



筑波大学

University of Tsukuba

Deep generative model for drug design from
protein target sequence

Yangyang Chen

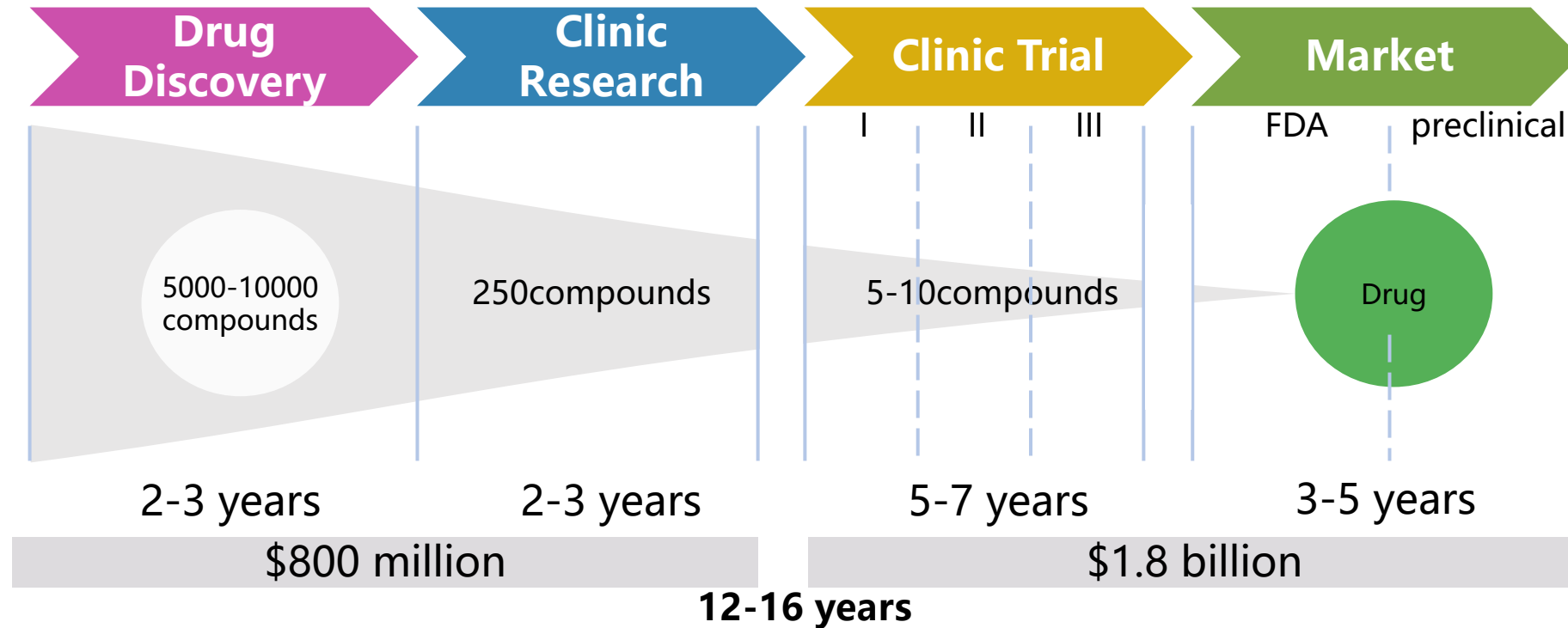


- **Background**
- **Existing research**
- **Our research**
- **Future outlook**



Background

Traditional R&D



Artificial Intelligence (AI)

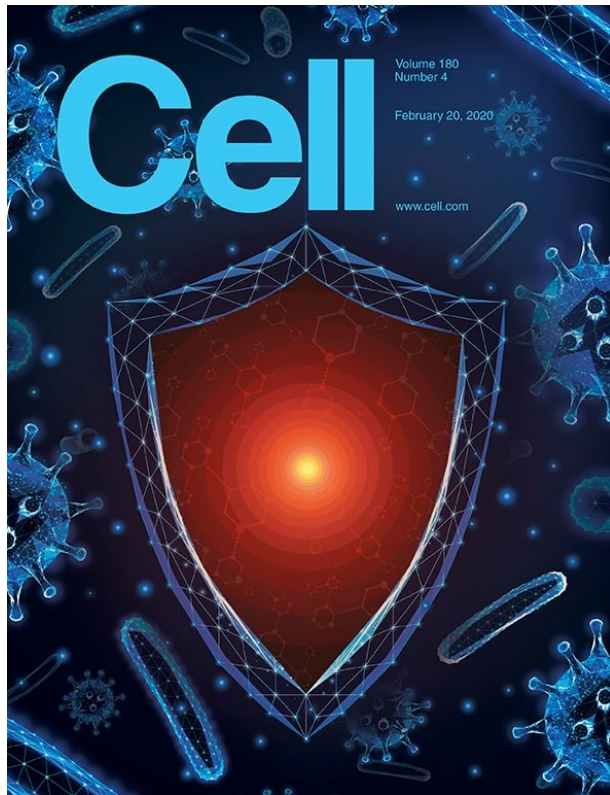
Target Discovery
Lead Compound Design
Lead Compound Screening
Lead Compound Optimization

Drug Repositioning
New indication discovery
Drug Property Prediction
Crystal shape prediction

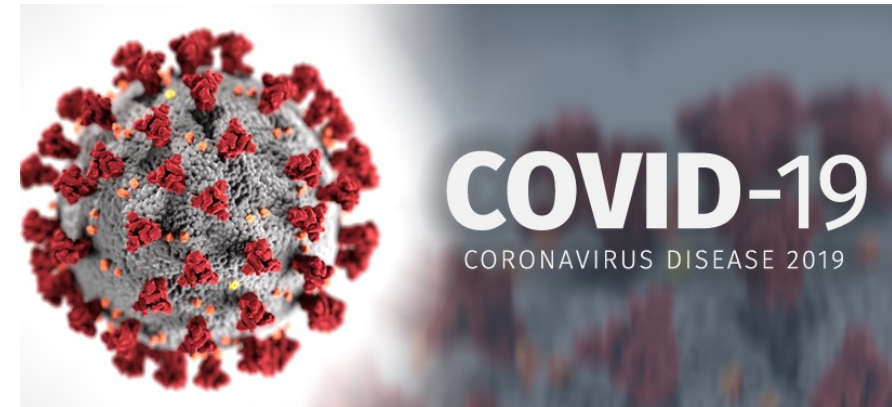
Clinical Trial
Patient Recruitment

Drug Interactions
Product Inspection
Drug Synthesis

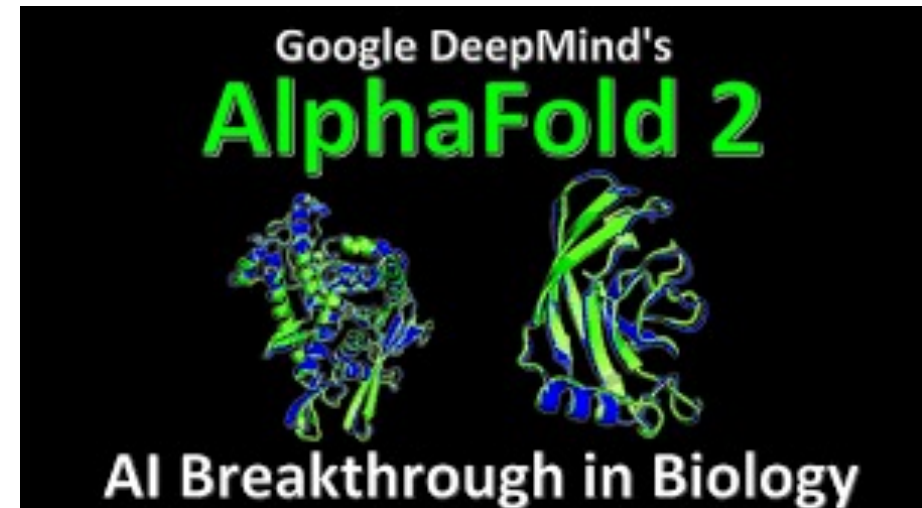
2020-The Year of Artificial Intelligence Drug Development



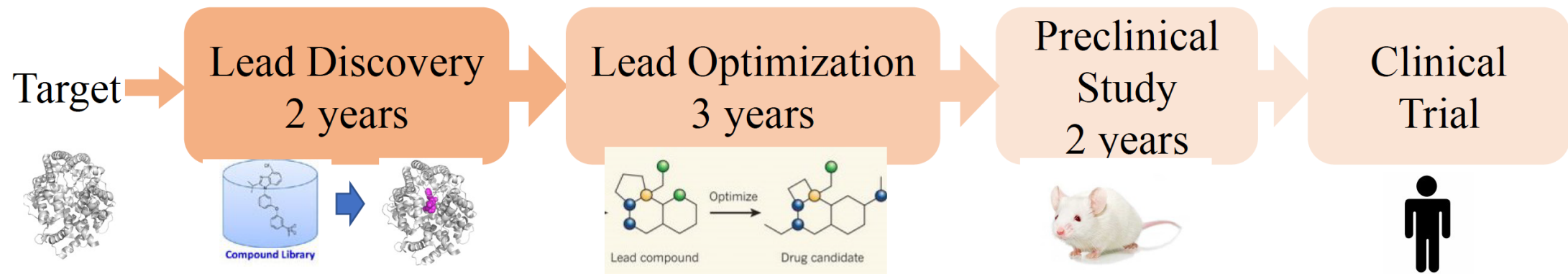
In 2020, Cell cover,
AI discovers the super-powerful antibiotic Halicin



In 2020, the COVID-19 epidemic Outbreak.



In 2020, the AlphaFold2 comes out.

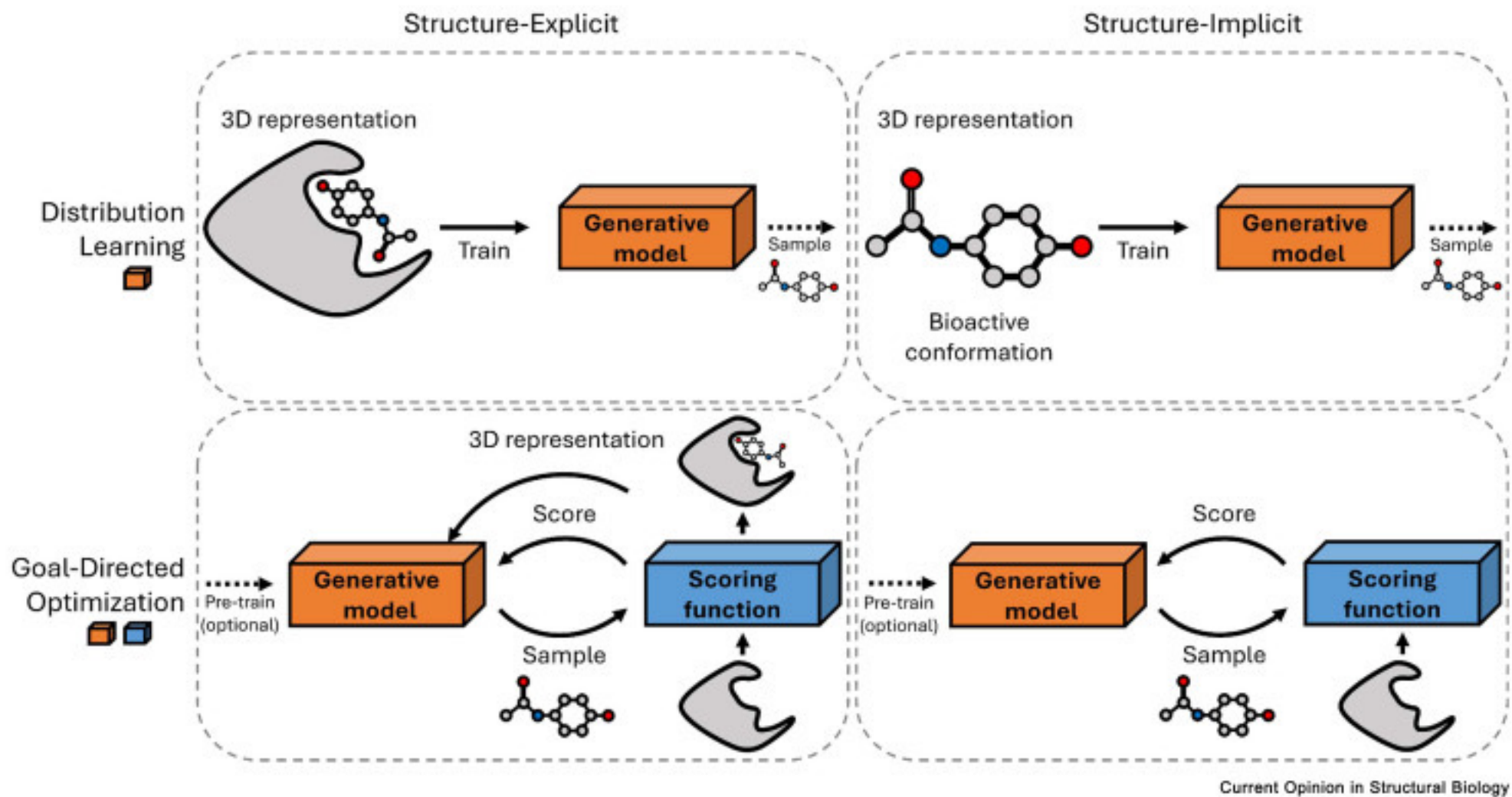


- Drug discovery: the process of finding new medicines to treat diseases.
- Finding the right drug ligand for a protein is a complex task.
- Traditional methods: biological experiments.
- **Intrinsic**: Finding the most suitable ligand from a large range of drug molecules.

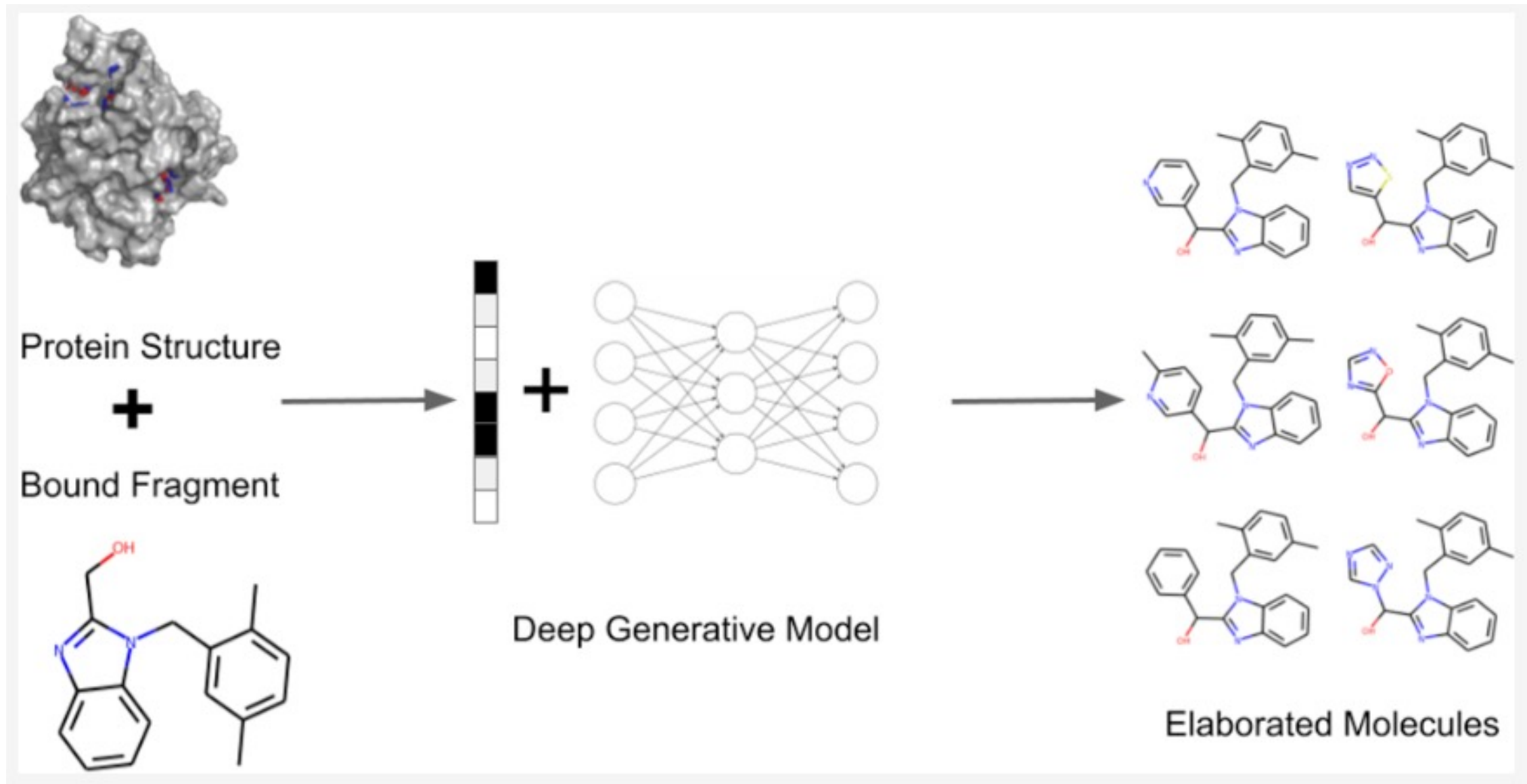


Existing research

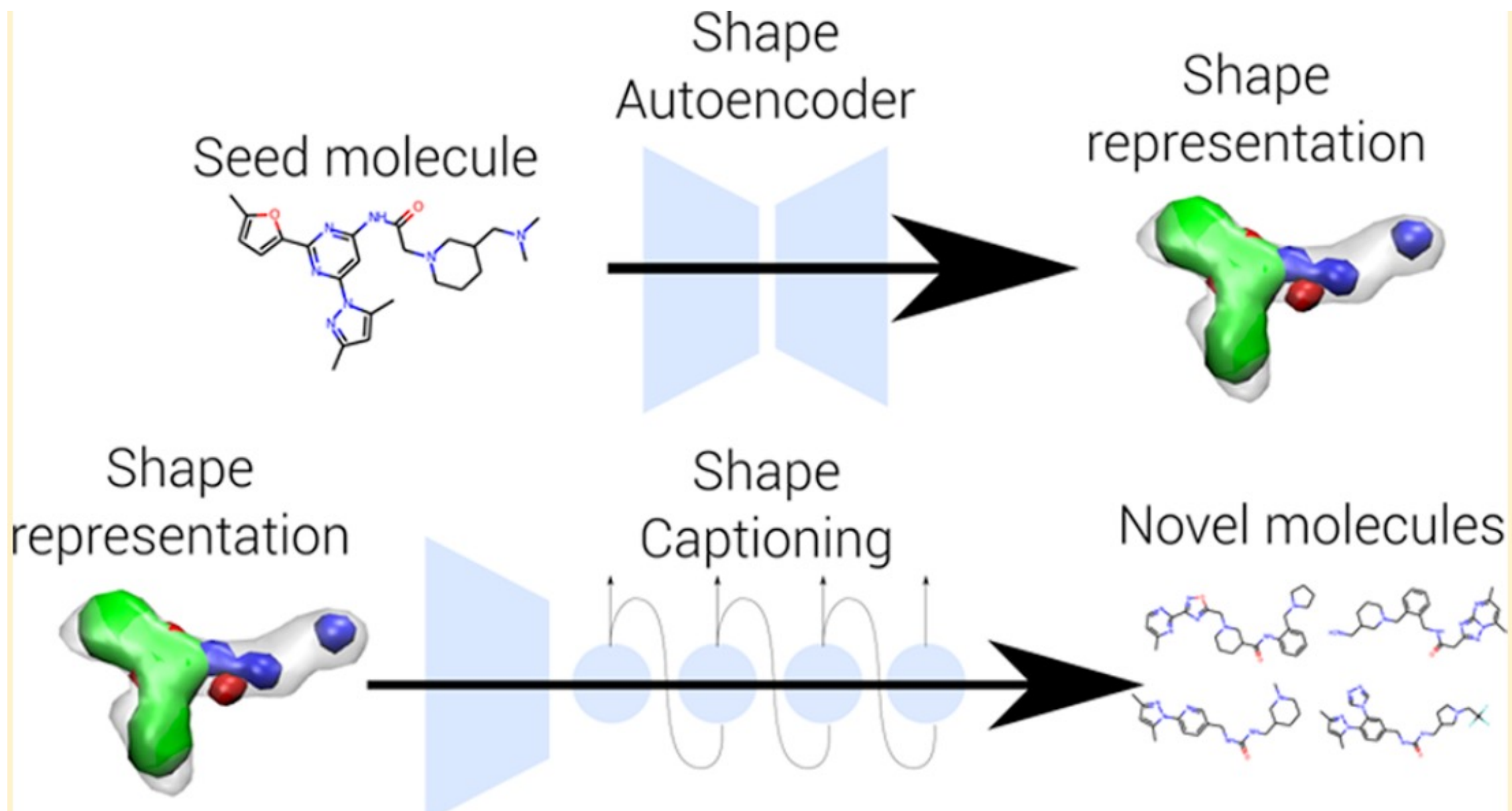
distribution learning or goal-directed optimization and structure-explicit/implicit



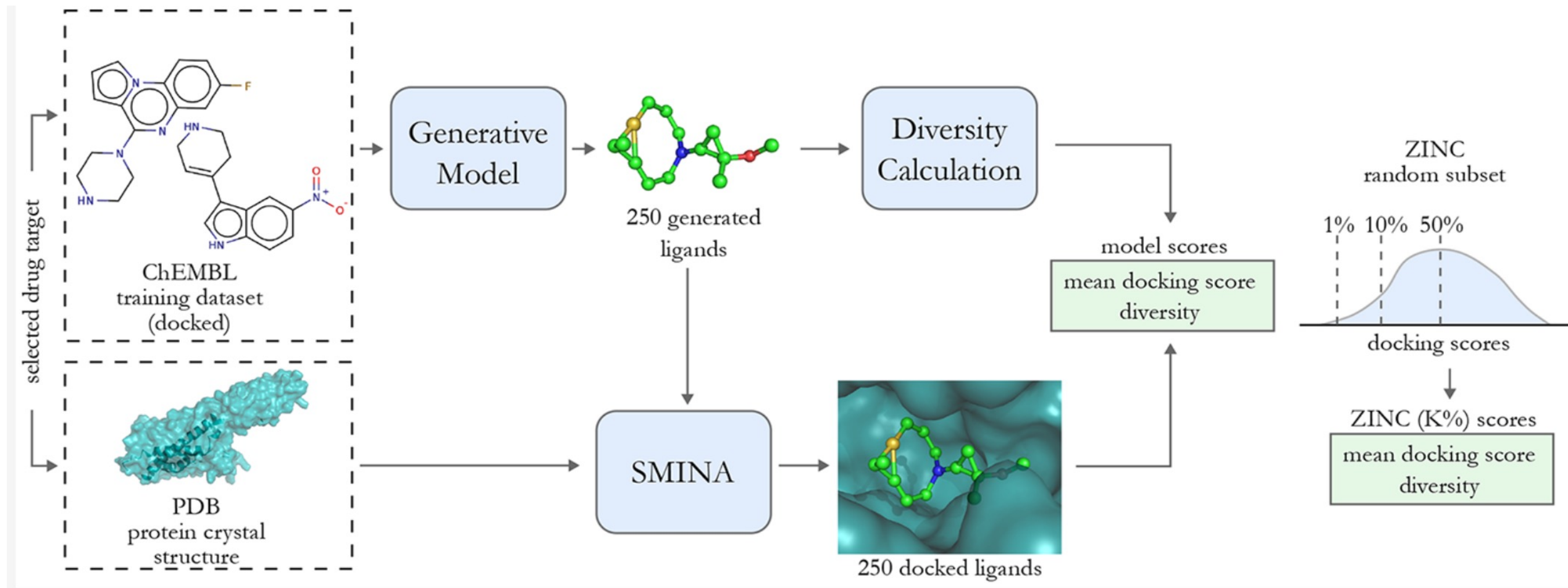
Distribution Learning & Structure-Explicit



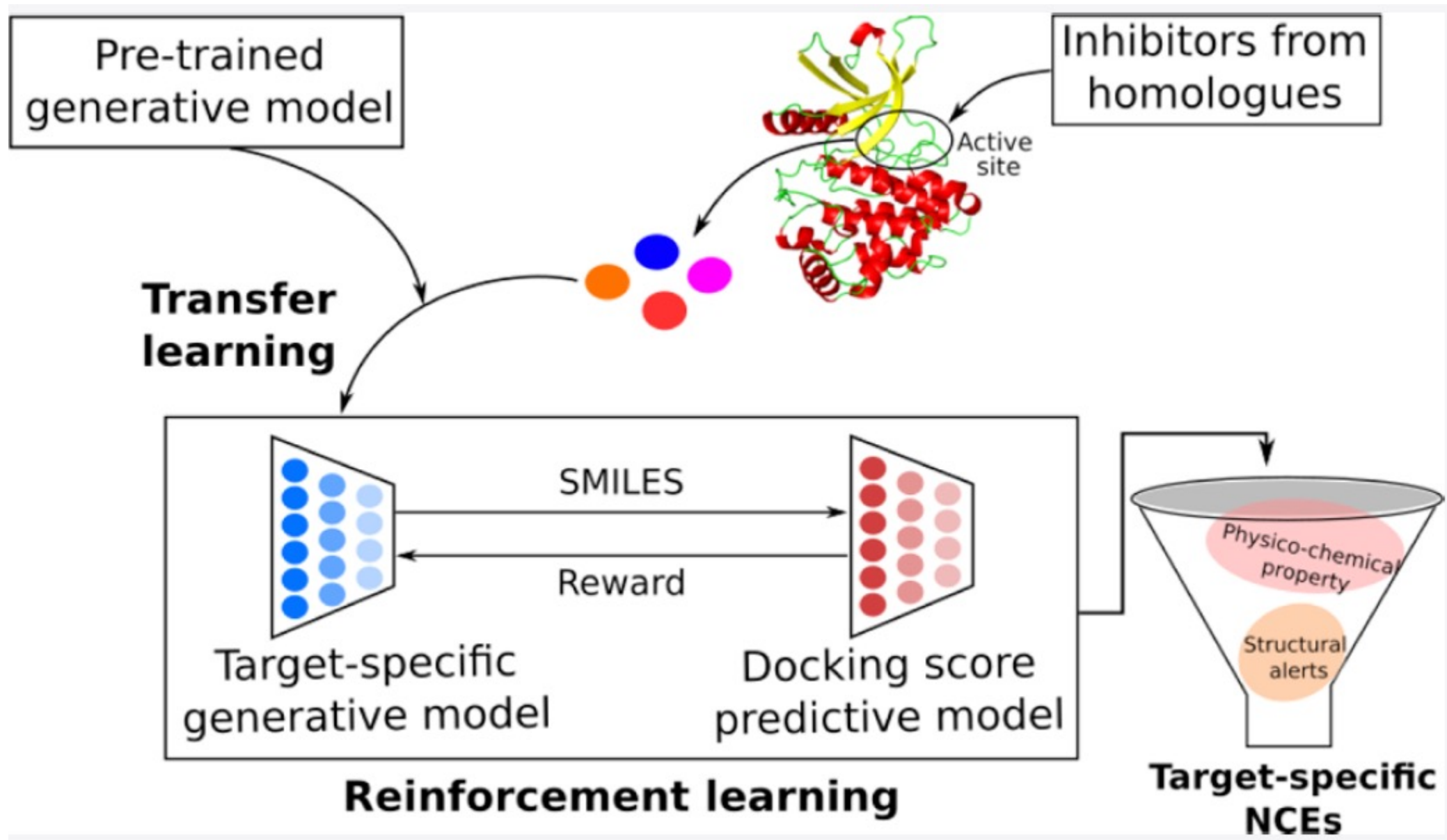
Distribution Learning & Structure-Implicit



Goal-directed Optimization & Structure-Explicit

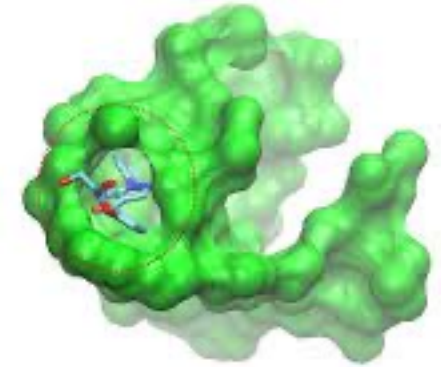
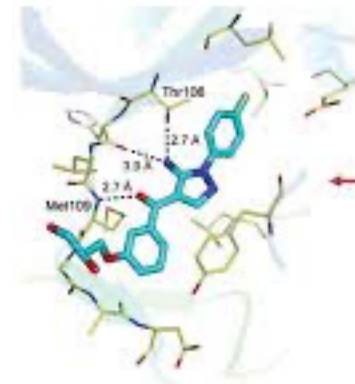
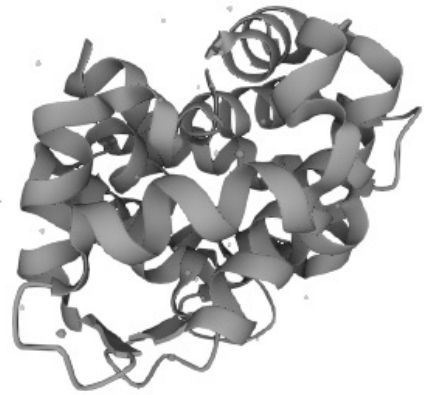


Goal-directed Optimization & Structure-Implicit

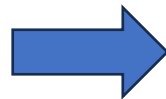


Problems and limitations

- Complexity of protein structure
- Structure of some proteins unknown (pocket)



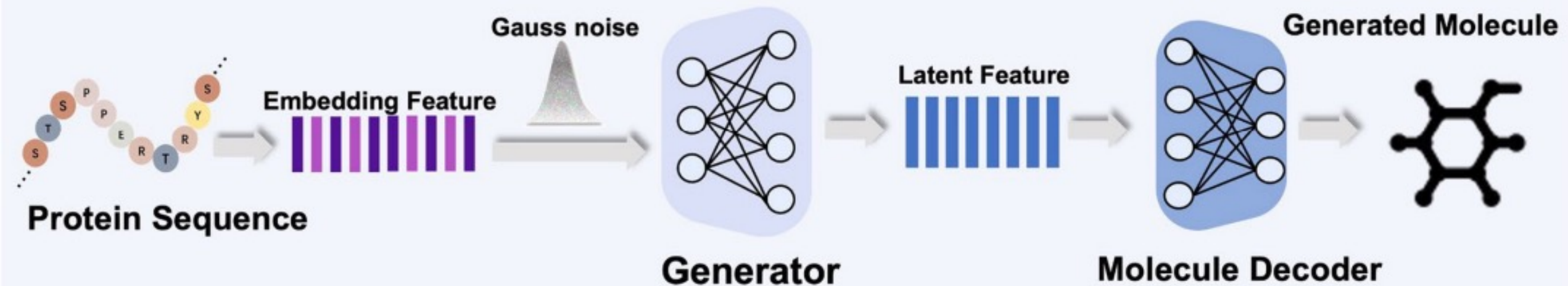
- **Simple form:**
protein sequence



```
>sp|P14416|DRD2_HUMAN D(2) dopamine receptor OS=Homo sapiens OX=9606 GN=DRD2 PE=1 SV=2
MDPLNLSWYD DDLERQNSR PFNGSDGKAD RPHYNYATL LTLIAVIVF GNVLVCMASV
REKALQTTN YLIVSLAVAD LLVATLVMPW VVYLEVVGEW KFSRIHCDIF VTLDVMMCTA
SILNLCAISI DRYTAVAMPM LYNTYSSKR RVTVMISIVW VLSFTISCPL LFGLNNADQN
ECIIANPAFV VYSSIVSFYV PFIVTLLVYI KIYIVLRRRR KRVNTRKSSR AFRAHLRAPL
KGNCTHPEDM KLCTVIMKSN GSFPVNRVV EAARRAQELE MEMLSSTSP ERTRYSPIPP
SHHQLTLPDP SHHGLHSTPD SPAKPEKNGH AKDHPKIAKI FEIQTMPNGK TRTSLKTMRSR
RKLSQQKEKK ATQMLAIVLG VFIICWLPFF ITHILNIHCD CNIPPVLYSA FTWLG YVNSA
VNPIIYTTFN IEFKAFKLI LHC
```



Our research



Model overflow

- Input: protein sequence Output: molecule SMILES
- Three modules:
 - AASE: Amino Acid Sequence Embedding module
 - SFI: Structural Feature Inference module
 - MG: Molecule Generation module

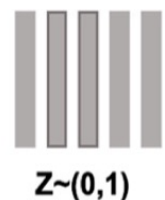
AASE: Amino Acid Sequence Embedding module

A Protein Sequence

```
MKSGSGGGSP TSLWGL  
LFLSAALSLWPTS GEIC  
GPGIDIRNDYQQL...
```

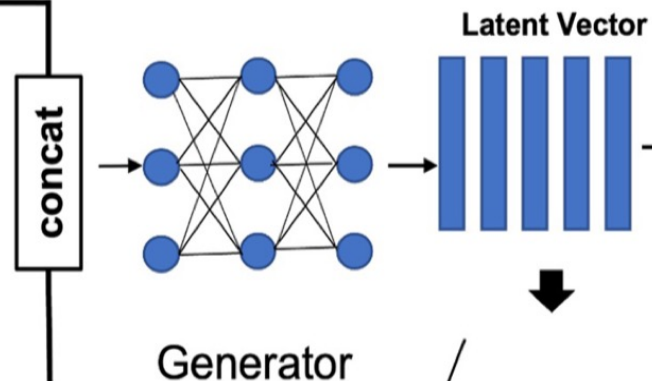
Protein Encoder

Protein Features



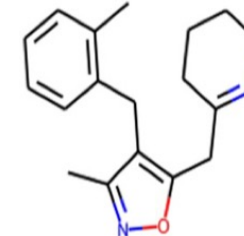
$Z \sim (0,1)$

SFI: Structural Feature Inference module



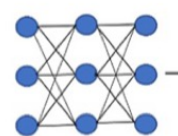
MG: Molecule Generation module

Molecule Decoder



```
Cc1noc(CC2=NCCCC2)c1Cc1  
c(C)cccc1
```

Generator



molecules

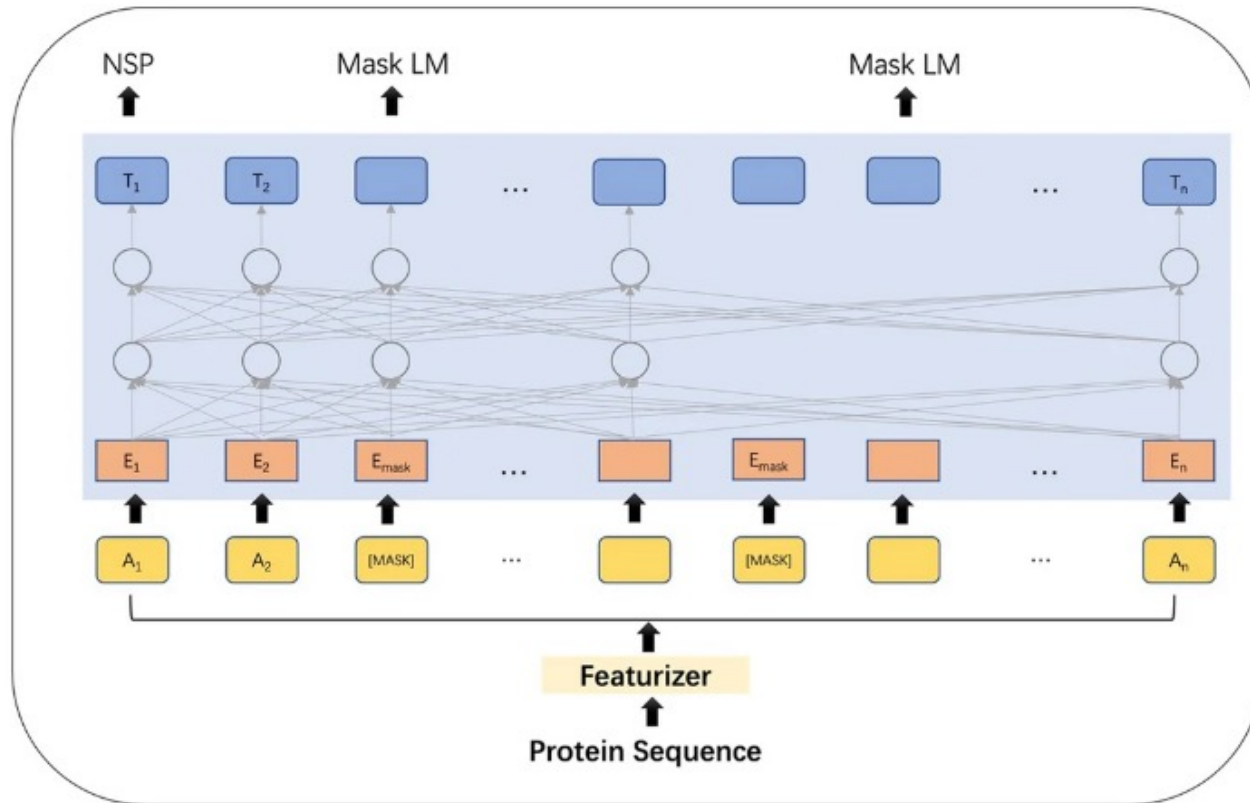


Contrastive Learning

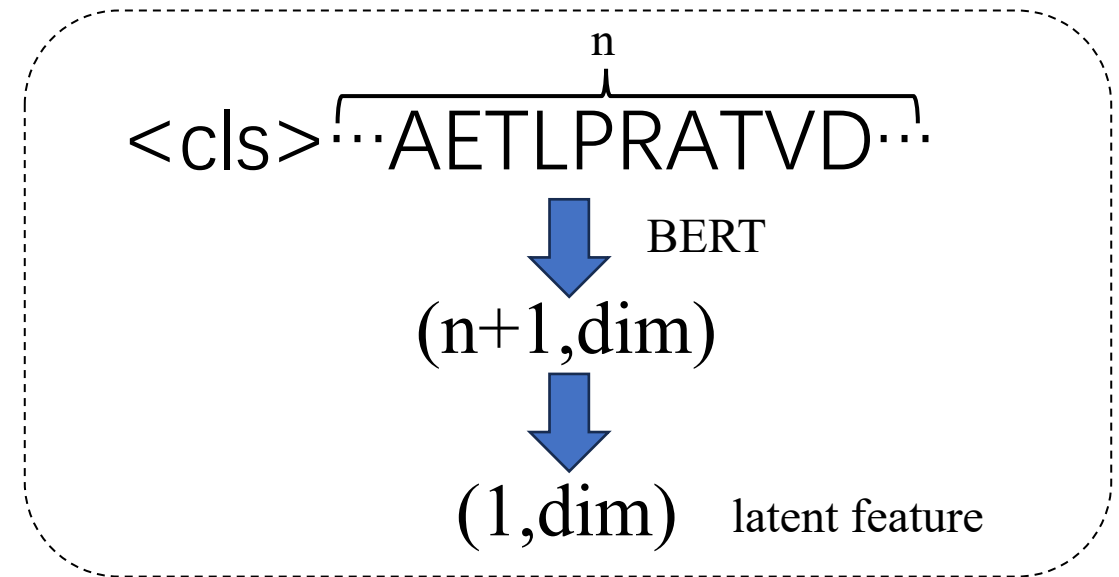
Optimize

- molecule of protein1
- molecule of protein2
- molecule of protein3

Amino Acid Sequence Embedding module

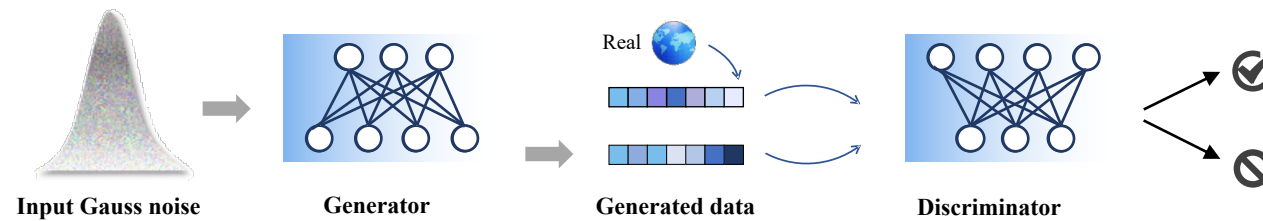


pre-trained protein encoder



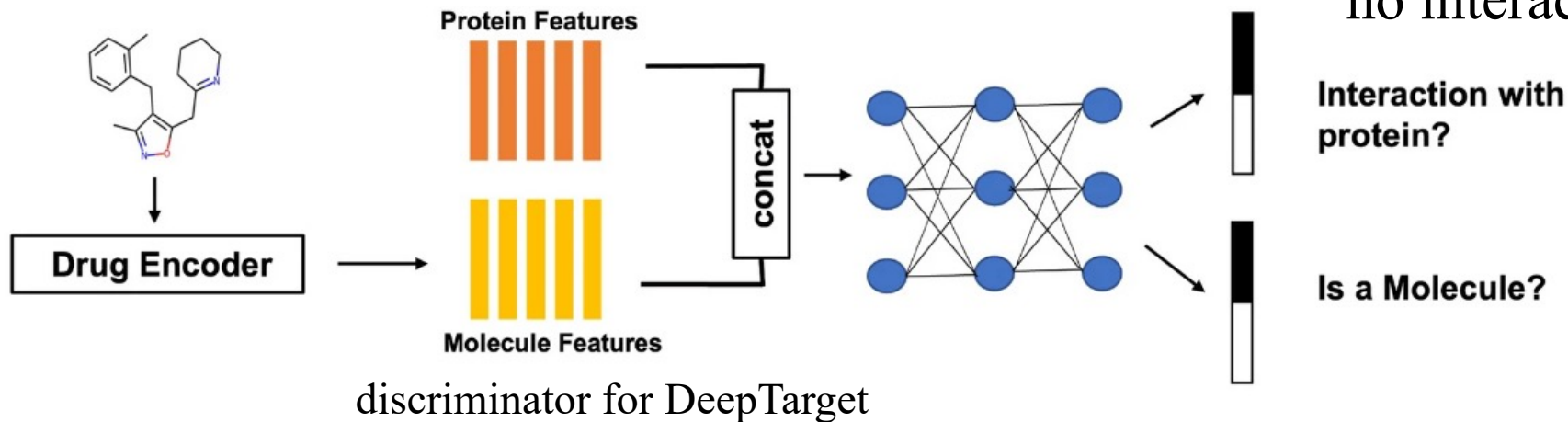
- Protein sequence is too long.
- BERT: Pre-training of deep bidirectional transformers for language understanding.
- Latent feature to represent the protein.

Structural Feature Inference module



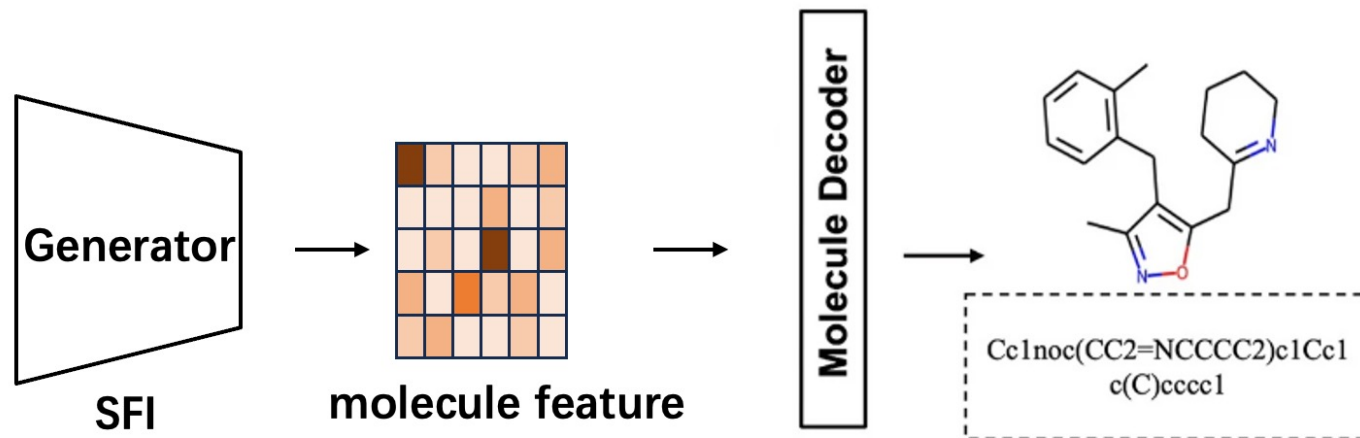
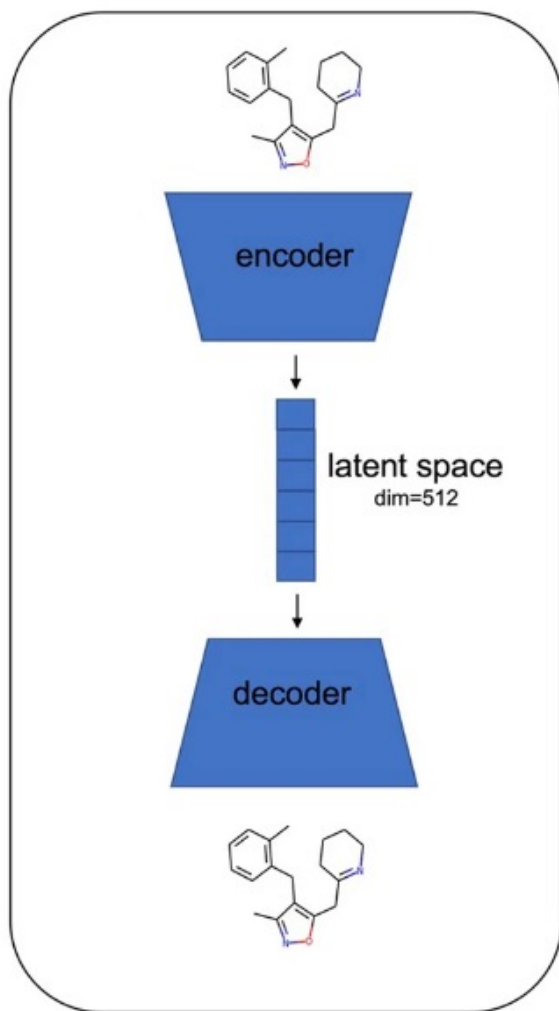
GANs: generative adversarial networks

- Generative Adversarial Networks
- Generator:
 - protein feature as input
- Discriminator:
 - one output->two outputs
- Real ligands: interact and valid
- Generated ligands:
 - no interaction and invalid



discriminator for DeepTarget

Molecule Generation module



- Autoencoder: pre-trained in ChEMBL
- Encoder:
embedding the molecules to latent feature
- Latent Feature:
represent the generated molecules
- Decoder
decoding feature to molecules

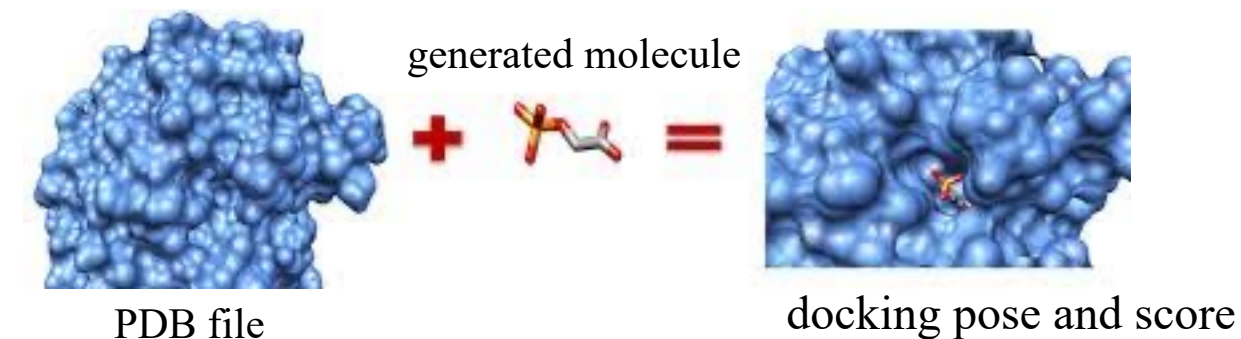
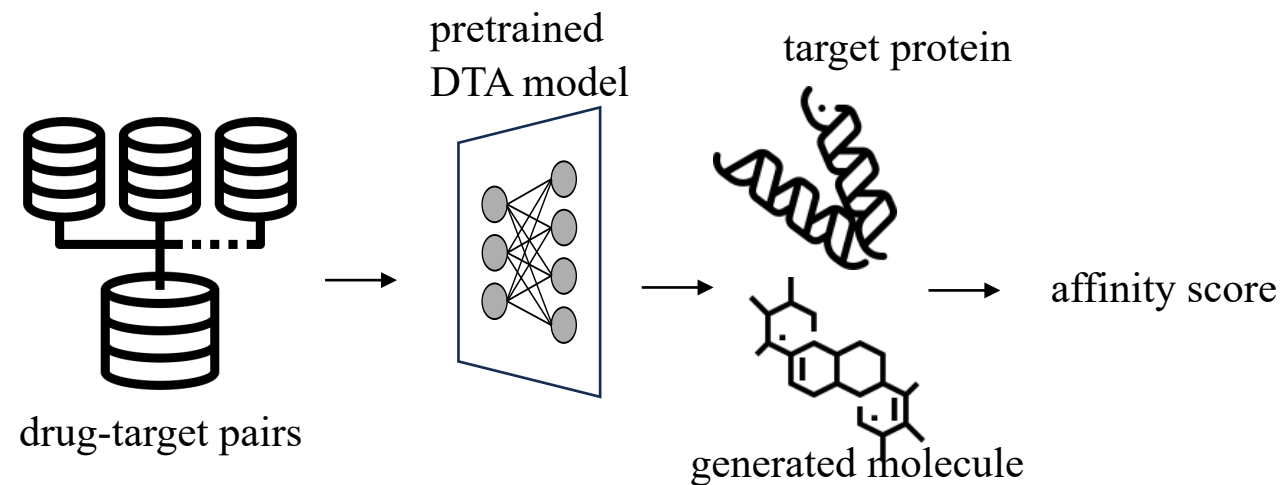
Evaluation Metrics

Basic Prosperities:

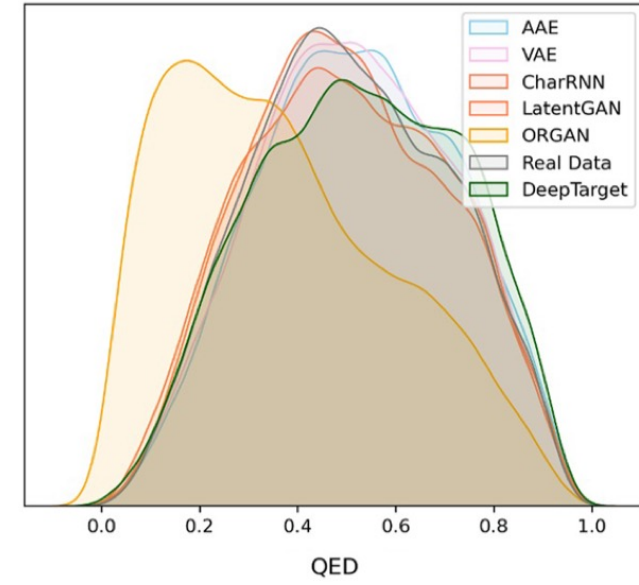
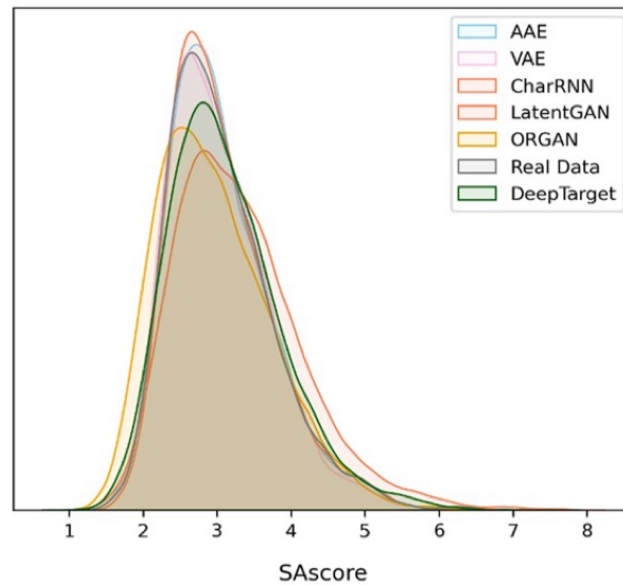
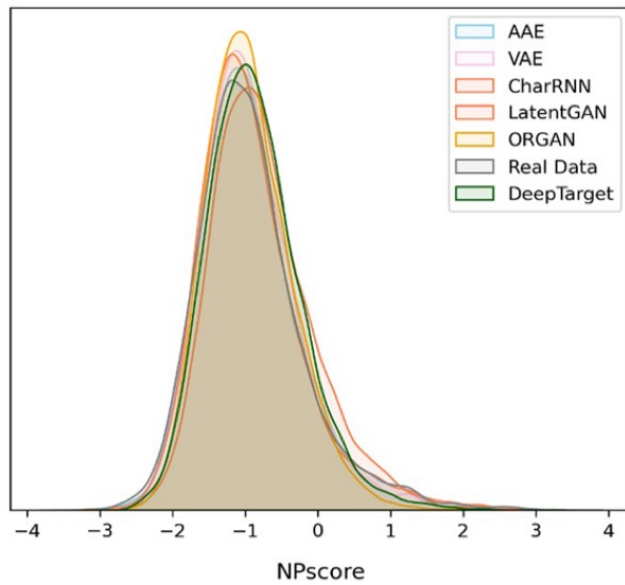
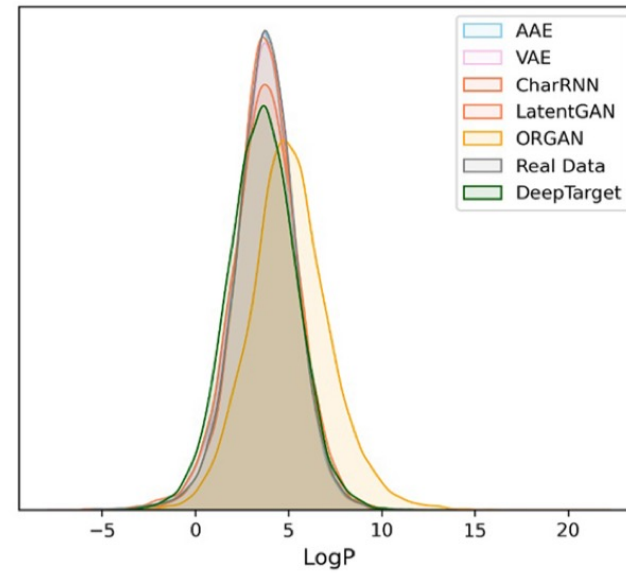
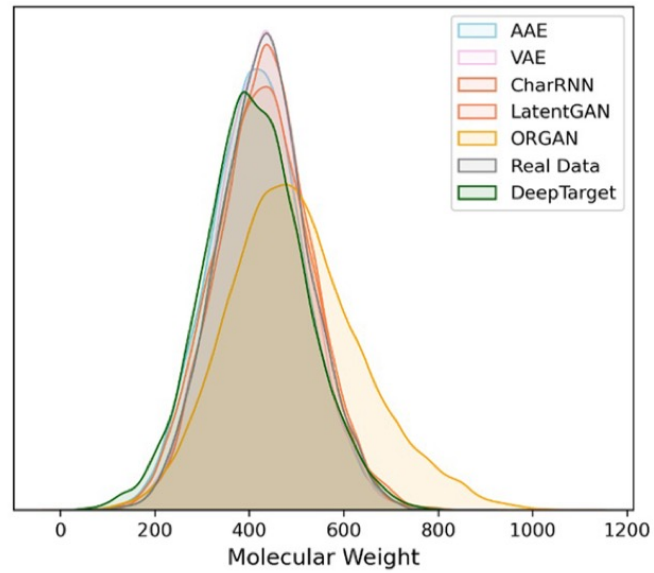
- The generated molecules should satisfy the basic property distribution
- Molecular weight (MW)
- Lipophilicity (LogP)
- Natural product-likeness (NP-likeness)
- Synthetic accessibility score (SAscore)
- Quantitative estimation of drug-likeness (QED)

Interaction with protein target:

- Affinity score
- Docking score



Results



Results

Two target examples:

📄 P14416 · DRD2_HUMAN

Proteinⁱ | D(2) dopamine receptor

Geneⁱ | DRD2

Statusⁱ | 📄 UniProtKB reviewed (Swiss-Prot)

Organismⁱ | [Homo sapiens \(Human\)](#)

<https://www.uniprot.org/uniprotkb/P14416/entry>

📄 P09874 · PARP1_HUMAN

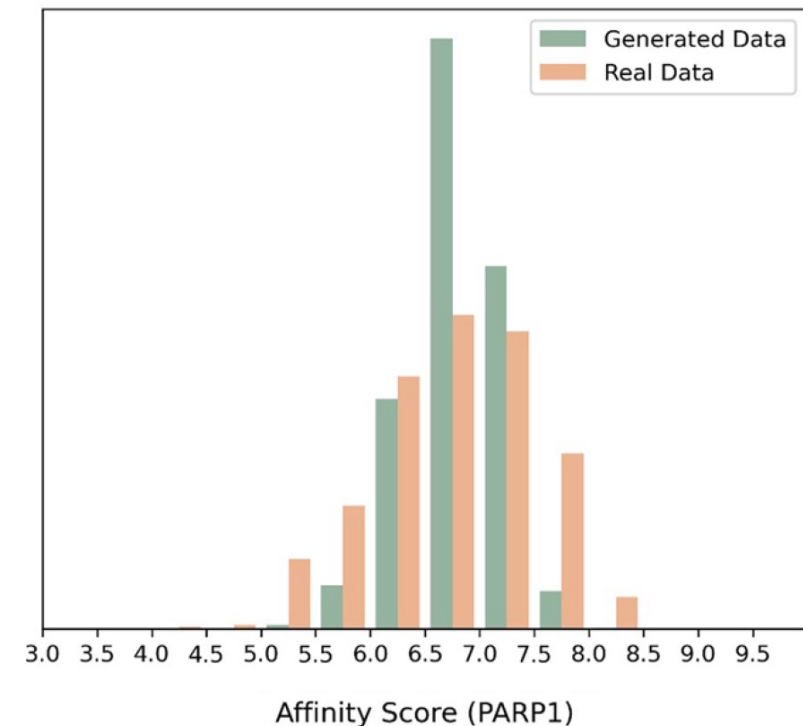
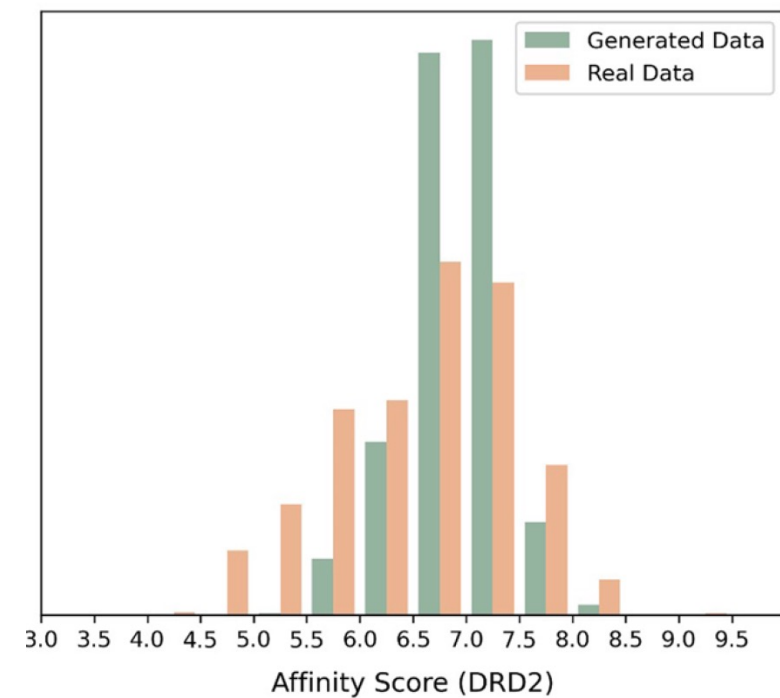
Proteinⁱ | Poly [ADP-ribose] polymerase 1

Geneⁱ | PARP1

Statusⁱ | 📄 UniProtKB reviewed (Swiss-Prot)

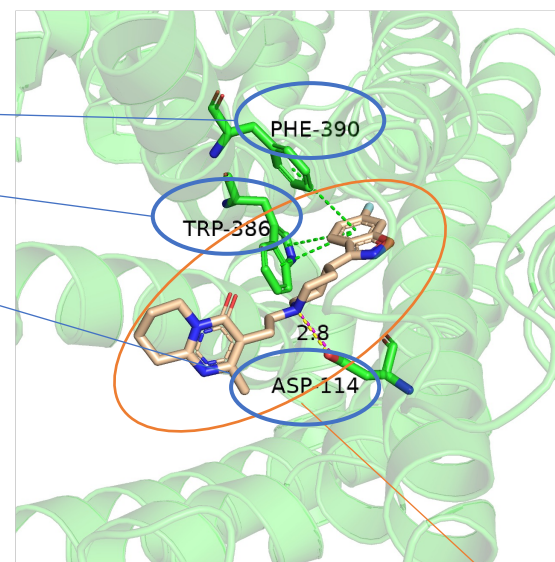
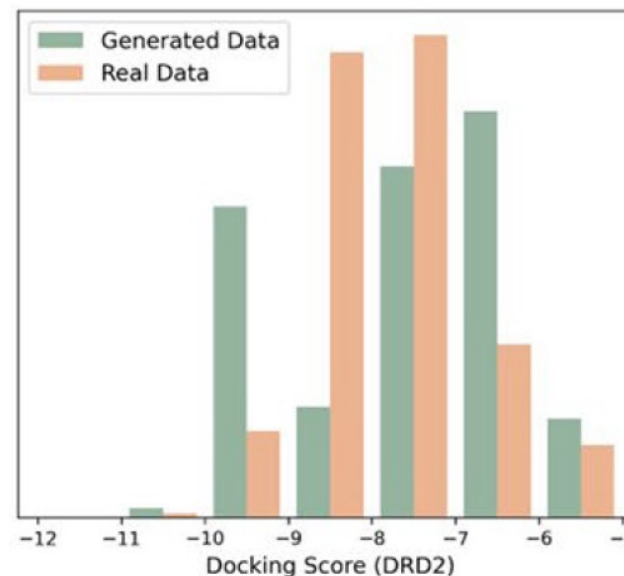
Organismⁱ | [Homo sapiens \(Human\)](#)

<https://www.uniprot.org/uniprotkb/P09874/entry>

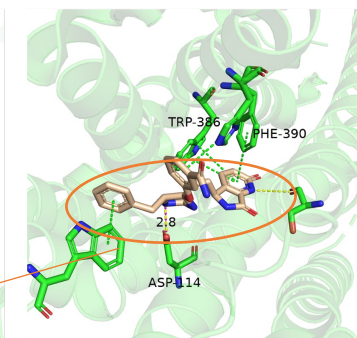
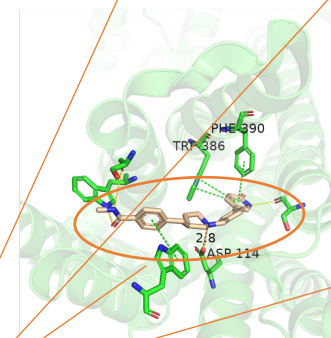
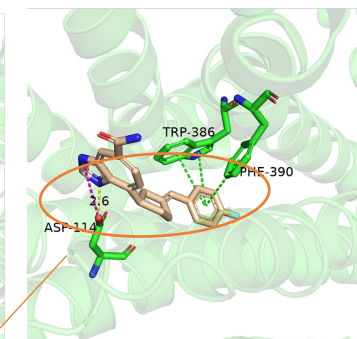
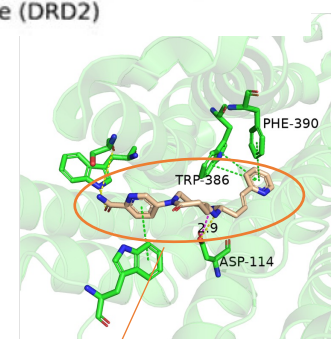


Result: Affinity

- Molecular docking is a computer technique that predicts how drugs and proteins interact to help discover new drugs.
- Three **key residues** interact with ligand.
- The generated molecules could interact with same key residues of real ligand.



Real data



Generated data

ligand molecules



Future outlook

Future outlook

- How to better represent features from protein and small molecules?
- How to further improve the affinity towards the target?
-

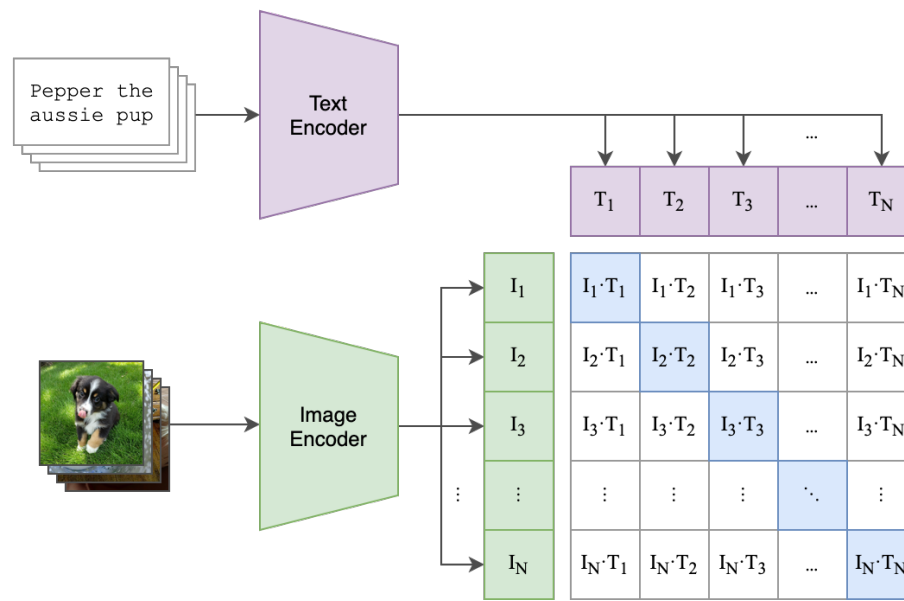
Future outlook

- How to better represent features from protein and small molecules?

Contrastive Learning

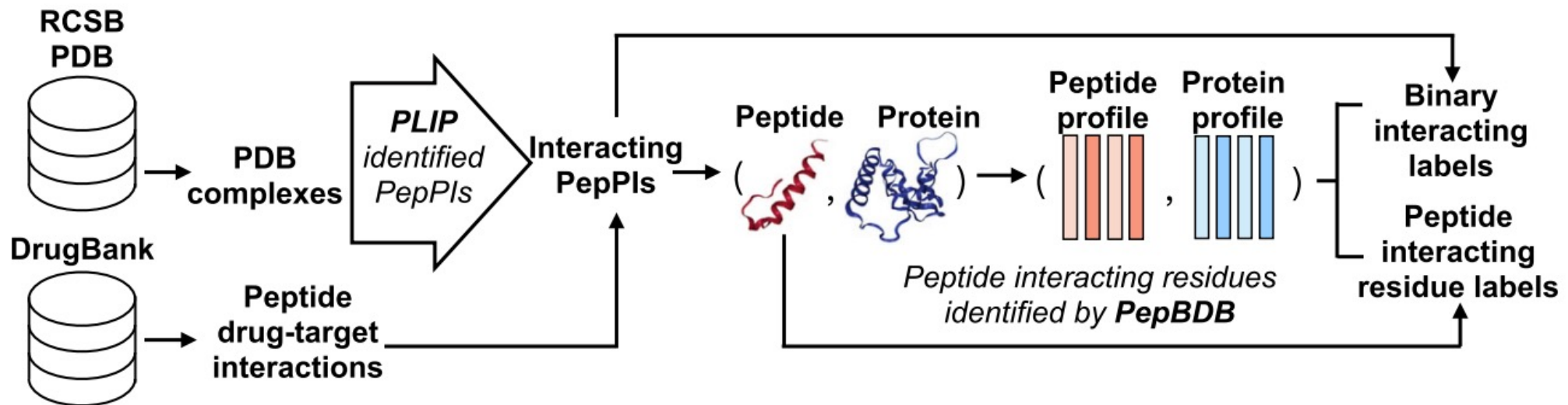
- multi-view
- close the feature distance of similar data

$$\text{score}(f(I), f(T^+)) \gg \text{score}(f(I), f(T^-))$$



Future outlook

- How to further improve the affinity towards the target?



multiple loss function constraints



Thanks