



You

Please design a drug molecule targeting APP, a drug target for Alzheimer's disease.



ChatGPT

As I mentioned earlier, I'm not capable of designing drug molecules, including those targeting APP (amyloid precursor protein) for Alzheimer's disease.



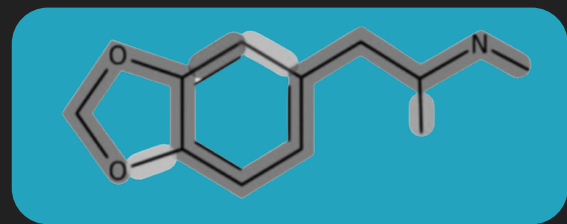
This Study's GPT Model

```
[C][ N][ C][ Branch3 ][ C][ C][ =C][ C][ =C][ C][ Branch3 ][ =O][ =C][ Ring ]  
[#N ][ O][ C][ O][ Ring][ #N ]
```



De Novo Drug Design as GPT Language Modeling

Gavin Ye



The drug discovery process is time consuming

Identify target protein



Drug (candidates) designing

Traditionally costly and time consuming



Drug selection



Entire process takes 10–15 years

...

Drug synthesis, testing, etc.

GPT = Generative ML that specializes in sequential data

Hello my...	name is ...
Q: 1+1=?	A:2
Q: What's ML?	A:ML is ...
...	...
Chat-GPT is...	a large ...

**Input text with
correct response**

Training



GPT model

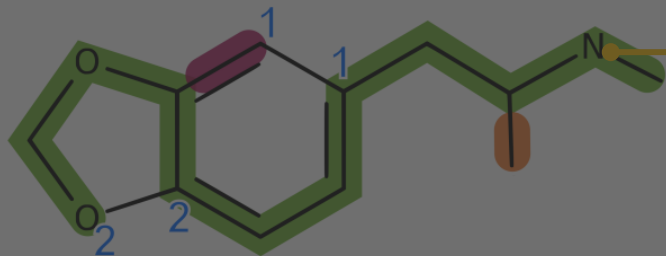


Perform certain task

**Such as:
-Text generation**

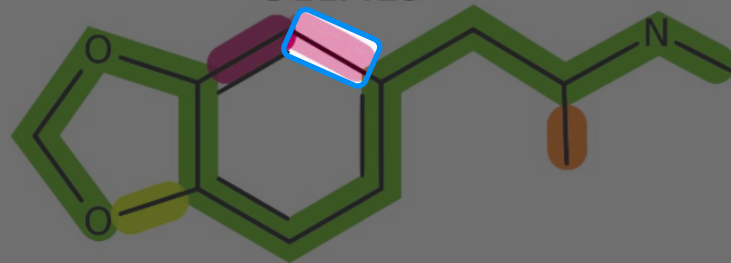
Molecules can be represented using sequences

SMILES



CNC (C) CC 1=CC=C 2 C (=C 1) OCO 2

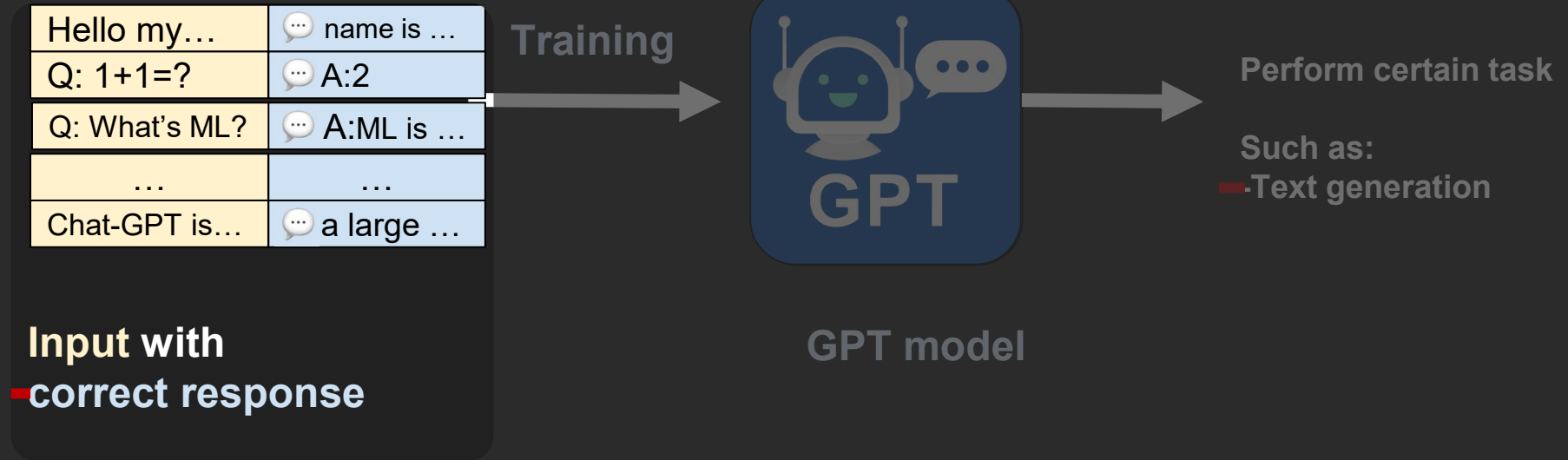
SELFIES



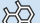
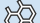
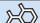
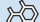
[C] [N] [C] [Branch3] [c] [C] [C] [C] [=C] [C] [=C] [C]
[Branch3] [=O] [=C] [Ring] [#N] [O] [C] [O] [Ring] [#N]

Token

GPT specializes in sequential data



GPT has not been used for designing effective drug molecules in previous studies

<empty>	 CCN...
C12H6...	 O12...
CO...	 CH4...
...	...
CH3C...	 OOH

Input with
~~correct response~~
desired molecule

Training



GPT model

Perform certain task

Such as:

~~Text generation~~

Molecule generation

Brief Recap of Problems: Non GPT models have low validity

Problems:

Sequential representations have been used for non -GPT models for different tasks

(Segler et al., 2018); (Abbasi et al., 2022);



Low Validity

(Abbasi et al., 2022); (Yasonik ., 2020);
(Popova et al., 2018)



Low Novelty or Efficacy

(Gao et al., 2020);



(Frey et al., 2022) 100% Validity
... **Not applied to any specific task
such as drug design.**

This Study:
GPT applied to Drug Design
by optimizing drug efficacy

My study: GPT applied to drug design

Goal:

Train GPT to generate drug-like **molecules with high efficacy** while **maintaining validity** towards treating a disease.



- Valid
- Effective

Methodology

Objectives

1. Evaluate Drug Efficacy

2. Design Drug -Like Molecules

3. Optimize Drug Efficacy

Steps

Reward Modeling



Drug Efficacy
Evaluation Model

Supervised Finetuning

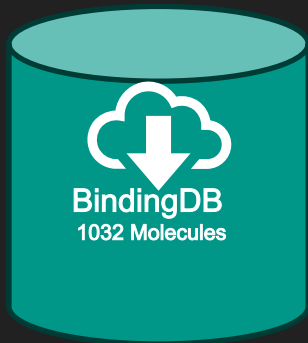


Drug Design
GPT Model

Proximal Policy
Optimization (PPO)



For case study, BindingDB dataset with molecules and experimentally determined drug efficacy values are used



(Liu et al., 2007)

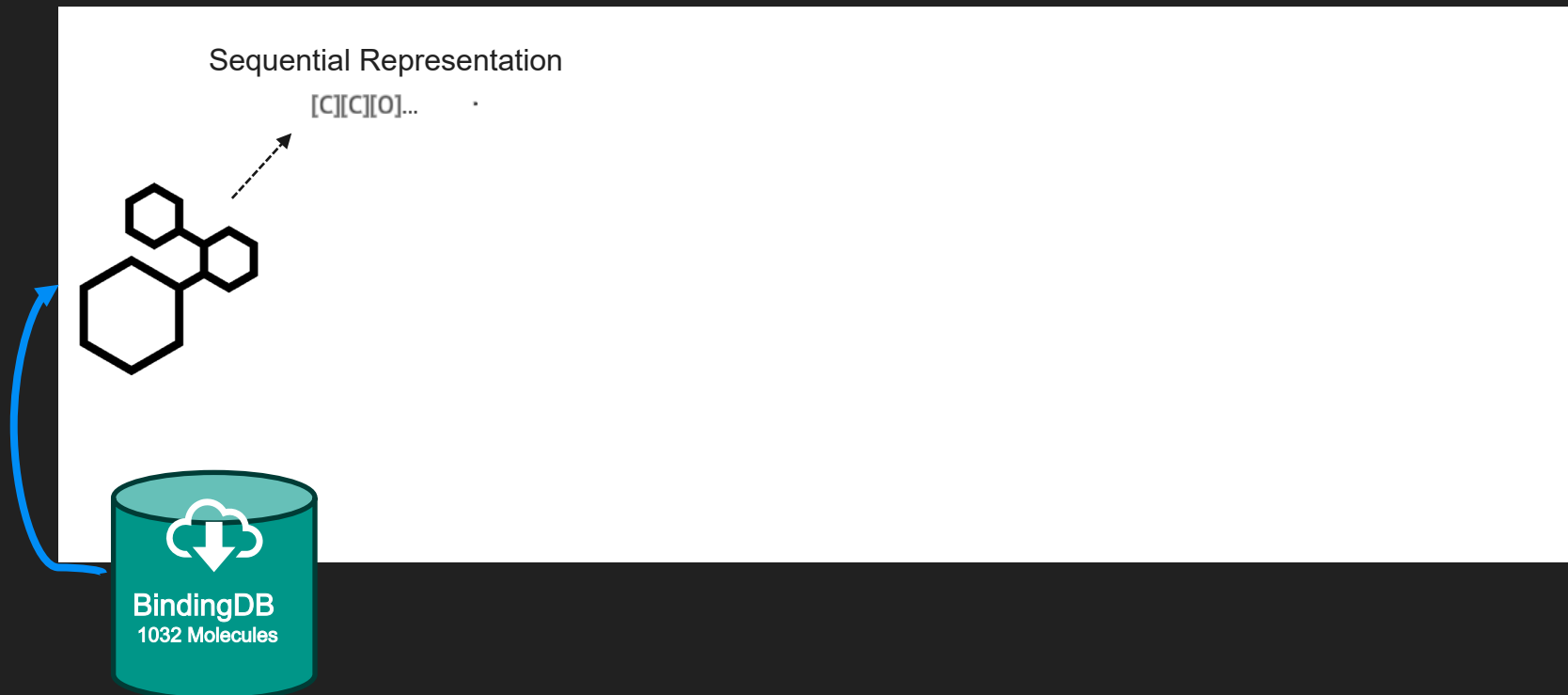
Training Data (bindingdb.org)

- drug candidates
- measured efficacy toward APP (in $^{\circ} p; \mathcal{H}_y$)

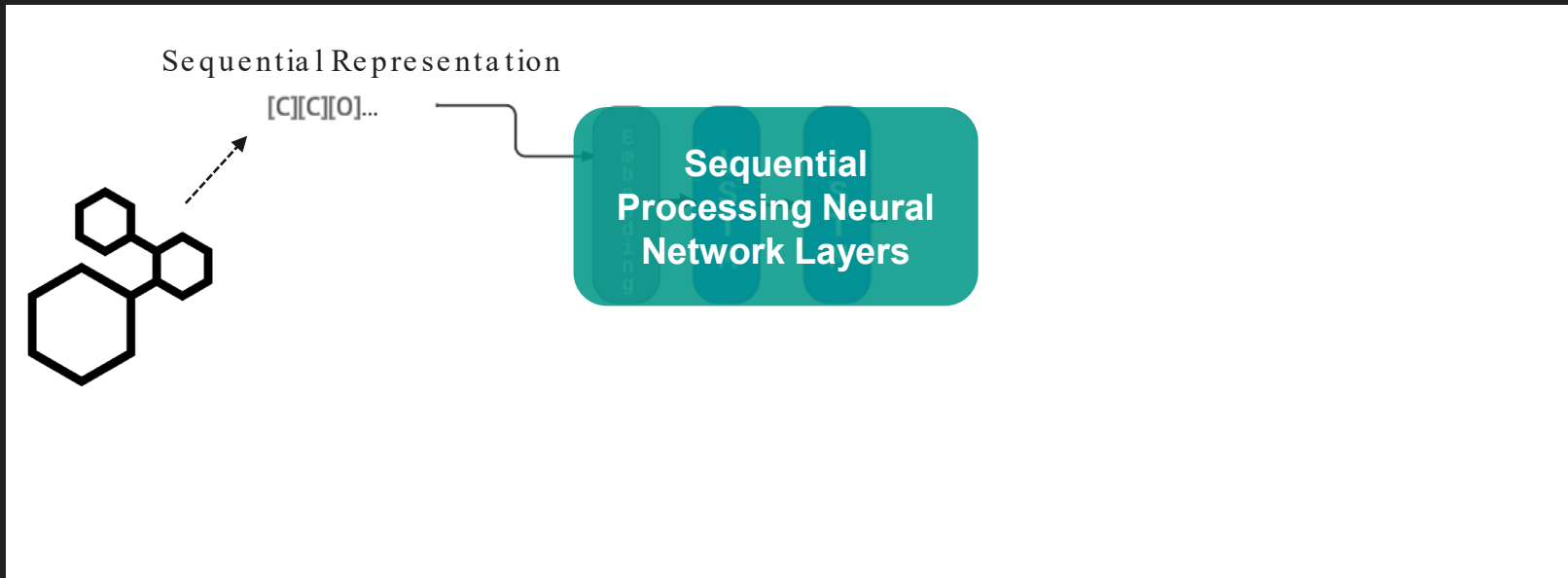
$^{\circ} p; \mathcal{H}_y > 7 \Rightarrow$ Highly Effective (Sydow et al., 2019)

↑ More Effective = **↑ pIC₅₀** = **↑ -log₁₀**

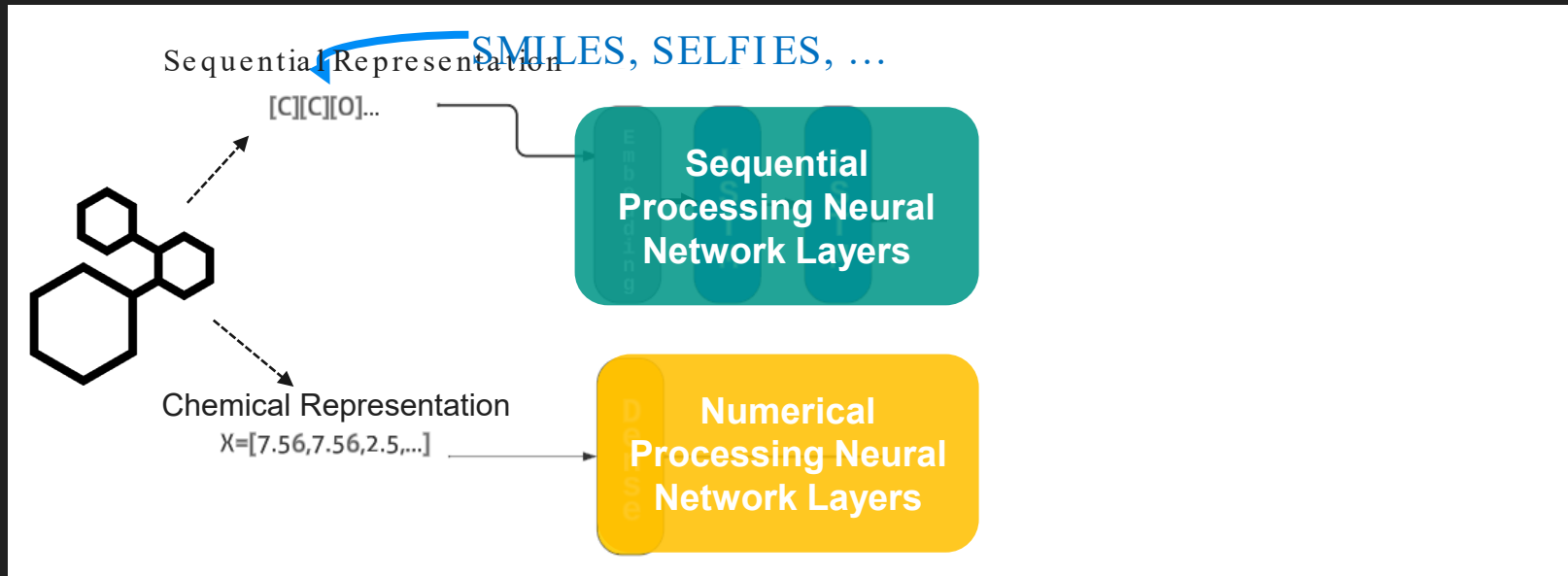
My novel drug efficacy evaluation model design combines both sequential and chemical representations



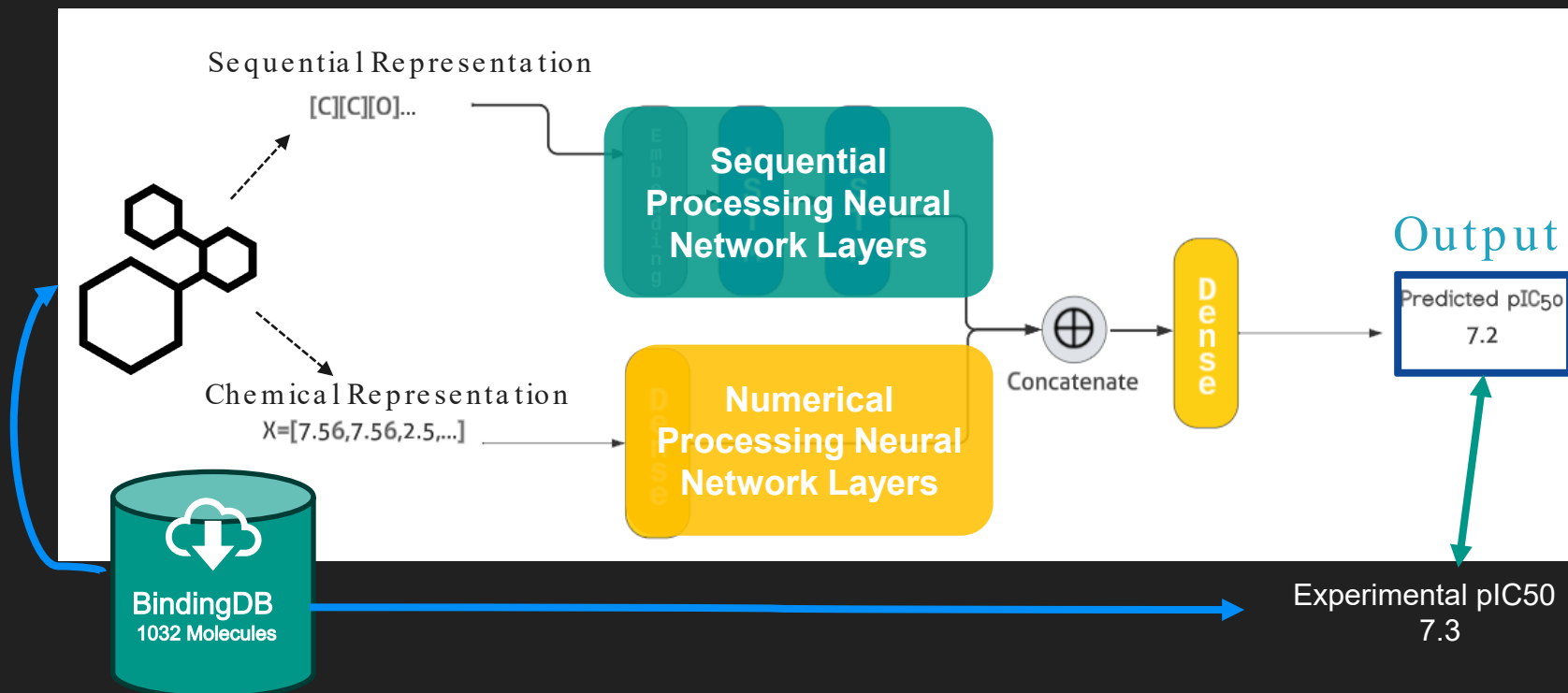
My novel drug efficacy evaluation model design combines both sequential and chemical representations



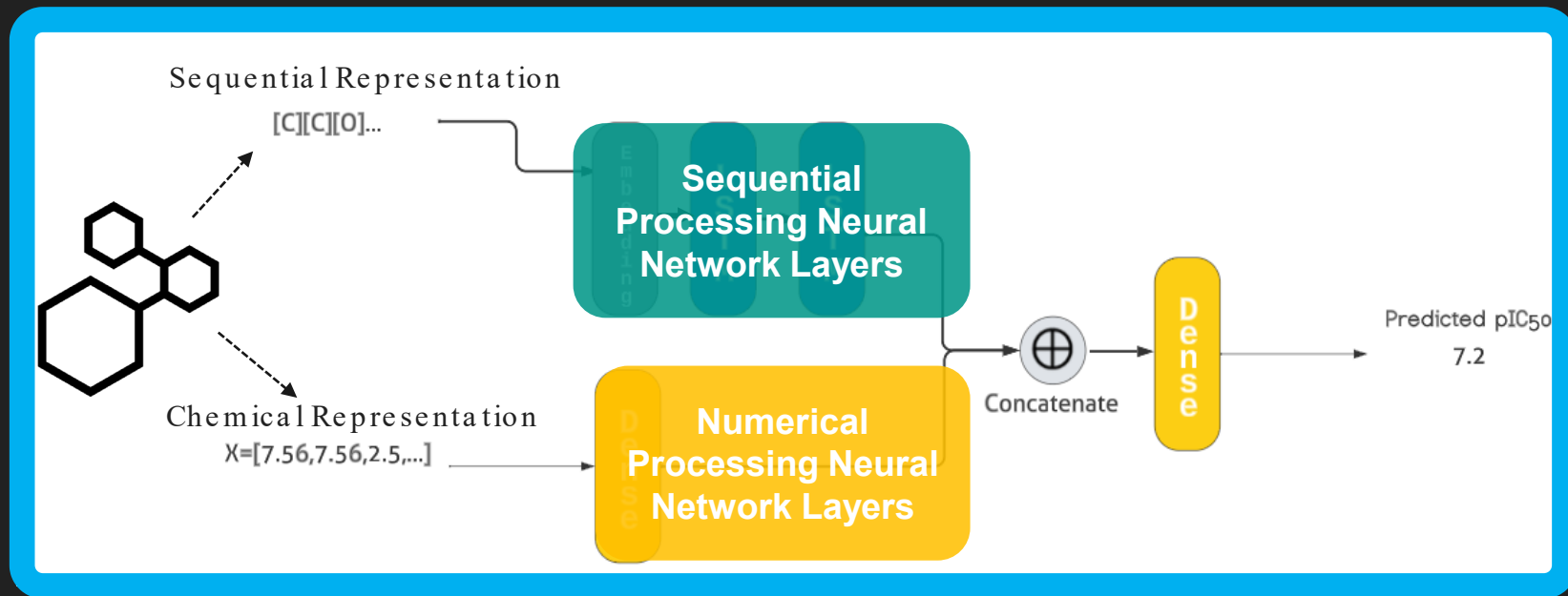
My novel drug efficacy evaluation model design combines both sequential and chemical representations



My novel drug efficacy evaluation model design combines both sequential and chemical representations



My novel drug efficacy evaluation model design combines both sequential and chemical representations



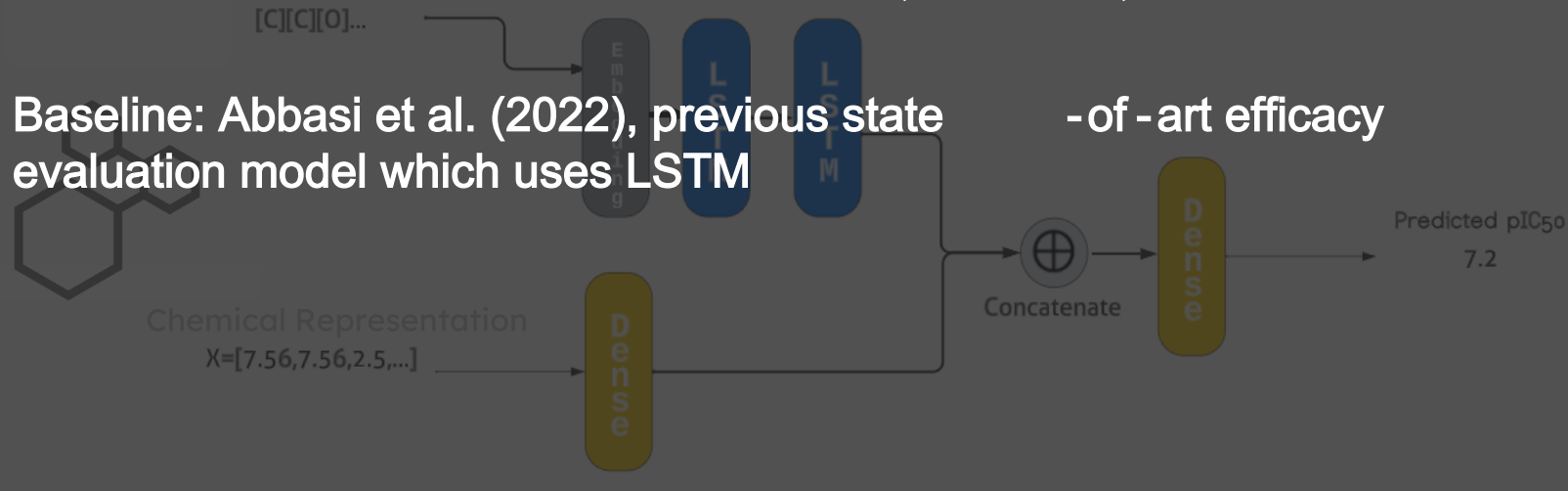
Novel Structure

My novel drug efficacy evaluation model design combines both sequential and chemical representations

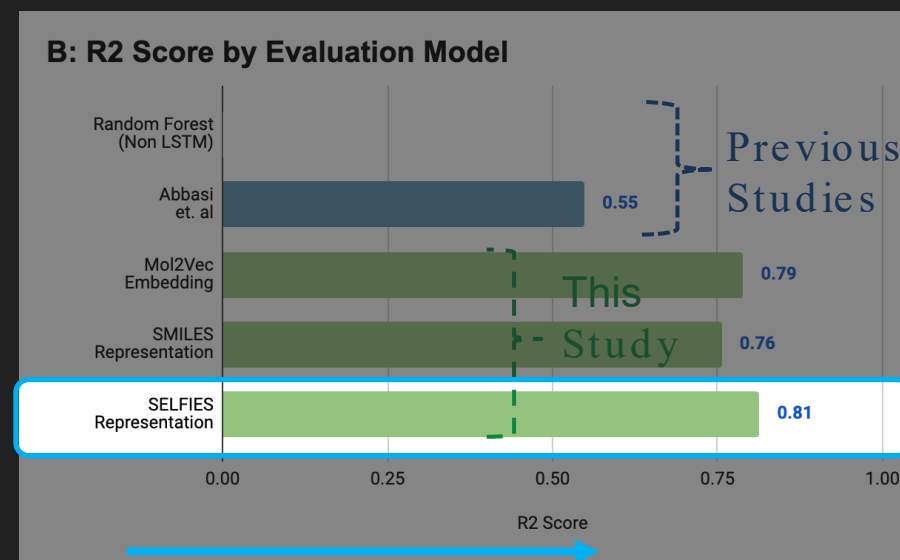
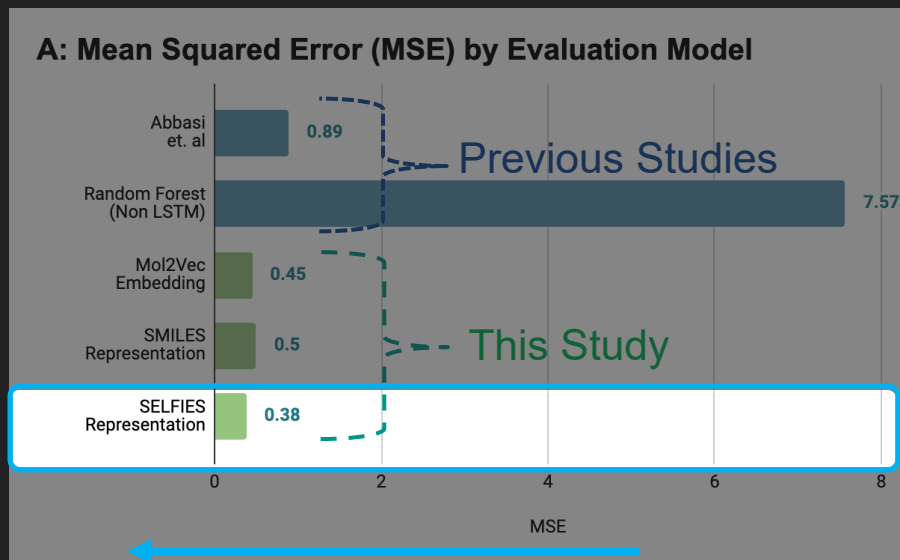
Same model structure is used for: SMILES, SELFIES, Mol2Vec

Baseline: Abbasi et al. (2022), previous state
evaluation model which uses LSTM

-of-art efficacy



Combining sequential representation with chemical descriptors **improves accuracy**

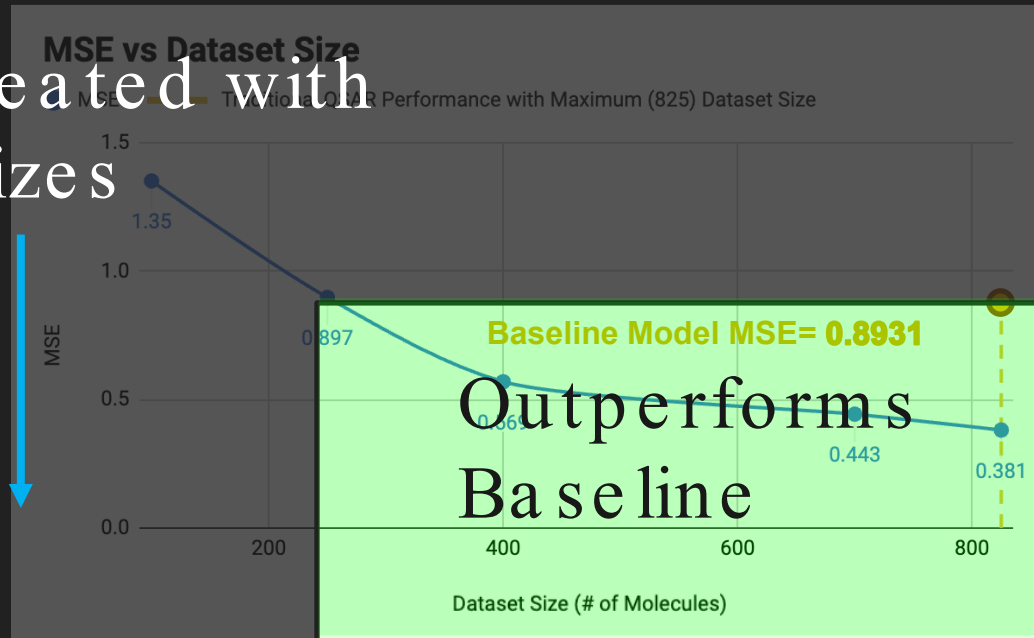



Effect of Dataset Size on Performance?

My efficacy evaluation model **outperforms**
baseline model **even with less data**

 Experiment repeated with
different dataset sizes

Performance
Increases



Methodology

Objectives

1. Evaluate Drug Efficacy

2. Design Drug -Like Molecules

3. Optimize Drug Efficacy

Steps

Reward Modeling



Drug Efficacy
Evaluation Model

Supervised Finetuning



Drug Design
GPT Model

Proximal Policy
Optimization (PPO)



Methodology

Objectives

1. Evaluate Drug Efficacy

2. Design Drug -Like Molecules

3. Optimize Drug Efficacy

Steps

Reward Modeling



Drug Efficacy
Evaluation Model



Supervised Finetuning



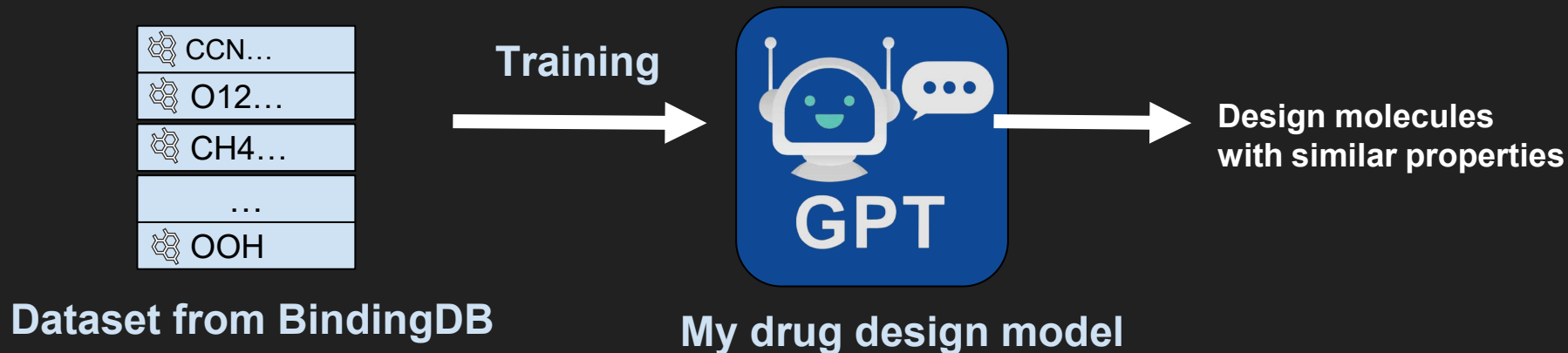
Drug Design
GPT Model



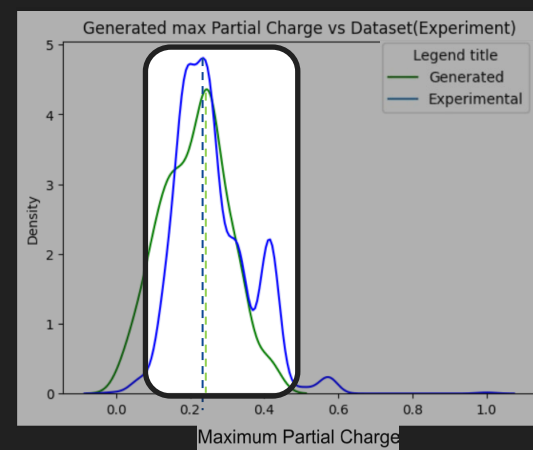
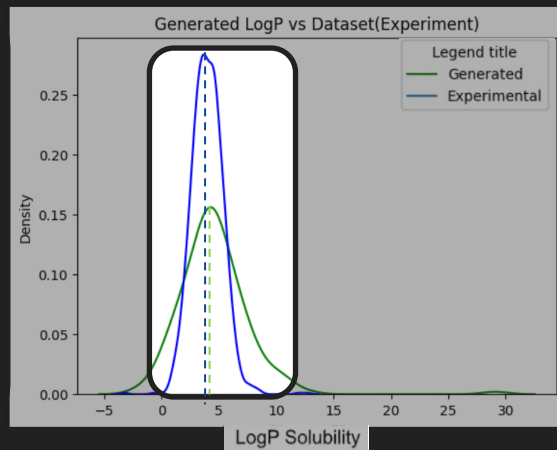
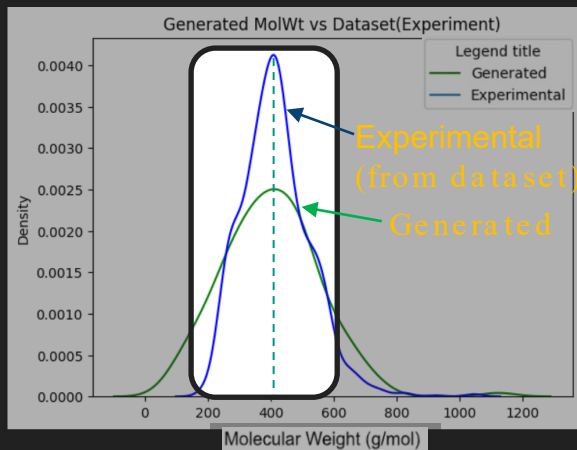
Proximal Policy
Optimization (PPO)



Supervised finetuning training uses the same dataset for designing drug-like molecules



Generated molecules exhibit similar properties as ones from the dataset using Supervised finetuning



Methodology

Objectives

1. Evaluate Drug Efficacy

2. Design Drug -Like Molecules

3. Optimize Drug Efficacy

Steps

Reward Modeling



Drug Efficacy
Evaluation Model



Supervised Finetuning



Drug Design
GPT Model



Proximal Policy
Optimization (PPO)



Methodology

Objectives

1. Evaluate Drug Efficacy

2. Design Drug -Like Molecules

3. Optimize Drug Efficacy

Steps

Reward Modeling



Drug Efficacy
Evaluation Model



Supervised Finetuning



Drug Design
GPT Model

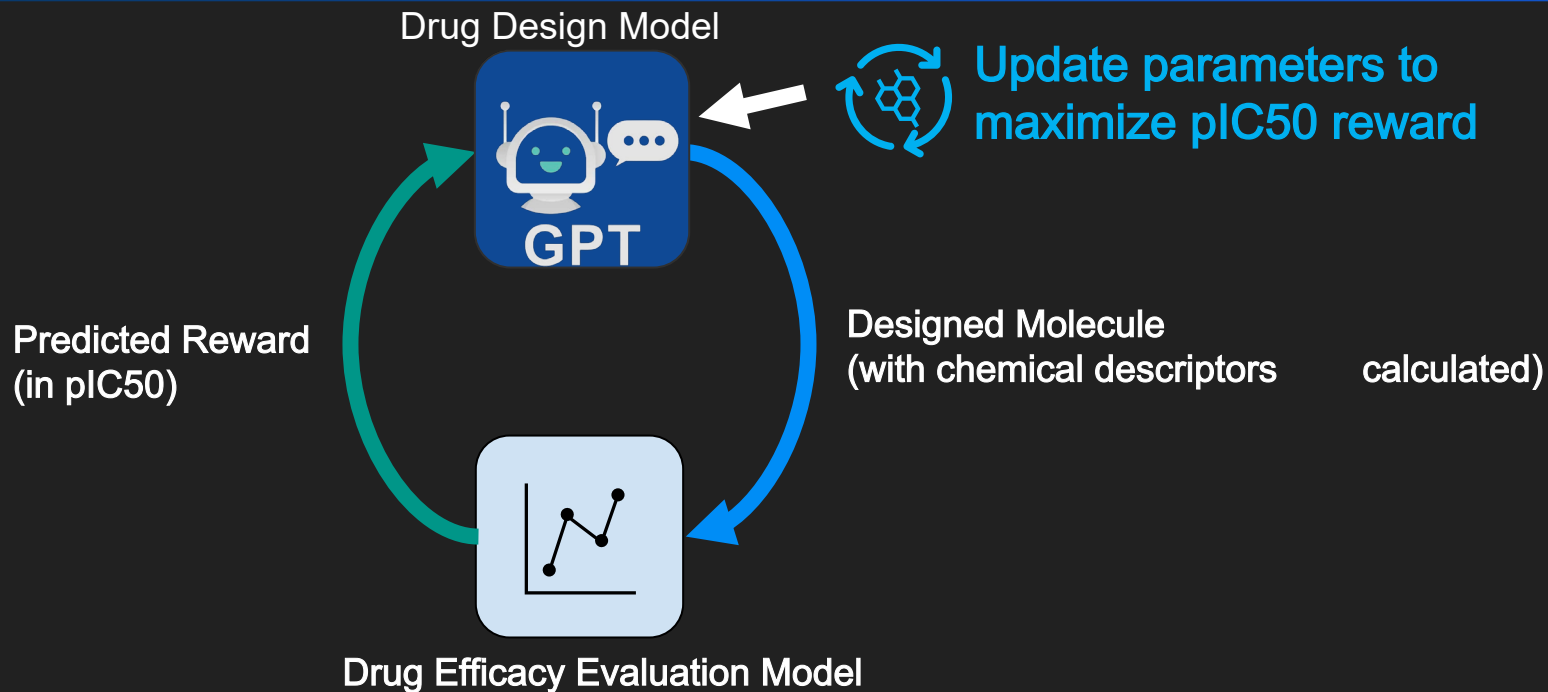


Proximal Policy
Optimization (PPO)



(first time used for drug design)

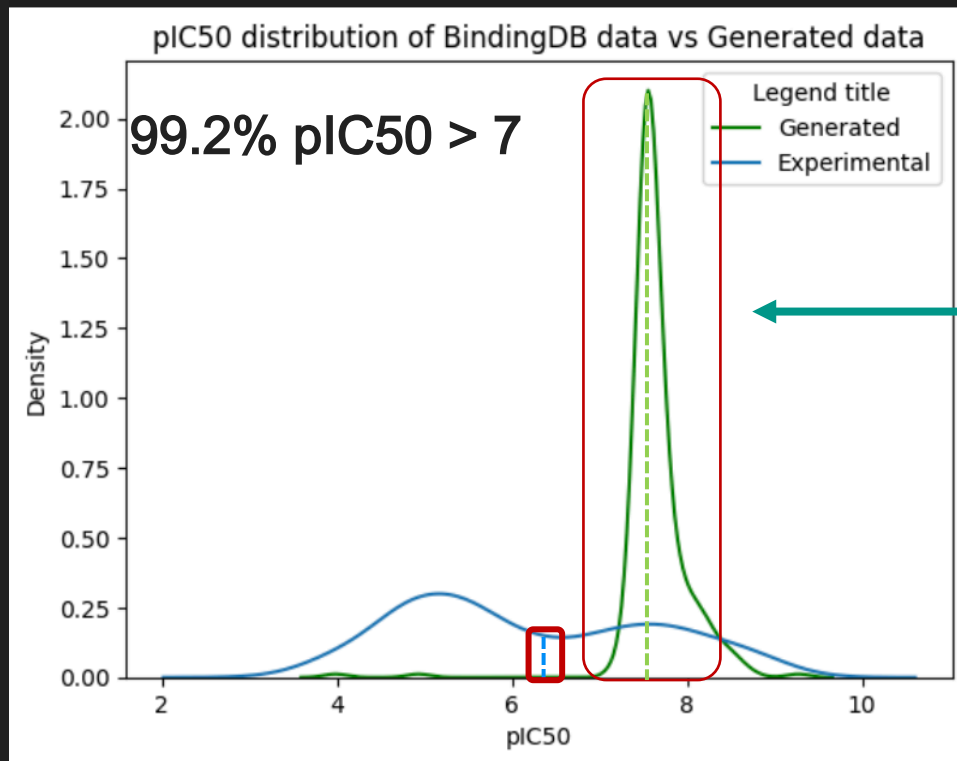
PPO Workflow



Result: PPO effectively
molecules for drug design

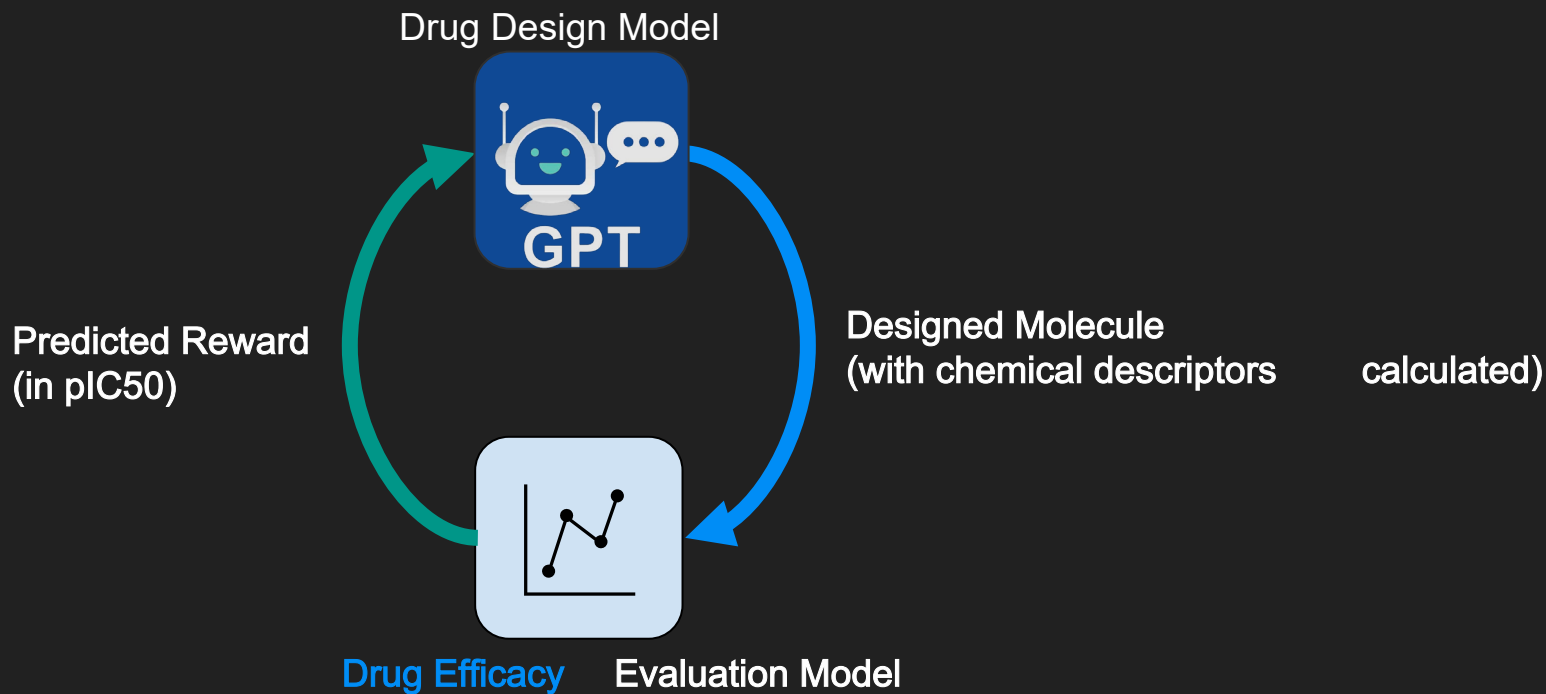
optimized efficacy of
for the first time

100 % Validity
(used RDKit
for validation)

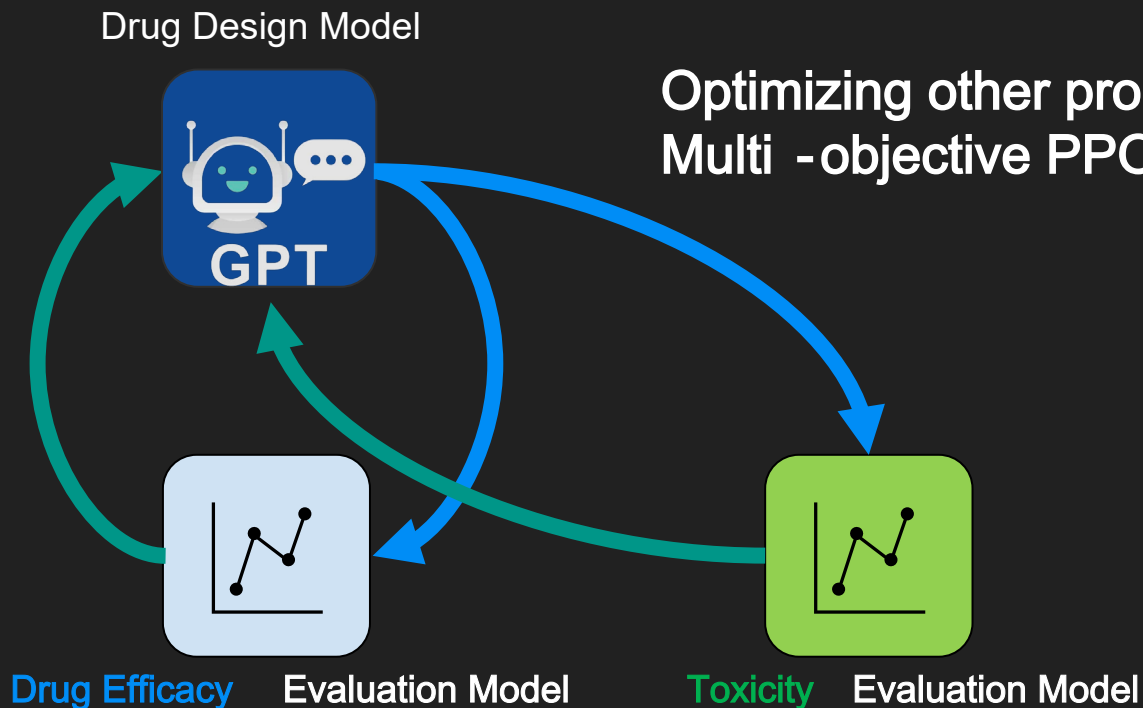


Generated

Limitations & Next Steps

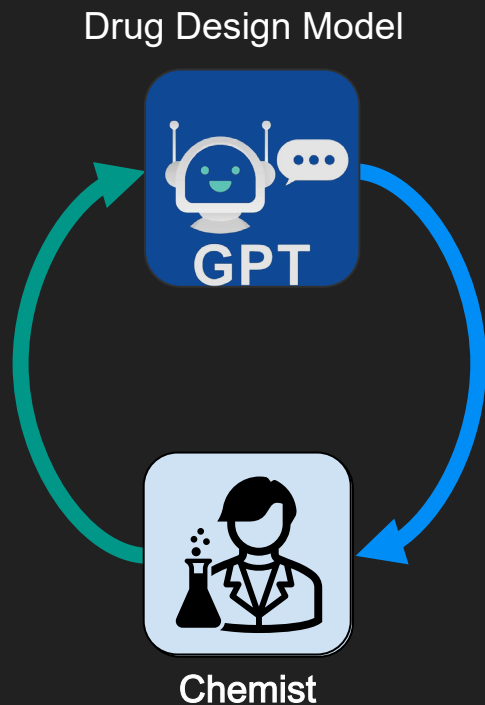


Future study can optimize multiple properties of the drug design model using similar methodology



(Khoi et al., 2021)

Future study can optimize multiple properties of the drug design model using similar methodology



Use human chemists
to provide feedback
(a.k.a RLHF with PPO)
(Ouyang et al., 2022)

Brief Recap of Problems: Non GPT models have low validity

Problems:

Traditional Drug Discovery
Costly and time consuming

(Abbasi et al., 2022) **Low Validity**

(Yasonik ., 2020) **Low Validity**

(Frey et al.,2022) **HIGH Validity**
...Not applied to any specific task

This Study
(High Efficacy + Validity)

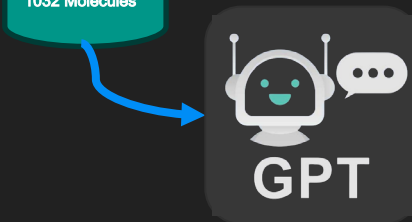
Novelty

Step 1, 2

Step 3

Drug Evaluation Model

Drug Design Model



Drug Design Model



Evaluate
Molecule

Design
Molecule



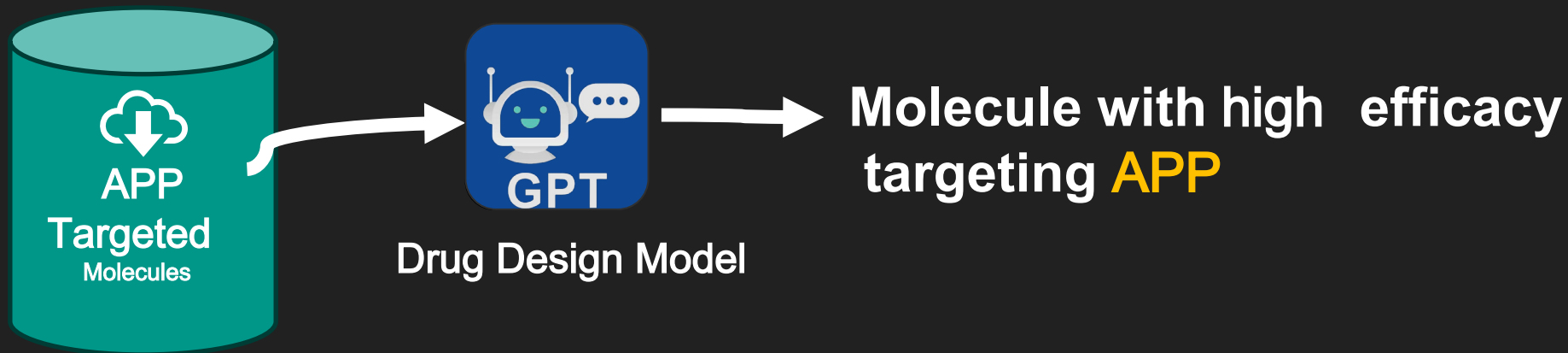
Drug Evaluation Model

Design model of
evaluation models

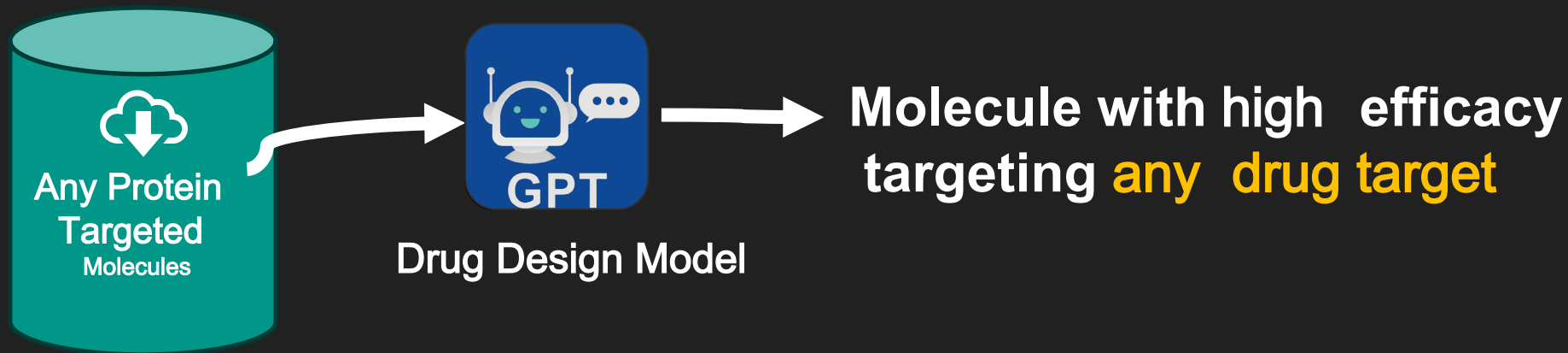
STEP 3:

GPT model successfully designed
molecules with high efficacy
(for the 1st time)

Training procedure is generalizable



Training procedure is generalizable



Significance: My drug design model can speed up drug discovery

Identify target protein



Traditionally costly and time consuming

Drug (candidates) Designing

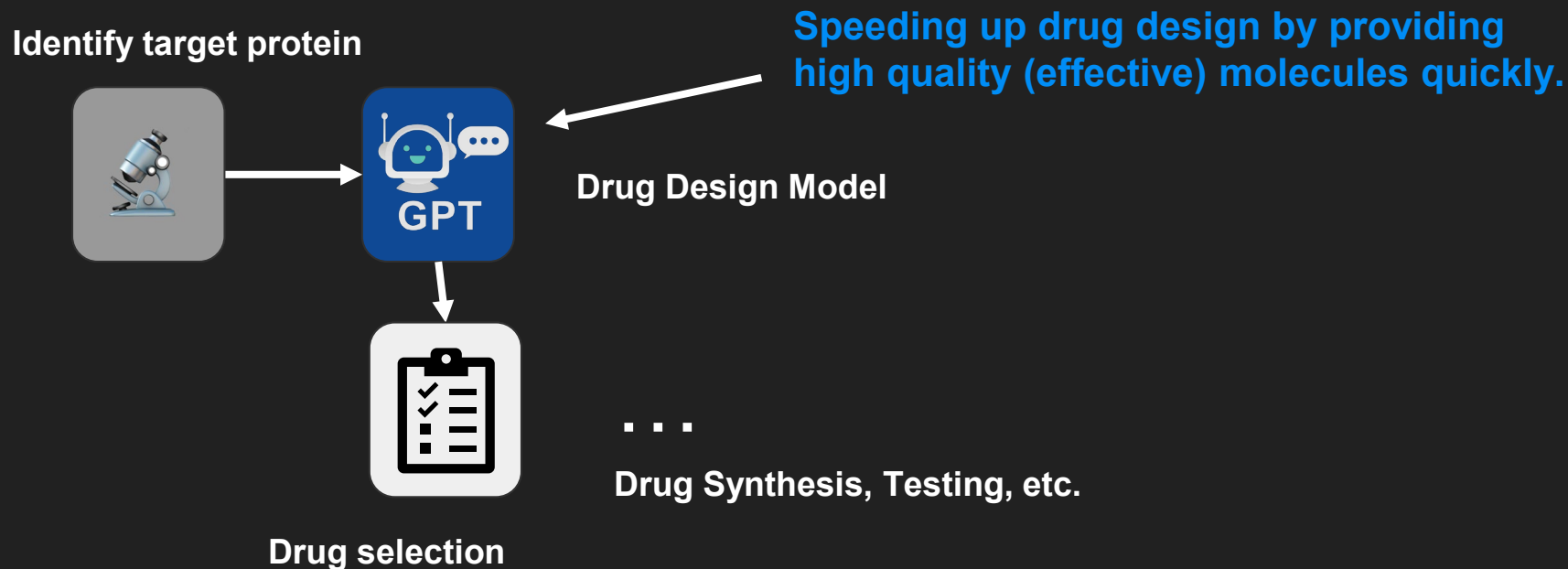


Drug selection

...

Drug Synthesis, Testing, etc.

Significance: My drug design model can speed up drug discovery



References

- Abbasi, M., Santos, B. P., Pereira, T. C., Sofia, R., Monteiro, N. R. C., Simões, C. J. V., Brito, R. M. M., Ribeiro, B., Oliveira, J. L., & Arrais, J. P. (2022). Designing optimized drug candidates with Generative Adversarial Network. *Journal of Cheminformatics*, 14(1), 40. <https://doi.org/10.1186/s13321-022-00623-6>
- DiMasi, J. A., Grabowski, H. G., & Hansen, R. W. (2016). Innovation in the pharmaceutical industry: New estimates of R&D costs. *Journal of Health Economics*, 47, 20–33. <https://doi.org/10.1016/j.jhealeco.2016.01.012>
- Frey, N., Soklaski, R., Axelrod, S., Samsi, S., Gomez-Bombarelli, R., Coley, C., & Gadepally, V. (2022). *Neural Scaling of Deep Chemical Models*. ChemRxiv. <https://doi.org/10.26434/chemrxiv-2022-3s512>
- Gao, K., Nguyen, D. D., Tu, M., & Wei, G.-W. (2020). Generative Network Complex for the Automated Generation of Drug-like Molecules. *Journal of Chemical Information and Modeling*, 60(12), 5682–5698. <https://doi.org/10.1021/acs.jcim.0c00599>
- Images (225×225)*. (n.d.). Retrieved January 23, 2024, from https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GeQWKtU0Mueg8u_SWPmmCfNhU0gYi0AVN3ZZvnH5k2szmB2XIaAd
- Khoi, N. D. H., Pham Van, C., Tran, H. V., & Truong, C. D. (2021). Multi-Objective Exploration for Proximal Policy Optimization. 2020 Applying New Technology in Green Buildings (ATiGB), 105–109. <https://doi.org/10.1109/ATiGB50996.2021.9423319>
- Krenn, M., Häse, F., Nigam, A., Friederich, P., & Aspuru-Guzik, A. (2020). Self-referencing embedded strings (SELFIES): A 100% robust molecular string representation. *Machine Learning: Science and Technology*, 1(4), 045024. <https://doi.org/10.1088/2632-2153/aba947>
- Liu, T., Lin, Y., Wen, X., Jorissen, R. N., & Gilson, M. K. (2007). BindingDB: A web-accessible database of experimentally determined protein–ligand binding affinities. *Nucleic Acids Research*, 35(Database issue), D198. <https://doi.org/10.1093/nar/gkl999>
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Aspell, A., Welinder, P., Christiano, P., Leike, J., & Lowe, R. (2022). *Training language models to follow instructions with human feedback* (arXiv:2203.02155). arXiv. <https://doi.org/10.48550/arXiv.2203.02155>
- Popova, M., Isayev, O., & Tropsha, A. (2018). Deep Reinforcement Learning for De-Novo Drug Design. *Science Advances*, 4(7), eaap7885. <https://doi.org/10.1126/sciadv.aap7885>
- Segler, M. H. S., Preuss, M., & Waller, M. P. (2018). Planning chemical syntheses with deep neural networks and symbolic AI. *Nature*, 555(7698), 604–610. <https://doi.org/10.1038/nature25978>
- Sydow, D., Morger, A., Driller, M., & Volkamer, A. (2019). TeachOpenCADD: A teaching platform for computer-aided drug design using open source packages and data. *Journal of Cheminformatics*, 11(1), 29. <https://doi.org/10.1186/s13321-019-0351-x>
- Weininger, D. (1988). SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *Journal of Chemical Information and Computer Sciences*, 28(1), 31–36. <https://doi.org/10.1021/ci00057a005>
- Yang, S., Wang, S., Zhao, M., Zhang, R., Zhou, W., Li, Y., Su, Y., Zhang, H., Yu, X., & Liu, R. (2012). A Peptide Binding to the β -Site of APP Improves Spatial Memory and Attenuates A β Burden in Alzheimer's Disease Transgenic Mice. *PLOS ONE*, 7(11), e48540. <https://doi.org/10.1371/journal.pone.0048540>
- Yasonik, J. (2020). Multiobjective de novo drug design with recurrent neural networks and nondominated sorting. *Journal of Cheminformatics*, 12(1), 14. <https://doi.org/10.1186/s13321-020-00419-6>

Summary

GPT For
Drug Design