Final Project
# MolDesigner UI and Model Implementation
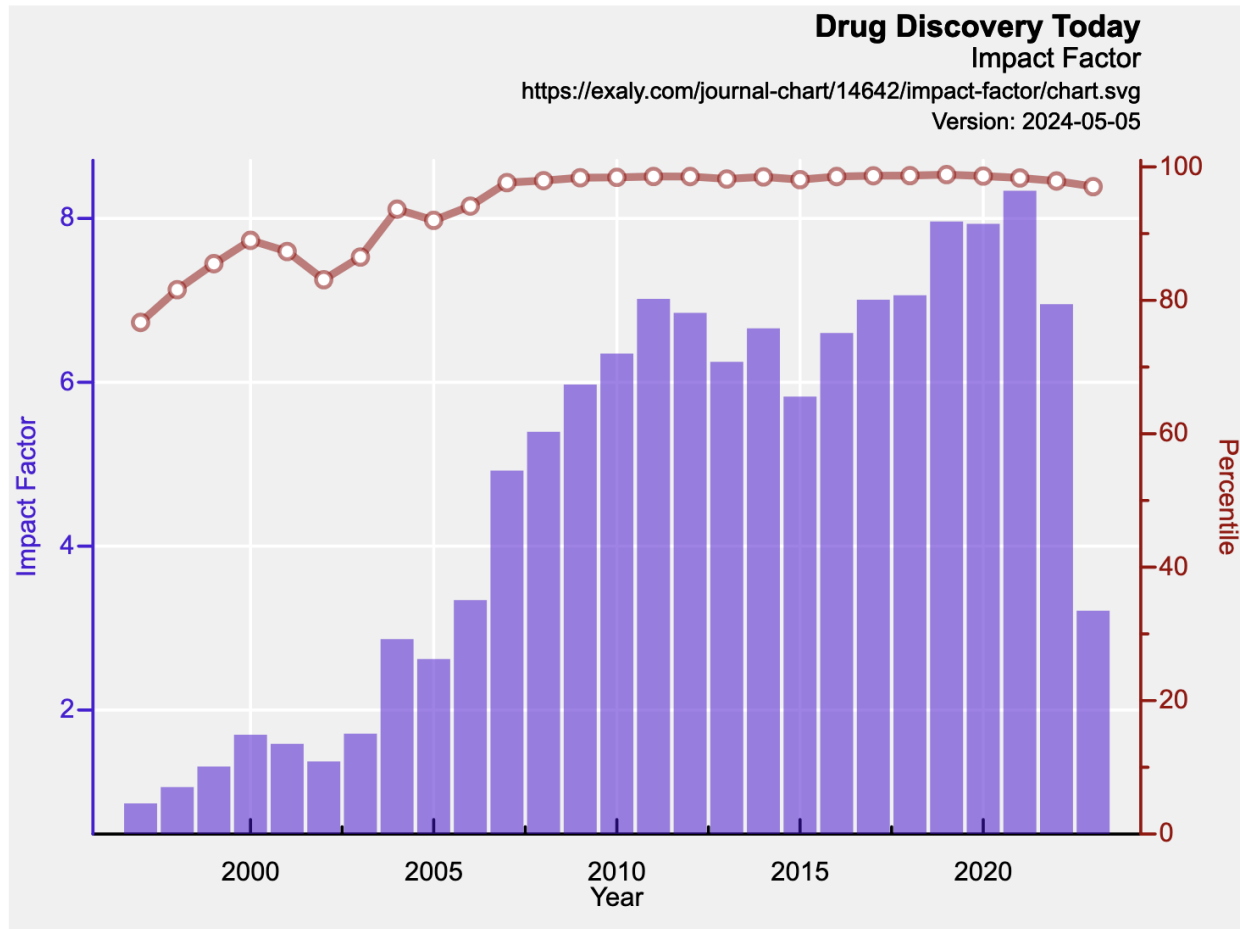
## Kotini, Jessica



CSCI E-104 Advanced Deep Learning, 2024
**Harvard University Extension School**
Prof. Zoran B. Djordjević

# Introduction

- Field of drug-discovery is rapidly evolving

- Need tools for fast visualization of drug-target interaction data

**Drug Discovery Today**
Impact Factor
https://exaly.com/journal-chart/14642/impact-factor/chart.svg
Version: 2024-05-05

The graph shows the changes in the impact factor of **Drug Discovery Today** and its corresponding percentile for the sake of comparison with the entire literature. Impact Factor is the most common scientometric index, which is defined by the number of citations of papers in two preceding years divided by the number of papers published in those years.
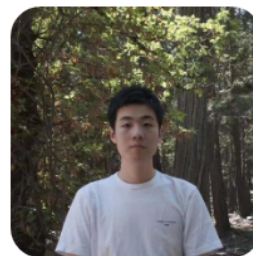
https://exaly.com/journal/14642/drug-discovery-today/

# Goal

1. Demonstrate utility of MolDesigner
2. Integrate my own models into MolDesigner

## kexinhuang12345/
## MolDesigner-Public

MolDesigner: Interactive Design of Efficacious
Drugs with Deep Learning (NeurIPS 2020 Demo)

1 Contributor    0 Issues    13 Stars    5 Forks

# Dataset

# Demo – Installation & Configuration

- Google Colab Pro+ highly recommended with A100 hardware accelerator

# Demo - Methodology

- Followed a SEMMA style framework
  - Loaded, processed, sampled data
  - Explored and viewed data
  - Modelled data
  - Viewed loss curves for each iteration within and across epochs
  - Displayed predictions with Imatinib and ABL1 receptor (common drug-target pair)

| Drug-Target Pairs | Unique Drugs | Unique Target Sequence |
|---|---|---|
| 2385 | 1300 | 438 |



Sample
Explore
Modify
Model
Assess

# Demo Methodology

- Separate convolutional neural network (CNN) architectures for processing drugs and proteins

- Each segment (drug and protein) features 3 convolutional layers
  - Increasing kernel sizes, filter counts
  - Linear layer, with dropout layer before final classification layers

CNN_CNN_BindingDB

```
Classifier(
  (model_drug): CNN(
    (conv): ModuleList(
      (0): Conv1d(63, 32, kernel_size=(4,), stride=(1,))
      (1): Conv1d(32, 64, kernel_size=(6,), stride=(1,))
      (2): Conv1d(64, 96, kernel_size=(8,), stride=(1,))
    )
    (fc1): Linear(in_features=96, out_features=256, bias=True)
  )
  (model_protein): CNN(
    (conv): ModuleList(
      (0): Conv1d(26, 32, kernel_size=(4,), stride=(1,))
      (1): Conv1d(32, 64, kernel_size=(8,), stride=(1,))
      (2): Conv1d(64, 96, kernel_size=(12,), stride=(1,))
    )
    (fc1): Linear(in_features=96, out_features=256, bias=True)
  )
  (dropout): Dropout(p=0.1, inplace=False)
  (predictor): ModuleList(
    (0): Linear(in_features=512, out_features=1024, bias=True)
    (1): Linear(in_features=1024, out_features=1024, bias=True)
    (2): Linear(in_features=1024, out_features=512, bias=True)
    (3): Linear(in_features=512, out_features=1, bias=True)
  )
)
```

# Results

# Results

```
Training at Epoch 10 iteration 0 with loss 2.25397. Total time 0.01055 hours
Validation at Epoch 10 with loss:3.03664, MSE: 2.59127 , Pearson Correlation: 0.77823 with p-value: 1.37E-49 , Concordance Index: 0.78491
--- Go for Testing ---
Testing MSE: 2.4322626316898774 , Pearson Correlation: 0.6884876546785907 with p-value: 2.85E-68 , Concordance Index: 0.7527270801708618
--- Training Finished ---
```

# Learnings/Future Work

- Learnings
  - ML moves quickly, expect deprecation/updates
  - Design flexible code for handling updated data structures
- Future Work
  - Model optimization, better results visualizations
    - train vs validation loss profiles
    - train vs validation MSE profiles
  - Model training with KIBA

# Conclusion

- Fun project, learned a lot about DeepPurpose API
    - Successfully trained models using DeepPurpose
    - Successfully integrated models into MolDesigner UI
- Need for fast responsive and user-friendly DTI tools is growing fast
    - Easy to use for users who are not programmers
    - Facilitate faster drug development cycles by aiding drug discovery

# YouTube URLs

- 2 minutes video URL (short):  https://www.youtube.com/watch?v=BhsMjA6aVAs
- 15 minutes video URL(long):  https://www.youtube.com/watch?v=87XkMBEzmhs