

MolecularRNN: Generating realistic molecular graphs with optimized properties

Mariya Popova, Olexandr Isayev

Carnegie Mellon University

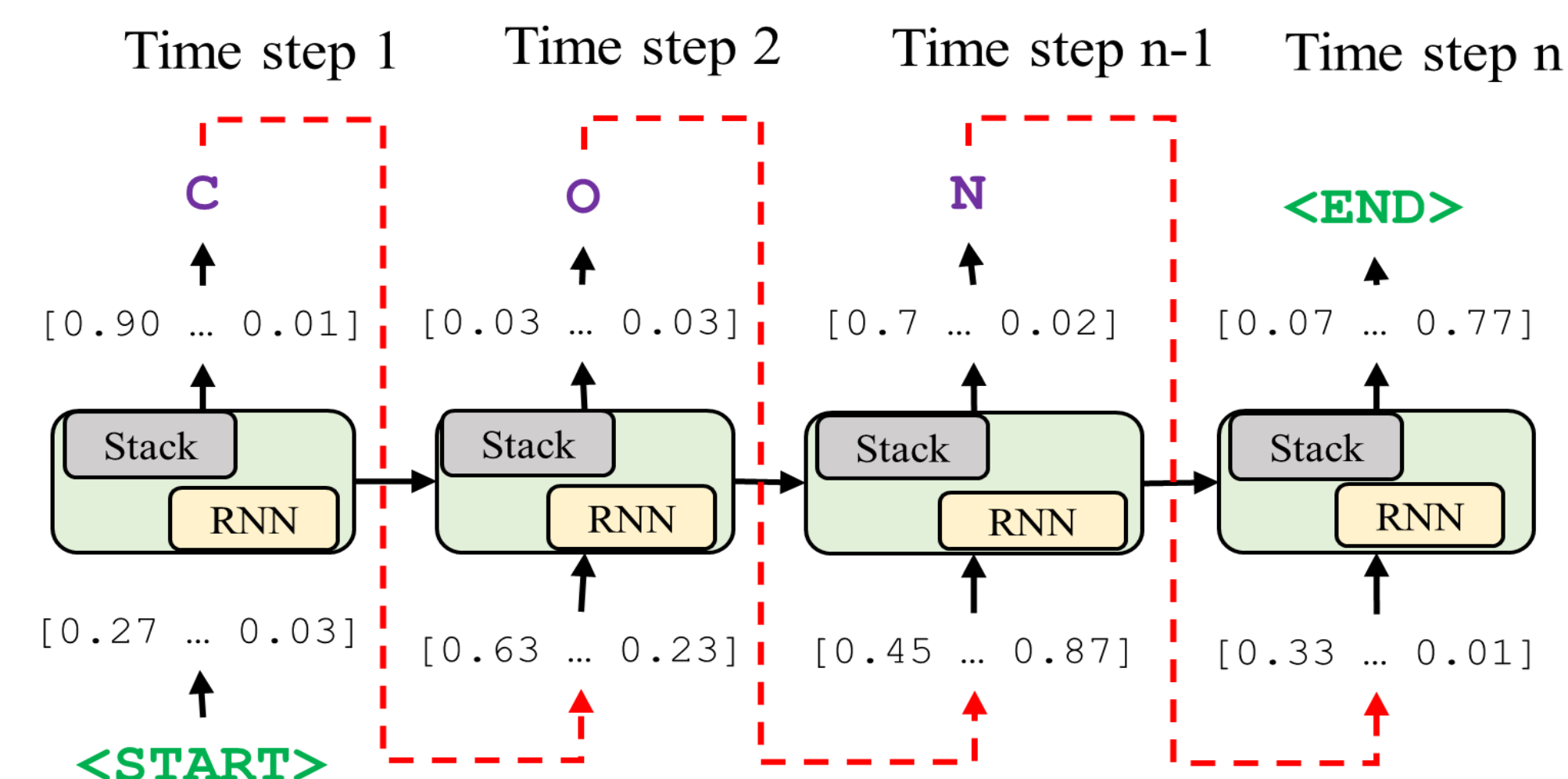
ABSTRACT

Designing new molecules with a set of predefined properties is a core problem in modern drug discovery and development. We present **MolecularRNN**, the graph recurrent generative model for molecular structures. Our model generates diverse realistic molecular graphs after likelihood pretraining on a big database of molecules. We perform an analysis of our pretrained models on large-scale generated datasets of 1 million samples. Further, the model is tuned with policy gradient algorithm, provided a critic that estimates the reward for the property of interest. We show a significant distribution shift to the desired range for lipophilicity, drug-likeness, and melting point outperforming state-of-the-art works. With the use of rejection sampling based on valency constraints, our model yields 100% validity. Moreover, we show that invalid molecules provide a rich signal to the model through the use of structure penalty in our reinforcement learning pipeline.

MODEL: SMILES vs Graphs

SMILES

$$p(s_t = a_i | s_{<t}; \theta) = f_i(s_{<t} | \theta)$$

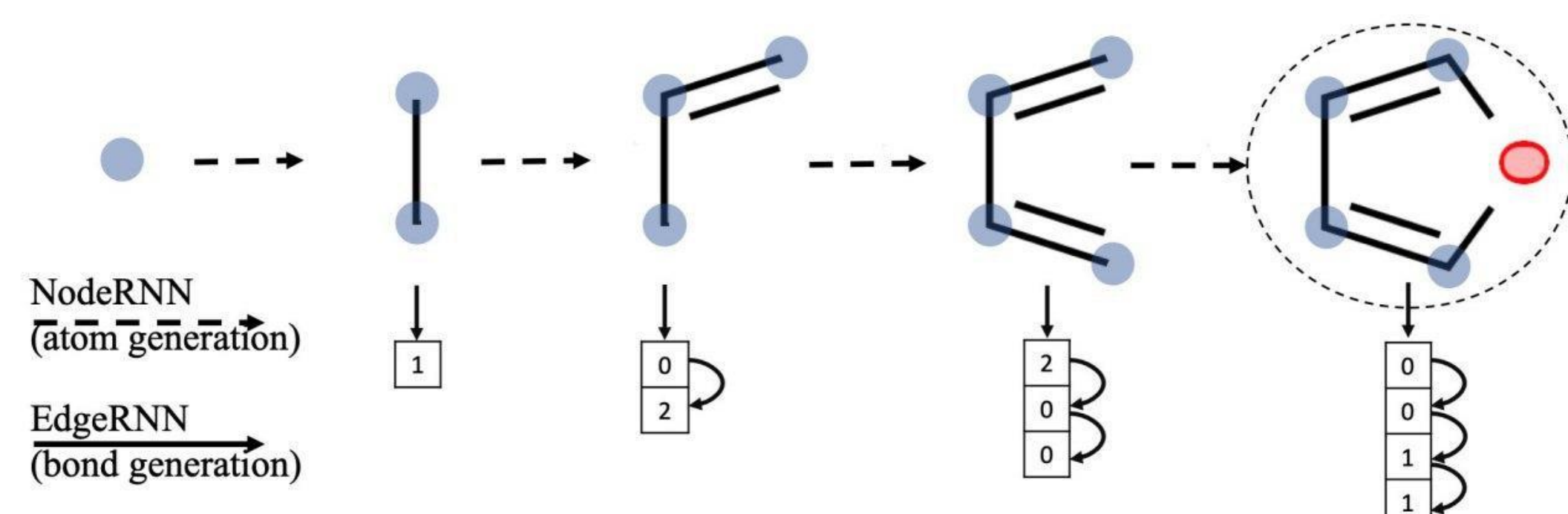


Graphs

- Adjacency matrix with bond types + node list
- g – graph, S – edges, C – vertices

$$p(g_t | g_{<t}; \theta) = p(S_t, C_t | S_{<t}, C_{<t}; \theta)$$

$$= p(S_t, C_t | S_{<t}, C_{<t}; \theta) \cdot \prod_{j=0}^t p(S_{tj} | S_{<t}, C_{<t}, C_t; \theta)$$



MolecularRNN model

- Bond types $S_{i,j}^\pi \in \{0, 1, 2, 3\}$ – no, single, double, triple bonds
- Atom types $C_i^\pi \in \{1, 2, \dots, K\}$ – oxygen, nitrogen, chlorine, etc.
- Nodes are ordered in BFS order permutation π starting from arbitrary Carbon atom

$$p(S^\pi, C^\pi) = \prod_{i=1}^{n+1} p(C_i^\pi | S_{<i}^\pi, C_{<i}^\pi) p(S_i^\pi | C_i^\pi, S_{<i}^\pi, C_{<i}^\pi)$$

- Model has 4 blocks – *NodeRNN*, *NodeMLP*, *EdgeRNN* and *EdgeMLP*
- Learnable embeddings dictionary for atom types and bond types

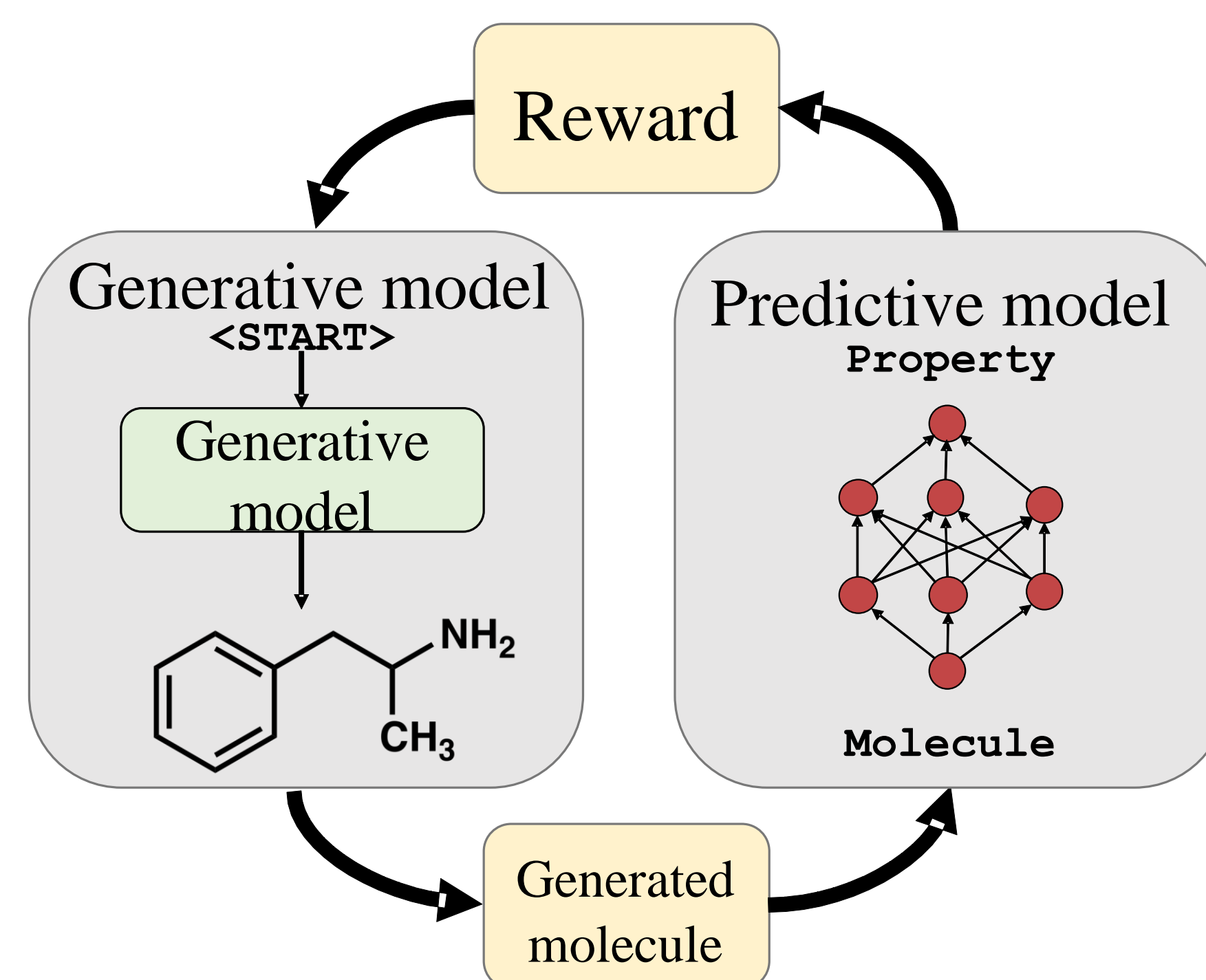
Valency-based rejection sampling

- Idea: reject bonds that violate per-atom valency constraints
- 100% of generated molecules are valid

$$\sum_j A_{i,j}^\pi + k \leq \text{valency}_{C_i^\pi} \ \& \ \sum_i A_{i,j}^\pi + k \leq \text{valency}_{C_j^\pi}$$

Property optimization

- Reward $r(s_N)$ is assigned when the molecule is generated



$$L(\theta) = - \sum_{i=1}^N r(s_N) \cdot \gamma^i \cdot \log p(s_i | s_{i-1}; \theta)$$

RESULTS

Unsupervised likelihood pretraining

- We considered 9 atom type – C, N, O, F, P, S, Cl, Br, I
- Graph size between 10 and 50 heavy atoms
- 3 training datasets:
 - ChEMBL (1.5M bioactive molecules)
 - ZINC 250k (250k molecules randomly selected from commercial database ZINC)
 - MOSES (2M molecules filtered from ZINC)

Table 1. Statistics for 1 million molecules generated by 3 models pretrained on 3 training datasets

Training set	Valid	Unique	Novel	IntDiv (p=1)	IntDiv (p=2)	SA score	QED
ChEMBL	100 %	99.2%	99.3 %	0.895	0.890	3.67 ± 1.20	0.56 ± 0.20
ZINC 250k	100 %	99.8 %	100 %	0.892	0.887	3.60 ± 1.01	0.68 ± 0.16
MOSES	100 %	99.4 %	100 %	0.881	0.876	3.24 ± 0.97	0.74 ± 0.14

Melting temperature optimization

- Model pretrained on ChEMBL
- Melting temperature predicted by NN trained on 40k samples
- Reward function

$$r(\text{mol}) = \exp(T_{\text{pred}}(\text{mol}) + 1)$$

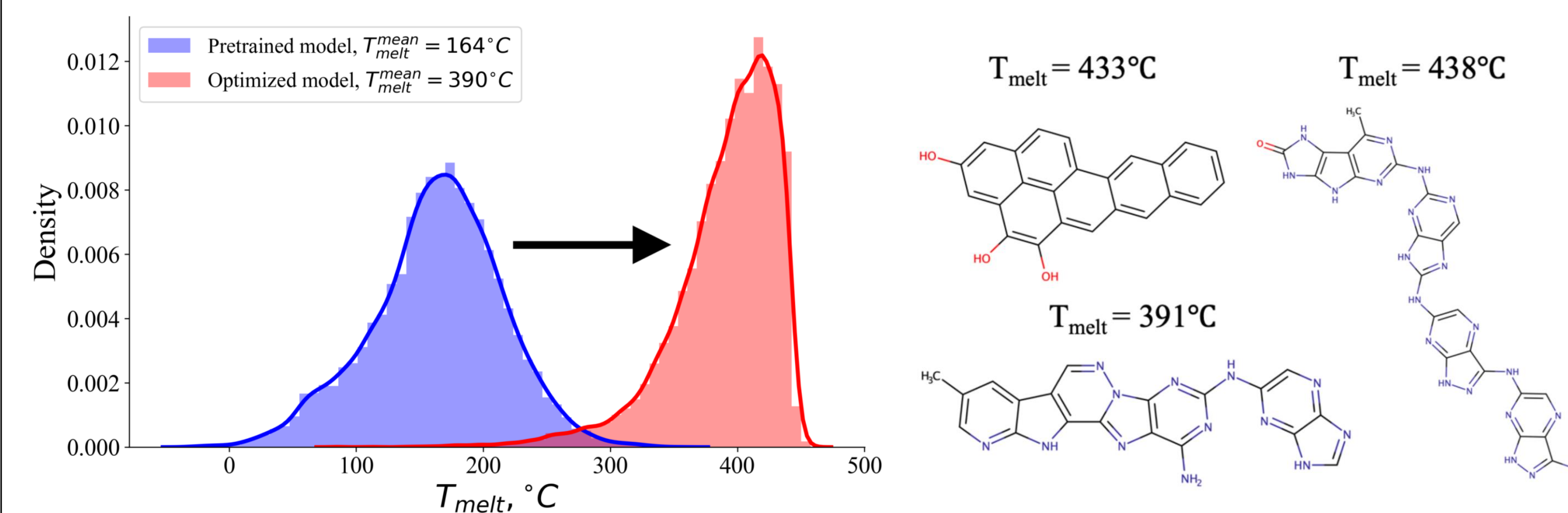


Figure 1. Melting temperature maximization. (A) Distribution of predicted melting temperature for base and optimized models (B) Examples of generated molecules with highest melting temperature.

CONCLUSIONS

- Molecular graph recurrent model, MolecularRNN, for direct generation of realistic molecular graph structures shows high validity/uniqueness/novelty
- Valency-based rejection sampling method during inference produces 100% valid molecules, and the structural penalty during training for atoms violating valency constraints
- Target property optimization with reinforcement learning for improves molecular properties

REFERENCES

- M Popova et al- arXiv preprint arXiv:1905.13372, 2019
- <https://github.com/Mariawelt/OpenChem>

ACKNOWLEDGEMENTS

O.I. acknowledges support from DOD-ONR (N00014-16-1-2311), National Science Foundation (NSFCHE-1802789), and Eshelman Institute for Innovation (EII) awards. M.P. acknowledges The Molecular Sciences Software Institute (MolSSI) Software Fellowship and NVIDIA Graduate Fellowship. We gratefully acknowledge the support and hardware donation from NVIDIA Corporation