

# Uncertainty-Aware Multi-Objective Reinforcement Learning-Guided Diffusion Models for 3D De Novo Molecular Design

---

**Lianghong Chen**, Dongkyu Kim, Mike Domaratzki, Pingzhao Hu\*

Western University, Ontario, Canada

\*Corresponding author (Email: [phu49@uwo.ca](mailto:phu49@uwo.ca))

# Background

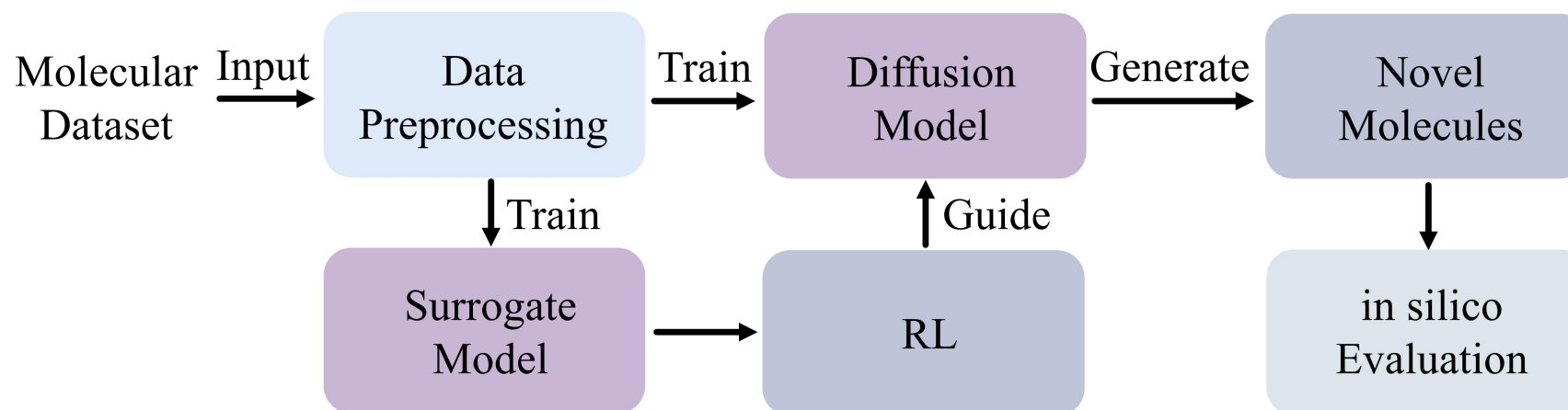
- Lack of a unified framework for generating molecules with multiple drug-relevant properties and strong protein binding
- Reinforcement Learning (RL) has been explored for 1D/2D molecular generative models, but 3D molecular diffusion models with RL guidance remain underexplored
- Previous RL methods for diffusion model optimization suffer from challenges such as reward sparsity and mode collapse
- When multiple conflicting optimization objectives exist, it difficult to balance these objectives, and algorithms tend to bias toward objectives that are easier to optimize

# Contributions

- **First unified framework** integrating RL, diffusion models, 3D molecular generation, and uncertainty-based multi-objective optimization
- **Uncertainty-aware reward** design enables **balanced** multi-objective optimization for drug-relevant properties
- Three auxiliary mechanisms address RL–diffusion challenges such as mode collapse and reward sparsity:
  - **Reward boosting** for validity, uniqueness, and novelty
  - **Diversity penalty** to enhance exploration
  - **Dynamic cutoff** strategy for stable, adaptive reward signals
- Consistent performance **gains** across datasets

# Study Design

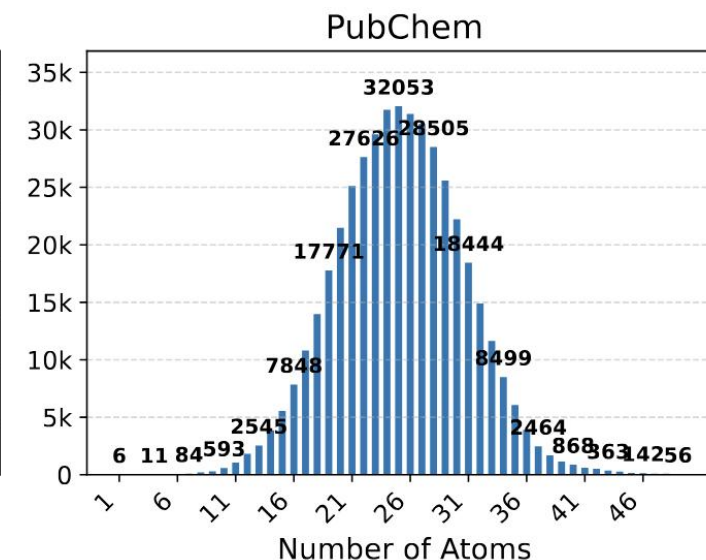
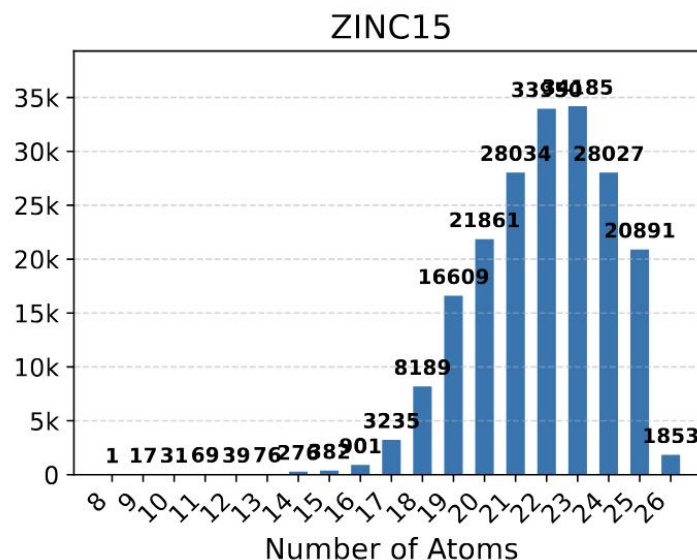
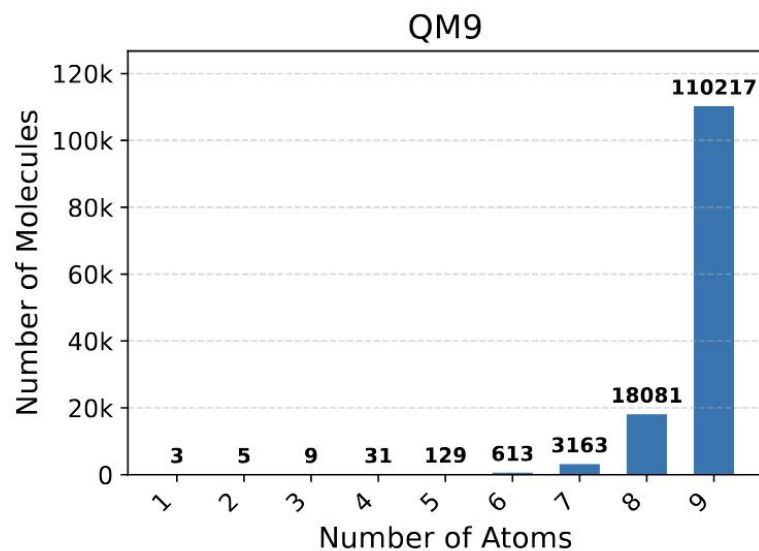
- Train surrogate models to predict uncertainty
- Pretrain the diffusion model to learn 3D molecular generation
- RL guides the diffusion model to generate molecules with desired drug-relevant objectives
- Perform in silico evaluation to further verify the drug potential and binding stability



# Datasets

- QM9: small-sized molecules used for basic 3D molecular generation tasks
- ZINC15: medium-sized, drug-like molecules
- PubChem: large-scale collection of diverse and complex molecules

| Dataset | Unique Atom Types           |
|---------|-----------------------------|
| QM9     | C, N, O, F                  |
| ZINC15  | C, N, O, F, P, S, Cl, Br, I |
| PubChem | C, N, O, F, P, S, Cl, Br, I |



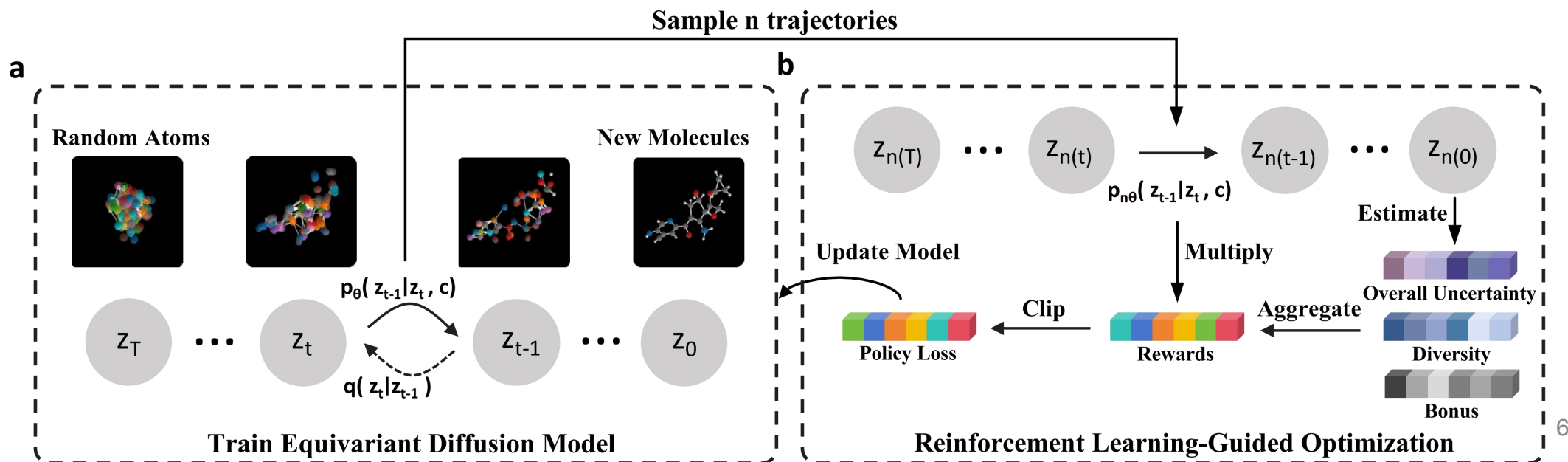
# Method

- Generate n candidate molecules and evaluate them to obtain RL rewards
- Reward design:

$$R_{\text{total}}(m; \delta_1, \dots, \delta_k, t_{\text{episode}}) = U_{\text{multi}}(m; \delta_1, \dots, \delta_k) \cdot R_{\text{bonus}}(m) - \lambda(t_{\text{episode}}) \cdot D(m) \quad (1)$$

Annotations for Equation (1):

- $U_{\text{multi}}(m; \delta_1, \dots, \delta_k)$ : generated molecule (points to  $m$ ), uncertainty (points to  $U$ ), # properties (points to  $\delta$ )
- $R_{\text{bonus}}(m)$ : quality-based bonus (points to  $R$ )
- $\lambda(t_{\text{episode}})$ : property threshold (points to  $\lambda$ )
- $D(m)$ : similarity-based penalty (points to  $D$ )



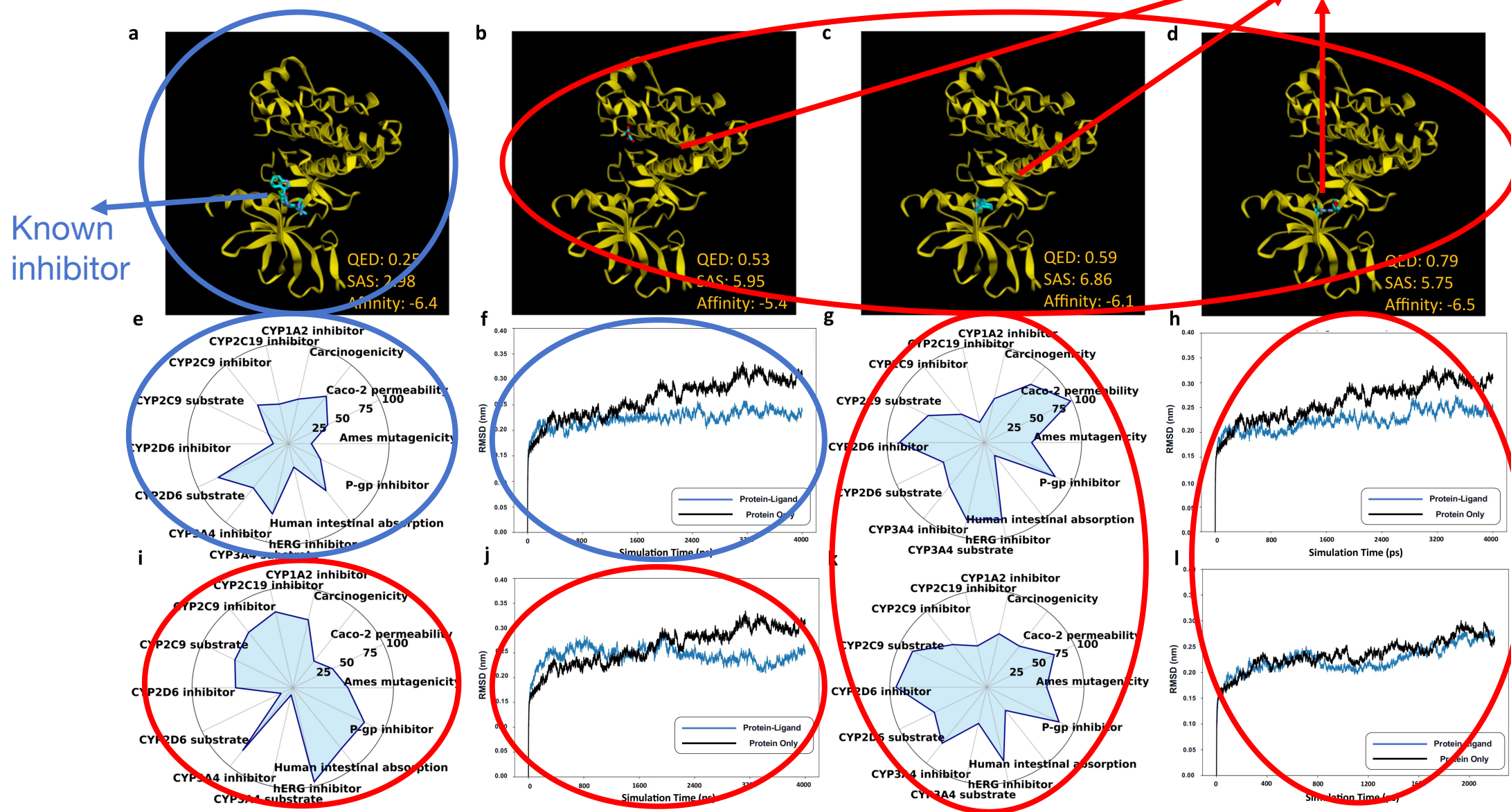


# Results

| Dataset | Method      | Val (%) ( $\uparrow$ )  | Uni (%) ( $\uparrow$ )   | Nov (%) ( $\uparrow$ )   | VUN (%) ( $\uparrow$ )  | ASta (%) ( $\uparrow$ ) | MSta (%) ( $\uparrow$ ) | Top (%) ( $\uparrow$ )  |
|---------|-------------|-------------------------|--------------------------|--------------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| QM9     | W/O RL      | 88.55 $\pm$ 0.65        | <b>97.57</b> $\pm$ 0.30  | 99.75 $\pm$ 0.15         | 86.19 $\pm$ 1.02        | 99.34 $\pm$ 0.11        | 95.90 $\pm$ 0.34        | 25.17 $\pm$ 1.17        |
|         | SFT-PG      | 88.57 $\pm$ 1.33        | 96.80 $\pm$ 0.33         | 99.80 $\pm$ 0.15         | 85.57 $\pm$ 1.60        | 99.24 $\pm$ 0.10        | 95.62 $\pm$ 0.58        | 25.58 $\pm$ 1.65        |
|         | DDPO-SF     | 88.65 $\pm$ 1.05        | 97.39 $\pm$ 0.48         | 99.79 $\pm$ 0.20         | 86.16 $\pm$ 1.59        | 99.25 $\pm$ 0.09        | 95.50 $\pm$ 0.34        | 25.65 $\pm$ 1.51        |
|         | DDPO-IS     | 88.82 $\pm$ 1.03        | 96.59 $\pm$ 0.71         | 99.39 $\pm$ 0.04         | 85.27 $\pm$ 1.30        | 97.31 $\pm$ 0.09        | 86.10 $\pm$ 0.64        | 25.77 $\pm$ 1.29        |
|         | DPOK        | 88.10 $\pm$ 0.37        | 97.52 $\pm$ 0.42         | <b>99.81</b> $\pm$ 0.15  | 85.75 $\pm$ 0.80        | 99.18 $\pm$ 0.03        | 95.28 $\pm$ 0.18        | 25.20 $\pm$ 1.44        |
|         | <b>Ours</b> | <b>98.17</b> $\pm$ 0.07 | 90.90 $\pm$ 0.72         | 99.63 $\pm$ 0.04         | <b>88.90</b> $\pm$ 0.68 | <b>99.87</b> $\pm$ 0.03 | <b>99.17</b> $\pm$ 0.27 | <b>28.33</b> $\pm$ 0.61 |
| ZINC15  | W/O RL      | 30.05 $\pm$ 1.34        | <b>100.00</b> $\pm$ 0.00 | <b>100.00</b> $\pm$ 0.00 | 30.05 $\pm$ 1.34        | 88.36 $\pm$ 0.40        | 12.00 $\pm$ 1.33        | 8.02 $\pm$ 0.46         |
|         | SFT-PG      | 41.25 $\pm$ 1.48        | <b>100.00</b> $\pm$ 0.00 | <b>100.00</b> $\pm$ 0.00 | 41.25 $\pm$ 1.48        | 91.71 $\pm$ 0.08        | 25.55 $\pm$ 0.65        | 10.43 $\pm$ 0.73        |
|         | DDPO-SF     | 30.25 $\pm$ 1.56        | <b>100.00</b> $\pm$ 0.00 | <b>100.00</b> $\pm$ 0.00 | 30.25 $\pm$ 1.56        | 88.37 $\pm$ 0.40        | 11.97 $\pm$ 1.41        | 8.05 $\pm$ 0.61         |
|         | DDPO-IS     | 30.47 $\pm$ 1.39        | <b>100.00</b> $\pm$ 0.00 | <b>100.00</b> $\pm$ 0.00 | 30.47 $\pm$ 1.39        | 88.35 $\pm$ 0.44        | 12.02 $\pm$ 1.57        | 8.13 $\pm$ 0.60         |
|         | DPOK        | 30.13 $\pm$ 2.28        | <b>100.00</b> $\pm$ 0.00 | <b>100.00</b> $\pm$ 0.00 | 30.13 $\pm$ 2.28        | 88.42 $\pm$ 0.61        | 12.01 $\pm$ 1.38        | 8.02 $\pm$ 0.78         |
|         | <b>Ours</b> | <b>99.02</b> $\pm$ 0.46 | 99.75 $\pm$ 0.06         | <b>100.00</b> $\pm$ 0.00 | <b>98.77</b> $\pm$ 0.49 | <b>99.86</b> $\pm$ 0.03 | <b>98.08</b> $\pm$ 0.63 | <b>33.40</b> $\pm$ 0.89 |
| PubChem | W/O RL      | 7.18 $\pm$ 4.78         | 99.67 $\pm$ 0.65         | <b>100.00</b> $\pm$ 0.00 | 7.17 $\pm$ 4.80         | 94.51 $\pm$ 0.16        | 38.18 $\pm$ 0.92        | 2.23 $\pm$ 1.65         |
|         | SFT-PG      | 7.47 $\pm$ 1.40         | 99.57 $\pm$ 0.85         | <b>100.00</b> $\pm$ 0.00 | 7.44 $\pm$ 1.37         | 82.99 $\pm$ 0.49        | 33.25 $\pm$ 0.34        | 2.03 $\pm$ 0.76         |
|         | DDPO-SF     | 7.98 $\pm$ 2.96         | <b>100.00</b> $\pm$ 0.00 | <b>100.00</b> $\pm$ 0.00 | 7.98 $\pm$ 2.96         | 94.49 $\pm$ 0.98        | 44.22 $\pm$ 0.32        | 2.40 $\pm$ 0.37         |
|         | DDPO-IS     | 10.50 $\pm$ 6.19        | 99.90 $\pm$ 0.20         | <b>100.00</b> $\pm$ 0.00 | 10.48 $\pm$ 6.16        | 95.36 $\pm$ 0.99        | 45.37 $\pm$ 1.85        | 2.52 $\pm$ 1.22         |
|         | DPOK        | 7.65 $\pm$ 1.75         | 99.67 $\pm$ 0.64         | <b>100.00</b> $\pm$ 0.00 | 7.62 $\pm$ 1.76         | 94.51 $\pm$ 0.20        | 38.17 $\pm$ 0.64        | 2.42 $\pm$ 0.46         |
|         | <b>Ours</b> | <b>16.23</b> $\pm$ 9.72 | <b>100.00</b> $\pm$ 0.00 | <b>100.00</b> $\pm$ 0.00 | <b>16.23</b> $\pm$ 9.72 | <b>99.04</b> $\pm$ 0.13 | <b>88.65</b> $\pm$ 0.59 | <b>2.97</b> $\pm$ 1.60  |

Val", "Uni", and "Nov" represent the percentages of valid, unique, and novel molecules, respectively. "VUN" is their joint metric computed as Val  $\times$  Uni  $\times$  Nov, representing the percentage of molecules that are simultaneously valid, unique, and novel. "ASta" and "MSta" denote atom-level and molecule-level stability. "Top" indicates the proportion of generated molecules that simultaneously satisfy all three property constraints

# Downstream Analysis





# Conclusion

- Proposed an uncertainty-aware RL framework to guide diffusion models for **3D drug-like molecular** generation
- Achieved superior performance across three benchmark datasets, demonstrating **strong stability** and **generalization**
- Molecular dynamic simulation and ADMET analysis **varified** the structural stability and drug-like potential of generated molecules

# Future Work

- Extend the framework to **broader drug-related property objectives**
- Incorporate chemical prior knowledge to tackle **more complex drug targets**

# Acknowledge

This work was supported in part by the Canada Research Chairs Tier II Program (CRC-2021-00482), the Canadian Institutes of Health Research (PLL 185683, PJT 190272, PJT204042), the Natural Sciences, Engineering Research Council of Canada (RGPIN-2021-04072) and The Canada Foundation for Innovation (CFI) John R. Evans Leaders Fund (JELF) program (\#43481), the Studentship/Fellowship funded by Breast Cancer Canada, and the Vector Scholarship in Artificial Intelligence provided through the Vector Institute.