# Auto-QChem: an automated workflow for the generation and storage of DFT calculations for organic molecules

Jingyi Liu

*29/09/2025*

# Outline

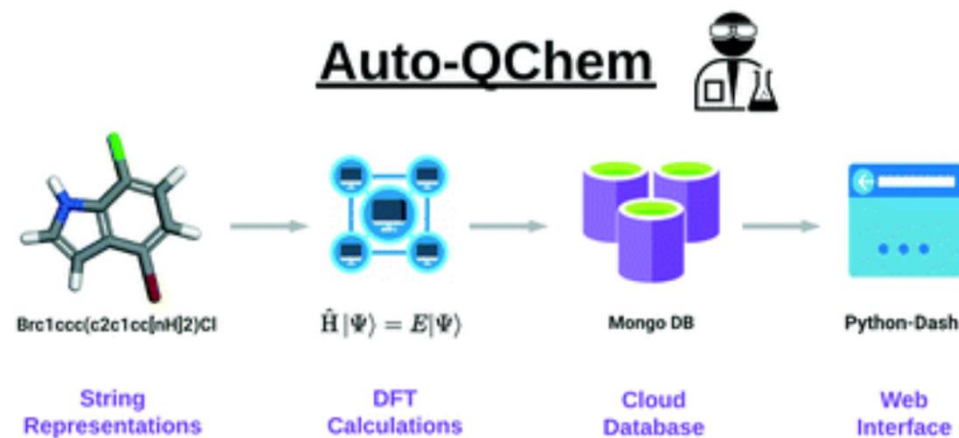**1. Background**

**2. Auto-Qchem workflow**

**3. Applications**

**4. Discussion**



Auto-QChem

Brc1ccc(c2c1cc[nH]2)Cl → $\hat{H}|\Psi\rangle = E|\Psi\rangle$ → Mongo DB → Python-Dash

String Representations | DFT Calculations | Cloud Database | Web Interface

# Research Background

# Research Background

- The application of machine learning models in organic chemistry requires effective representations of chemical structures.

- ML models trained with chemical descriptors are more interpretable than molecular fingerprints and other representations.

- High throughput DFT presents a significant barrier to experimental chemists.

- Previous DFT automation packages are mainly designed for material science rather than small organic molecules.

# Auto-Qchem

- An automatic, high-throughput and end-to-end DFT calculation workflow

- Computes chemical descriptors for organic molecules

[1] https://github.com/PrincetonUniversity/auto-qchem;
[2] https://princetonuniversity.github.io/auto-qchem
[3] https://autoqchem.org

# Auto-Qchem workflow

# ❖ Auto-Qchem

## ➤ Features

1. Generate user-specified input.

2. HPC interface.

3. Automatic information extraction from DFT results.
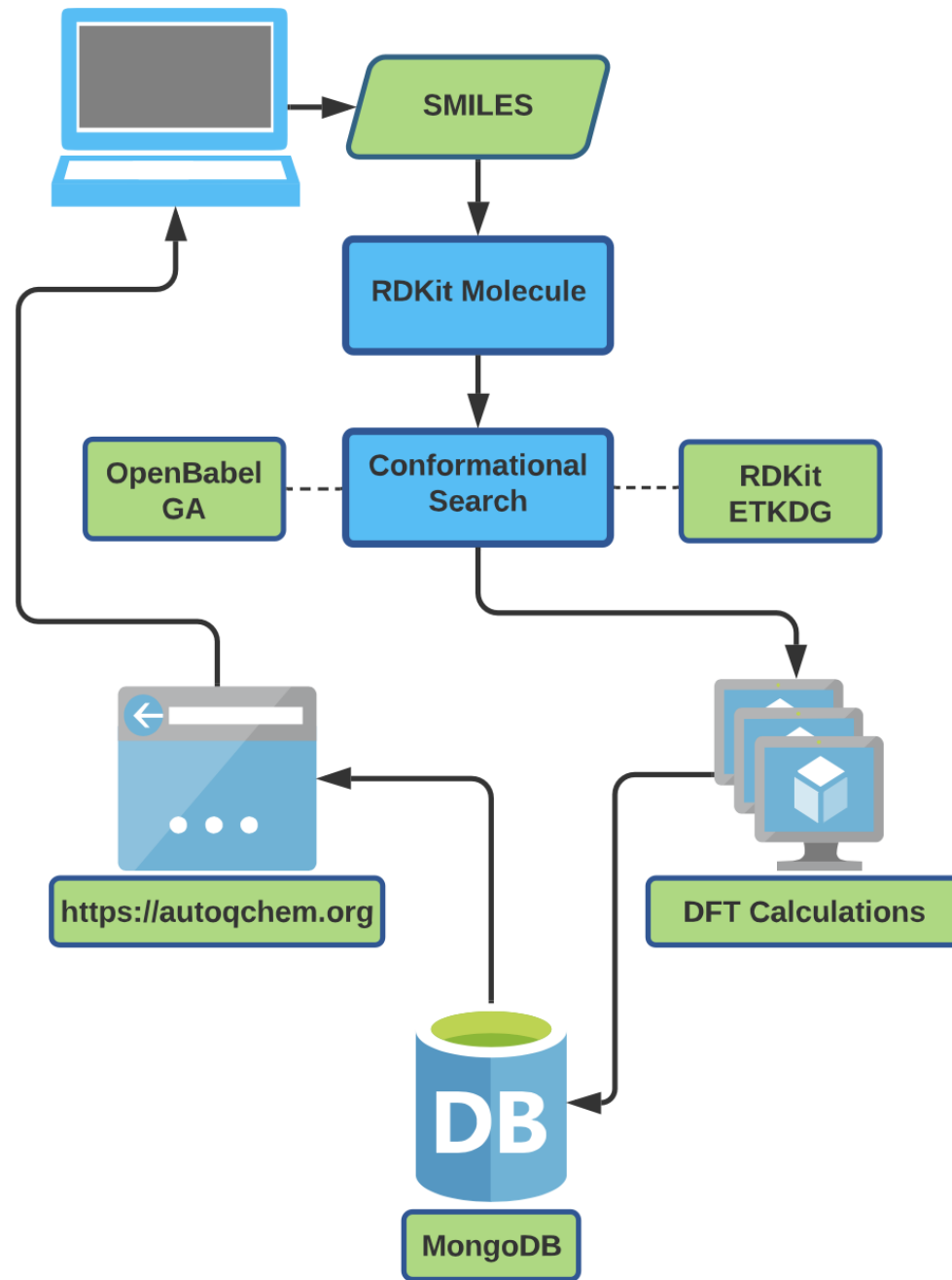
4. Convenient storage and data accessing.



**Fig. 1** Computational workflow of Auto-QChem.

## ❖ Auto-Qchem



**Fig. 1** Computational workflow of Auto-QChem.

SMILES

RDKit Molecule

OpenBabel GA · · · Conformational Search · · · RDKit ETKDG

https://autoqchem.org

DFT Calculations

- ❖ **Functional, basis, ..**
- ❖ **Geometry optimization**
- ❖ **Frequency and thermochemistry**
- ❖ **Excited state TDDFT**

DB

MongoDB

- ❖ **Molecules**
- ❖ **Matadata**
- ❖ **Log_files**
- ❖ **Qchem descriptors**
- ❖ **tags**

8

# ❖ Auto-Qchem

➢ **Features**

  • Avoid duplicate calculations



**Fig. 2** Collection schema of Auto-QChem database.

# ❖ Auto-Qchem

➢ **Features:** Queries and data retrieval

➢ 1692 molecules in total



**Fig. 3** Query view (left) and the molecule view (right) of the web interface. The molecule view is a snapshot while viewing the second lowest energy conformation in 3D.

## 1. Substrate scope design in Ni/ photoredox methodology development

- Ni/ photoredox catalyzed alkylation reaction of aryl halides using acetals as alcohol-derived aliphatic radical sources.
- Data set: 2683 aryl bromides, with 168 electronic and steric features,
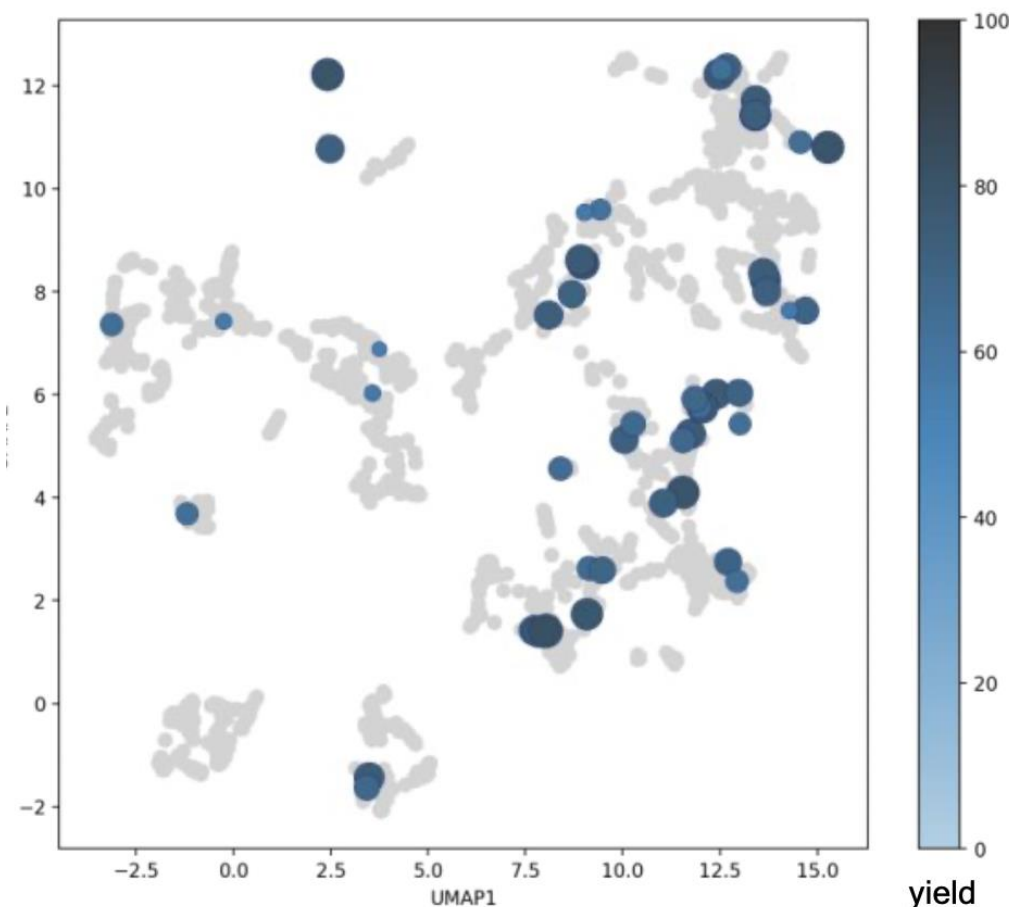- 95 of 168 features were employed for clustering.



Reaction scheme conditions:

[Ir[dF(CF$_3$)ppy]$_2$(dtbbpy)]PF$_6$ (1 mol%)
NiBr$_2$•glyme (2 mol%), dtbbpy (3 mol%)

K$_3$PO$_4$ (1 equiv)
PhH:MeCN (1:1) (0.1 M)
34W blue LEDs, 24 h

1.1 equiv

A (R' = i-Amyl)
22% yield

D (R' = CH$_2$CH$_2$OMe)
98% yield

E (R' = Et)
87% yield

K (R' = Me)
0% yield
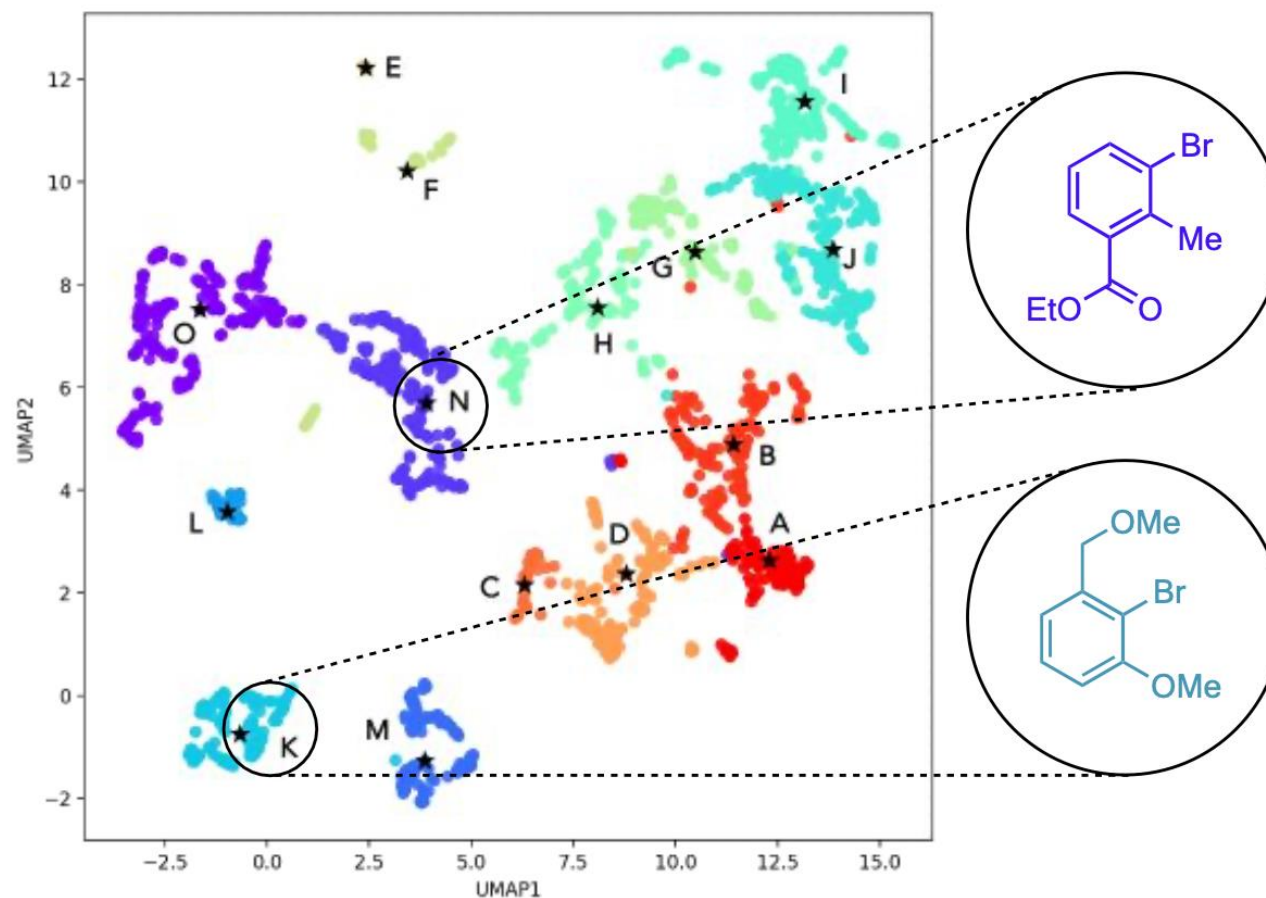
N (R' = Me)
70% yield

## 1. Substrate scope design in Ni/ photoredox methodology development

- Electronegativity of the aryl bromides was highly correlated with yield.
- Generalized additive model (GAM) trained on electronegativity outperformed ML models trained on literatures
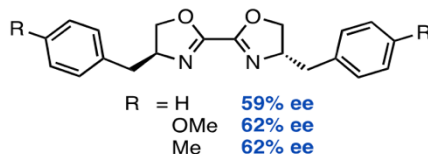


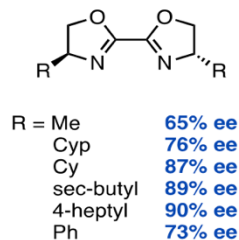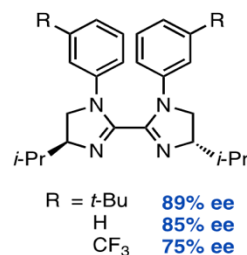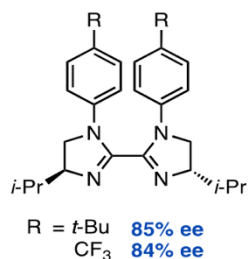(a). Substrates from literature and yields

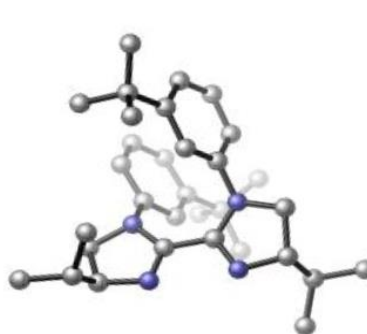(b). Clustering results and selected aryl bromides

## 2. Ligand parametrization and enantioselectivity prediction in nickel catalysis

- Previous used ligand: Bioxazoline (BiOx)
  - Good enantioselectivity, low yield
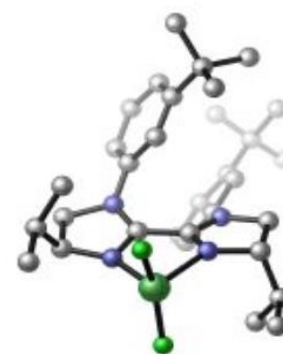- Optimal ligand: chiral biimidazoline (BiIm) : time-consuming to discover



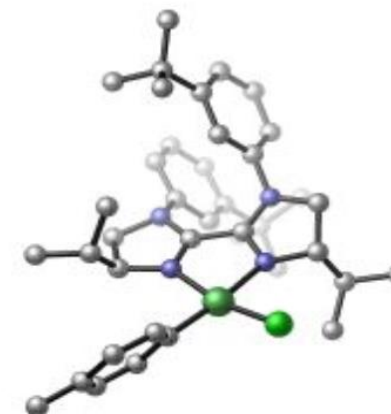(a). Model reaction system and representative examples of ligands tested.
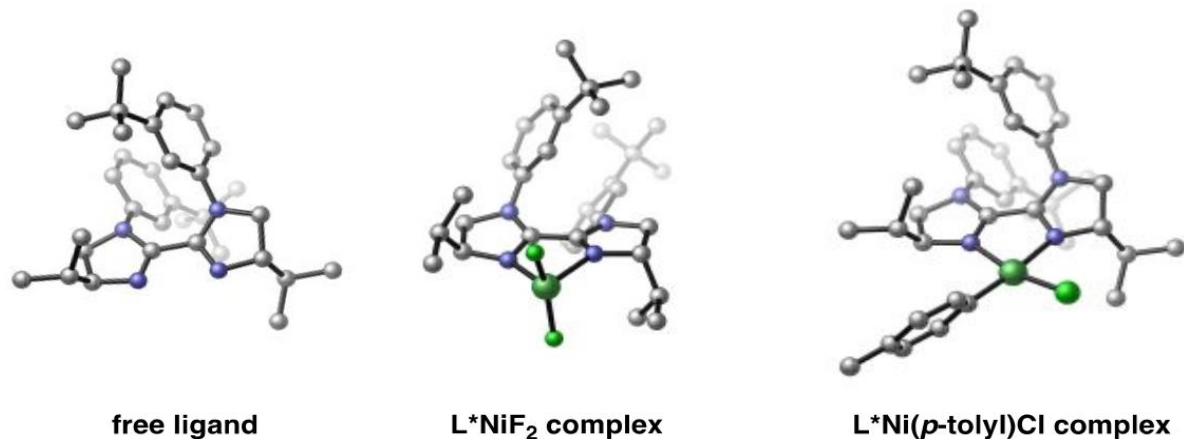
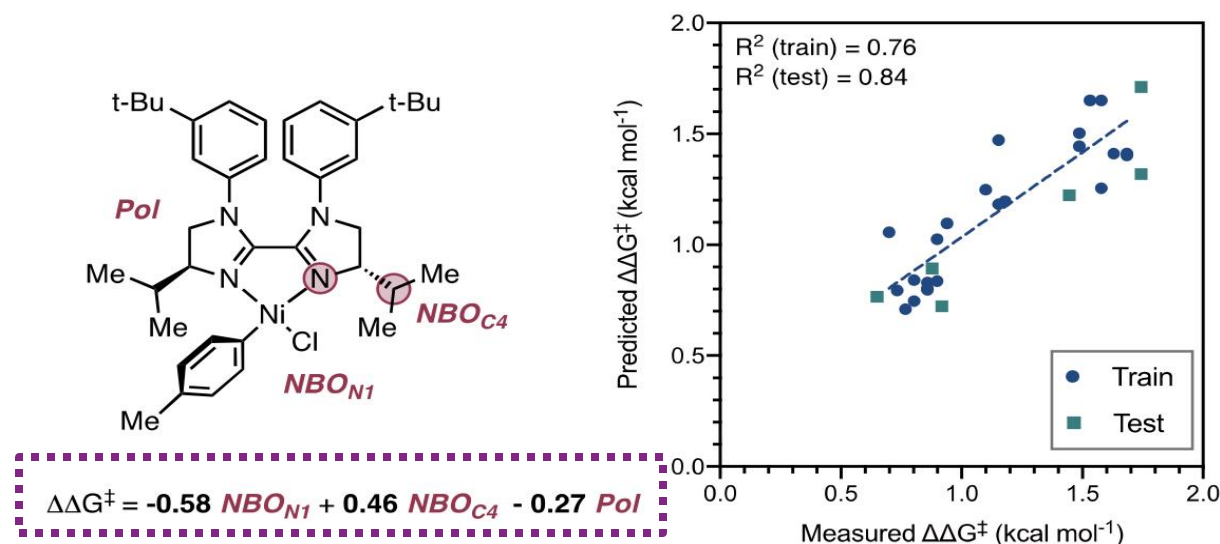(b). Ligand environment models.

free ligand          L*NiF₂ complex          L*Ni(p-tolyl)Cl complex

## 2. Ligand parametrization and enantioselectivity prediction in nickel catalysis

**(b). Ligand environment models.**



free ligand     L*NiF₂ complex     L*Ni(*p*-tolyl)Cl complex

**(c). Regression modeling for L*Ni(*p*-tolyl)Cl with DFT-derived features.**



$$\Delta\Delta G^{\ddagger} = -0.58\ NBO_{N1} + 0.46\ NBO_{C4} - 0.27\ Pol$$

- Hypothesis: computed features depend on ligand environment.

- Manually generated conformers ← Auto-Qchem cannot handle transition metal complexes

- Conclusion: electronic, rather than steric attributes of BiIm ligands govern the enantioselectivity of this reaction

14

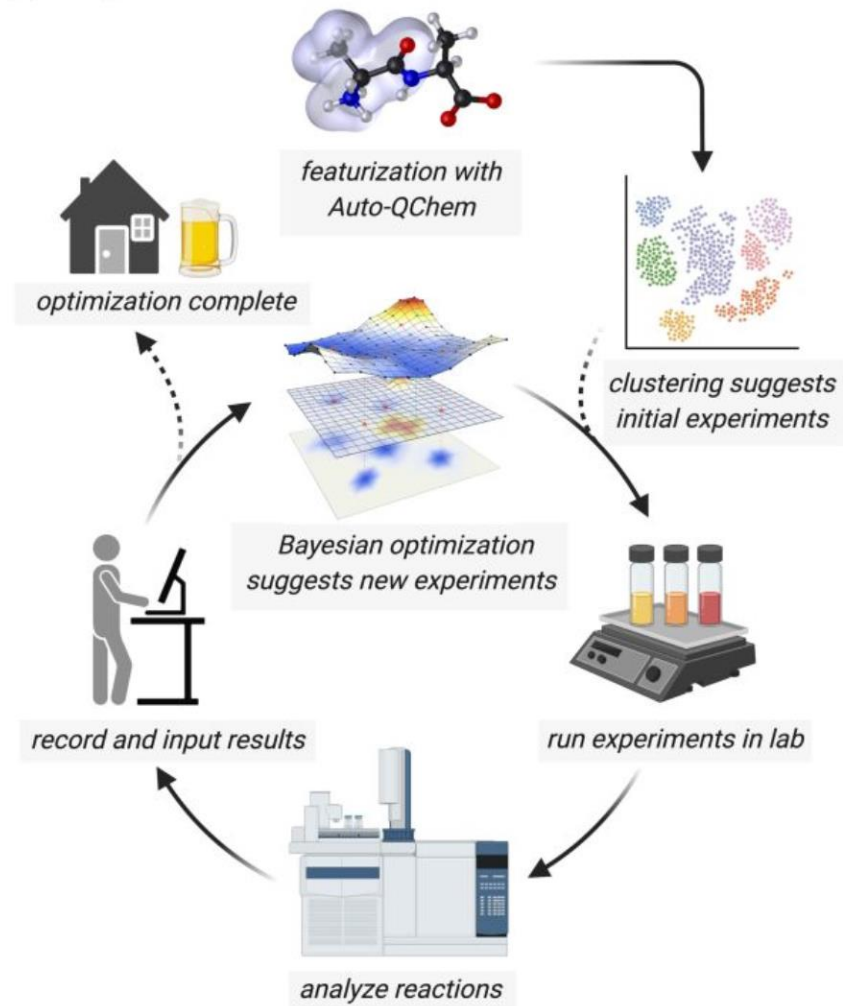## 3. Reaction condition optimization via Bayesian optimization

- Traditional reaction optimization is slow and resource-intensive, relying on trial-and-error or costly methods, with limited exploration of chemical space.

- Bayesian optimization, a sequential design algorithm for global optimization of black-box functions, in efficient reaction condition optimization.

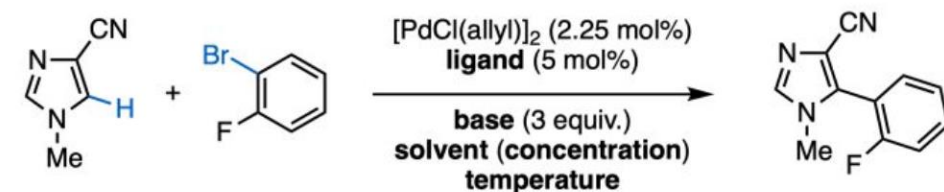- EDBO (Experimental Design via Bayesian Optimization)

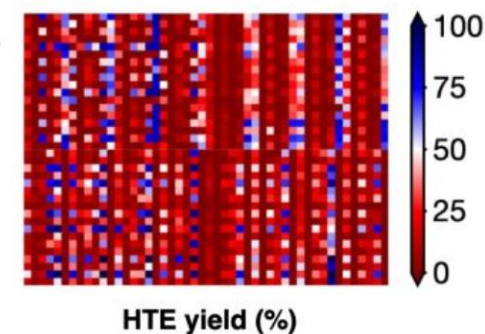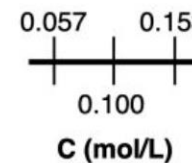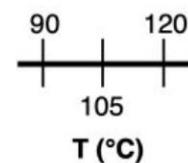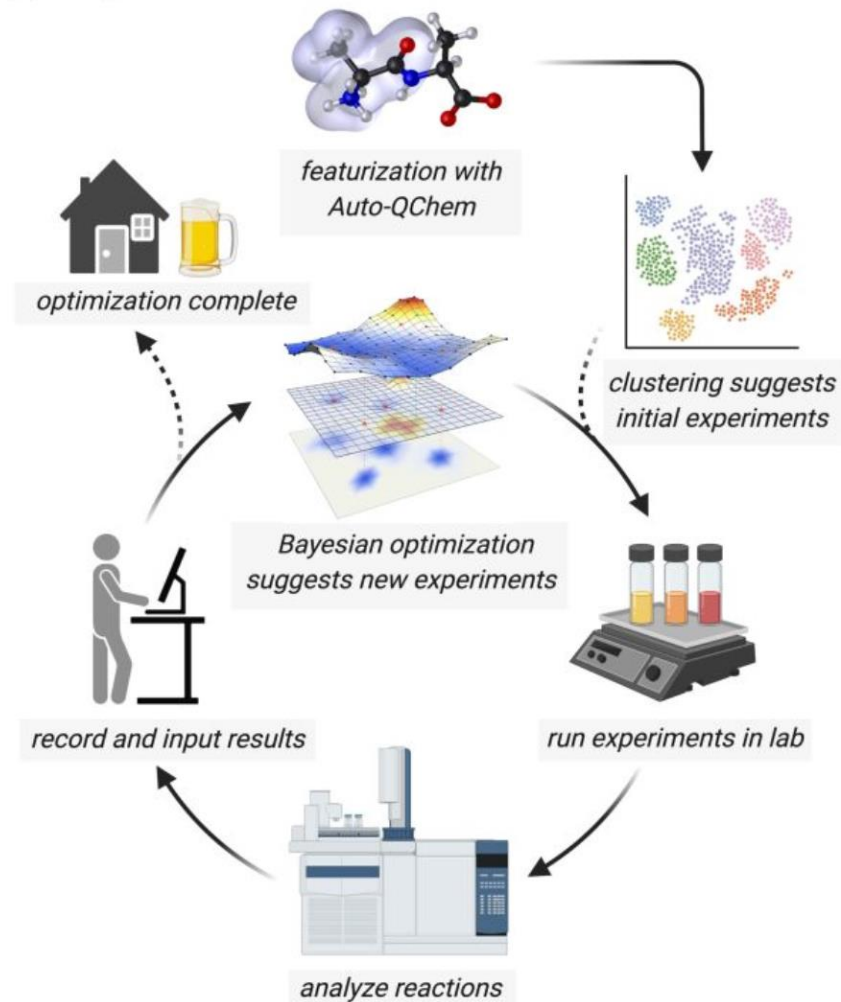➢ 3. Reaction condition optimization via Bayesian optimization



**Fig. 6** Use case 3: reaction condition optimization via Bayesian optimization.

3. Reaction condition optimization via Bayesian optimization



(a). the general optimization workflow of EDBO.

featurization with Auto-QChem

clustering suggests initial experiments

optimization complete

Bayesian optimization suggests new experiments

record and input results

run experiments in lab

analyze reactions

EDBO (simulated 50 times) achieved a higher average performance within the first 15 experiments even with random initialization and found conditions with >99% yield 100% of the time!

# Discussion

## ❖ Limits

- Cannot generate accurate conformers for transition metal complexes and molecules with non-canonical bonds.

- Lack of supports for other cluster schedulers

## ❖ Future perspective

- External packages support

- Data insufficiency

# Questions? Comments?

# Thank You