

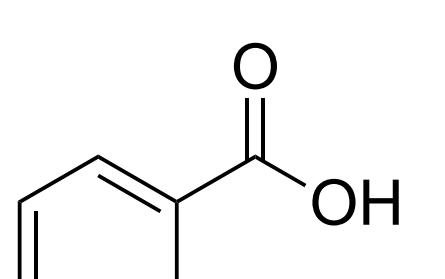
Bio-Isostere Guided Molecular Property Prediction

Anatol Ehrlich, Nils M. Kriege, Christoph Flamm

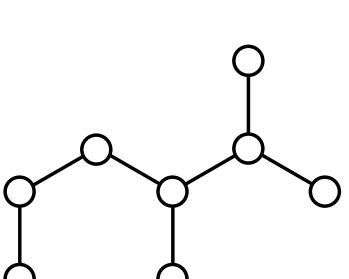
University of Vienna



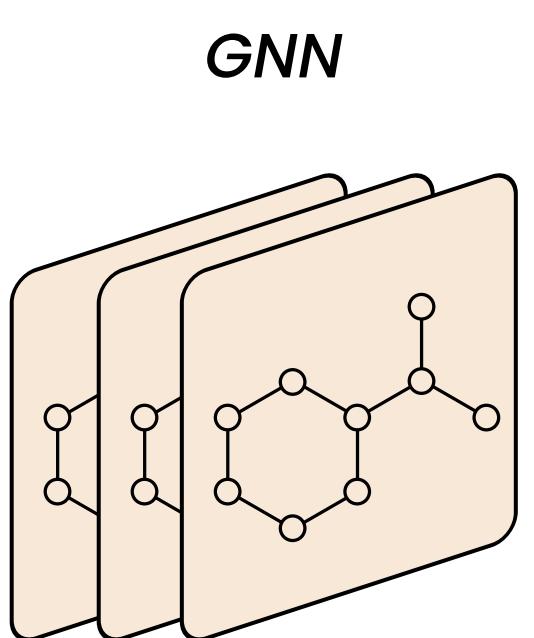
Molecule



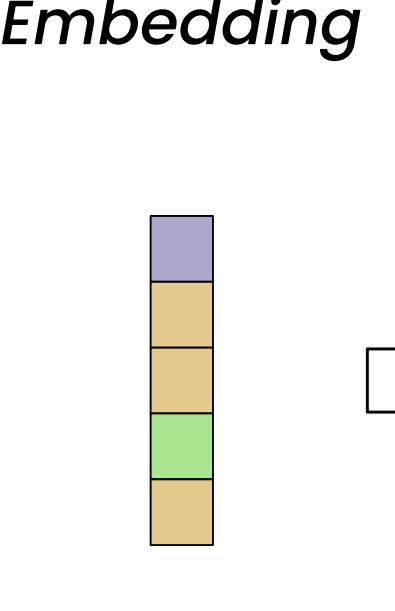
Graph



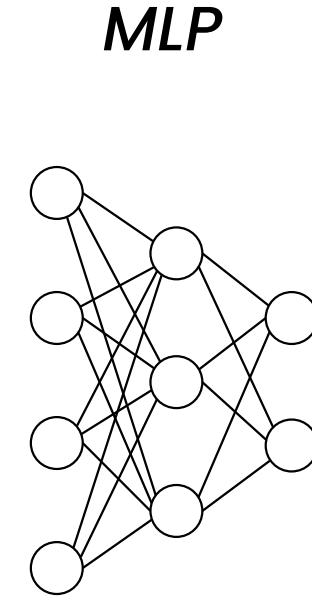
GNN



Embedding



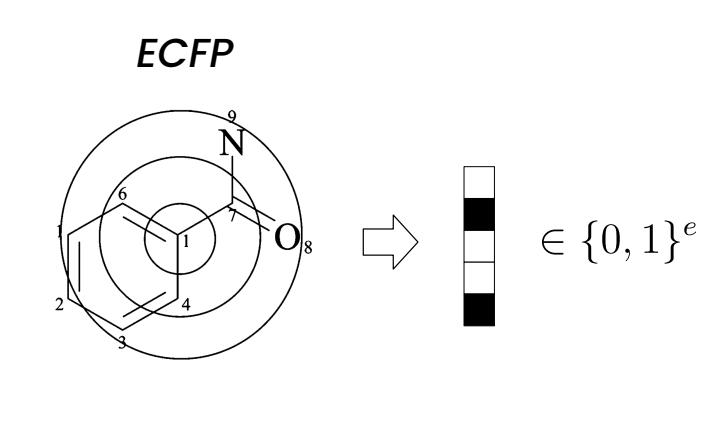
MLP



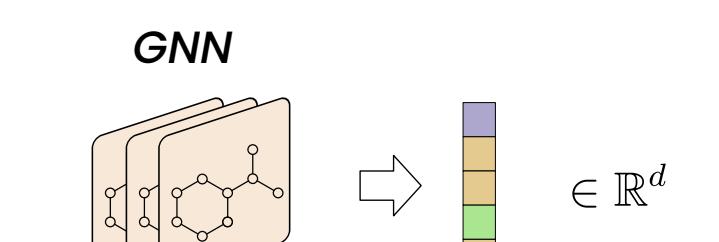
Property

We want to improve this!

ECFP



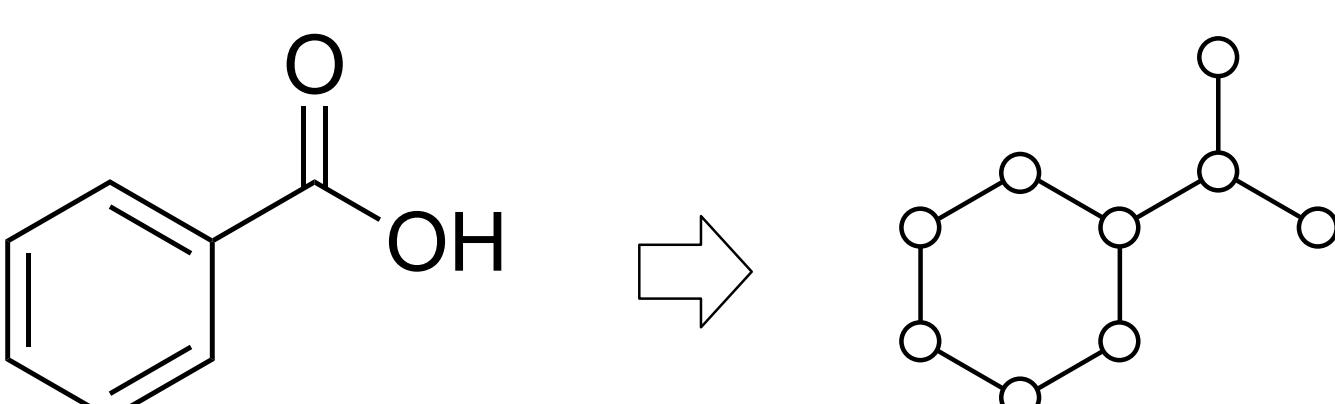
GNN



0. Motivation

- Improve **Molecular Property Predictions**
- Create better **embeddings** of molecules using **Graph Neural Networks (GNNs)**
- Make use of large amounts of **unlabeled data**

1. Molecules as Graphs



- Abstract **Molecules as Graphs**
- Enrich with **Node & Edge Features**

Node Features

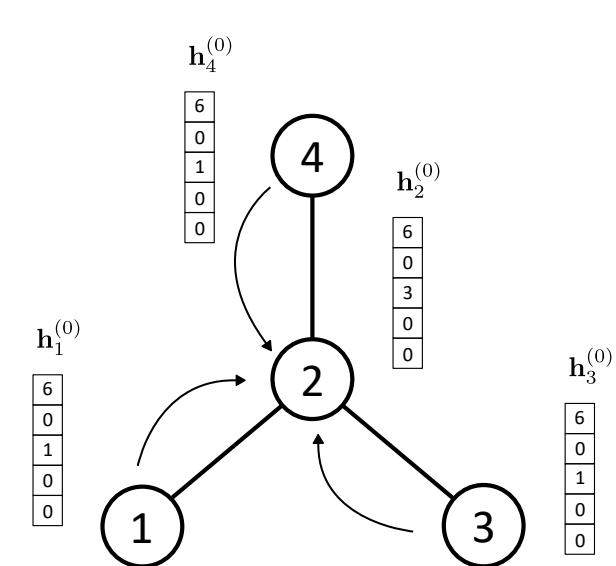
atomic_num,
chirality,
degree,
formal charge
num_hs
num_radical_e
hybridization,
is_aromatic,
is_in_ring

Edge Features

bond_type,
stereoinfo,
is_conjugated

The Approach

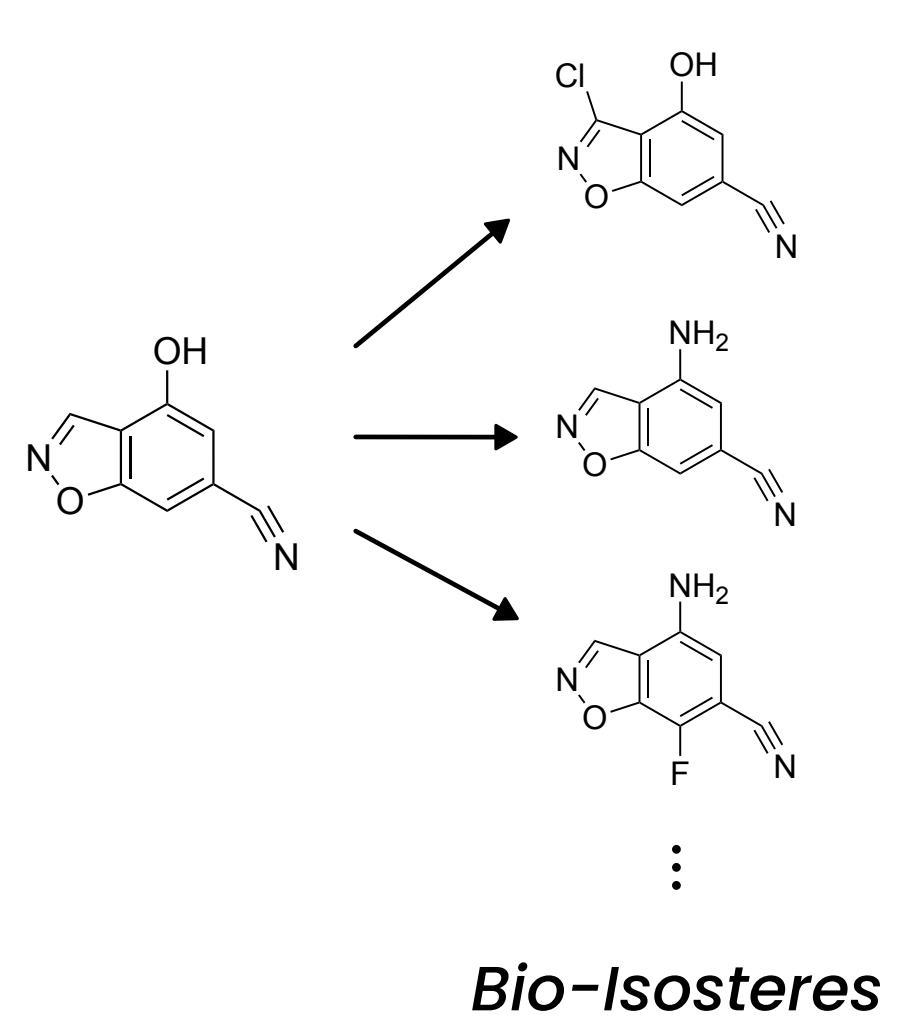
2. Graph Neural Networks (GNN)



Message Passing Framework

- $\mathbf{m}_u^{(i)} = \text{Aggregate}^{(i)} (\{\mathbf{h}_v^{(i)} \mid v \in \mathcal{N}(u)\})$
- $\mathbf{h}_u^{(i+1)} = \text{Combine}^{(i)} (\mathbf{h}_u^{(i)}, \mathbf{m}_u^{(i)})$

3. Data Augmentation



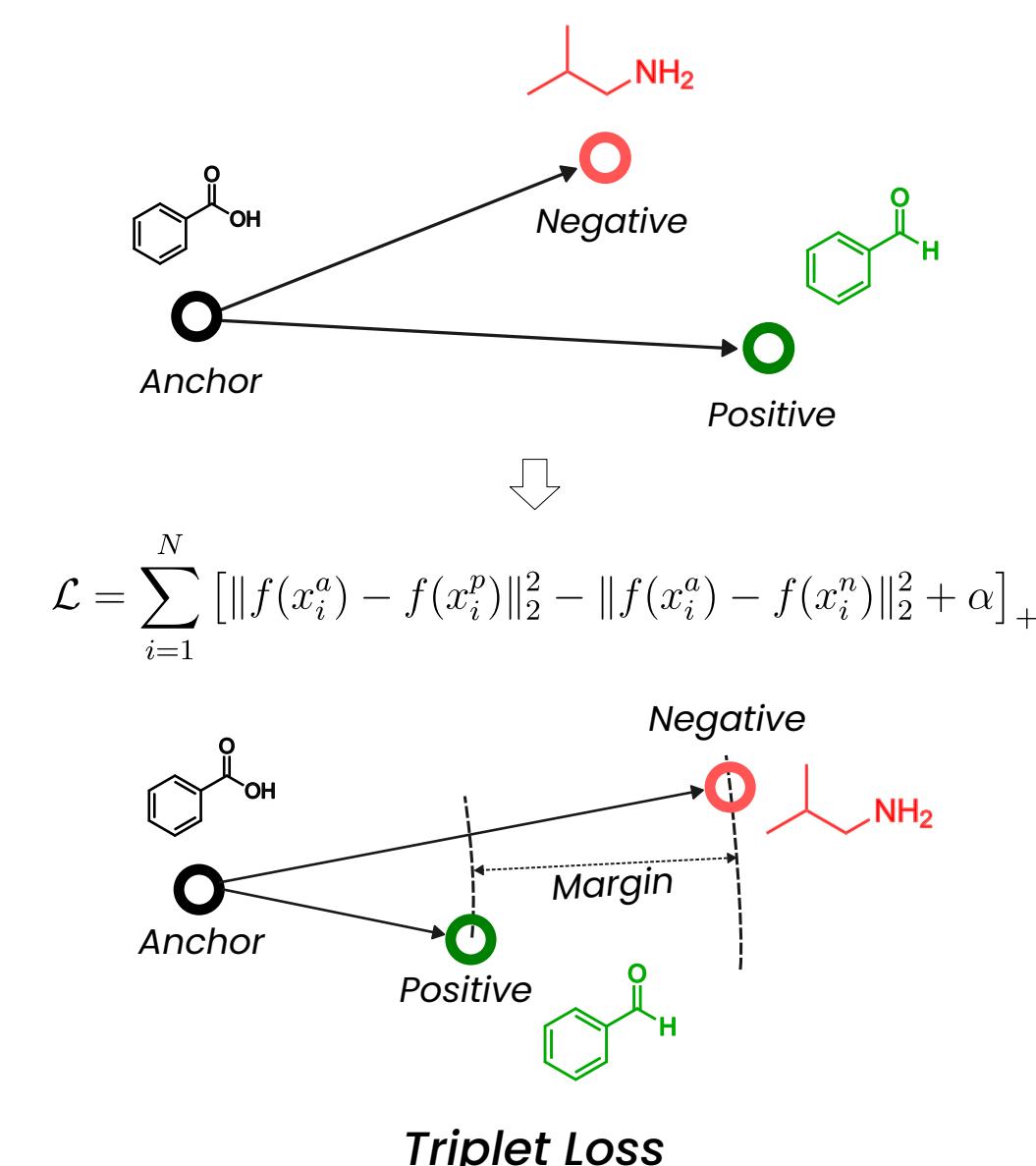
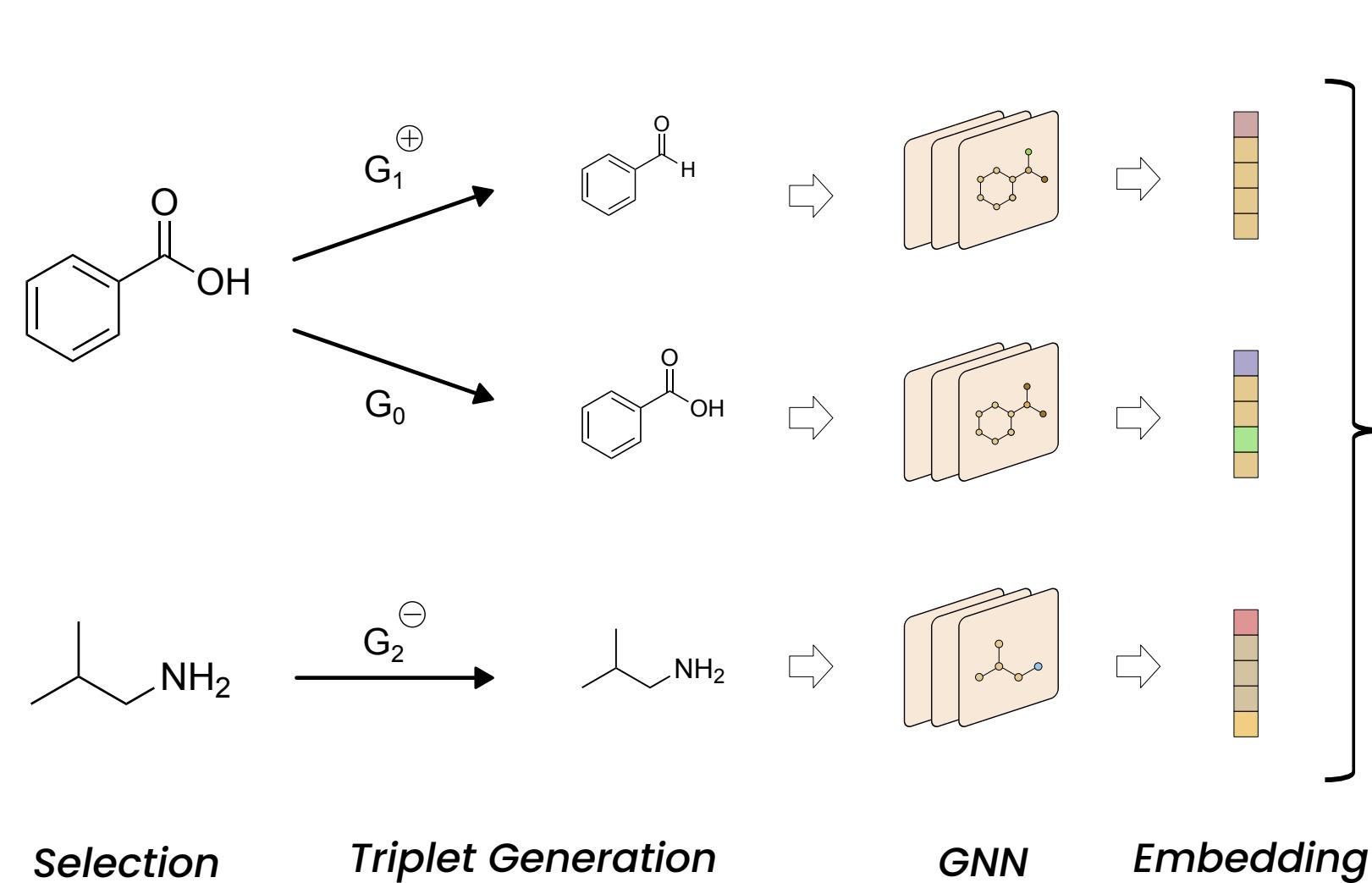
10M Molecules	
label	smiles
0	N#Cc1cc(O)c2oncc2c1
1	Cc1nccc(C=CC(C)Cl)n1
2	O=PCOCC(COC(=O)OP=O)c1SC
3	COc1ccc(Br)cc1SC
4	CSC1(=O)N=NC=C1Cl
5	COCCNc1nccc2oncc2c1F
6	O=C1CCCCC(OCCO)c1
7	Fc1cccc(Oc2cccn2)c1
8	CCOc1c(CC)ccc1OC
9	CCOCCCO[Si](C1)C1CCCCO1
...	

Augmentations

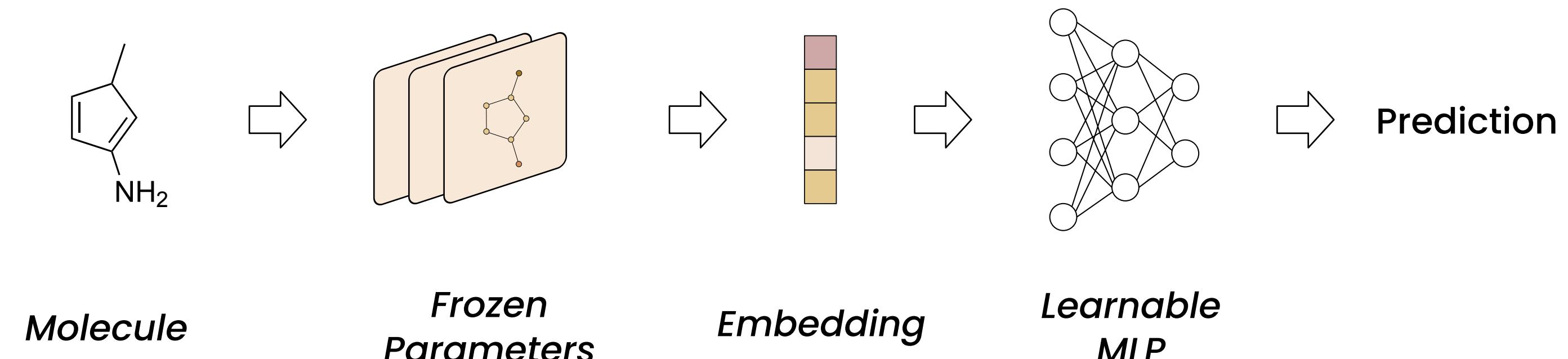
120M Molecules	
label	smiles
0	N#Cc1cc(O)c2oncc2c1
0	N#Cc1cc(O)c2oncc(Cl)c2c1
0	N#Cc1cc(N)c2oncc2c1
0	N#Cc1cc(O)c2oncc2c1F
1	Cc1nccc(C=CCC(C)Cl)n1
1	Cc1nccc(C=CCC(Cl)F)n1
2	O=C1CCCCC(OCCO)c1
2	O=PCOCC(COC(=O)OP=O)c1SC
2	O=PCOCC(COC(=O)OP=O)c1SC
2	O=PCOCC(COC(=O)OP=O)c1SC
...	

- Input: **10M Molecules**
- 130 **Augmentation Rules**
- Out: **120M Molecules**

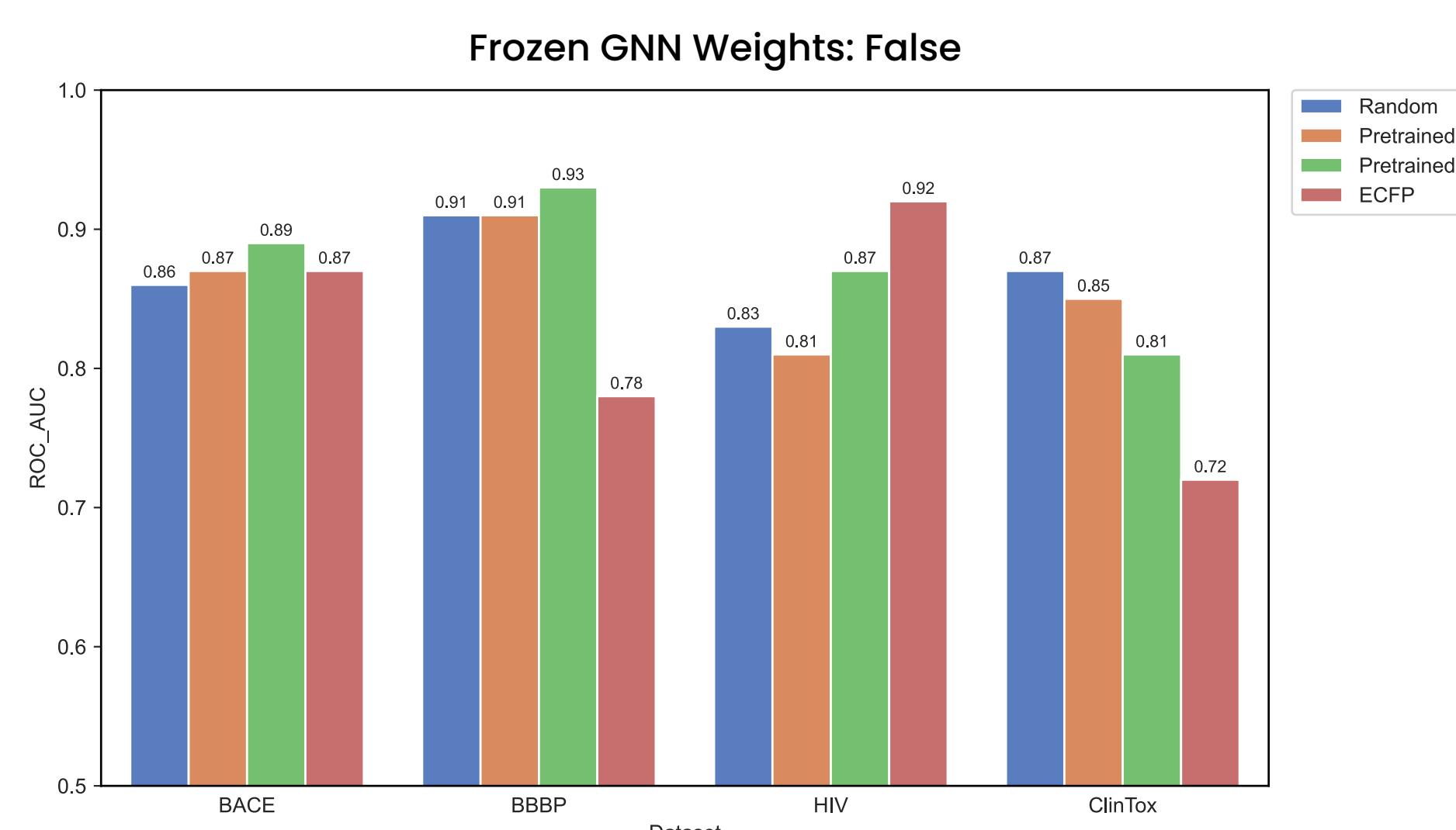
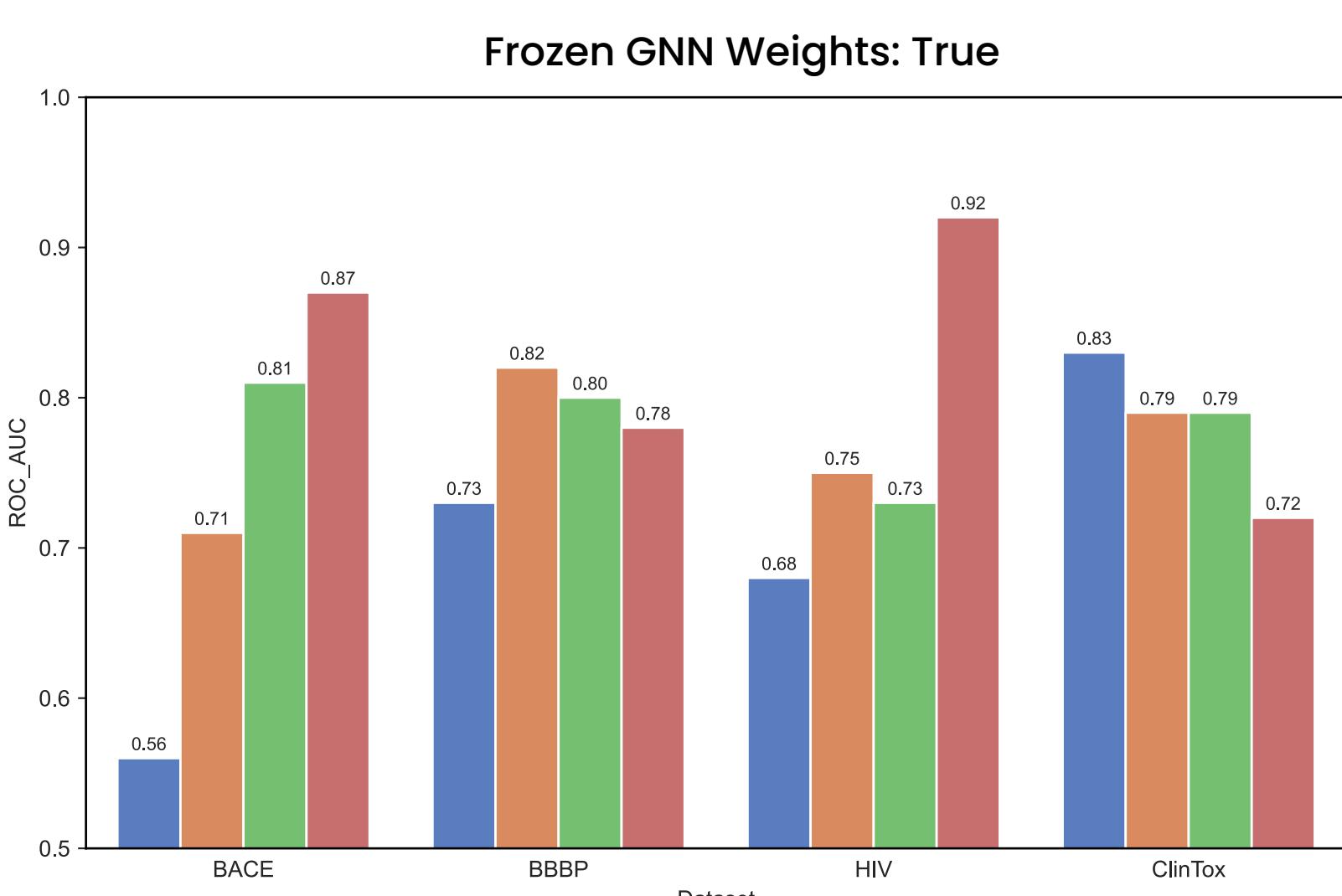
4. Self-Supervised Pre-Training



5. Fine-Tuning



The Results



6. Conclusion

- Pre-Training **improves embeddings** vs. randomly initialized GNNs
- ECFPs still have **comparatively good embeddings**
- Further investigation of embeddings is required