Drug repurposing significantly reduces drug discovery time.



- Development cycle
 - Traditional Drug
 Development: ~ 13 years[1].
 - Drug repurposing: ~6 years.

- Orange physical or *in-silico* experiments.
- Green animal and human experiments.

Data from Ref [1], illustration by the presenter

[1] S. Pushpakom et al, Nat. Rev. Drug Dis., vol. 18, no. 1, Jan. 2019

DTI (Drug-Target Interaction)

DTI Prediction – (*in silico*-base approach)



Traditional Models Rely on Complex and Rare Drug/Protein Spatial information

Fast but inaccurate

Traditional Machine Learning Methods: Using human selecting features to do the prediction, the precision of prediction are not sufficient for finding potential drugs.

Accurate but limited/rare

3D-Structure-Based Models: High accuracy but the model may cannot be deployed into real-life situation.

Research Question: Can we use simpler input (drug/protein) information to make the model both fast and accurate?

	FNN	SVM	RF	KNN
StaticF	0.687 ± 0.131	0.668 ± 0.128	0.665 ± 0.125	0.624 ± 0.120
SemiF	0.743 ± 0.124	0.704 ± 0.128	0.701 ± 0.119	0.660 ± 0.119
ECFP6	0.724 ± 0.125	0.715 ± 0.127	0.679 ± 0.128	0.669 ± 0.121
DFS8	0.707 ± 0.129	0.693 ± 0.128	0.689 ± 0.120	0.648 ± 0.120
ECFP6 + ToxF	0.731 ± 0.126	0.722 ± 0.126	0.711 ± 0.131	0.675 ± 0.122

Tranditional ML model prediction precision on binding affinity MSE Data from references [1] https://doi.org/10.1186/s13321-017-0209-z

[2] https://doi.org/10.1039/C8SC00148K



Images taken from references [3] https://doi.org/10.3389/fgene.2020.607824 [4] https://arxiv.org/abs/1510.02855

Our proposed model DeepLPI (Ligand-Protein Interaction)



Illustrations by the presenter

Our Model: Treat drug and protein as language and adapt NLP techniques for DTI.

Traditional	 use drug/protein spatial information, complex
Our model	• use drug formula/protein sequence, simple

DeepLPI model overview (best after 9 versions)

Common Setup

Dropout	0.3	
Weight initialization	Kaiming	
Optimizer	Adam	
Batch size	256	
Learning rate (LR)	0.001	
	0.0001	
LR decay rate	0.8	

Input

Drug Molecule -- SMILES format



Target Protein -- FASTA format

;LCBO - Prolactin precursor - Bovine ; a sample sequence in FASTA format MDSKGSSQKGSRLLLLLVVSNLLLCQGVVSTPVCPNGPGNCQVSLRDLFDRAVMVSHYIHDLSS EMFNEFDKRYAQCKGFITMALNSCHTSSLPTPEDKEQAQQTHHEVLMSLILGLLRSWNDPLYHL VTEVRGMKGAPDAILSRAIEIEEENKRLLEGMEMIFGQVIPGAKETEPYPVWSGLPSLQTKDED ARYSAFYNLLHCLRRDSSKIDTYLKLLNCRIIYNNNC*

Image and Data from [1] https://en.wikipedia.org/wiki/FASTA_format [2] https://en.wikipedia.org/wiki/Simplified_molecular-input_line-entry_system

Loss Function = Binary Cross Entropy + L2 regularization

$$\text{Loss} = -\frac{1}{N} \sum_{i=1}^{N} [y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)] + \underbrace{\alpha \|W\|_2^2}_{\text{L2-norm regularization}}$$



Molecule Embedding by Mol2Vec



Protein Embedding by ProSE



ProSE maximizes global similary and residue contacts between proteins

Target sequence

0

sequence

С

QLE

Predicte

contacts

We tested a few embedding methods, including AllenNLP, SeqVec, and etc, and finally decided to use **ProSE** to embed the protein sequence in FASTA format and obtained both a 6165-dimension embedded vector, and a 100-dimension embedded vector.

We test both embedding for our model because the longer vecter might contain more useful information but may potentially lead to overfitting, while the shorter vector might lose some information but could suppress overfitting and run at faster speed.

Therefore, we have two models: DeepLPI-6165 and DeepLPI-100.



Image from Ref. [1]

S N P

в

[1] T. Bepler and B. Berger, "Learning the protein language: Evolution, structure, and function," Cell Systems, vol. 12, no. 6, (2021)

Train Data Selection

No duplicates, High Confidence Experiment, Balanced Label

All Images on this page created by the presenter

Data stats and processing for BindingDB dataset. Davis dataset follow a similar pre-processing and stats.







Independent Testing Results On *BindingDB* dataset with DeepLPI-6165

All Images on this page created by the presenter



Independent Testing Results

On Davis dataset with DeepLPI-6165

All Images on this page created by the presenter



Performance Comparison

BindingDB	AUROC	Sensitivity	Specificity	PPV	NPV	Remark		
Our 6165	0.790	0.684	0.773	0.671	0.783			
Our 100	0.751	0.541	0.818	0.668	0.725			
DeepCDA	0.448	0.000	1.000	Nan	0.596	All nonbinding		
Transfer to COVID Data								
Our 6165	0.610	0.538	0.576	0.110	0.928			
Our 100	0.475	0.692	0.332	0.092	0.912			
DeepCDA	0.400	0.000	1.0	nan	0.911	All nonbinding		
Davis	AUROC	Sensitivity	Specificity	PPV	NPV	Remark		
Davis Our 6165	AUROC 0.791	Sensitivity 0.661	Specificity 0.789	PPV 0.132	NPV 0.980	Remark		
Davis Our 6165 Our 100	AUROC 0.791 0.673	Sensitivity 0.661 0.395	Specificity 0.789 0.820	PPV 0.132 0.439	NPV 0.980 0.791	Remark		
Davis Our 6165 Our 100 DeepCDA	AUROC 0.791 0.673 0.741	Sensitivity 0.661 0.395 0.511	Specificity 0.789 0.820 0.813	PPV 0.132 0.439 0.495	NPV 0.980 0.791 0.823	Remark		
Davis Our 6165 Our 100 DeepCDA	AUROC 0.791 0.673 0.741	Sensitivity 0.661 0.395 0.511 Transfe	Specificity 0.789 0.820 0.813 r to COVID Data	PPV 0.132 0.439 0.495	NPV 0.980 0.791 0.823	Remark		
Davis Our 6165 Our 100 DeepCDA Our 6165	AUROC 0.791 0.673 0.741 0.534	Sensitivity 0.661 0.395 0.511 Transfe 0.000	Specificity 0.789 0.820 0.813 r to COVID Data 1.000	PPV 0.132 0.439 0.495 nan	NPV 0.980 0.791 0.823 0.911	Remark All nonbinding		
Davis Our 6165 Our 100 DeepCDA Our 6165 Our 100	AUROC 0.791 0.673 0.741 0.534 0.482	Sensitivity 0.661 0.395 0.511 Transfe 0.000 0.040	Specificity 0.789 0.820 0.813 r to COVID Data 1.000 1.000	PPV 0.132 0.439 0.495 nan 1	NPV 0.980 0.791 0.823 0.911 0.914	Remark All nonbinding		

Result Summary

Use 1-dimension drug SMILES and protein sequence as input

Use NLP technique to treat drug and protein

Model DeepLPI-6165 performance in classification on BindingDB dataset is **76% better** than the state-of-the-art DeepCDA model

Model DeepLPI-6165 performance in classification on transferability is 25% (Davis to Covid) and 50% (BindingDB to Covid) better than DeepCDA