

# Molecular Generation with Chemical Feedback

<https://openreview.net/forum?id=9rPyHyjfwP>

Yin Fang, Ningyu Zhang<sup>†</sup> , Zhuo Chen, Lingbing Guo,

Xiaohui Fan, Huajun Chen<sup>†</sup> 





01 Introduction & Background

02 Model

03 Experiments

04 Conclusion & Future Work



01 Introduction & Background

02 Model

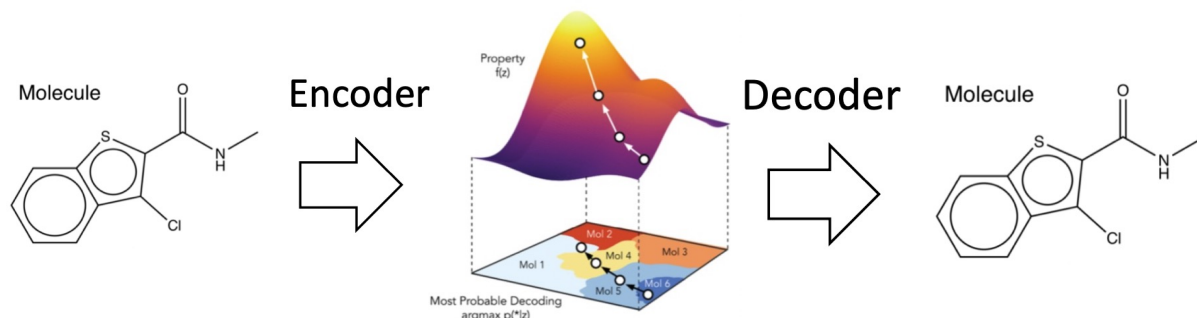
03 Experiments

04 Conclusion & Future Work

SEEKING TRUTH  
PURSUING INNOVATION

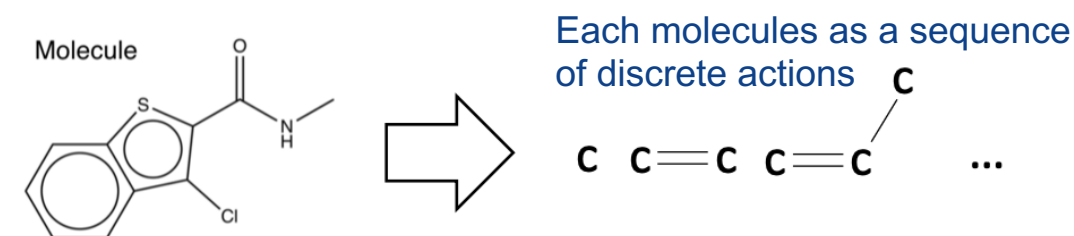
## ❑ Molecule Generation: Finding novel molecular structures with desired properties

### ❑ Search in *continuous* hidden space



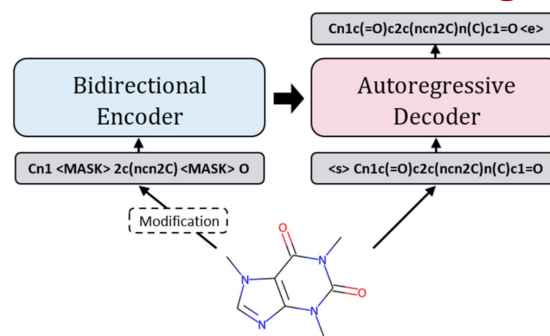
Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules (2018)

### ❑ Search in *discrete* chemical space



Graph Convolutional Policy Network for Goal-Directed Molecular Graph Generation (2019)

### ❑ Search in *molecular language* space



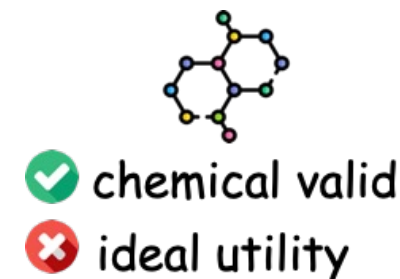
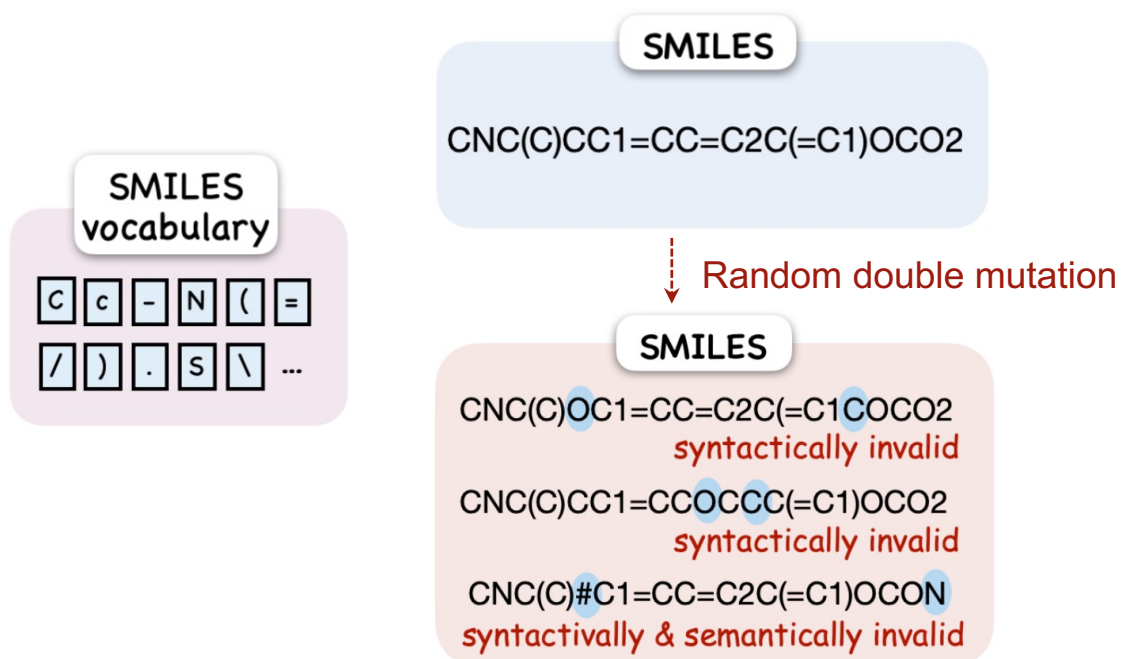
Chemformer: a pre-trained transformer for computational chemistry (2022)

# Challenges in Molecular Language Models

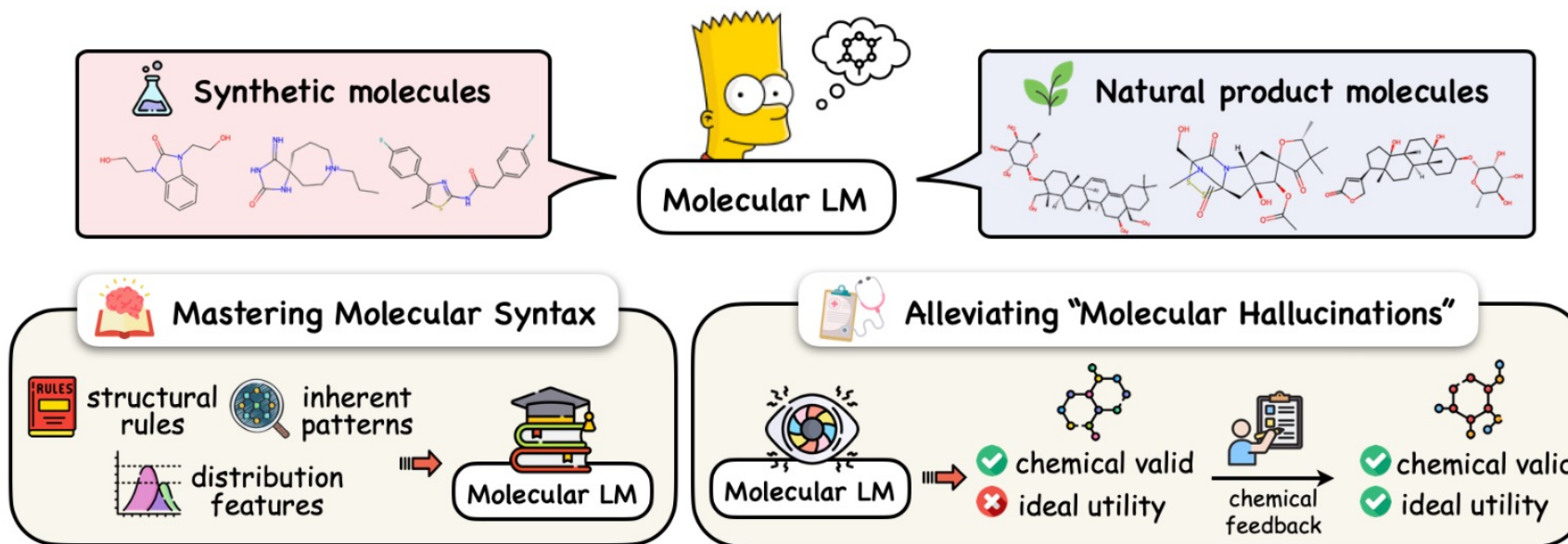


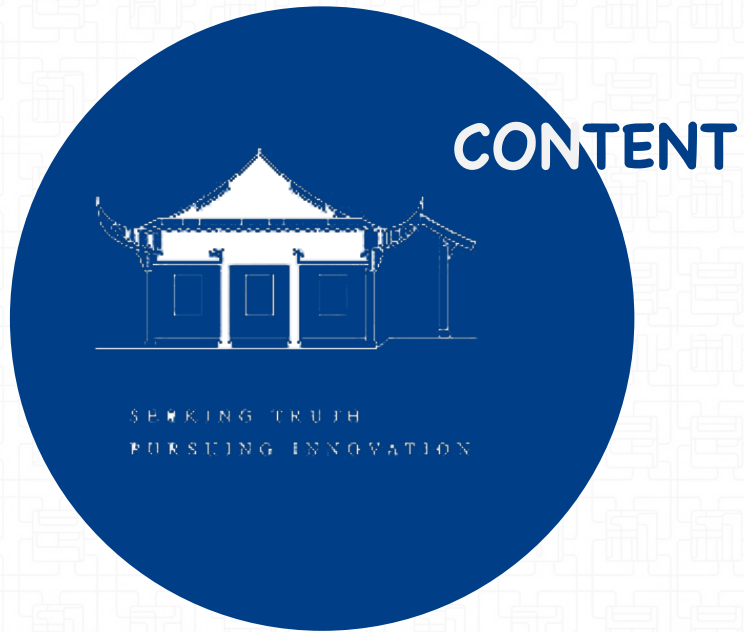
- SMILES-based language models have a certain probability of producing **invalid** molecules

- Molecular language models often suffer from “**molecular hallucinations**”



## 💡 **Aligning** pre-trained molecular language model with **chemical preferences**





01 Introduction & Background

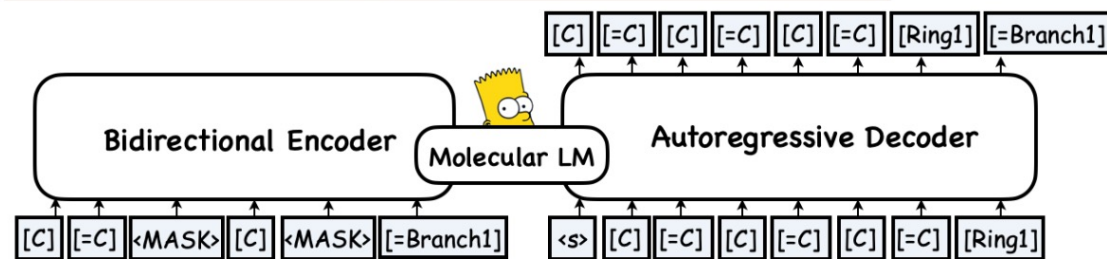
02 Model

03 Experiments

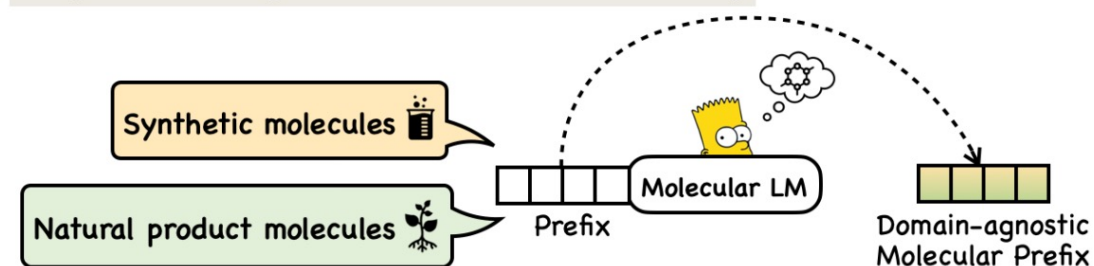
04 Conclusion & Future Work

SEEKING TRUTH  
PURSUING INNOVATION

## Step1: Molecular Language Syntax and Semantic Learning



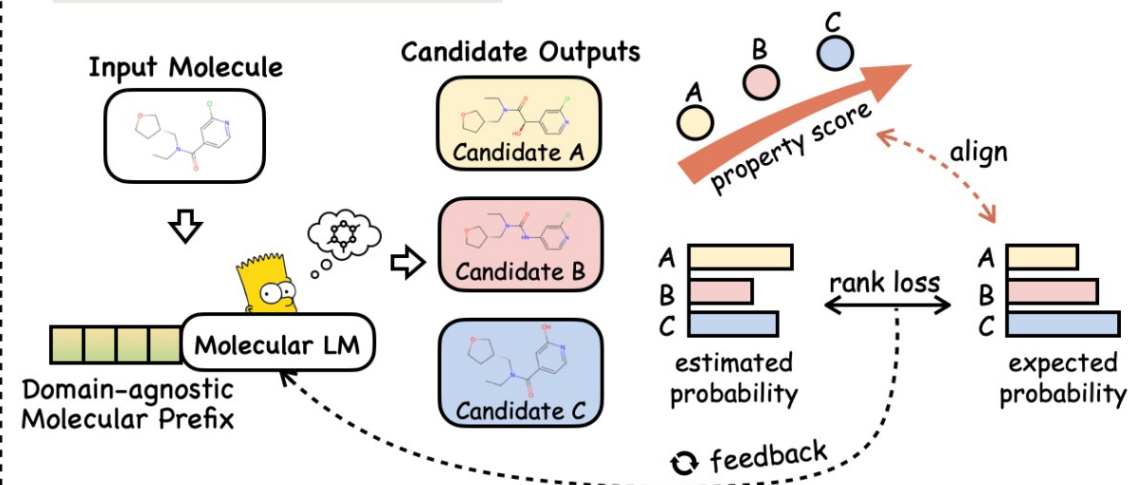
## Step2: Domain-agnostic Molecular Prefix Tuning



### Domain-agnostic Molecular Pre-training:

- Stage 1: Understand the molecular structure, grammar, and intrinsic semantics.
- Stage 2: Harness knowledge transferable across diverse domains.

## Chemical Feedback Paradigm



### Self-feedback Paradigm - align PLM with chemical preference:

- Align the probabilistic rankings with chemical preference rankings.
- Learn to evaluate and rectify its molecular outputs.





01 Introduction & Background

02 Model

03 Experiments

04 Conclusion & Future Work

PURSUING INNOVATION

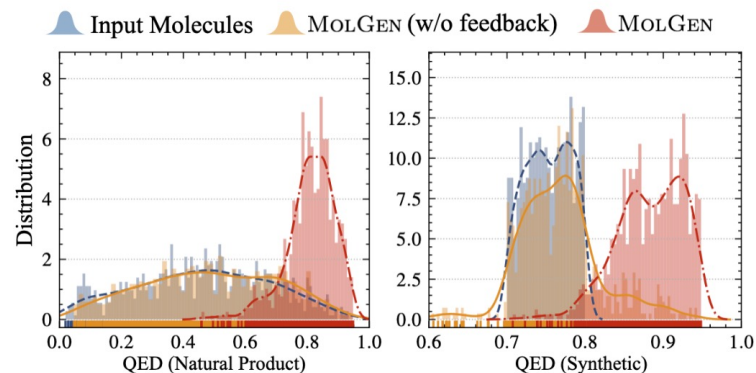
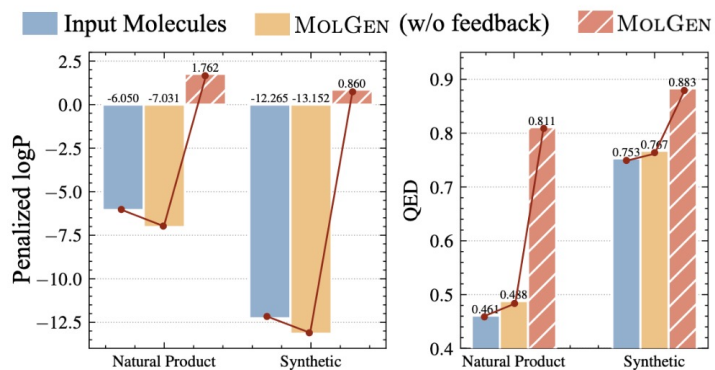
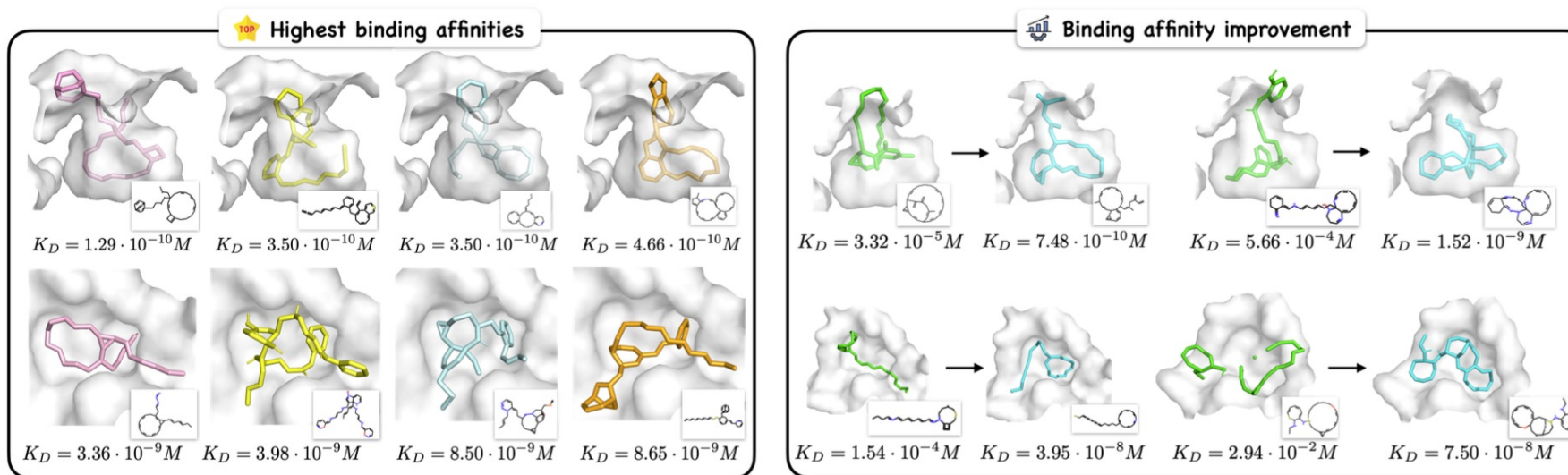
Reflects **real-world** molecular **distributions**

MODEL	SYNTHETIC MOLECULES							NATURAL PRODUCT MOLECULES						
	Validity↑	Frag↑	Scaf↑	SNN↑	IntDiv↑	FCD↓	Novelty↑	Validity↑	Frag↑	Scaf↑	SNN↑	IntDiv↑	FCD↓	Novelty↑
AAE	.9368	.9910	.9022	.6081	.8557	.5555	.7931	.0082	.9687	.2638	.3680	.8704	4.109	.9943
LATENTGAN	.8966	.9986	.8867	.5132	.8565	.2968	.9498	.9225	.2771	.0884	.5321	.6009	45.53	.9949
CHARRNN	.9748	.9998	.9242	.6015	.8562	.0732	.8419	.7351	.8816	.5212	.4179	.8756	2.212	.9792
VAE	.9767	.9994	.9386	.6257	.8558	.0990	.6949	.2627	.8840	.4563	.3950	.8719	4.318	.9912
JT-VAE	<b>1.000</b>	.9965	.8964	.5477	.8551	.3954	.9143	<b>1.000</b>	.8798	.5012	.3748	.8743	12.03	.9957
LIMO	<b>1.000</b>	.9562	.1073	.6125	.8544	.1532	.8956	<b>1.000</b>	.7242	.0005	.3416	.7726	31.84	.9962
CHEMFORMER	.9843	.9889	.9248	.5622	.8553	.0061	.9581	.9825	.9826	.4126	.5875	.8650	.8346	.9947
MOLGEN	<b>1.000</b>	<b>.9999</b>	<b>.9999</b>	<b>.9996</b>	<b>.8567</b>	<b>.0015</b>	<b>1.000</b>	<b>1.000</b>	<b>.9994</b>	<b>.8404</b>	<b>.8148</b>	<b>.8878</b>	<b>.6519</b>	<b>.9987</b>

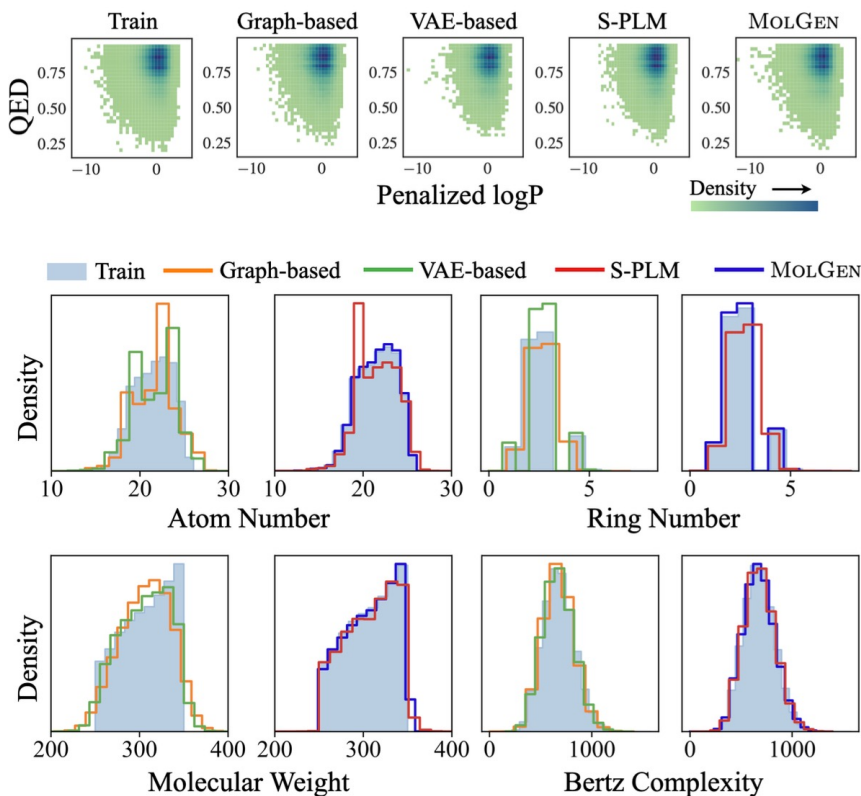
# Molecular Optimization



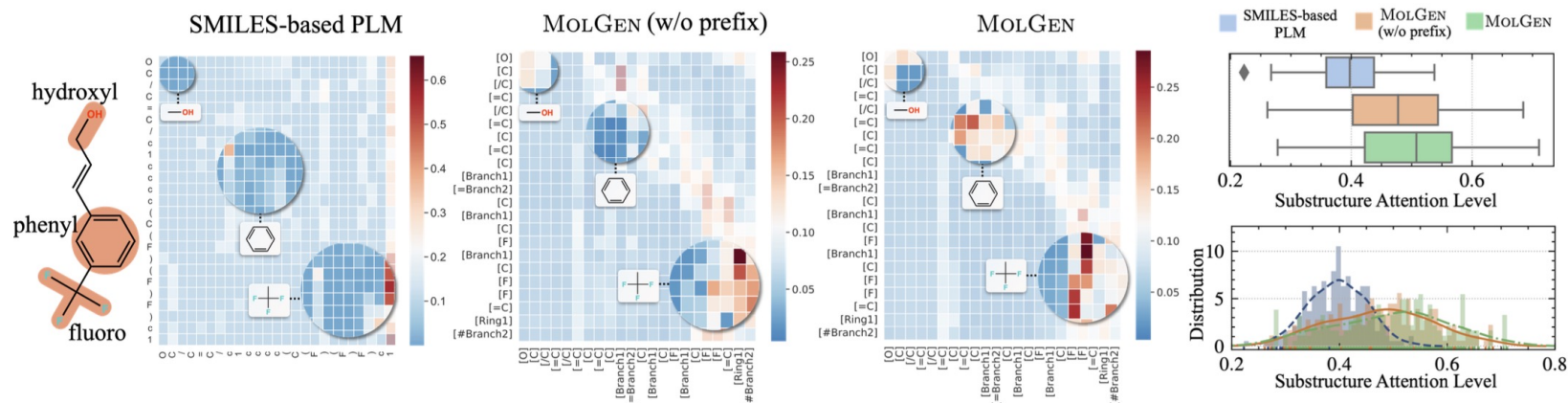
## Alleviate molecular hallucinations



## Captures molecular **characteristics**



## Recognizes **meaningful substructures**





01 Introduction & Background

02 Model

03 Experiments

04 Conclusion & Future Work

- ❑ This study proposes a pre-trained molecular language model tailored for molecule generation:
  - ❑ generating valid molecules while avoiding “molecular hallucinations”
  - ❑ identifying essential molecular substructures

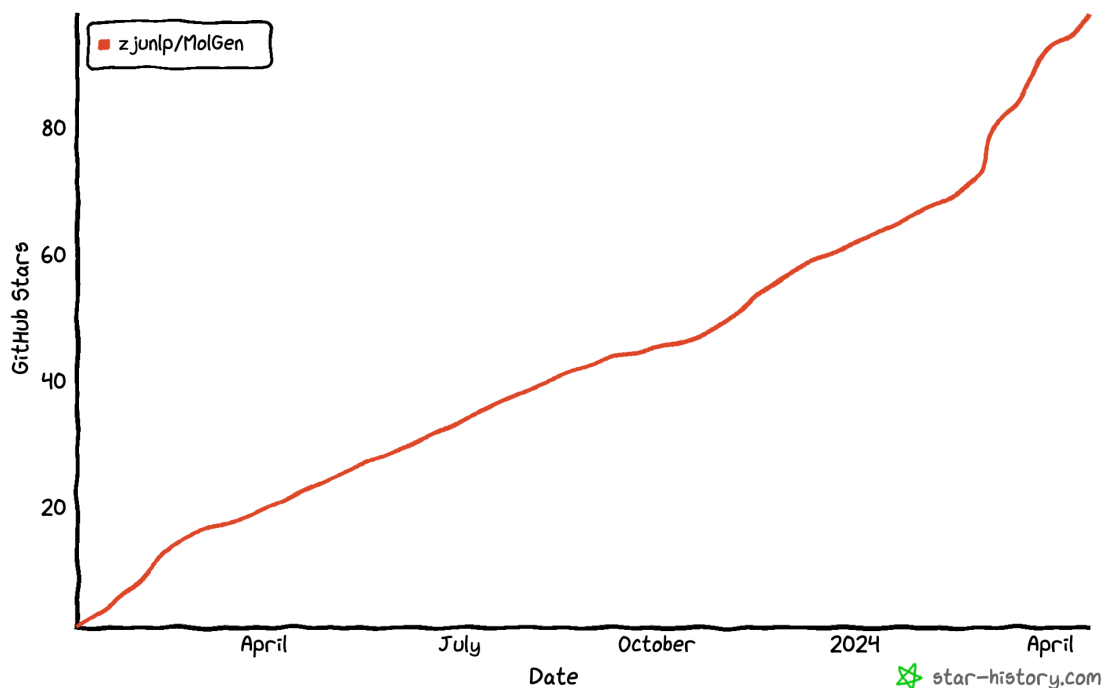
## Future Work

- ❑ Apply to other tasks such as retrosynthesis and reaction prediction
- ❑ Explore multimodal pre-training
- ❑ Incorporate additional sources of knowledge



GitHub

[github.com/zjunlp/MolGen](https://github.com/zjunlp/MolGen)



star-history.com



Hugging Face

[zjunlp/MolGen-large](https://huggingface.co/zjunlp/MolGen-large)

↓ **Total downloads**

12,632 (all time, tracked internally since January 2021)

[zjunlp/MolGen-large-opt](https://huggingface.co/zjunlp/MolGen-large-opt)

↓ **Total downloads**

1,726 (all time, tracked internally since January 2021)

[zjunlp/MolGen-7B](https://huggingface.co/zjunlp/MolGen-7B)

↓ **Total downloads**

1,568 (all time, tracked internally since January 2021)

# Thank you!



**Code**



**Model**



**浙江大学**  
ZHEJIANG UNIVERSITY