performance on low-
207 coverage and zero-shot interactions.** A key advance of Con-
208 PLex is the use of pre-trained protein language models (PLMs)
209 for protein representation. As foreshadowed by Scaiewicz and
210

Table 1. ConPLex is highly accurate and generalizes broadly in low coverage settings. ConPLex outperforms several state-of-the-art methods, including EnzPred-CPI (25), MolTrans (13), GNN-CPI (28), and DeepConv-DTI (12), as well as a simple single-target Ridge regression model, on several low- and zero- coverage benchmark data sets. We report the average and standard deviation of the area under the precision-recall curve (AUPR) for 5 random initializations of each model. Metrics for models with † are taken from (13). Ridge regression cannot be applied for the Unseen Drugs data set, since a separate model is trained for each drug in the training set.

Data Set	ConPLex	EnzPred-CPI	MolTrans	GNN-CPI†	DeepConv-DTI†	Ridge
BIOSNAP	0.897 ± 0.001	0.866 ± 0.003	0.885 ± 0.005	0.890 ± 0.004	0.889 ± 0.005	0.641 ± 0.000
BindingDB	0.628 ± 0.012	0.602 ± 0.006	0.598 ± 0.013	0.578 ± 0.015	0.611 ± 0.015	0.516 ± 0.000
DAVIS	0.458 ± 0.016	0.277 ± 0.009	0.335 ± 0.017	0.269 ± 0.020	0.299 ± 0.039	0.320 ± 0.000
Unseen Drugs	0.874 ± 0.002	0.844 ± 0.005	0.863 ± 0.005	-	0.847 ± 0.009	N/A
Unseen Targets	0.842 ± 0.006	0.795 ± 0.004	0.668 ± 0.045	-	0.766 ± 0.022	0.617 ± 0.000

193 Levitt (29), PLMs have repeatedly been shown to encode evolutionary and structural information (30–32), and to enable broad
194 generalization in low-coverage scenarios (33, 34). Here, we show
195 that ConPLex achieves state-of-the-art performance on three low-
196 coverage benchmark data sets – **BIOSNAP**, **BindingDB**, and
197 **DAVIS** – where it is important to learn the broad strokes of the
198 DTI landscape. In Table 1 we show the average area under the
199 precision-recall curve (AUPR) over 3 random initializations of
200 each model evaluated on a held-out test set (**Methods**). Here,
201 we compare with several methods which use non-PLM protein
202 features: MolTrans (13), GNN-CPI (28), and DeepConv-DTI
203 (12). In addition, we compare to the EnzPred-CPI model from
204 Goldman et al. (25) (developed simultaneously and indepen-
205 dently), which uses a PLM for protein featurization but does not
206 perform a co-embedding or utilize a contrastive training step.
207 Finally, we compare with the single-task Ridge regression model
208 described in (25), which trains a different model per-drug rather
209 than a single model for the entire benchmark.
210

211 Observing the strength of ConPLex to generalize on low-
212 coverage data, we sought to evaluate its performance on fully
213 zero-shot prediction. **Unseen drugs** and **Unseen targets** are
214 variants of the BIOSNAP data set where drugs/targets in the
215 test set do not appear in any interactions in the training set
216 (**Methods**). Note that for the unseen drugs setting, the Ridge
217 model cannot be applied since a different model must be trained
218 for each drug that appears in the training set. We show that
219 ConPLex achieves the best zero-shot prediction performance
220 (Table 1), further demonstrating the applicability of the model to
221 large-scale, very low-coverage prediction tasks.

222 **Contrastive learning enables high-specificity DTI map-
223 ping.** Another key advance of our method is the use of con-
224 trastive learning to fine-tune model predictions on high-coverage
225 data to achieve high specificity. Recently, Heinzinger et al. (36)
226 demonstrated the use of semi-supervised contrastive learning
227 for effective protein embedding-based annotation transfer. Here,
228 we adapt contrastive learning to a fully-supervised setting and
229 demonstrate that the contrastive training is essential to achieving
230 specificity using DTI pairs from the Database of Useful Decoys
231 (**DUD-E**) (35). The DUD-E data set contains 57 protein targets
232 and drugs which are known to interact with each target. However,
233 it also contains 50 negative “decoy” small molecules for each

drug, which have similar physicochemical properties to the truly interacting small molecule, but are known to not bind the target. Thus, accurate prediction on DUD-E requires a model to achieve high-specificity and to accurately differentiate between highly similar compounds. Additionally, DUD-E contains four different classes of targets (GPCRs, kinases, proteases, and nucleases), so models must generalize across target classes (note that single task models don’t have this generalization requirement, since a different model is trained per-target).

We derive evaluation sets from DUD-E by holding out 50% of proteins in each target class for testing and using the remaining targets for training (full splits are specified in Supplementary S1). Here, we evaluate a ConPLex model trained on BIOSNAP, both with and without contrastive training on DUD-E, and show that contrastive training is essential to achieving specificity on decoys.

For each target in the DUD-E test set, we use t-SNE to visualize the target alongside all drugs and decoys using embeddings learned by both versions of the model. Figure 3a,b shows one such example, the tyrosine kinase *VGFR2*. We also show the distribution of distances in the latent space between the target embedding and the embeddings of the drugs and decoys for each model (Figure 3c, d) (*p*-values from one-sided t-test). Without contrastive training, drugs are interspersed with decoys and are far away in space from the target, while ConPLex clusters most true drugs very close to both each other and the *VGFR2* embedding.

In Figure 3e, we show a quantitative analysis of all 31 test-set targets. We compute the effect size (Cohen’s *d*) of the difference between predicted drug and decoy scores. We plot these effect sizes for ConPLex trained with and without contrastive training. An increase in the effect size indicates that the co-embedding distances learned by the model better represent binding specificity. The effect size increases for every target, and the median effect size between predicted true and decoy compound scores was 0.730 prior to contrastive training compared to 4.716 after. For each class of targets, we also report the median *p*-value (one-sided t-test) between drug and decoy scores predicted by ConPLex. While contrastive training has an extremely large impact on specificity in high-coverage domains, we also show that this additional training does not significantly decrease the

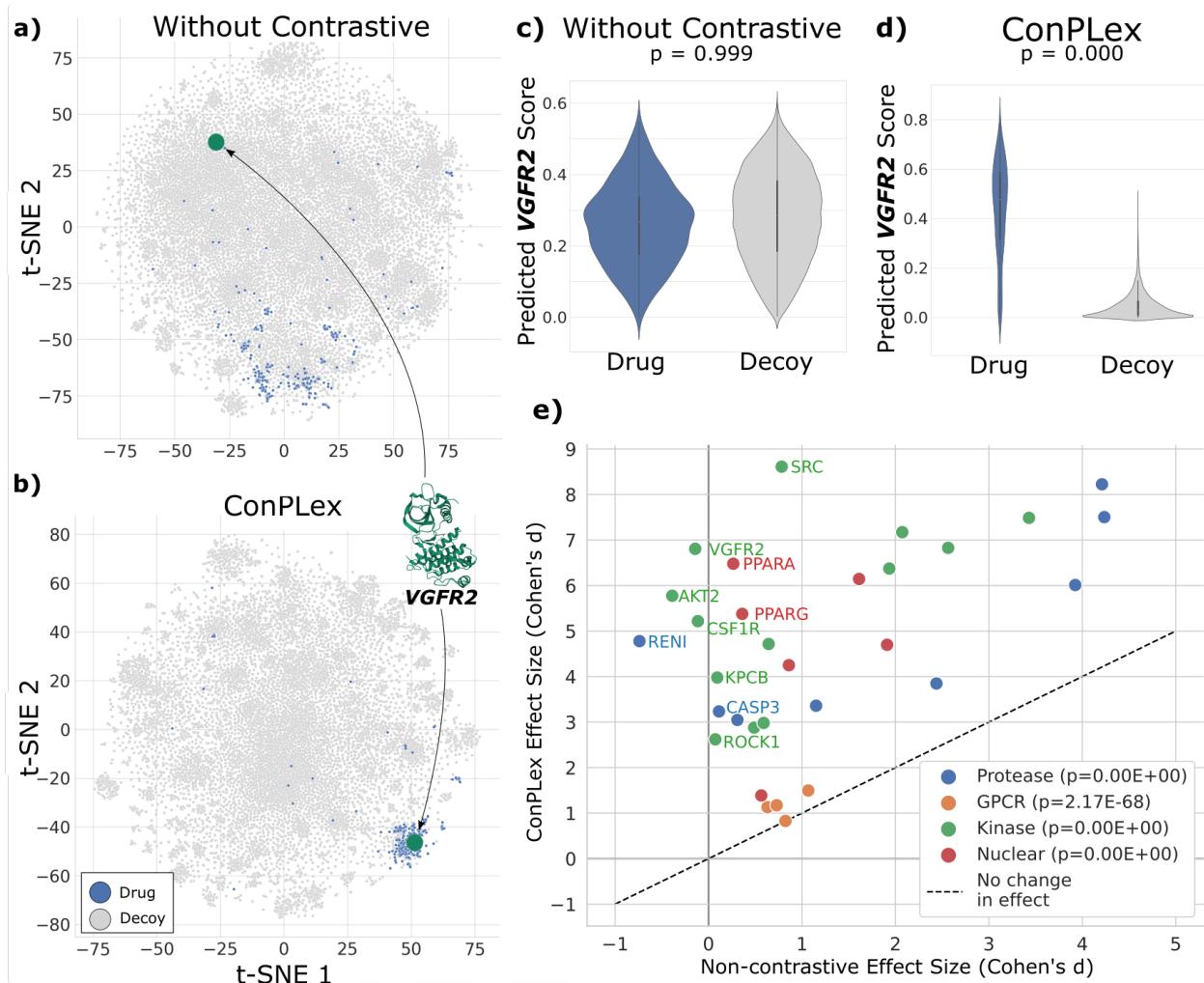


Fig. 3. Contrastive training enables high-specificity in discriminating drugs from decoys. We demonstrate that contrastive learning is essential for ConPlex to achieve high specificity using the DUD-E (35) data set of drugs and decoys (non-binding small molecules with similar physicochemical properties to the true drugs). (a, b) Using t-SNE, we show the learned ConPlex latent space for *VGFR2* (green) and known drugs (blue) and decoys (grey). Without contrastive training, drugs and decoys representations do not separate and true drugs are far from their target. With contrastive training, *VGFR2* and drugs cluster very tightly compared to decoys. (c, d) ConPlex predictions significantly differentiate between drugs and decoys after contrastive training ($p = 0.000$ paired t-test), but do not differ at all without such training ($p = 0.999$). (e) We compute the effect size between drug and decoy predictions using Cohen's d for all 31 targets in the test set. Targets are classified as proteases (blue), GPCRs (orange), kinases (green), and nuclear proteins (red). This effect is computed for ConPlex both with and without contrastive training. Contrastive training increases the effect size for every target (median 0.730 vs 4.716). For each class, we report the median p -value for ConPlex drug vs. decoy predictions. ConPlex performs particularly well for kinases and nuclear proteins, and more poorly for GPCRs.

275 model performance on low-coverage benchmarks via an ablation
276 study in Supplementary S3.

277 In addition to evaluation on DUD-E, we also evaluate Con-
278 PLEX on 5 benchmark data sets derived from family-specific
279 enzyme-substrate screens (**Methods**). These data sets are ex-
280 tremely high coverage, generally including data points for all
281 possible pairs of drugs and targets. We find that in this regime,
282 ConPlex and other PLM based models like EnzPred-CPI have
283 strong but highly variable performance, and are still generally
284 outperformed by a Ridge regression model (Supplementary S5)
285 as shown previously in (25). However, a fine scale single-task
286 model is limited in its generalizability beyond the enzyme family

on which it was trained (**Discussion**).

287

Interpreting the DTI landscape. One of the advantages of the co-embedding approach that our model takes is the ability to visualize and interpret the shared embedding space, and to use proximity within this space to make additional predictions beyond DTI. We apply a ConPlex model trained on BindingDB and fine-tuned on DUD-E to co-embed cell surface proteins from the Surfaceome database (17). In Figure 4a, we show the projections of 2716 proteins from Surfaceome, colored by their classification into one of five functional categories (from Almén et al. (37)) – transporters, receptors, enzymes, miscella-

288

289

290

291

292

293

294

295

296

297

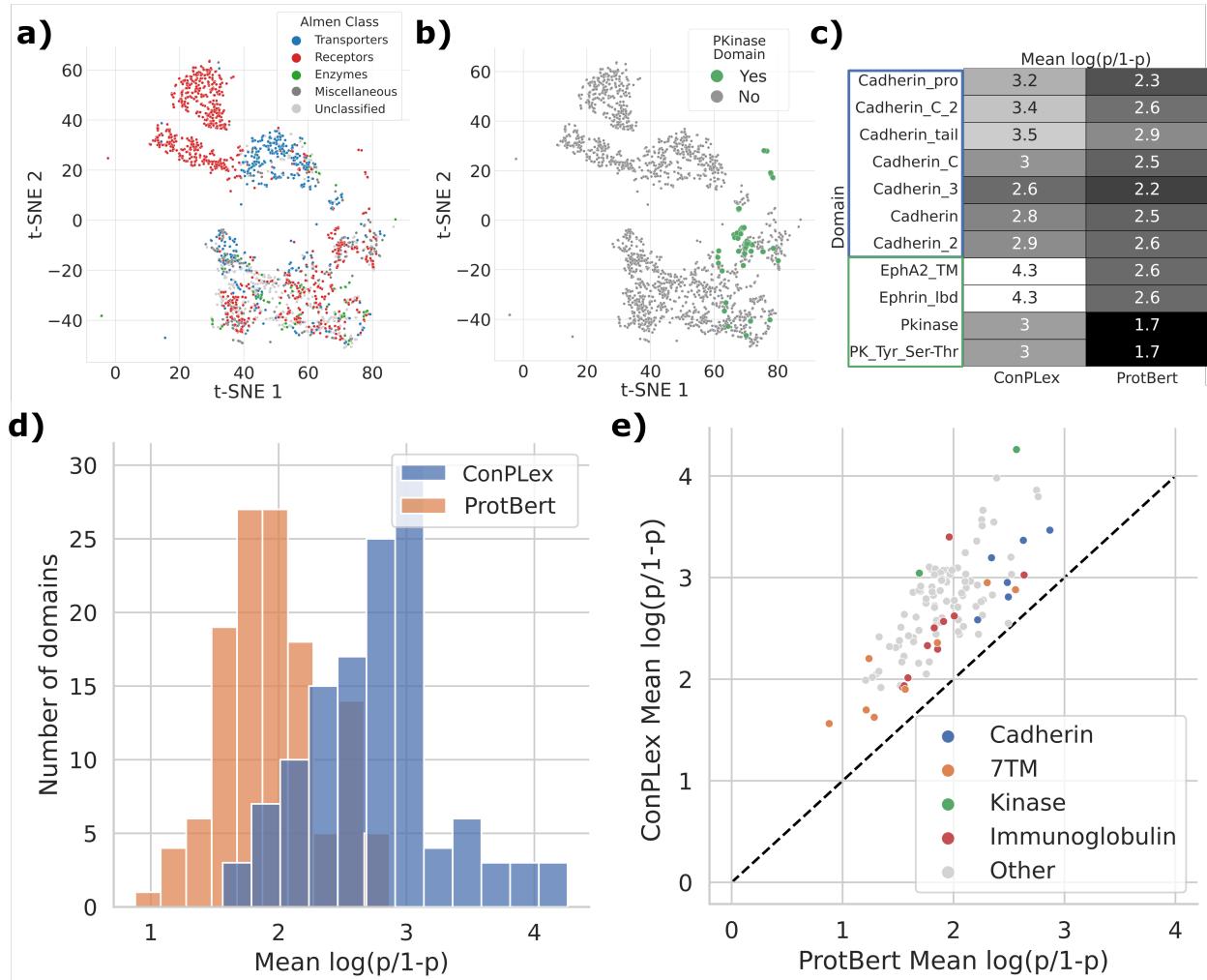


Fig. 4. The shared representation space learned by ConPlex is interpretable and captures protein function. (a) ConPlex representations of cell surface proteins from the Surfaceome (17) cluster by functional class as assigned in Almén et al. (37). (b) These representations also cluster by several functional Pfam domains (38), like the PKinase domain (PF00069) shown in green. (c) We evaluated the coherence of representations for each domain by training a logistic regression classifier and report the model's average confidence for proteins containing that domain as $\log(\frac{p}{1-p})$. While ConPlex separates all domains better than the untransformed ProtBert embeddings, it differentiates kinase domains (green) especially well, while doing more poorly on cadherin domains (blue). (d) Distribution of domain separation for all 126 domains represented in ≥ 10 proteins ($p = 4.85 \times 10^{-54}$, paired t-test). (e) Domain separation using ProtBert vs. ConPlex representations for all 126 domains. ConPlex improves for every domain. We have highlighted several classes of domains, including cadherins (blue), 7-transmembrane proteins (orange), kinases (green), and immunoglobulins (red).

neous, and those that are unclassified. ConPlex projections of surface proteins cluster in embedding space by functional type, with transporters and receptors especially separating from other classes.

However, the Almén functional classification is quite broad and may group proteins with vastly different functions and binding properties. We further demonstrate the link between ConPlex projections and protein function, by evaluating how the learned DTI embedding space separates proteins by domains contained therein. We identified Pfam domains (38) for each protein in the Surfaceome database using HMMscan (39) and compared the projections of proteins that share the same domains. We identified 780 unique domains across all proteins, of

which 126 domains were represented in at least 10 proteins. To quantitatively evaluate the coherence of ConPlex embeddings, we trained separate logistic regression classifiers for each domain to separate proteins with that domain from others, and used the model's confidence ($\log(\frac{p}{1-p})$) for in-sample proteins as a measure of separation for the domain. We find that for all 126 domains, the model more confidently discriminated domains when trained on ConPlex representations than the baseline ProtBert embeddings.

Figure 4d shows the distributions of all such scores for each model, while Figure 4e shows the change in confidence scores for all 126 domains, where the dotted line represents equal confidence using either ConPlex or ProtBert. We find that while

324 prediction of all domains were improved using ConPlex, kinase
325 domains (PF14575, PF01404, PF00069, PF07714) separated
326 especially well, while Cadherin domains (PF08785, PF16492,
327 PF15974, PF01049, PF16184, PF00028, PF08266) showed more
328 modest improvement (Figure 4c). As discussed previously, Con-
329 PLex was trained contrastively with several kinase targets and
330 excels at kinase prediction on DUD-E (Figure 3e), so it is un-
331 surprising that proteins with these domains separate well. In fact,
332 one of the top differentiated domains is the Ephrin ligand binding
333 domain (PF01404), which is responsible for binding to the ephrin
334 ligand (40). In Figure 4b, we show the same visualization of pro-
335 jections as in Figure 4a, but colored by another top-differentiated
336 domain, PKinase (PF00069). While the model was not explicitly
337 trained on targets with Cadherin domains, cadherins typically are
338 responsible for cell-cell adhesion and structural stability, and are
339 less involved in ligand binding (41). Thus, we would not expect
340 a model trained to represent small molecule binding to explicitly
341 differentiate cadherins. This demonstrates that the latent space
342 learned by ConPlex is useful not only for predicting drug-target
343 interaction, but for broadly classifying protein function as it re-
344 lates to ligand binding. Future work in this area might adjust
345 distances in this landscape to account for the low metric entropy
346 of biological sequences, as demonstrated in Berger, Waterman,
347 and Yu (42).

348 **Adapting ConPlex for affinity prediction.** While we have to
349 this point been using the model to predict probabilities of inter-
350 action and perform binary classification, we show that ConPlex
351 can be easily adapted to perform binding affinity prediction, and
352 that this model too achieves state-of-the-art performance. The
353 final step of our binary interaction predictor is converting the
354 cosine distance between the projections in the DTI space to a
355 probability using a sigmoid activation (**Methods**). However, it is
356 completely natural to replace this activation with a dot product
357 between the two projections, which enables the model to make
358 real-valued predictions, which can then be interpreted as a bind-
359 ing affinity. We evaluated ConPlex trained for affinity prediction
360 on the Therapeutics Data Commons (TDC) DTI Domain Gen-
361 eralization (**TDC-DG**) benchmark. The TDC-DG benchmark
362 contains binding affinity (K_d) data from interactions patented
363 between 2013 and 2018, with the test set drawn from interactions
364 patented in 2019–2021 (**Methods**). Thus, this data requires out-
365 of-domain generalization and corresponds with the real-life sce-
366 nario of training on interactions up to a known point, and predict-
367 ing interactions which are yet to be documented. We trained Con-
368 PLex to predict binding affinity with 5 random train/validation
369 splits, and achieve an average Pearson correlation coefficient
370 between the true and predicted affinity of $0.538(\pm 0.008)$ on the
371 held-out test set. At the time of submission, ConPlex is the top
372 performing method on the TDC-DG benchmark on TDC (Table
373 2).

374 Discussion

375 Much previous work has recognized the value of meaningful
376 drug representations (44, 45) for DTI prediction, yet relatively
377 little work has focused on the target protein representation. As
378 the first method to use pre-trained protein language models

379 **Table 2. ConPlex can be adapted for state-of-the-art K_d predic-
380 tion.** By replacing the cosine distance in the final step of Con-
381 PLex with a dot product between the projections, ConPlex can
382 be used for affinity prediction rather than binary classifica-
383 tion. The Therapeutics Data Commons-Domain Generalization
384 (**TDC-DG**) data set contains K_d values for patented drug-target pairs,
385 where training/testing data are split from before/after 2018. We
386 report the average and standard deviation of the Pearson Corre-
387 lation Coefficient between true and predicted values across five
388 train/validation splits. Metrics for all methods other than Con-
389 PLex come from the TDC leaderboard (19, 43), where at the time
390 of submission ConPlex is the best performing method.
391

Model	PCC
ConPlex	0.538 ± 0.008
MMD	0.433 ± 0.010
CORAL	0.432 ± 0.010
ERM	0.427 ± 0.012
MTL	0.425 ± 0.010
GroupDRO	0.384 ± 0.006
AndMASK	0.288 ± 0.019
IRM	0.284 ± 0.021

(PLMs) for DTI prediction, ConPlex is yet another example
379 of the power of transferring learned representations for biology
380 (13, 32, 33, 46, 47). This approach enables broad generalization
381 to unseen proteins, as well as extremely fast model inference
382 ($>10x$ speed-up even over other sequence-based approaches,
383 Supplementary S4). This speed is particularly valuable for drug
384 re-purposing and iterative screening, where large compound li-
385 braries are evaluated against hitherto-uncharacterized proteins
386 implicated in a disease of interest. The co-embedding approach
387 which enables this speedup could also be effective for integrative
388 multi-structure models (e.g., the IMP framework (48)) where
389 efficient scanning of possible combinations is important. Re-
390 cent methods have also demonstrated the power of PLMs for
391 transferring knowledge between species (33), and our framework
392 may enable more accurate transfer of DTI from the model or-
393 ganisms on which drugs are initially tested, to their eventual use
394 in human patients. Skolnick and Zhou (49) have reported the
395 importance of considering small molecule binding pockets for
396 protein-protein interaction prediction; thus our DTI-informed
397 protein representations may also be useful in that context. While
398 structural similarity is often implicitly learned by PLMs, future
399 work could explicitly incorporate structure where such data is
400 available, perhaps by incorporating a more advanced projection
401 architecture like the Geoformer (3).
402

403 It has been shown in previous work that the performance
404 of different PLMs vary on different tasks, and that there is not
405 one clearly “best” language model (14, 50, 51). While we have
406 chosen to use ProtBert here, it is likely that other existing or
407 newly developed language models may yield better performance
408 for certain types of drugs or targets. Likewise, advancements in
409 drug representation may improve performance—the ConPlex
410 framework is flexible to different input features, and it remains
411 important to experiment with different feature choices for the
412 task at hand (Supplementary S2).

413 ConPlex approaches the DTI decoy problem from the per-
414 spective of adversarial machine learning, where the model must
415 act as a discriminator for adversarial examples from the decoy
416 database. This approach is directly enabled by the co-embedding
417 architecture—to compute the triplet distance loss, the protein
418 and drugs must be co-embedded, and the distance between them
419 must be meaningful and simply computed. Such an approach
420 would not be feasible using a model which concatenates features
421 up front, nor for a model which has significant computation
422 defining the probability of interaction after the co-embedding.
423 Thus the shared lexicographic space in which we embed the
424 proteins, targets, and decoys is key. Future work could explore
425 adapting molecular generation methods such as JT-VAE or Hi-
426 erG2G (52, 53) to directly act as a generator for decoys. High-
427 specificity DTI prediction is valuable beyond decoy detection—
428 greater specificity of inference can help improve personalized
429 medicine or the modeling of drug effects against rare variants
430 from under-represented populations.

431 It is also important to consider the coverage of the problem
432 to select an appropriate method. If enough data is available, for
433 specific enzyme-family prediction tasks we still recommend the
434 use of single-task models (25). To verify individual interactions,
435 energy-based molecular docking will likely be more accurate,
436 although at the cost of being substantially slower (4). Different
437 classes of computational tools for DTI prediction each have
438 varying strengths, and the highest quality predictions can be
439 achieved by leveraging all of these methods together where each
440 is most fit.

441 Drug discovery is a fundamental task for human health, yet
442 remains both extremely coarse expensive and time-consuming,
443 with the median drug requiring about 1 billion dollars (54) and 10
444 years (55) from development to approval and distribution. While
445 experimental results will remain the gold-standard for validating
446 drug functionality, *in silico* prediction of drug-target binding
447 remains much faster and cheaper and so will continue to play
448 an important role in early screening of therapeutic candidates
449 (56). To address this step in the drug design pipeline, we have
450 introduced ConPlex. DTI prediction methods should be able
451 to generalize to unseen types of drugs and targets, while also
452 discriminating between highly similar molecules with different
453 binding properties. ConPlex tackles both of these challenges
454 through its dual use of protein language models and contrastive
455 learning. We hope that its broad applicability, specificity on
456 decoys, and ability to scale to massive data will allow ConPlex
457 to be a critical step in this pipeline and contribute to the efficient
458 discovery of effective therapeutics.

459 Materials and Methods

460

461 **Computing Data Set Coverage.** Let $1_{(i,j)}$ be the indicator variable
462 meaning there exists an observation of drug i and target j . For a
463 data set with m unique drugs and n unique targets, we can define the
464 coverage for drug d as $C_d = \frac{1}{n} \sum_{j=0}^n 1_{(d,j)}$ and for a target t as
465 $C_t = \frac{1}{m} \sum_{i=0}^m 1_{(i,t)}$. Then, for a given data set we can evaluate the
466 median drug and target coverage. A data set with maximum coverage
467 would have a single data point for each drug-target pair, and thus a me-
468 dian coverage of 1 for both drugs and targets. Conversely, each drug and

target would only be represented a single time in a minimum coverage
469 data set, resulting in drug and target coverages of $\frac{1}{n}$ and $\frac{1}{m}$ respectively.
470 We report the median drug and target coverage for each benchmark data
471 set in Table 3. Since the DUD-E data set is separated out by targets, we
472 instead report the median number of drugs against each target.
473

474 Benchmarks Overview.

475 **Low coverage benchmarks** We evaluate our framework on three
476 broad-scale, low-coverage benchmark data sets. Two data sets, **DAVIS**
477 (59) and **BindingDB** (58), consist of pairs of drugs and targets with
478 experimentally determined dissociation constants (K_d). Following (13),
479 we treat pairs with $K_d < 30$ as positive DTIs, while larger K_d values
480 are negative. The third data set, ChG-Miner from **BIOSNAP** (57),
481 consists of only positive DTIs. We create negative DTIs by randomly
482 sampling an equal number of protein-drug pairs, making the assumption
483 that a random pair is unlikely to be positively interacting. The DAVIS
484 data set represents a few-shot learning setting: it contains only 2,086
485 training interactions, compared to 12,668 for BindingDB and 19,238 for
486 BIOSNAP. The rest of the data preparation follows (13). The data sets
487 are split into 70% for training, 10% for validation, and the remaining
488 20% for testing. Training data are artificially sub-sampled to have an
489 equal number of positive and negative interactions, while validation and
490 test data is left at the ratio originally in the data set.

491 **Zero-shot benchmarks.** We evaluate our framework on two zero-
492 shot prediction modifications of BIOSNAP. Following (13), the **Unseen**
493 **proteins** set was created by selecting 20% of proteins from the full set,
494 and selecting any interactions including these proteins for the test set.
495 Thus, there are no proteins which appear in both the training and test set.
496 The corresponding process was used to create the **Unseen drugs** data
497 set. The training set was then further split using 7/8 of the interactions
498 for training and 1/8 of the interactions for testing. As above, data are
499 sub-sampled so that training is balanced.

500 **Continuous benchmarks.** Continuous affinity prediction data
501 come from the Therapeutics Data Commons DTI Domain Generalization
502 benchmark (**TDC-DG**) (19). The TDC-DG consists of 140,746 unique
503 drugs and 477 unique targets derived from BindingDB (58) interactions
504 that have patent information. Each interaction is labeled with an ex-
505 perimentally determined dissociation constant (K_d). Interactions are
506 temporally split, so that training pairs are from patents filed between
507 2013 and 2018, and test pairs are from between 2019 and 2021. 20% of
508 the training pairs are randomly set aside as a validation set. We train 5
509 different models with the train/validation splits determined by the TDC
510 benchmarking framework, and report the average Pearson correlation
511 coefficient of predictions on the test set.

512 **High coverage benchmarks.** The Database of Useful Decoys: En-
513 hanced (**DUD-E**) (35) consists of 102 protein targets and known binding
514 partners (average 224 molecules per target). For each binding partner,
515 there are 50 “decoys”, or physio-chemically similar compounds that are
516 known not to bind with the target. 57 of the targets are classified as either
517 GPCRs, kinases, nuclear proteins, or proteases. We generate train-test
518 splits by splitting targets within classes, so that there are representative
519 members of each class in both the training and test set, but no target
520 appears in both the training and test set (26 train, 31 test). These data
521 are by definition high-coverage, since there are several true and decoy
522 compounds available for each target. We provide the full target splits in
523 Supplementary S1.

524 We also evaluate on several protein-family specific data sets from
525 various different sources and compiled by Goldman et al. (25). These
526 include DTI data on β -ketoacid cleavage (**BKACE**) (64), **Esterase** (61),
527 **Glycosyltransferases** (62), **Halogenase** (63), and **Phosphatase** (60)
528 enzymes. These data are uniformly very high coverage, with a known data
529 point for nearly every drug-target pair. Following (25), we performed a
530 10-fold cross validation where the data were split into train-test sets by
531 target, so that all drugs appear in both the training and test set, but no
532 target does.

533 ConPlex Model.

Table 3. Full specification of benchmark data sets. We report the number of unique drugs and targets, the median (drug/target) coverage, and the number of training, validation, and test samples in each data set. Number of pairs are shown as (positive/negative), except for TDC-DG (19, 43), which is a regression task, thus total number of pairs is shown. We consider BIOSNAP (57), BindingDB (58), DAVIS (59), and TDC-DG as low-coverage, while Phosphatase (60), Esterase (61), Glycosyltransferase (62), Halogenase (63), BKACE (64), and DUD-E (35) are considered high-coverage. † Because DUD-E is a decoy data set, we report as coverage the median number of true drugs or decoys for each target.

Data Set	Drugs	Targets	Median Coverage	# Training	# Validation	# Test
BIOSNAP	4510	2181	0.0023 / 0.0020	9670 / 9568	1396 / 1352	2770 / 2727
Unseen Drugs				9535 / 9616	1383 / 1353	2918 / 2675
Unseen Targets				9876 / 9499	1382 / 1386	2578 / 2762
BindingDB	7165	1254	0.0008 / 0.0010	6334 / 6334	927 / 5717	1905 / 11384
DAVIS	68	379	0.3707 / 0.3676	1043 / 1043	160 / 2846	303 / 5708
TDC-DG	140746	477	0.0021 / 0.0005	146891	36539	49028
Phosphatase	165	218	1.0 / 1.0	5054 / 27286	—	370 / 3260
Esterase	96	146	1.0 / 1.0	2150 / 10426	—	926 / 514
Glycosyltransferase	89	54	0.9259 / 0.9778	725 / 3042	—	113 / 417
Halogenase	62	42	1.0 / 1.0	303 / 1991	—	20 / 290
BKACE	17	161	1.0 / 1.0	255 / 2193	—	19 / 270
DUD-E †				8996 / 406208	—	11430 / 521132
GPCR	99671	5	18563			
Kinase	315399	26	15409			
Protease	286089	15	9271			
Nuclear	151133	11	16257			

534 **Target featurization.** We generate protein target features using pre-
 535 trained protein language models (PLM): These models generate a protein
 536 embedding $E_{full} \in \mathbb{R}^{n \times d_t}$ for a protein of length n , which is then
 537 mean-pooled along the length of the protein resulting in a vector $E \in$
 538 \mathbb{R}^{d_t} . Specifically, we investigate the pre-trained models Prose (31), ESM
 539 (65), and ProtBert (27), with default dimensions $d_t = 6165, 1280, 1024$
 540 respectively (Supplementary S2). Elnaggar et al. recommend the use of
 541 ProtT5XLUniref50, but we found that it did not perform as well as
 542 ProtBert for the DTI prediction task. We emphasize that the language and
 543 projection models are used exclusively to generate input features— their
 544 weights are kept unchanged and are not updated during DTI training.

545 **Drug featurization.** We featurize the drug molecule by its Morgan
 546 fingerprint (26), an encoding of the SMILES string of the molecular
 547 graph as a fixed-dimension embedding $M \in \mathbb{R}^{d_m}$ (we chose $d_m =$
 548 2048) by considering the local neighborhood around each atom. The
 549 utility of the Morgan fingerprint for small molecule representation has
 550 been demonstrated in (25, 66). We additionally investigated the use of
 551 molecule embeddings from Mol2Vec (67) and MolR (68) and found they
 552 failed to perform as well as the Morgan fingerprint (Supplementary S2).

553 **Transformation into a shared latent space and prediction.** Given a
 554 target embedding $T \in \mathbb{R}^{d_t}$ and small molecule embedding $M \in \mathbb{R}^{d_m}$,
 555 we transform them separately into $T^*, M^* \in \mathbb{R}^h$ using a single fully-
 556 connected layer with a ReLU activation. These layers are parameterized
 557 with weight matrices $W_t \in \mathbb{R}^{h \times d_t}, W_m \in \mathbb{R}^{h \times d_m}$ and bias vectors
 558 $b_t, b_m \in \mathbb{R}^h$.

$$T^* = \text{ReLU}(W_t T + b_t) \quad [1]$$

$$M^* = \text{ReLU}(W_m M + b_m) \quad [2]$$

562 Given the latent embeddings T^*, M^* , we compute the probability of
 563 a drug-target interaction $\hat{p}(T^*, M^*)$ as the cosine similarity between the
 564 embedding vectors, followed by a sigmoid activation. Thus, we compute
 565 the predicted probability as

$$\hat{p}(T^*, M^*) = \sigma\left(\frac{T^* \cdot M^*}{\|T^*\|_2 \cdot \|M^*\|_2}\right) \quad [3]$$

When predicting compound binding affinity $\hat{y}(T^*, M^*)$, we substitute the sigmoid and cosine similarity (Equation 3) with a dot product followed by a ReLU activation, which gives a non-negative distance in the embedding space (Equation 4).

$$\hat{y}(T^*, M^*) = \text{ReLU}(T^* \cdot M^*) \quad [4]$$

572 **Training.** The model is trained both for broad and fine prediction, with
 573 the loss computed depending on the training data set. Broad-scale training
 574 data uses the binary cross-entropy loss (L_{BCE}) between the true
 575 labels y and the predicted interaction probabilities \hat{p} . When the model
 576 was trained to predict binding affinity, we substitute the binary cross-
 577 entropy loss for the mean squared error loss (L_{MSE}) is used during
 578 supervision.

579 Training on fine-scale data (DUD-E) was performed using contrastive
 580 learning. Contrastive learning uses triplets of training points rather than
 581 pairs, denoted the **anchor**, **positive**, and **negative**, and aims to minimize
 582 the distance between the anchor and positive examples while maximizing
 583 the distance between the anchor and the negative examples. In the
 584 DTI setting, the natural choice for a triplet is the protein target as the
 585 anchor, the true drug as the positive and decoy as the negative example,
 586 respectively. We derive a training set of triplets in the following manner:
 587 for each known interacting drug-target pair (T, M^+) , we randomly
 588 sample $k = 50$ non-interacting pairs (T, M^-) and generate the triplets
 589 (T, M^+, M^-) , where M^- is drawn from the set of all decoys against
 590 T . We map these to latent space embeddings as described above. Since
 591 all the entities are now comparable to each other, we can compute the
 592 triplet margin-distance loss (L_{TRM}).

$$L_{TRM}(a, p, n) = \frac{1}{N} \sum_{i=1}^N \max(D(a, p) - D(a, n) + m, 0) \quad [5]$$

593 where

$$D(u, v) = 1 - \hat{p}(u, v) \quad [6]$$

596 The margin m sets the maximum required delta between distances,
 597 above which the loss is zero.

598 **Margin annealing.** The margin m sets the maximum required delta
599 between distances, above which the loss is zero. Initially, a large margin
600 requires the decoy to be much further from the target than the drug to
601 avoid a penalty, resulting in larger weight updates. As training progresses,
602 lower margins relax this constraint, requiring only that the drug be
603 closer than the decoy as $m \rightarrow 0$. Here, the margin is initialized at
604 $M_{max} = 0.25$ according to a tanh decay with restarts decay schedule.
605 Every $E_{max} = 10$ contrastive epochs, the margin is reset to the initial
606 M_{max} , for a total of 50 epochs. At epoch i , the margin is set to

$$m(i) = M_{max}(1 - \tanh(\frac{2(i \bmod E_{max})}{E_{max}})) \quad [7]$$

607 **Implementation.** Model weights were initialized using the Xavier
608 method from a normal distribution (69). Weights were updated with
609 error back-propagation using the AdamW optimizer (70) for a total of 50
610 epochs. For the binary classification task, the learning rate was initially
611 set to 10^{-4} and adjusted according to a cosine annealing schedule with
612 warm restarts (71) every 10 epochs. For the contrastive task, the learning
613 rate was initially set to 10^{-5} and the same annealing schedule was fol-
614 lowed. The margin for the contrastive loss was initially set to 0.25 and
615 decreased to a minimum of 0 over 50 epochs according to a tanh decay
616 schedule with restarts every 10 epochs. We used a latent dimension
617 $d = 1024$ (results were robust to even with lower dimensions, and much
618 higher dimensions may over-fit or be subject to topological restrictions)
619 and a batch size of 32. The model was implemented in PyTorch version
620 1.11. Model training and inference was performed on a single NVIDIA A100 GPU.
621

622 **Surfaceome analysis.** We evaluate the interpretability and functional
623 use of ConPLex embeddings using data from the Surfaceome database
624 (17), which contains 2,886 cell-surface proteins. We identified Pfam
625 domains using HMMscan from HMMER3 (39) with default settings. We
626 analyzed domains hit in > 10 proteins. For each domain, we trained a lo-
627 gistic regression classifier from sklearn with balanced class weights. We
628 also evaluated domain coherence using spectral clustering with $k = 10$
629 clusters, and evaluated the adjusted mutual information (AMI) between
630 true clusters (protein has/doesn't have domain) and predicted clusters
631 (Supplementary S6).

632 **Genome wide ChEMBL scan.** We trained a ConPLex model using
633 BindingDB and DUD-E, and used it to make predictions for all pairs of
634 human proteins against all drugs in ChEMBL. Human protein sequences
635 were taken from the STRING database and processed following (33),
636 resulting in 15,816 proteins between 50 and 800 amino acids long. Small
637 molecule structures were downloaded from ChEMBL 30 (18), resulting
638 in 1,533,652 compounds. Prediction took just under a day, accounting
639 for embedding time.

641 **ACKNOWLEDGMENTS.** RS and BB were supported by the NIH
642 grant R35GM141861. SS was supported by the National Science Foun-
643 dation Graduate Research Fellowship under Grant No. 2141064. LC was
644 supported by CCF-1934553. The authors thank Kapil Devkota, Tristan
645 Bepler, and Tim Truong for helpful discussions.

- 646 1. Jumper J, et al. (2021) Highly accurate protein structure prediction with AlphaFold.
Nature 596(7873):583–589.
- 647 2. Baek M, et al. (2021) Accurate prediction of protein structures and interactions using
648 a three-track neural network. Science 373(6557):871–876.
- 649 3. Wu R, et al. (2022) High-resolution de novo structure prediction from primary se-
650 quence. *bioRxiv*.
- 651 4. Pinzi L, Rastelli G (2019) Molecular docking: shifting paradigms in drug discovery.
International journal of molecular sciences 20(18):4331.
- 652 5. Bonk BM, Tarasova Y, Hicks MA, Tidor B, Prather KL (2018) Rational design of thio-
653 lase substrate specificity for metabolic engineering applications. Biotechnology and
654 bioengineering 115(9):2167–2182.
- 655 6. de Melo-Minardi RC, Bastard K, Artiguenave F (2010) Identification of subfamily-
656 specific sites based on active sites modeling and clustering. Bioinformatics 26(24):3075–3082.
- 657 7. Trudeau SJ, et al. (2022) Prepc: A structure-and chemical similarity-informed
658 database of predicted protein compound interactions. *bioRxiv*.

- 659 8. Singh R, Park D, Xu J, Hosur R, Berger B (2010) Struct2Net: a web service to predict
660 protein–protein interactions using a structure-based approach. Nucleic acids research
661 38(suppl_2):W508–W515. publisher: Oxford University Press.
- 662 9. Anderson E, Veith GD, Weininger D (1987) SMILES, a line notation and computerized
663 interpreter for chemical structures. (US Environmental Protection Agency, Environ-
664 mental Research Laboratory).
- 665 10. Bagherian M, et al. (2021) Machine learning approaches and databases for prediction
666 of drug–target interaction: a survey paper. *Briefings in Bioinformatics* p. 23.
- 667 11. Hie B, Cho H, Berger B (2018) Realizing private and practical pharmacological collab-
668 oration. Science 362(6412):347–350.
- 669 12. Lee I, Keum J, Nam H (2019) DeepConv-DTI: Prediction of drug-target interactions
670 via deep learning with convolution on protein sequences. PLoS computational biology
671 15(6):e1007129.
- 672 13. Huang K, Xiao C, Glass LM, Sun J (2021) MolTrans: Molecular Interaction Transformer
673 for drug–target interaction prediction. Bioinformatics 37(6):830–836.
- 674 14. Sledzieski S, Singh R, Cowen L, Berger B (2021) Adapting protein language models
675 for rapid DTI prediction. Machine Learning for Structural Biology Workshop (MLSB) at
676 NeurIPS.
- 677 15. Bommasani R, et al. (2021) On the opportunities and risks of foundation models. arXiv
678 preprint arXiv:2108.07258.
- 679 16. Gururangan S, et al. (2020) Don't stop pretraining: adapt language models to domains
680 and tasks. arXiv preprint arXiv:2004.10964.
- 681 17. Bausch-Fluck D, et al. (2018) The *in silico* human surfaceome. Proceedings of the
682 National Academy of Sciences 115(46):E10988–E10997.
- 683 18. Mendez D, et al. (2019) ChEMBL: towards direct deposition of bioassay data. Nucleic
684 acids research 47(D1):D930–D940.
- 685 19. Huang K, et al. (2021) Therapeutics data commons: Machine learning datasets and
686 tasks for drug discovery and development. arXiv preprint arXiv:2102.09548.
- 687 20. Zong N, et al. (2022) Beta: a comprehensive benchmark for computational drug-
688 target prediction. *Briefings in Bioinformatics*.
- 689 21. Dönertaş HM, Fuentealba Valenzuela M, Partridge L, Thornton JM (2018) Gene
690 expression-based drug repurposing to target aging. Aging cell 17(5):e12819.
- 691 22. Morselli Gysi D, et al. (2021) Network medicine framework for identifying drug-
692 repurposing opportunities for covid-19. Proceedings of the National Academy of Sciences
693 118(19):e2025581118.
- 694 23. Huang K, et al. (2020) DeepPurpose: a deep learning library for drug–target interac-
695 tion prediction. Bioinformatics 36(22–23):5545–5547.
- 696 24. Huang Y, et al. (2019) A framework for identification of on-and off-target transcriptional
697 responses to drug treatment. Scientific reports 9(1):1–9.
- 698 25. Goldman S, Das R, Yang KK, Coley CW (2022) Machine learning modeling of
699 family wide enzyme-substrate specificity screens. PLoS computational biology
700 18(2):e1009853.
- 701 26. Morgan HL (1965) The generation of a unique machine description for chemical
702 structures—a technique developed at chemical abstracts service. Journal of Chemical
703 Documentation 5(2):107–113.
- 704 27. Elnaggar A, et al. (2020) ProtTrans: towards cracking the language of life's code
705 through self-supervised deep learning and high performance computing. arXiv
706 preprint arXiv:2007.06225.
- 707 28. Tsubaki M, Tomi K, Sese J (2019) Compound–protein interaction prediction with
708 end-to-end learning of neural networks for graphs and sequences. Bioinformatics
709 35(2):309–318.
- 710 29. Scaiiewicz A, Levitt M (2015) The language of the protein universe. Current opinion in
711 genetics & development 35:50–56.
- 712 30. Bepler T, Berger B (2019) Learning protein sequence embeddings using information
713 from structure in 7th International Conference on Learning Representations, ICLR
714 2019.
- 715 31. Bepler T, Berger B (2021) Learning the protein language: Evolution, structure, and
716 function. Cell Systems 12(6):654–669.e3. Publisher: Elsevier.
- 717 32. Heinzinger M, et al. (2019) Modeling aspects of the language of life through transfer-
718 learning protein sequences. BMC bioinformatics 20(1):1–17.
- 719 33. Sledzieski S, Singh R, Cowen L, Berger B (2021) D-ScripT translates genome to
720 phenotype with sequence-based, structure-aware, genome-scale predictions of protein-
721 protein interactions. Cell Systems 12:1–14.
- 722 34. Singh R, Devkota K, Sledzieski S, Berger B, Cowen L (2022) Topsy-Turvy:
723 integrating a global view into sequence-based ppi prediction. Bioinformatics
724 38(Supplement_1):i264–i272.
- 725 35. Mysinger MM, Carchia M, Irwin JJ, Shoichet BK (2012) Directory of useful decoys,
726 enhanced (DUD-E): better ligands and decoys for better benchmarking. Journal of
727 medicinal chemistry 55(14):6582–6594.
- 728 36. Heinzinger M, et al. (2022) Contrastive learning on protein embeddings enlightens
729 midnight zone. NAR genomics and bioinformatics 4(2):lqac043.
- 730 37. Almén MS, Nordström KJ, Fredriksson R, Schiöth HB (2009) Mapping the human
731 membrane proteome: a majority of the human membrane proteins can be classified
732 according to function and evolutionary origin. BMC biology 7(1):1–14.
- 733 38. El-Gebali S, et al. (2019) The Pfam protein families database in 2019. Nucleic Acids
734 Research 47(D1):D427–D432. Publisher: Oxford University Press.
- 735 39. Eddy SR (2011) Accelerated profile HMM searches. PLoS computational biology
736 7(10):e1002195.
- 737 40. 662–663
- 738 41. 664–665
- 739 42. 666–667
- 740 43. 668–669
- 741 44. 670–671
- 742 45. 672–673
- 743 46. 674–675
- 744 47. 676–677
- 745 48. 678–679
- 746 49. 680–681
- 747 50. 682–683
- 748 51. 684–685
- 749 52. 686–687
- 750 53. 688–689
- 751 54. 690–691
- 752 55. 692–693
- 753 56. 694–695
- 754 57. 696–697
- 755 58. 698–699
- 756 59. 700–701
- 757 60. 702–703
- 758 61. 704–705
- 759 62. 706–707
- 760 63. 708–709
- 761 64. 710–711
- 762 65. 712–713
- 763 66. 714–715
- 764 67. 716–717
- 765 68. 718–719
- 766 69. 720–721
- 767 70. 722–723
- 768 71. 724–725
- 769 72. 726–727
- 770 73. 728–729
- 771 74. 730–731
- 772 75. 732–733
- 773 76. 734–735
- 774 77. 736–737
- 775 78. 738–739

- 740 40. Himanen JP, Henkemeyer M, Nikolov DB (1998) Crystal structure of the ligand-binding
741 domain of the receptor tyrosine kinase ephb2. *Nature* 396(6710):486–491.
742 41. Maitre JL, Heisenberg CP (2013) Three functions of cadherins in cell adhesion. *Current Biology* 23(14):R626–R633.
743 42. Berger B, Waterman MS, Yu YW (2020) Levenshtein distance, sequence comparison
744 and biological database search. *IEEE transactions on information theory* 67(6):3287–
745 3294.
746 43. Gulrajani I, Lopez-Paz D (2020) In search of lost domain generalization. *arXiv preprint
747 arXiv:2007.01434*.
748 44. Huang K, et al. (2021) DeepPurpose: a deep learning library for drug–target interaction
749 prediction. *Bioinformatics* 36(22–23):5545–5547.
750 45. Ramsundar B (2018) Ph.D. thesis (Stanford University).
751 46. Hie B, Zhong ED, Berger B, Bryson B (2021) Learning the language of viral evolution
752 and escape. *Science* 371(6526):284–288.
753 47. Littmann M, Heinzinger M, Dallago C, Weissenow K, Rost B (2021) Protein embed-
754 dings and deep learning predict binding residues for various ligand classes. *Scientific
755 Reports* 11(1):1–15.
756 48. Russel D, et al. (2012) Putting the pieces together: integrative modeling platform
757 software for structure determination of macromolecular assemblies. *PLoS biology*
758 10(1):e1001244.
759 49. Skolnick J, Zhou H (2022) Implications of the essential role of small molecule ligand
760 binding pockets in protein–protein interactions. *The Journal of Physical Chemistry B*
761 126(36):6853–6867.
762 50. Hie BL, Yang KK, Kim PS (2021) Evolutionary velocity with protein language models.
763 *bioRxiv*.
764 51. Hsu C, Nisonoff H, Fannjiang C, Listgarten J (2021) Combining evolutionary and
765 assay-labelled data for protein fitness prediction. *bioRxiv*.
766 52. Jin W, Barzilay R, Jaakkola T (2018) Junction tree variational autoencoder for molec-
767 ular graph generation in *International Conference on Machine Learning*. (PMLR), pp.
768 2323–2332.
769 53. Jin W, Barzilay R, Jaakkola T (2020) Hierarchical generation of molecular graphs using
770 structural motifs in *International Conference on Machine Learning*. (PMLR), pp. 4839–
771 4848.
772 54. Wouters OJ, McKee M, Luyten J (2020) Estimated research and development invest-
773 ment needed to bring a new medicine to market, 2009–2018. *Jama* 323(9):844–853.
774 55. Van Norman GA (2016) Drugs, devices, and the fda: part 1: an overview of approval
775 processes for drugs. *JACC: Basic to Translational Science* 1(3):170–179.
776 56. Sabe VT, et al. (2021) Current trends in computer aided drug design and a highlight
777 of drugs discovered via computational techniques: A review. *European Journal of
778 Medicinal Chemistry* 224:113705.
779 57. Zitnik M, Sosić R, Maheshwari S, Leskovec J (2018) BioSNAP Datasets: Stanford
780 biomedical network dataset collection (<http://snap.stanford.edu/biodata>).
781 58. Liu T, Lin Y, Wen X, Jorissen RN, Gilson MK (2007) BindingDB: a web-accessible
782 database of experimentally determined protein–ligand binding affinities. *Nucleic acids
783 research* 35(suppl_1):D198–D201.
784 59. Davis MI, et al. (2011) Comprehensive analysis of kinase inhibitor selectivity. *Nature
785 biotechnology* 29(11):1046–1051.
786 60. Huang H, et al. (2015) Panoramic view of a superfamily of phosphatases through
787 substrate profiling. *Proceedings of the National Academy of Sciences* 112(16):E1974–
788 E1983.
789 61. Martínez-Martínez M, et al. (2017) Determinants and prediction of esterase substrate
790 promiscuity patterns. *ACS chemical biology* 13(1):225–234.
791 62. Yang M, et al. (2018) Functional and informatics analysis enables glycosyltransferase
792 activity prediction. *Nature chemical biology* 14(12):1109–1117.
793 63. Fisher BF, Snodgrass HM, Jones KA, Andorf MC, Lewis JC (2019) Site-selective c–
794 h halogenation using flavin-dependent halogenases identified via family-wide activity
795 profiling. *ACS central science* 5(11):1844–1856.
796 64. Bastard K, et al. (2014) Revealing the hidden functional diversity of an enzyme family.
797 *Nature chemical biology* 10(1):42–49.
798 65. Rives A, et al. (2021) Biological structure and function emerge from scaling unsuper-
799 vised learning to 250 million protein sequences. *Proceedings of the National Academy
800 of Sciences* 118(15).
801 66. Rogers D, Hahn M (2010) Extended-connectivity fingerprints. *Journal of chemical
802 information and modeling* 50(5):742–754.
803 67. Jaeger S, Fulle S, Turk S (2018) Mol2vec: unsupervised machine learning approach
804 with chemical intuition. *Journal of chemical information and modeling* 58(1):27–35.
805 68. Wang H, et al. (2021) Chemical-reaction-aware molecule representation learning.
806 *arXiv preprint arXiv:2109.09888*.
807 69. Glorot X, Bengio Y (2010) Understanding the difficulty of training deep feedforward
808 neural networks in *Proceedings of the thirteenth international conference on artificial
809 intelligence and statistics*. (JMLR Workshop and Conference Proceedings), pp. 249–
810 256.
811 70. Loshchilov I, Hutter F (2017) Decoupled weight decay regularization. *arXiv preprint
812 arXiv:1711.05101*.
813 71. Loshchilov I, Hutter F (2016) Sgdr: Stochastic gradient descent with warm restarts.
814 *arXiv preprint arXiv:1608.03983*.