

ChemicalX: A Deep Learning Library for Drug Pair Scoring

Turing Workshop on Open-Source AI Software for Healthcare

Benedek Rozemberczki, Isomorphic Laboratories

The presentation covers ideas and results from:

A Unified View of Relational Deep Learning for Drug Pair Scoring. Rozemberczki et al. IJCAI 2022.

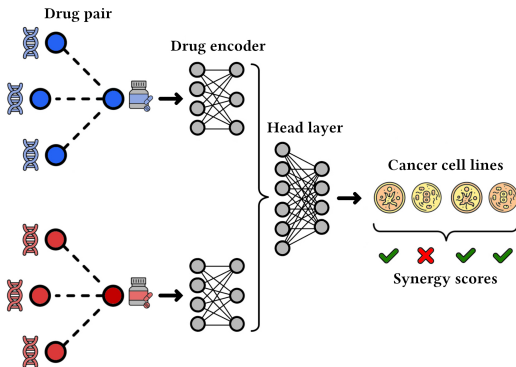
ChemicalX: A Deep Learning Library for Drug Pair Scoring. Rozemberczki et al. KDD 2022.

MOOMIN: Deep Molecular Omics Network for Anti-Cancer Drug Combination Therapy. Rozemberczki et al. CIKM 2022.

1. Drug Pair Scoring
2. A System for Repurposing
3. Experimental Results
4. Conclusions

Drug Pair Scoring

What is drug pair scoring?



- ▶ The drug pair scoring task
- ▶ Motivation
 - ▶ Timelines
 - ▶ Material costs (assays)
 - ▶ Labour costs
 - ▶ Tractability
- ▶ Application domains - tasks
 - ▶ Interaction
 - ▶ Polypharmacy side effect
 - ▶ Synergy
- ▶ Multi-objective optimization

How do we represent drugs in drug pair scoring?

Let $G = (\mathcal{V}, \mathcal{R}, \mathcal{E})$ be a heterogeneous biological graph with drug entities $\mathcal{D} \subset \mathcal{V}$.

1. Molecular - low level encoders:

$$\mathbf{h}^d = f_{\Theta_D}(\mathcal{M}^d)$$

2. Systems biology based - high level encoders:

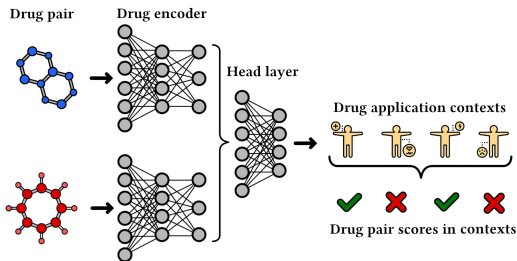
$$\mathbf{h}^d = \text{AGGREGATE}(\{\Theta^u, \forall u \in \mathcal{N}(d)\})$$

3. Hierarchical encoders (structure and systems view):

$$\mathbf{h}^d = \text{AGGREGATE}(\{f_{\Theta_D}(\mathcal{M}^u), \forall u \in \mathcal{N}(d)\})$$

| Drug Pair Scoring ○○○●○○○○○ | | A System for Repurposing ○○○○○○○○○ | | | Experimental Results ○○○○○ | | | Conclusions ○○○ | | References |
|--------------------------------|------------------|---------------------------------------|-----------|--------------|-------------------------------|---------|---------------|--------------------|----------|------------|
| Task | Method | Model view | Induction | Entity Types | | | Drug Features | | | Generic |
| | | | | Drug | Protein | Disease | SMILES | Graph | Geometry | |
| Polypharmacy | DECAGON | Higher | | ● | ● | | | | | |
| | KBLRN | Higher | | ● | ● | | | | | |
| | SDHINE | Higher | | ● | ● | ● | | | | |
| | ESP | Higher | | ● | ● | | | | | |
| | MHCADDI | Lower | ● | | | | | ● | | |
| | TIP | Higher | | ● | ● | | | | | |
| Interaction | DeepCCI | Lower | ● | | | | ● | | | ● |
| | MVGAE | Hierarchical | ● | ● | ● | ● | ● | | | |
| | DeepDDI | Lower | ● | | | | ● | | | |
| | D ³ I | Hierarchical | | ● | | | | | | ● |
| | MR-GNN | Lower | ● | | | | | ● | | |
| | SkipGNN | Higher | | ● | ● | ● | | | | |
| | CASTER | Lower | ● | | | | ● | | | |
| | DeepDrug | Lower | ● | | | | ● | | ● | |
| | GoGNN | Hierarchical | | ● | | | | ● | | |
| | DPDDI | Lower | ● | | | | | | | ● |
| | KGNN | Higher | | ● | ● | ● | | | | |
| | BiGNN | Hierarchical | ● | ● | ● | | | ● | | |
| | MIRACLE | Hierarchical | ● | ● | | | | ● | | |
| | EPGCN-DS | Lower | ● | | | | | | | |
| | SumGNN | Higher | | ● | ● | ● | | | | |
| | DANN-DDI | Higher | | ● | ● | | | | | |
| | SSI-DDI | Lower | | | | | | ● | | |
| | MTDDI | Hierarchical | | ● | ● | ● | | | | ● |
| | MUFFIN | Hierarchical | | ● | ● | ● | ● | | | |
| | DDIAAE | Higher | | ● | ● | ● | | | | |
| RWGCN | Higher | | ● | ● | ● | | | | | |
| SmileGNN | Hierarchical | | ● | ● | | ● | | | | |
| GCN-BMP | Lower | ● | ● | | | ● | | | | |

Why molecule level models?



- Induction
- Attribution (explanation)
- Pre-training
- Transfer learning
- Structural ablation

How do we represent biological contexts and score pairs?

Given \mathcal{C} a set of biological contexts the representation of $c \in \mathcal{C}$ is:

$$\mathbf{h}_c = f_{\Theta_C}(\mathbf{x}_c).$$

We score the pair $d, d' \in \mathcal{D}$ in the context $c \in \mathcal{C}$ with:

$$\hat{\mathbf{y}}^{d,d',c} = f_{\Theta_H}(\mathbf{h}^d, \mathbf{h}^{d'}, \mathbf{h}^c).$$

The loss for the pair in the context is defined as:

$$\ell_{d,d',c} = \ell(\hat{\mathbf{y}}^{d,d',c}, \mathbf{y}^{d,d',c}).$$

The parametric functions $f_{\Theta_D}(\cdot)$, $f_{\Theta_C}(\cdot)$, and $f_{\Theta_H}(\cdot)$ can be trained jointly by minimizing the cost accumulated from data point level losses.

How do we conceptualize these models?

Data: $\mathcal{X}_{\mathcal{D}}$ - Drug feature set.

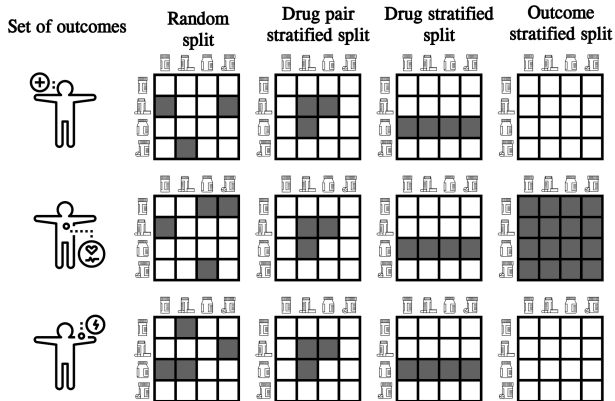
$\mathcal{X}_{\mathcal{C}}$ - Context feature set.

\mathcal{B} - Labeled drug pair - context batch.

Result: \mathcal{L} - The cost for the batch.

```
1  $\mathcal{L} \leftarrow 0$ 
2 for  $(d, d', c, y^{d,d',c}) \in \mathcal{B}$  do
3    $\mathbf{h}^d \leftarrow f_{\Theta_D}(\mathbf{x}^d, \mathcal{G}^d, \mathbf{X}_N^d, \mathbf{X}_E^d)$  // Compute drug representation for  $d \in \mathcal{D}$ .
4    $\mathbf{h}^{d'} \leftarrow f_{\Theta_D}(\mathbf{x}^{d'}, \mathcal{G}^{d'}, \mathbf{X}_N^{d'}, \mathbf{X}_E^{d'})$  // Compute drug representation for  $d' \in \mathcal{D}$ .
5    $\mathbf{h}^c \leftarrow f_{\Theta_C}(\mathbf{x}^c)$  // Compute context representation for  $c \in \mathcal{C}$ .
6    $\hat{y}^{d,d',c} \leftarrow f_{\Theta_H}(\mathbf{h}^d, \mathbf{h}^{d'}, \mathbf{h}^c)$  // Score based on the representations.
7    $\mathcal{L} \leftarrow \mathcal{L} + \ell(y^{d,d',c}, \hat{y}^{d,d',c})$  // Add loss to the accumulated cost.
8 end
```

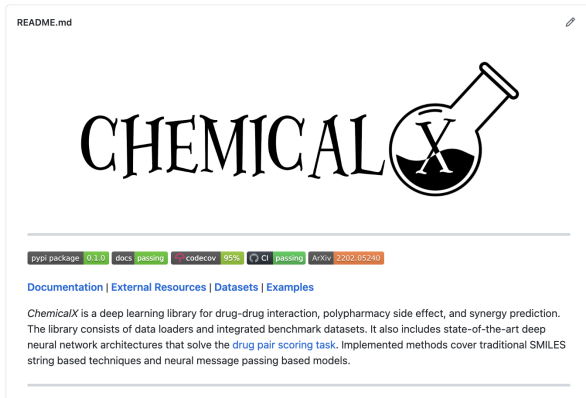
How do we design evaluation regimes?



| Model | Evaluation metric | | | | | |
|------------------|-------------------|-------|-----------|--------|----------|----------------|
| | AUPRC | AUROC | Precision | Recall | Accuracy | F ₁ |
| DECAGON | ● | ● | ● | | | |
| KBLRN | ● | ● | ● | | | |
| SDHINE | | ● | ● | | | |
| ESP | ● | ● | ● | | | |
| MHCADDI | | ● | | | | |
| TIP | ● | ● | ● | | | |
| DeepCCI | | ● | | | ● | ● |
| MVGAE | ● | ● | | | | |
| DeepDDI | | | ● | ● | ● | ● |
| D ³ I | | | ● | ● | ● | ● |
| MR-GNN | | ● | | ● | | ● |
| SkipGNN | ● | | | | | ● |
| CASTER | ● | ● | | | | ● |
| DeepDrug | ● | ● | | | | |
| GoGNN | | ● | ● | | | |
| DPDDI | ● | ● | ● | ● | ● | ● |
| KGNN | ● | ● | | | ● | ● |
| BiGNN | | ● | ● | | | ● |
| MIRACLE | ● | ● | | | | ● |
| EPGCN-DS | | ● | | | | ● |
| SumGNN | ● | ● | ● | | | ● |
| DANN-DDI | ● | ● | ● | ● | ● | ● |
| SSI-DDI | ● | ● | | | ● | |
| MTDDI | ● | ● | ● | ● | ● | ● |
| MUFFIN | | ● | ● | ● | ● | ● |
| DDIAAE | ● | ● | ● | | | |
| RWGCN | ● | ● | | | | |
| SmileGNN | | ● | | | ● | ● |

A System for Repurposing

What is ChemicalX?



<https://github.com/AstraZeneca/chemicalx/>

How did you engineer it?

- ▶ Documentation
- ▶ Unit and integration tests with coverage reports
- ▶ Tutorials
- ▶ Example datasets
- ▶ Continuous integration
- ▶ Linting, type hinting and docstrings

What are the competitors of ChemicalX?

| Library | Year | Backend | Drug Domain | Pair Scoring |
|------------------|------|----------|-------------|--------------|
| PyG [7] | 2018 | PT | ✗ | ✗ |
| DGL [24] | 2019 | PT/TF/MX | ✗ | ✗ |
| StellarGraph [5] | 2019 | TF | ✗ | ✗ |
| DeepChem [19] | 2019 | TF | ✓ | ✗ |
| CHChem [16] | 2019 | CH | ✓ | ✗ |
| Jraph [8] | 2020 | JAX | ✗ | ✗ |
| Spektral [9] | 2020 | TF | ✗ | ✗ |
| DIG [15] | 2021 | PT | ✗ | ✗ |
| TorchDrug [28] | 2021 | PT | ✓ | ✗ |
| CogDL [3] | 2021 | PT | ✗ | ✗ |
| TFG [10] | 2021 | TF | ✗ | ✗ |
| DGL-LS [13] | 2021 | PT | ✓ | ✗ |
| Our Work | 2022 | PT | ✓ | ✓ |

What is included in ChemicalX?

| Model | Year | Domain | Encoder |
|------------------|------|--------------|-------------|
| DeepDDI [20] | 2018 | Interaction | Feedforward |
| DeepSynergy [18] | 2018 | Synergy | Feedforward |
| MHCADDI [6] | 2019 | Polypharmacy | GAT |
| MR-GNN [25] | 2019 | Interaction | GCN |
| CASTER [12] | 2019 | Interaction | Feedforward |
| SSI-DDI [17] | 2020 | Interaction | GAT |
| EPGCN-DS [21] | 2020 | Interaction | GCN |
| DeepDrug [2] | 2020 | Interaction | GCN |
| GCN-BMP [4] | 2020 | Interaction | GCN |
| DeepDDS [23] | 2021 | Synergy | GCN or GAT |
| MatchMaker [1] | 2021 | Synergy | Feedforward |

How do we load the dataset?

```
1 from chemicalx.data import DrugCombDB, BatchGenerator
2
3 loader = DrugCombDB()
4
5 context_set = loader.get_context_features()
6 drug_set = loader.get_drug_features()
7 triples = loader.get_labeled_triples()
8
9 generator = BatchGenerator(batch_size=1024,
10                             context_features=True,
11                             drug_features=True,
12                             drug_molecules=False,
13                             context_feature_set=context_set,
14                             drug_feature_set=drug_set,
15                             labeled_triples=triples)
```

How do we train a model?

```
1 import torch
2 from chemicalx.models import DeepSynergy
3
4 model = DeepSynergy(context_channels=112,
5                     drug_channels=256)
6
7 optimizer = torch.optim.Adam(model.parameters())
8 model.train()
9 loss = torch.nn.BCELoss()
10
11 for epoch in range(200):
12     for batch in generator:
13         optimizer.zero_grad()
14         prediction = model(batch.context_features,
15                           batch.drug_features_left,
16                           batch.drug_features_right)
17         loss_value = loss(prediction, batch.labels)
18         loss_value.backward()
19         optimizer.step()
```

How do we score a dataset with the model?

```
1 import pandas as pd
2 from chemicalx.data import HAEM
3
4 model.eval()
5
6 loader = HAEM()
7
8 generator.labeled_triples = loader.get_labeled_triples()
9
10 predictions = []
11 for batch in generator:
12     prediction = model(batch.context_features,
13                        batch.drug_features_left,
14                        batch.drug_features_right)
15     prediction = prediction.detach().cpu().numpy()
16     identifiers = batch.identifiers
17     identifiers["prediction"] = prediction
18     predictions.append(identifiers)
19 predictions = pd.concat(predictions)
```

How can you BYOD (Bring Your Own Data)?

Training and scoring for specific pairs would only need these things:

- ▶ Drug pairs with biological/chemical contexts and labels.
- ▶ Context set with context identifier keys and context feature vector values.
- ▶ Drug set with SMILES strings and molecule level features.

Checkout the following link for an example:

<https://github.com/AstraZeneca/chemicalx/tree/main/dataset/drugbankddi>

Experimental Results

What are the datasets integrated?

Table 1: Datasets available in ChemicalX in the domain of the pair scoring task and the number of drugs ($|\mathcal{D}|$), administration contexts ($|\mathcal{C}|$), and labeled triples ($|\mathcal{Y}|$).

| Dataset | Task | $ \mathcal{D} $ | $ \mathcal{C} $ | $ \mathcal{Y} $ |
|--------------------------|--------------|-----------------|-----------------|-----------------|
| TWOSIDES [22] | Polypharmacy | 644 | 10 | 499,582 |
| Drugbank DDI [20] | Interaction | 1,706 | 86 | 383,496 |
| DrugComb [26, 27] | Synergy | 4,146 | 288 | 659,333 |
| DrugCombDB [14] | Synergy | 2,956 | 112 | 191,391 |
| OncolyPharm [11] | Synergy | 38 | 39 | 23,052 |

How about predictive performance?

Table 2: The predictive performance of (some) models in *ChemicalX* on TWOSIDES [22]. Feed-forward encoder based architectures are noted with a ■.

| | | AUROC | AUPR | F ₁ |
|------------------|---|-------------|-------------|----------------|
| DeepDDI [20] | ■ | .929 ± .001 | .907 ± .001 | .848 ± .009 |
| DeepSynergy [18] | ■ | .940 ± .001 | .919 ± .001 | .887 ± .001 |
| MR-GNN [25] | | .937 ± .002 | .917 ± .001 | .875 ± .002 |
| SSI-DDI [17] | | .823 ± .002 | .800 ± .003 | .756 ± .001 |
| EPGCN-DS [21] | | .855 ± .003 | .834 ± .002 | .785 ± .004 |
| DeepDrug [2] | ■ | .923 ± .004 | .904 ± .002 | .857 ± .002 |
| GCN-BMP [4] | | .709 ± .003 | .694 ± .002 | .592 ± .003 |
| DeepDDS [23] | | .915 ± .002 | .898 ± .002 | .839 ± .003 |
| MatchMaker [1] | ■ | .912 ± .002 | .892 ± .001 | .849 ± .001 |

How about training time?

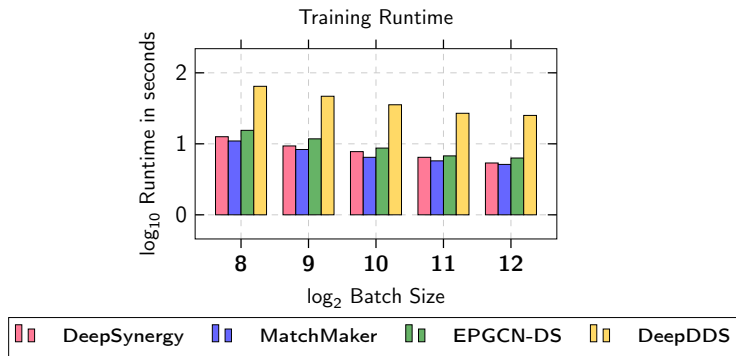


Figure 1: The average runtime of doing an epoch on DrugBankDDI [11].

How about inference time?

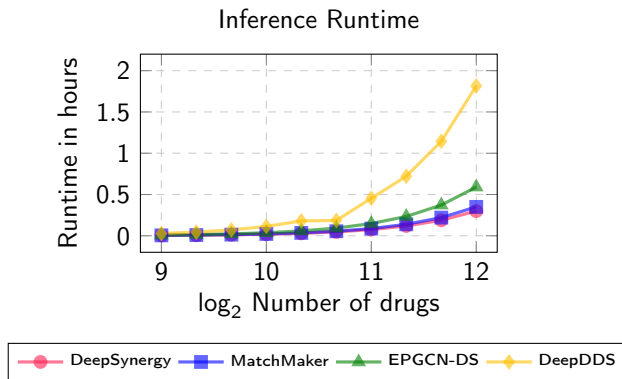


Figure 2: The average runtime of doing a scoring pass for all combinations in DrugBankDDI [11].

Conclusions

What are the main takeaways?

Having impact!

- ▶ ChemicalX is used for targeting haematological malignancies.
- ▶ Early oncology scientists in AstraZeneca are using ChemicalX self-service.

Having fun!

- ▶ Internal and external collaborations.
- ▶ Skills elevated for people in AZ who are not core machine learning.

What could be improved?

- ▶ Splits that take scaffolds into account.
- ▶ Geometric graph encoders.
- ▶ Scaling with locality sensitive hashing.

Thank you for the kind attention!

- [1] Kuru Halil Brahim, Oznur Tastan, and Ercument Cicek. 2021. MatchMaker: A Deep Learning Framework for Drug Synergy Prediction. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* (2021).
- [2] Xusheng Cao, Rui Fan, and Wanwen Zeng. 2020. DeepDrug: A General Graph-Based Deep Learning Framework for Drug Relation Prediction. *bioRxiv* (2020).
- [3] Yukuo Cen, Zhenyu Hou, Yan Wang, Qibin Chen, et al. 2021. CogDL: An Extensive Toolkit for Deep Learning on Graphs. (2021).
- [4] Xin Chen, Xien Liu, and Ji Wu. 2020. GCN-BMP: Investigating Graph Representation Learning for DDI Prediction Task. *Methods* 179 (2020), 47–54. Interpretable machine learning in bioinformatics.
- [5] CSIRO's Data61. 2018. StellarGraph Machine Learning Library. <https://github.com/stellargraph/stellargraph>.

- [6] Andreea Deac, Yu-Hsiang Huang, Petar Velickovic, Pietro Liò, and Jian Tang. 2019. Drug-Drug Adverse Effect Prediction with Graph Co-Attention. *ICML Workshop on Computational Biology* (2019).
- [7] Matthias Fey and Jan E. Lenssen. 2019. Fast Graph Representation Learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*.
- [8] Jonathan Godwin, Thomas Keck, Peter Battaglia, Victor Bapst, Thomas Kipf, et al. 2020. *Jraph: A Library for Graph Neural Networks in Jax*.
- [9] Daniele Grattarola and Cesare Alippi. 2021. Graph Neural Networks in TensorFlow and Keras with Spektral. *IEEE Computational Intelligence Magazine* 16, 1 (2021), 99–106.
- [10] Jun Hu, Shengsheng Qian, Quan Fang, et al. 2021. Efficient Graph Deep Learning in TensorFlow with TF Geometric. *arXiv preprint 2101.11552* (2021).

- [11] Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, et al. 2021. Therapeutics Data Commons: Machine Learning Datasets and Tasks for Drug Discovery and Development. In *35th Conference on Neural Information Processing Systems*.
- [12] Kexin Huang, Cao Xiao, Trong Nghia Hoang, Lucas M Glass, and Jimeng Sun. 2020. CASTER: Predicting Drug Interactions with Chemical Substructure Representation. *AAAI* (2020).
- [13] Mufei Li, Jinjing Zhou, Jiajing Hu, Wenxuan Fan, Yangkang Zhang, Yaxin Gu, and George Karypis. 2021. DGL-LifeSci: An Open-Source Toolkit for Deep Learning on Graphs in Life Science. *ACS Omega* 6, 41 (2021), 27233–27238.
- [14] Hui Liu, Wenhao Zhang, Bo Zou, Jinxian Wang, and Yuanyuan Deng. 2020. Drug-CombDB: A Comprehensive Database of Drug Combinations Toward the Discovery of Combinatorial Therapy. *Nucleic acids research* 48 (2020), 871–881.
- [15] Meng Liu, Youzhi Luo, Limei Wang, et al. 2021. DIG: A Turnkey Library for Diving into Graph Deep Learning Research. *Journal of Machine Learning Research* 22, 240 (2021), 1–9.

- [16] Abe Motoki, Mihai Mororiu, Tomoya Otabi, Kenshin Abe, and Others. 2017. *Chainer Chemistry: A Library for Deep Learning in Biology and Chemistry*. <https://github.com/chainer/chainer-chemistry>
- [17] Arnold K Nyamabo, Hui Yu, and Jian-Yu Shi. 2021. SSI-DDI: Substructure–Substructure Interactions for Drug–Drug Interaction Prediction. *Briefings in Bioinformatics* (2021).
- [18] Kristina Preuer, Richard PI Lewis, Sepp Hochreiter, Andreas Bender, Krishna C Bulusu, and Günter Klambauer. 2018. DeepSynergy: Predicting Anti-Cancer Drug Synergy with Deep Learning. *Bioinformatics* 34, 9 (2018), 1538–1546.
- [19] Bharath Ramsundar, Peter Eastman, Patrick Walters, Vijay Pande, et al. 2019. *Deep Learning for the Life Sciences*. O'Reilly Media.
- [20] Jae Yong Ryu, Hyun Uk Kim, and Sang Yup Lee. 2018. Deep Learning Improves Prediction of Drug–Drug and Drug–Food Interactions. *Proceedings of the National Academy of Sciences* 115, 18 (2018), E4304–E4311.

- [21] Mengying Sun, Fei Wang, Olivier Elemento, and Jiayu Zhou. 2020. Structure-Based Drug-Drug Interaction Detection via Expressive Graph Convolutional Networks and Deep Sets. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 13927–13928.
- [22] Nicholas P Tatonetti, P Ye Patrick, Roxana Daneshjou, and Russ B Altman. 2012. Data-Driven Prediction of Drug Effects and Interactions. *Science translational medicine* 4, 125 (2012).
- [23] Jinxian Wang, Wenhao Zhang, Siyuan Shen, Lei Deng, and Hui Liu. 2021. Deep-DDS: Deep Graph Neural Network with Attention Mechanism to Predict Synergistic Drug Combinations. *bioRxiv* (2021).
- [24] Minjie Wang, Lingfan Yu, Da Zheng, et al. 2019. Deep Graph Library: Towards Efficient and Scalable Deep Learning on Graphs. (2019).
- [25] Nuo Xu, Pinghui Wang, Long Chen, Jing Tao, and Junzhou Zhao. 2019. MR-GNN: Multi-Resolution and Dual Graph Neural Network for Predicting Structured Entity Interactions. *Proceedings of IJCAI* (2019).

- [26] Bulat Zagidullin, Jehad Aldahdooh, Shuyu Zheng, Wenyu Wang, Yinyin Wang, et al. 2019. DrugComb: An Integrative Cancer Drug Combination Data Portal. *Nucleic acids research* 47, W1 (2019), W43–W51.
- [27] Shuyu Zheng, Jehad Aldahdooh, Tolou Shadbahr, Yinyin Wang, Dalal Aldahdooh, et al. 2021. DrugComb Update: A More Comprehensive Drug Sensitivity Data Repository and Analysis Portal. *Nucleic Acids Research* (2021).
- [28] Zhaocheng Zhu, Shengchao Liu, Chence Shi, et al. 2021. *TorchDrug: A Powerful and Flexible Machine Learning Platform for Drug Discovery*.