



# Enabling late-stage drug diversification by high-throughput experimentation with geometric deep learning

David F. Nippa, Kenneth Atz, Remo Hohler, Alex T. Müller, Andreas Marx, Christian Bartelmus, Georg Wuitschik, Irene Marzuoli, Vera Jost, Jens Wolfard, Martin Binder, Antonia F. Stepan, David B. Konrad\*, Uwe Grether\*, Rainer E. Martin\*& Gisbert Schneider\*

Nature Chemistry

<https://doi.org/10.1038/s41557-023-01360-5>

Bin Yang

2024.01.13



**Prof. Gisbert Schneider**

**Research interest:**

**Integration of artificial intelligence into  
pharmaceutical research and chemical biology**

- **Full professor, Department of Chemistry and Applied Bioscience, ETH Zurich;**
- **Director, Singapore-ETH Centre (SEC)**
  
- **Graduate from the Freie Universität Berlin, Germany;**
- **Pharmaceuticals Division at Roche, Switzerland;**
- **Goethe-University in Frankfurt, Germany (Beilstein Endowed Chair for Chem and Bioinformatics);**
- **Current position at the ETH**



**David B. Konrad**

**Research interest:**

- Medicinal Chemistry
- Chemical Methodology
- Cancer Cell Biology
- Chemo Proteomics

**ORGANIC CHEMISTRY**

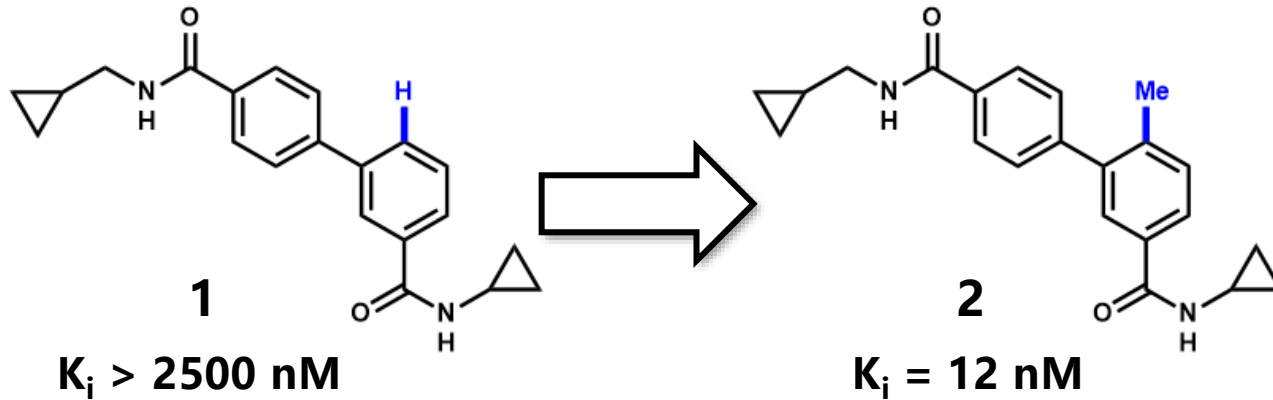
## A concise synthesis of tetrodotoxin

**David B. Konrad**<sup>1</sup> ‡, Klaus-Peter Rühmann<sup>2</sup> ‡, Hiroyasu Ando<sup>3</sup>, Belinda E. Hetzler<sup>2</sup>, Nina Strassner<sup>4</sup>, Kendall N. Houk<sup>4</sup>, Bryan S. Matsuura<sup>2\*</sup> ‡, Dirk Trauner<sup>2\*</sup> §

Tetrodotoxin (TTX) is a neurotoxic natural product that is an indispensable probe in neuroscience, a biosynthetic and ecological enigma, and a celebrated target of synthetic chemistry. Here, we present a stereoselective synthesis of TTX that proceeds in 22 steps from a glucose derivative. The central cyclohexane ring of TTX and its  $\alpha$ -tertiary amine moiety were established by the intramolecular 1,3-dipolar cycloaddition of a nitrile oxide, followed by alkynyl addition to the resultant isoxazoline. A ruthenium-catalyzed hydroxylactonization set the stage for the formation of the dioxo-adamantane core. Installation of the guanidine, oxidation of a primary alcohol, and a late-stage epimerization gave a mixture of TTX and anhydro-TTX. This synthetic approach could give ready access to biologically active derivatives.

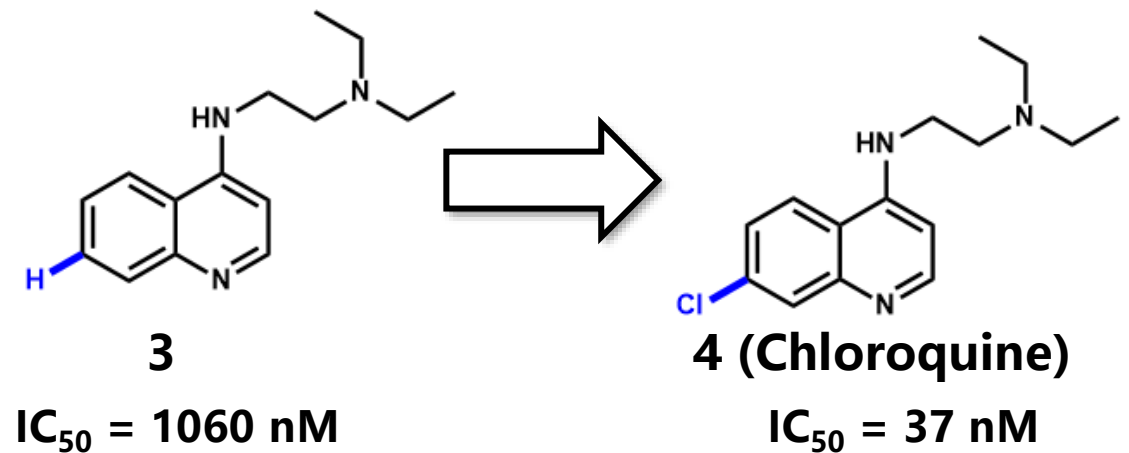
# Late-stage Functionalization (LSF)

## p38 $\alpha$ MAP Kinase Inhibitor



“Methyl Effect”

## Antimalarial



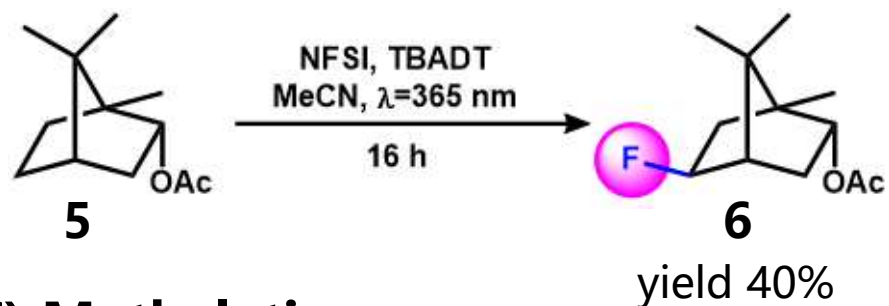
“Chloro Effect”

*J. Med. Chem.* 2012, 55, 9, 4489

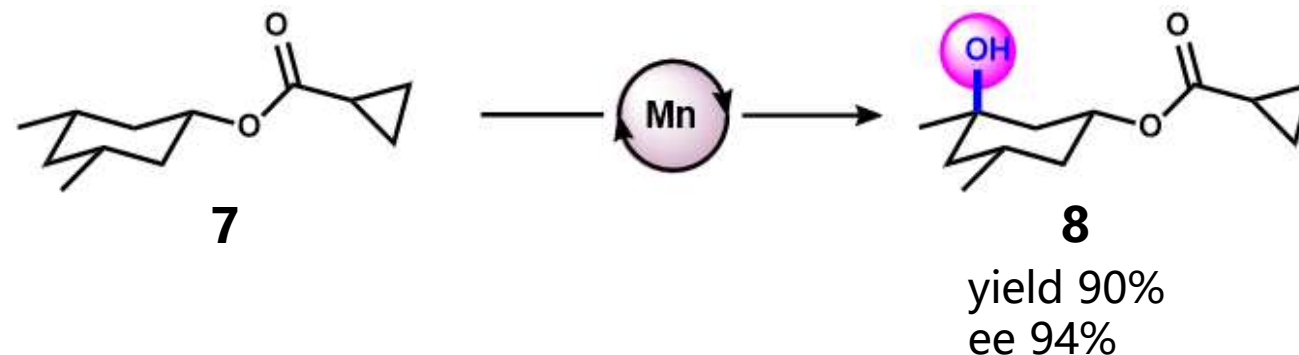
*J. Med. Chem.* 2023, 66, 8, 5305

# Late-stage C-H Functionalization

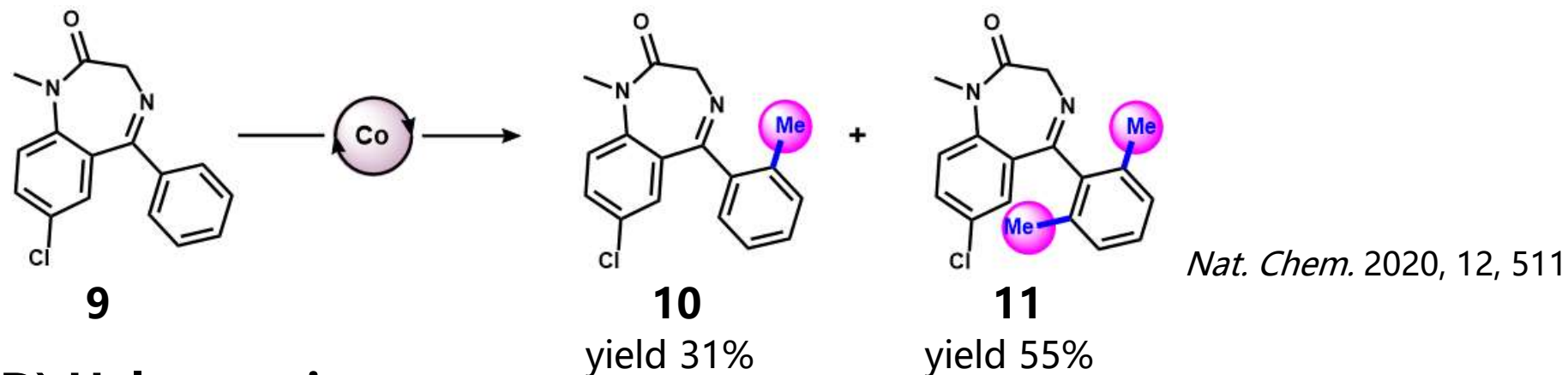
## (A) Fluorination *Angew. Chem.* 2014,126, 4778



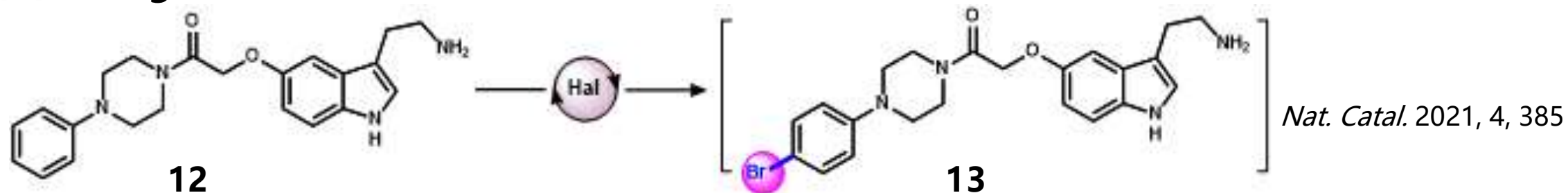
## (B) Oxidation *J. Am. Chem. Soc.* 2023, 145, 29, 15742



## (C) Methylation



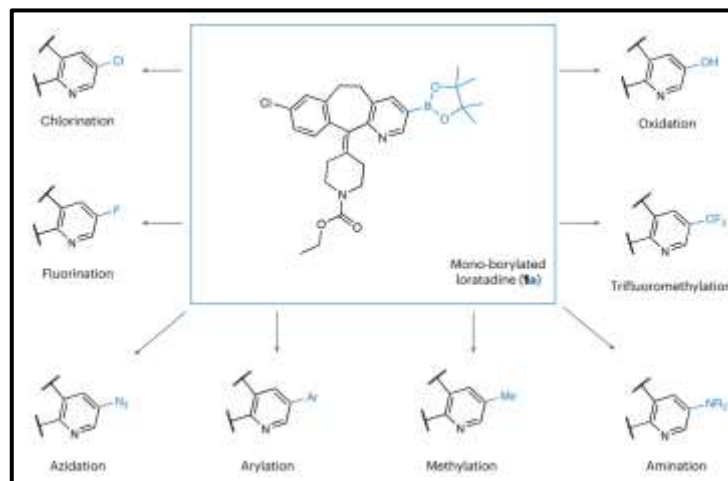
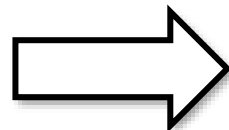
## (D) Halogenation



# Late-stage C-H Functionalization

- Rare reports
- One single reaction type

## Limitations

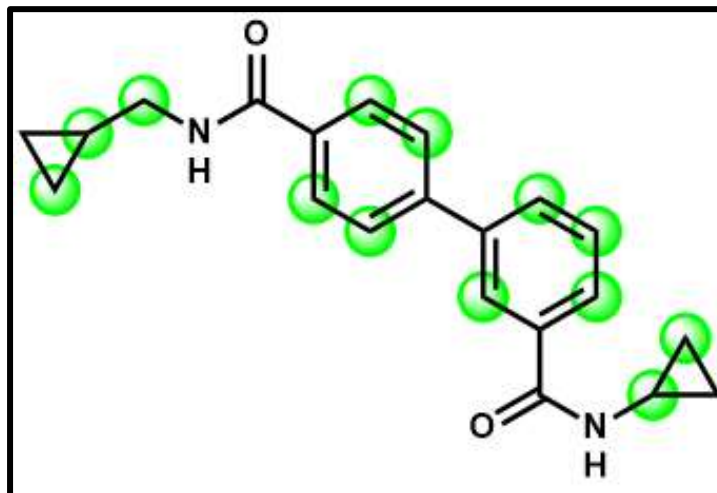
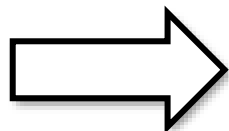


## C-H Borylation

Guidelines?

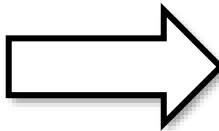
- Various types of C-H bonds
- Different bond strengths
- Electronic properties
- Steric properties
- FG environment

## Reaction features



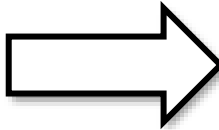
# Deep Learning & High-throughput Experimentation (HTE)

Guidelines

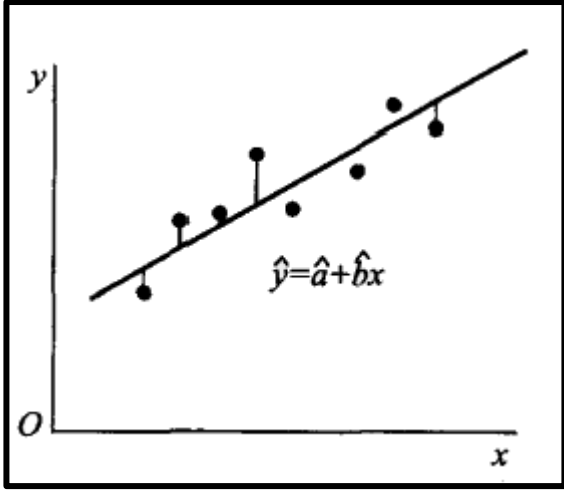


Deep Learning

Regression & Prediction

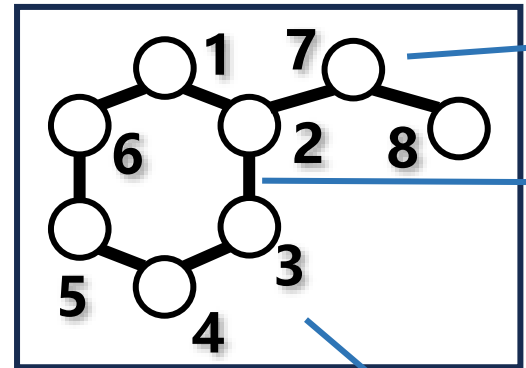


Samples HTE



Learning Method — Graph Neuro Network (GNN)

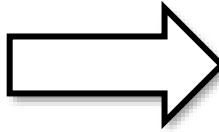
Adjacency Matrix



Nodes

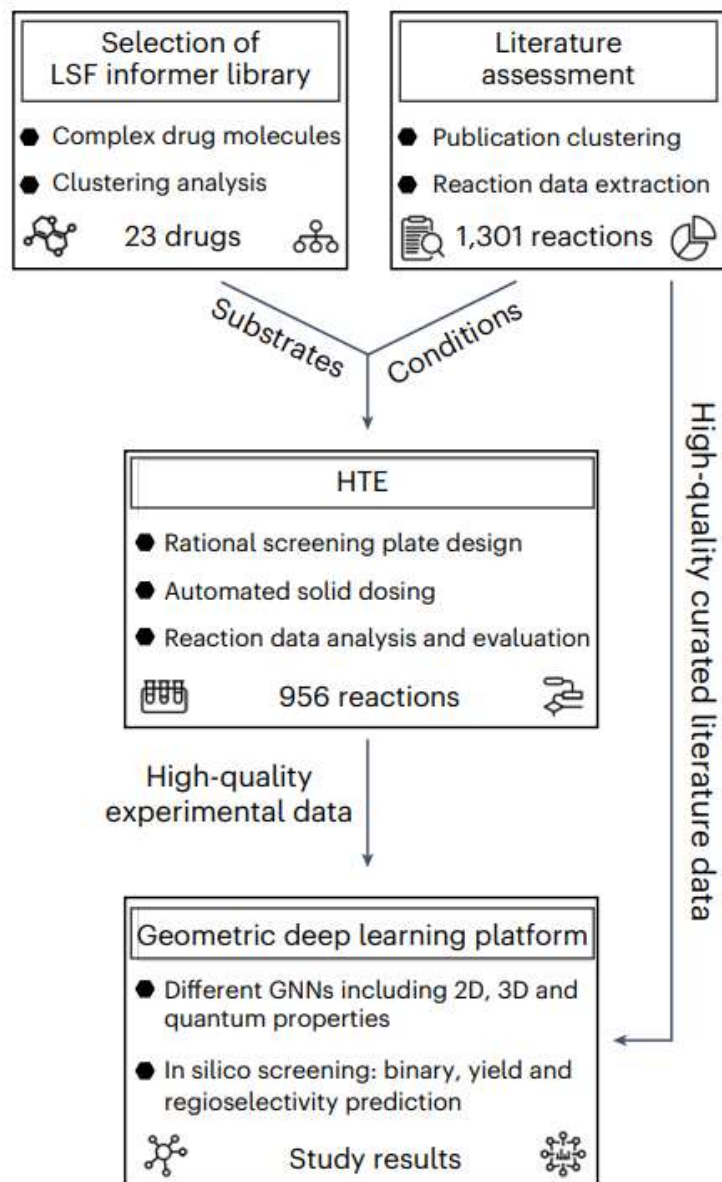
Edges

Graph



	1	2	3	4	5	6	7	8
1	0	1	0	0	0	1	0	0
2	1	0	1	0	0	0	1	0
3	0	1	0	1	0	0	0	0
4	0	0	1	0	1	0	0	0
5	0	0	0	1	0	1	0	0
6	1	0	0	0	1	0	0	0
7	0	1	0	0	0	0	0	1
8	0	0	0	0	0	0	1	0

# Overview of this study



## Process:

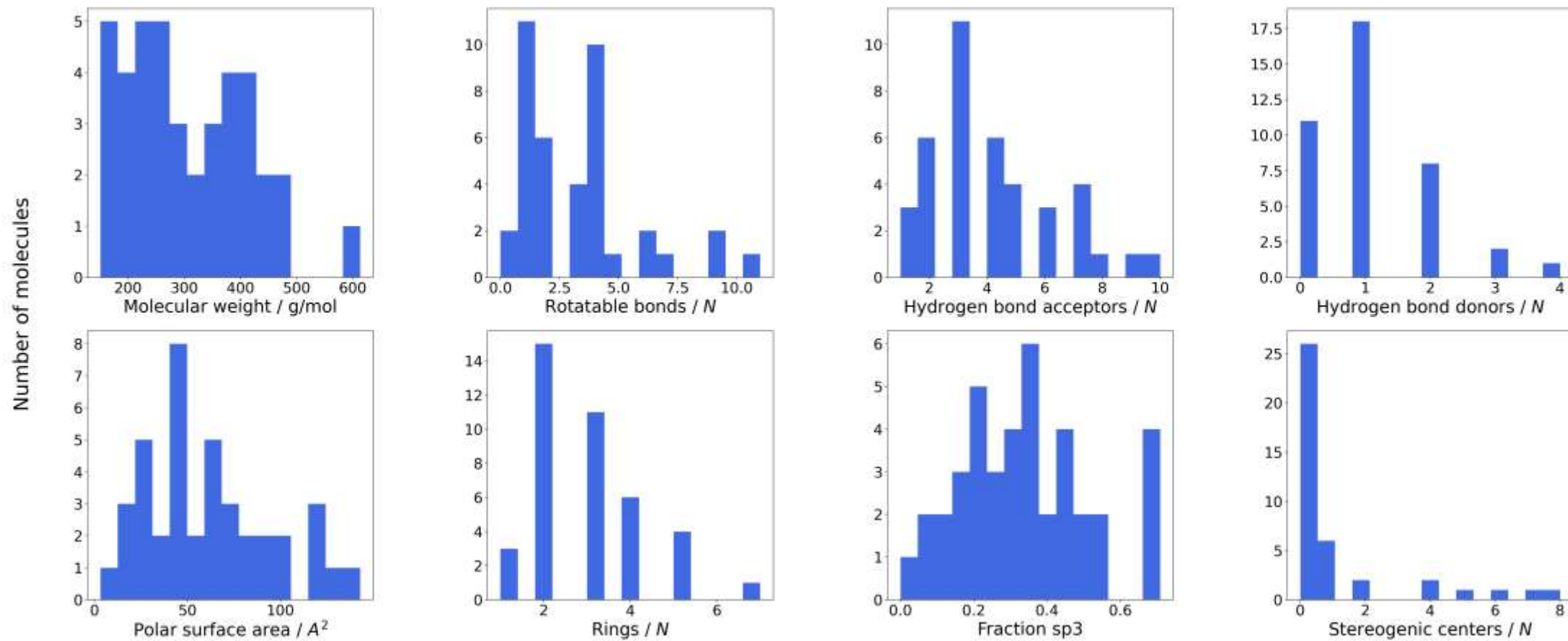
- Select conditions;
- Select substrates;
- HTE;
- Model fitting

## Three questions:

- Reaction or not?
- Yield?
- Regioselectivity?



# Method—Substrate & Condition selection



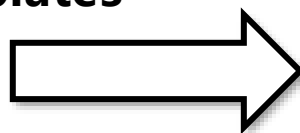
# Method—HTE

	1	2	3	4	5	6
Ligand (5.0 mol%)						
Solvent (0.2 M)						
A CyHex (10)	<div style="display: flex; justify-content: space-around;"> <div style="text-align: center;"> <p>Catalyst (2.5 mol%)</p> <p>(2) [Ir(OMe)(1,5-cod)]<sub>2</sub></p> </div> <div style="text-align: center;"> <p>B source (1.0 equiv.)</p> <p>(3) B<sub>2</sub>Pin<sub>2</sub></p> </div> <div style="text-align: center;"> </div> </div>					
B Me-THF (11)						
C CPME (12)						
D MeCN (13)						

➤ Operation in a glove-box

➤ Operation in 24- or 96-well plates

➤ Detection by LC-MS



Reaction or not?

Yield?

Reference products?

New products?

# Deep learning

## (A) Outcome & yield

	Reaction yield <i>r</i> value	Reaction yield m.a.e. (%)	Binary reaction outcome (random split), AUC (%)	Binary reaction outcome (substrate split), AUC (%)
GTNN2D	0.896±0.006	4.53±0.09	<b>91.8±2.1</b>	52±2
GNN2D	0.866±0.005	5.61±0.06	87.5±1.0	51±2
GTNN3D	0.884±0.01	4.51±0.11	91.4±0.7	58±4
GNN3D	0.877±0.001	5.33±0.34	89.4±0.8	65±5
GTNN2DQM	<b>0.898±0.003</b>	4.41±0.17	90.9±1.5	53±5
GNN2DQM	0.876±0.01	5.41±0.10	89.0±1.1	59±5
GTNN3DQM	0.890±0.01	<b>4.23±0.08</b>	<b>91.8±0.9</b>	<b>67±2</b>
GNN3DQM	0.890±0.006	4.88±0.24	89.1±0.9	64±4
ECFP4NN	0.885±0.0006	4.55±0.14	89.3±1.3	52±3

## 2 models

GNN using sum pooling (GNN)  
Graph Transformer Neural Network (GTNN)

+ 1 baseline model:

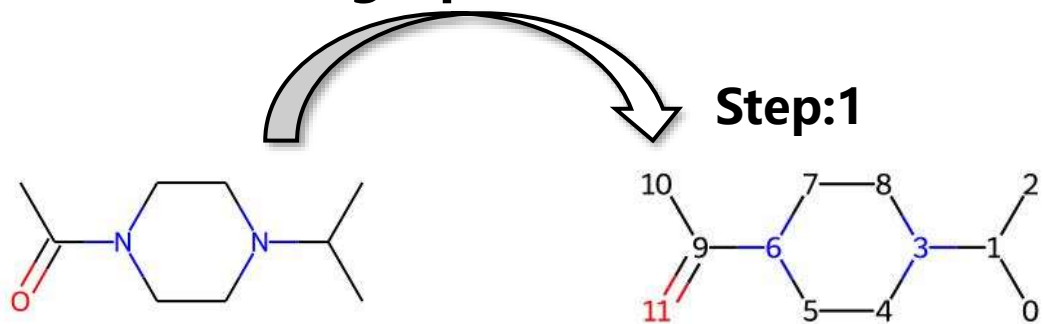
ECFP4-based Neuro Network (ECFP4NN)

## 4 input types

2D  
3D  
2DQM  
3DQM

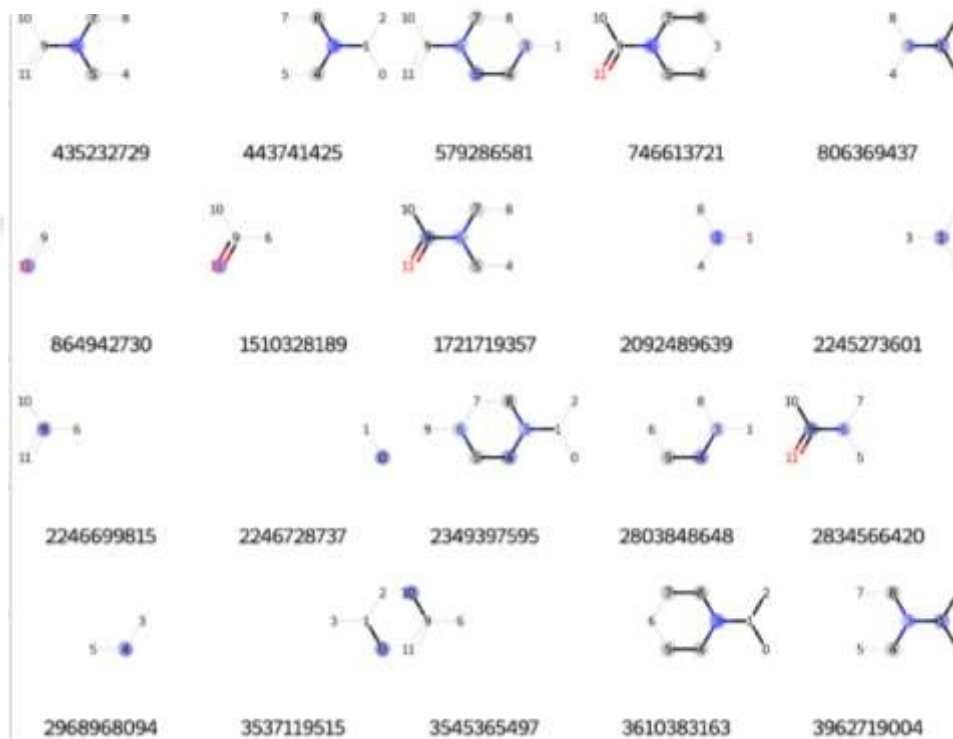
# Deep learning

## Molecular fingerprint (ECFP4)



Step:1

Step:2

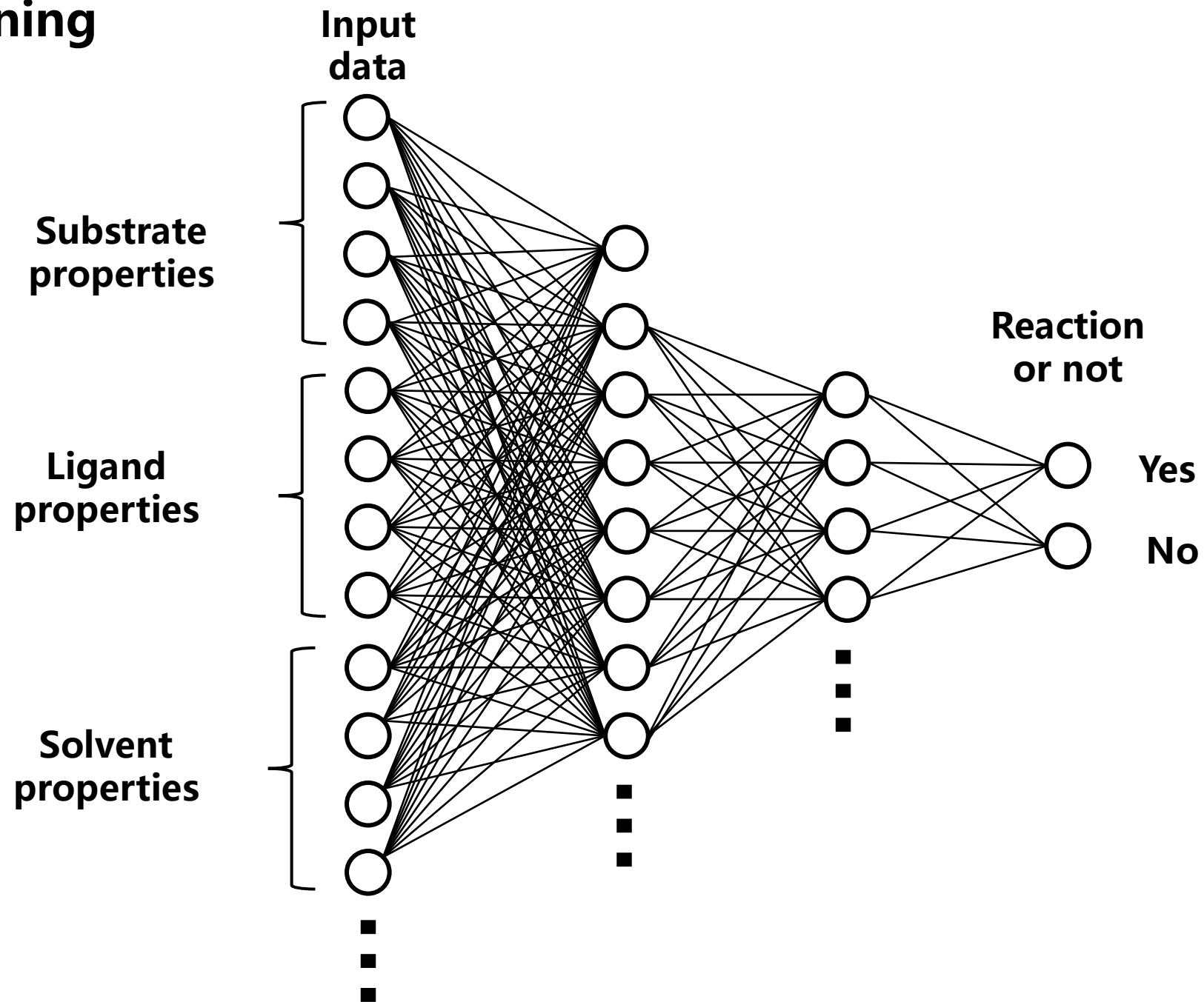


```
16
17 '''Step3: 计算ECFP4分子指纹'''
18 ECFP_bitinfo = {}
19 ECFP = AllChem.GetMorganFingerprint(mol, radius=2, bitInfo=ECFP_bitinfo, useFeatures=False)
20 # print(ECFP)
21 '''一个分子指纹是一个特殊的bit数组，其中囊括的信息就是一个邻接矩阵'''
22
23
24 '''查看每个有效信息的bit, '''
25 for v in ECFP_bitinfo.values():
26     print(v)
27 '''每个bit的信息(x, y)代表，第x个原子的半径y个化学键的信息，平行代表信息一致'''
```

Run: 1 x

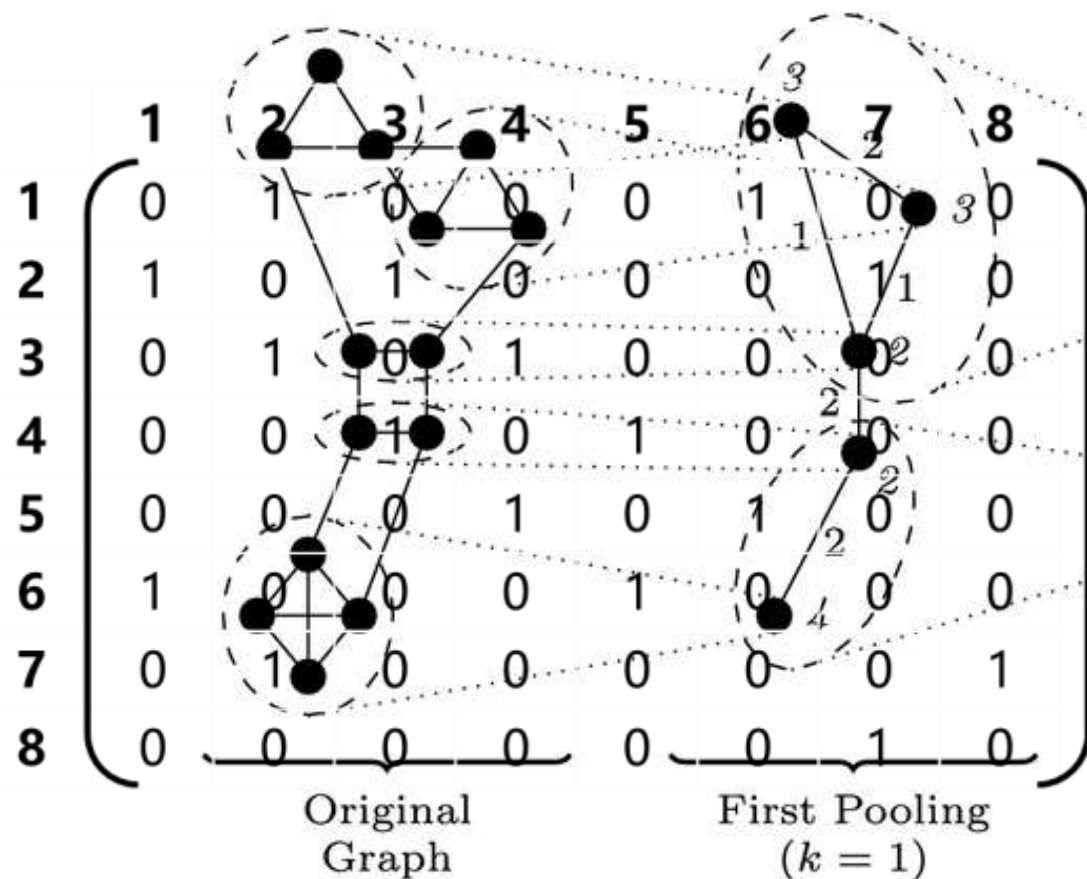
```
D:\Download\envs\rdkit-env-2\python.exe C:/Users/Lenovo/PycharmProjects/untitled/1.py
[(6, 1),]
[(3, 1),]
[(5, 2), (7, 2)]
[(6, 2),]
[(1, 1),]
[(11, 8),]
[(11, 1),]
```

# Deep learning

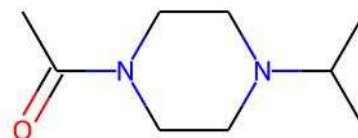


# Deep learning

## Sum pooling



Why pooling ?



```
''' 邻接矩阵中有很多元素是0, 是无用信息. '''  
print("number of non-zero bit in this ecfp: ", len(ECFP.GetNonzeroElements()))  
  
1 x  
D:\Download\envs\rdkit-env-2\python.exe C:/Users/lenovo/PycharmProjects/untitled/1.py  
number of non-zero bit in this ecfp: 28
```

[https://www.researchgate.net/figure/Two-examples-sum-pooling-on-the-same-graph-where-phv1documentclass12ptminimal\\_fig2\\_353839563](https://www.researchgate.net/figure/Two-examples-sum-pooling-on-the-same-graph-where-phv1documentclass12ptminimal_fig2_353839563)

# Deep learning

## GNN input in this article:

### Nodes:

```
atoms = []
qml_atoms = []
is_ring = []
hyb = []
arom = []
crds_3d = []

AllChem.EmbedMolecule(mol, randomSeed=randomseed)
AllChem.UFFOptimizeMolecule(mol)

for idx, i in enumerate(mol.GetAtoms()):
    atoms.append(ATOMTYPE_DICT[i.GetSymbol()])
    qml_atoms.append(QML_ATOMTYPE_DICT[i.GetSymbol()])
    is_ring.append(IS_RING_DICT[str(i.IsInRing())])
    hyb.append(HYBRIDISATION_DICT[str(i.GetHybridization())])
    arom.append(AROMATICITY_DICT[str(i.GetIsAromatic())])
    crds_3d.append(list(mol.GetConformer().GetAtomPosition(idx)))
```

- Atom ID
- Hybridization
- Coordinates (for 3D conformation)
- Ring status
- Aromatic status

### Edges:

```
edge_dir1 = []
edge_dir2 = []
for idx, bond in enumerate(mol.GetBonds()):
    a2 = bond.GetEndAtomIdx()
    a1 = bond.GetBeginAtomIdx()
    edge_dir1.append(a1)
    edge_dir1.append(a2)
    edge_dir2.append(a2)
    edge_dir2.append(a1)

edge_2d = torch.from_numpy(np.array([edge_dir1, edge_dir2]))
```

## 2D: Start atom and end atom

```
# 3D graph for qml and qml prediction
qml_atoms = torch.LongTensor(qml_atoms)
xyzs = torch.FloatTensor(crds_3d)
edge_index = np.array(nx.complete_graph(qml_atoms.size(0)).edges())
edge_index = to_undirected(torch.from_numpy(edge_index).t().contiguous())
edge_index, _ = add_self_loops(edge_index, num_nodes=crds_3d.shape[0])

qml_graph = Data(
    atoms=qml_atoms,
    coords=xyzs,
    edge_index=edge_index,
    num_nodes=qml_atoms.size(0),
)

charges = QMLMODEL(qml_graph).unsqueeze(1).detach().numpy()

# Get edges for 3d graph
distance_matrix = squareform(pdist(crds_3d))
np.fill_diagonal(distance_matrix, float("inf")) # to remove self-loops
edge_3d = torch.from_numpy(np.vstack(np.where(distance_matrix <= radius)))
```

## 3D: Tensor

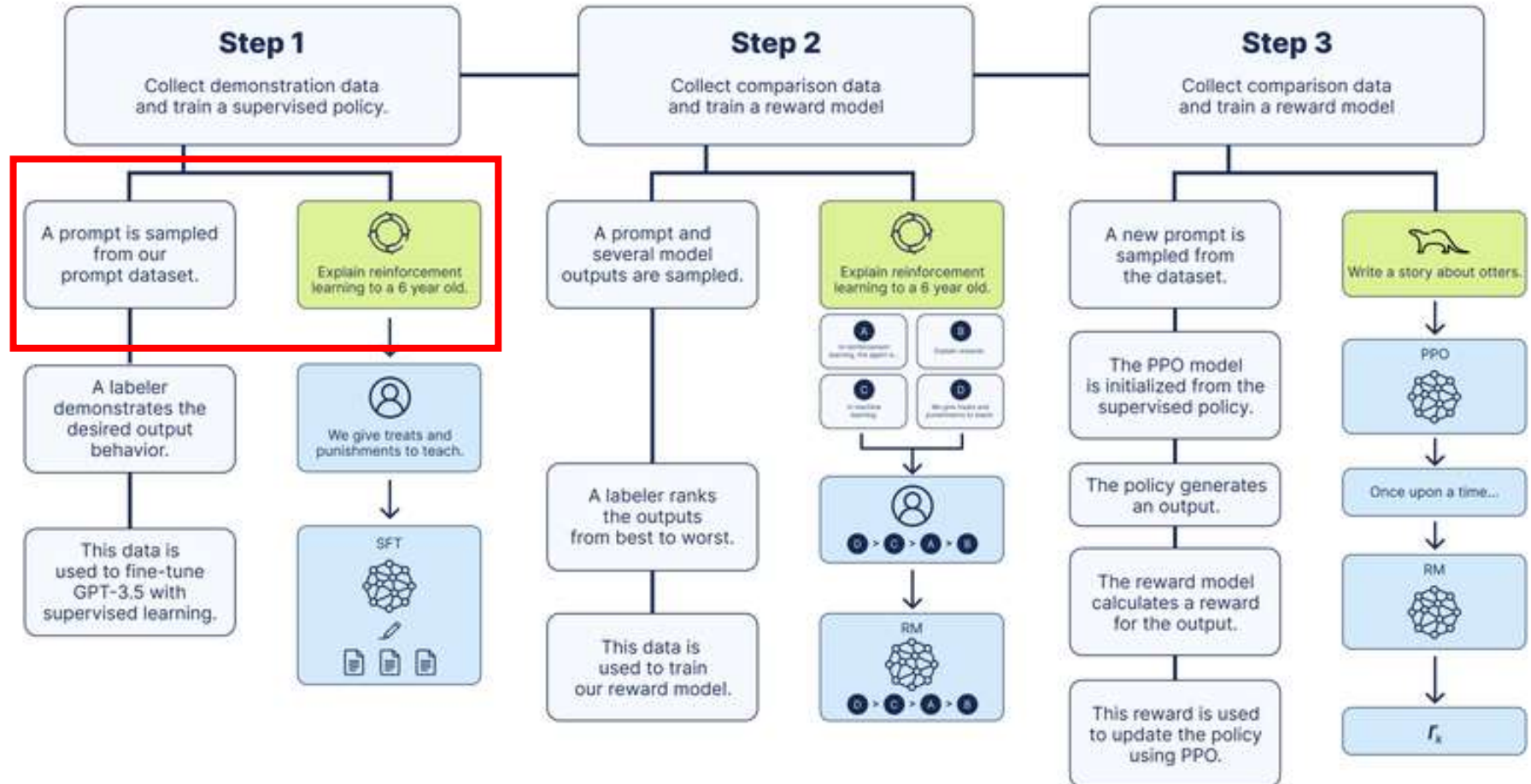
# Deep learning

## Graph transformer



ChatGPT

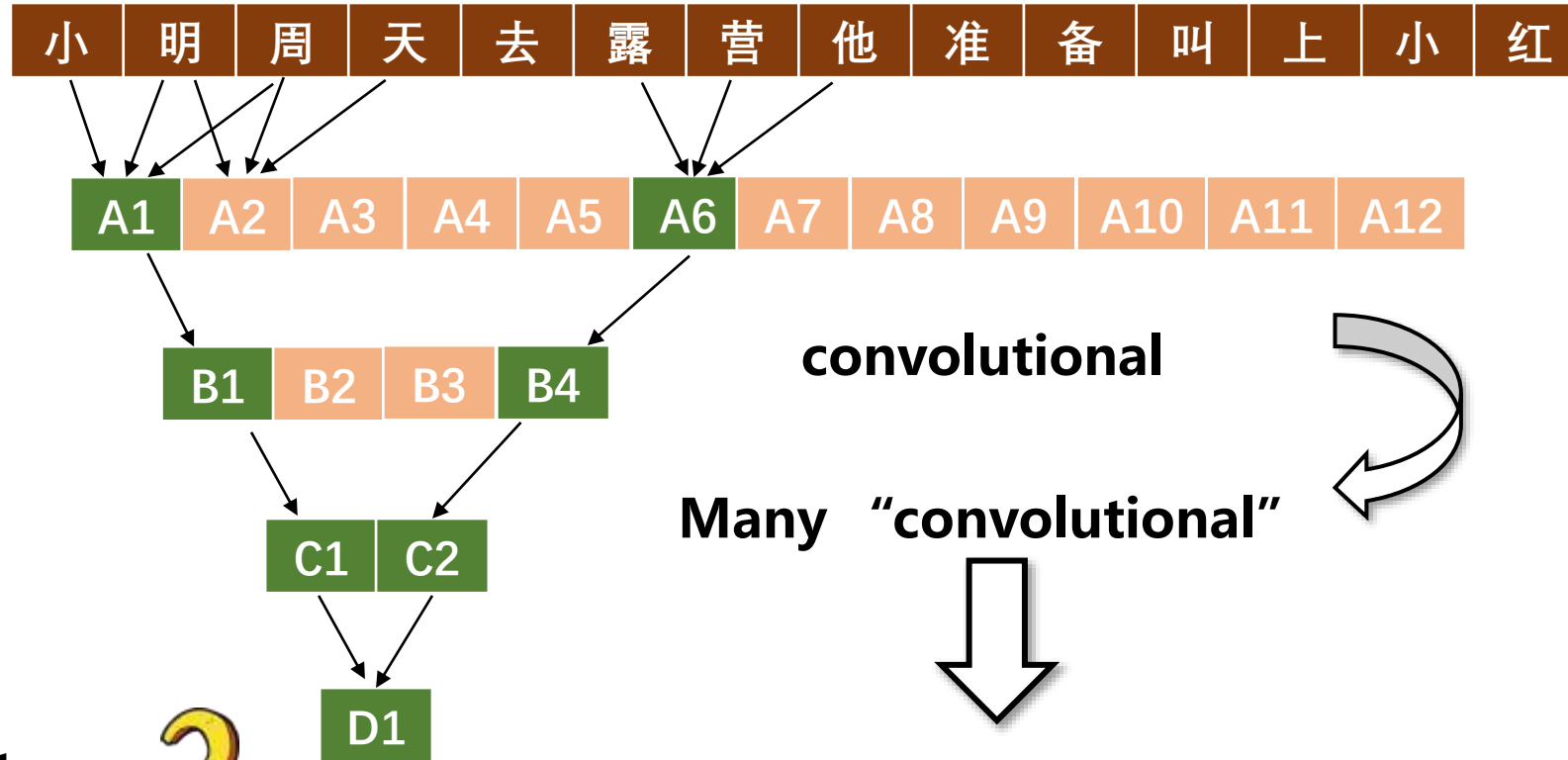
What is “transformer” ?





# Deep learning

## CNN VS Transformer in NLP

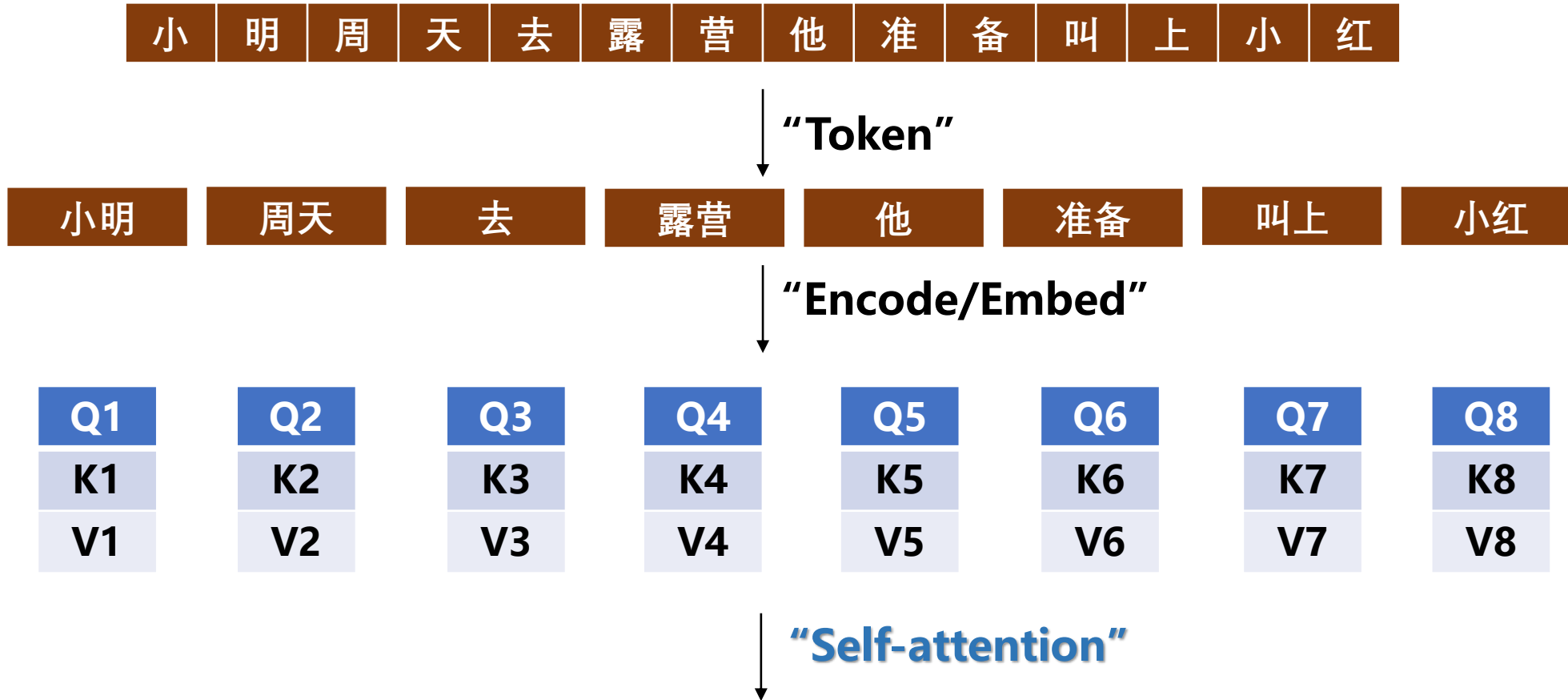


Disadvantages ?

Indirect information capture !

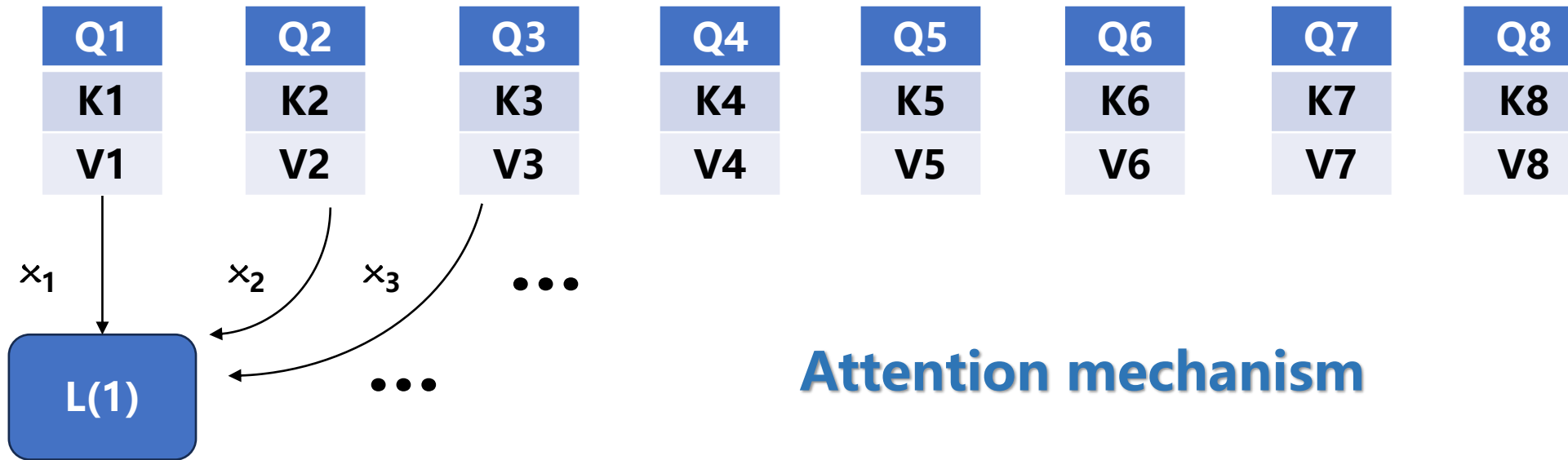
# Deep learning

## CNN VS Transformer in NLP



# Deep learning

## CNN VS Transformer in NLP

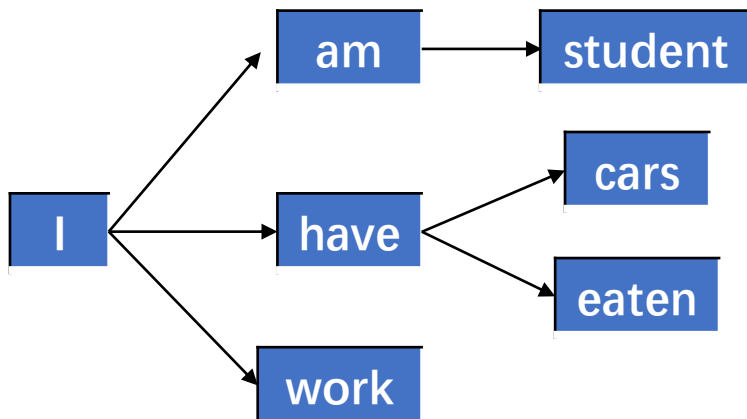
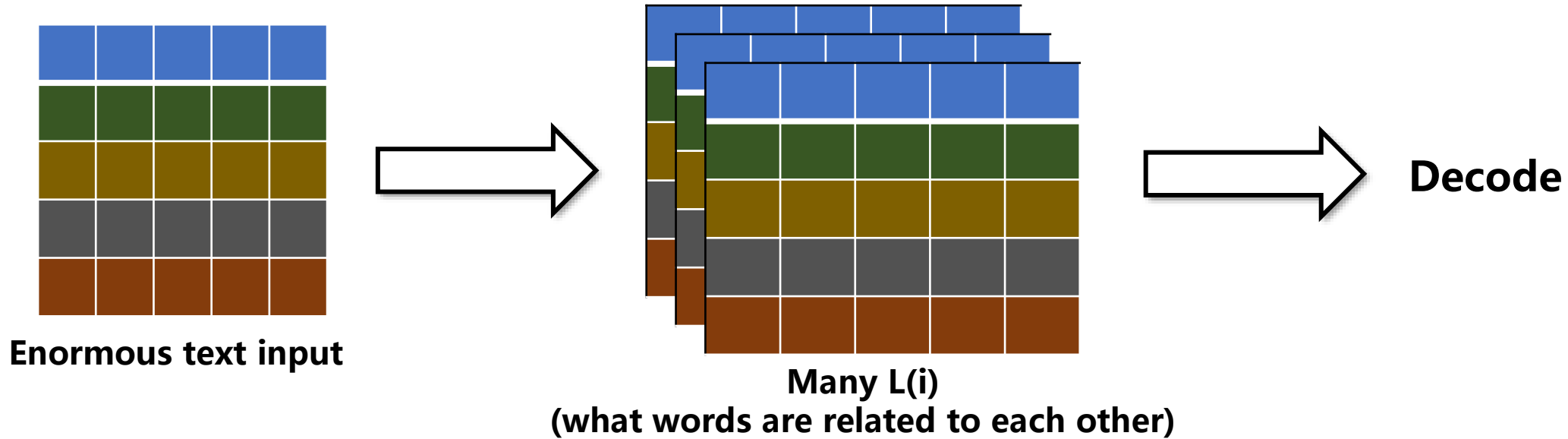


This framework can be applied to many other tasks



# Deep learning

## How was ChatGPT trained?



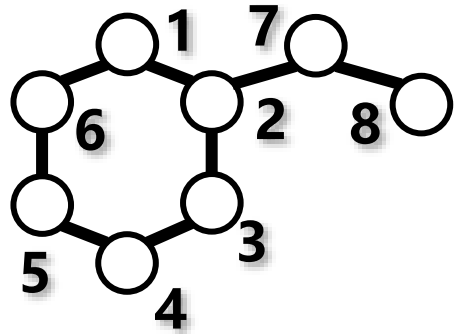
**Speak what humans speak**

**Speak what humans understand**

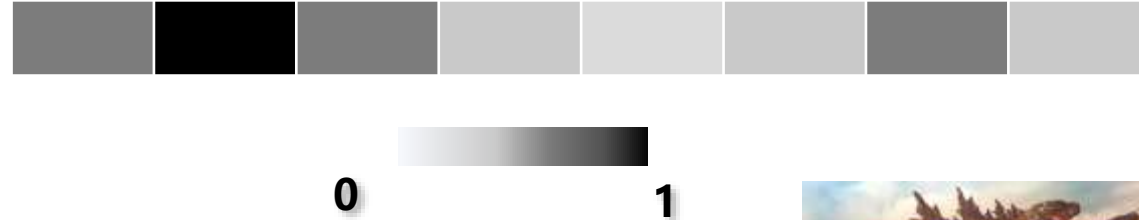
<https://www.stackbuilders.com/blog/inside-the-brain-of-chatgpt/>  
[https://www.nolibox.com/creator\\_articles/principle\\_of\\_ChatGPT.html](https://www.nolibox.com/creator_articles/principle_of_ChatGPT.html)

# Deep learning

## What is Graph Transformer?



$$L(A2) = x1 * H(A1) + x2 * H(A2) + x3 * H(A3) + \dots$$



## Comparison to normal GNN?

Normal GNN:

Manually identification of Nodes & Edges before learning

GTNN:

trap key features of nodes & Edges directly for learning



# Deep learning

## (A) Outcome & yield

	Reaction yield <i>r</i> value	Reaction yield m.a.e. (%)	Binary reaction outcome (random split), AUC (%)	Binary reaction outcome (substrate split), AUC (%)
GTNN2D	0.896±0.006	4.53±0.09	<b>91.8±2.1</b>	52±2
GNN2D	0.866±0.005	5.61±0.06	87.5±1.0	51±2
GTNN3D	0.884±0.01	4.51±0.11	91.4±0.7	58±4
GNN3D	0.877±0.001	5.33±0.34	89.4±0.8	65±5
GTNN2DQM	<b>0.898±0.003</b>	4.41±0.17	90.9±1.5	53±5
GNN2DQM	0.876±0.01	5.41±0.10	89.0±1.1	59±5
GTNN3DQM	0.890±0.01	<b>4.23±0.08</b>	<b>91.8±0.9</b>	<b>67±2</b>
GNN3DQM	0.890±0.006	4.88±0.24	89.1±0.9	64±4
ECFP4NN	0.885±0.0006	4.55±0.14	89.3±1.3	52±3

Testing data from the experimental data

***r*** (Pearson coefficient):  
预测值和观测值的相关系数。  
*r*越接近1, 相关性越强。

**m.a.e.** (Mean absolute error):  
预测值和观测值之间绝对误差的平均值。  
m.a.e.越小, 预测越准确。

**AUC** (Area under curve):  
在二分类问题中, AUC越大, 性能越好。

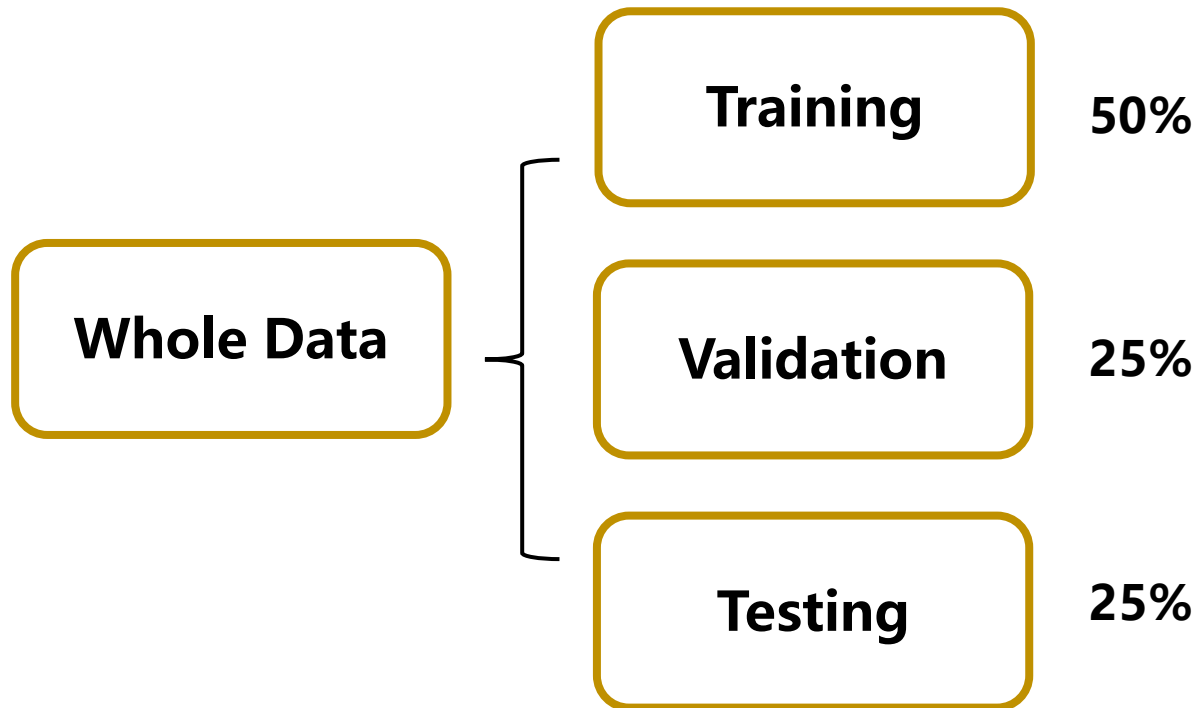
Reaction  
yield

Reaction  
outcome

# Deep learning

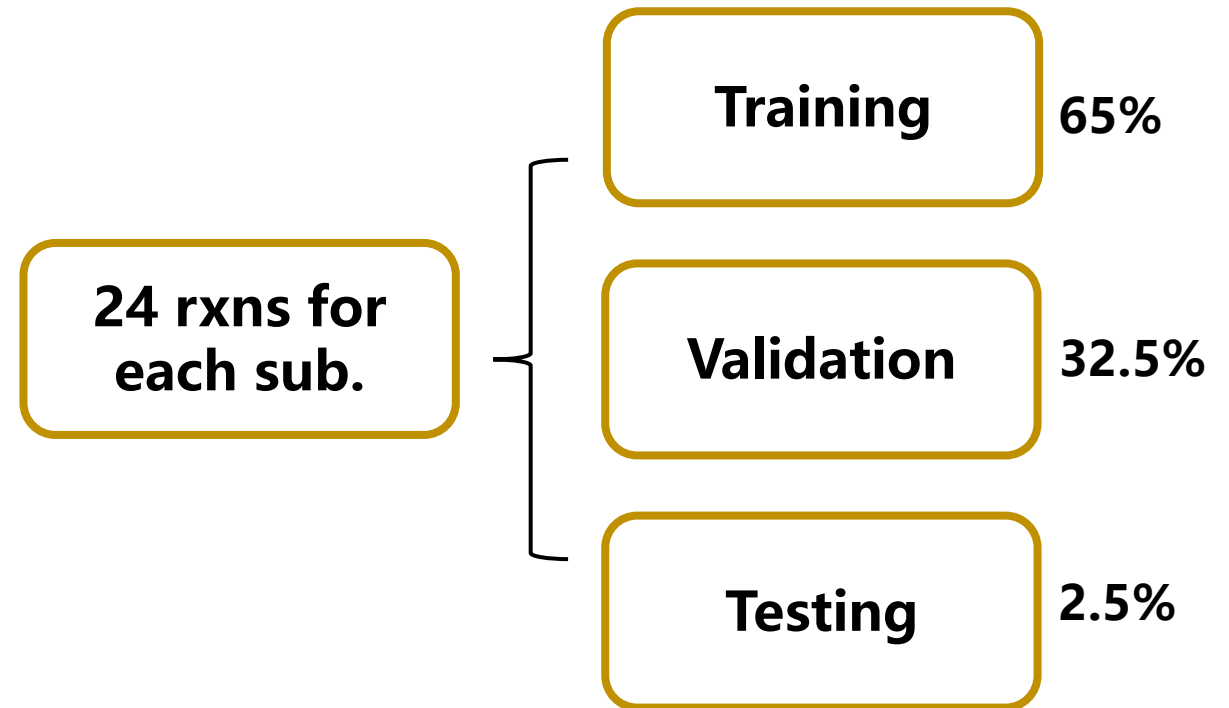
## (A) Outcome & yield

### Random split



Investigate the performance on **new conditions** for **known substrates**

### Substrate-based split



Investigate the performance on **unknown substrates** for **known conditions**

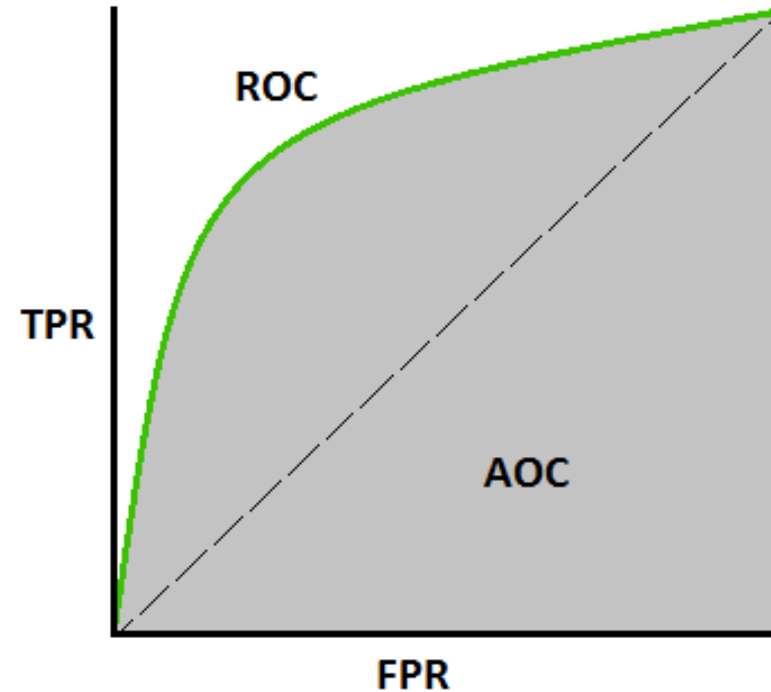
# Deep learning

## (A) Outcome & yield

### What is AUC in DL?

**TPR: 真正率, 正确预测为正样本的样本数占真正样本的比例**

**FPR: 假正率, 错误预测为正样本的样本数占真负样本的比例**



	预测		
		正样本	负样本
实际	正样本	True Positive (TP)	False Negative (FN)
	负样本	False Positive (FP)	True Negative (TN)

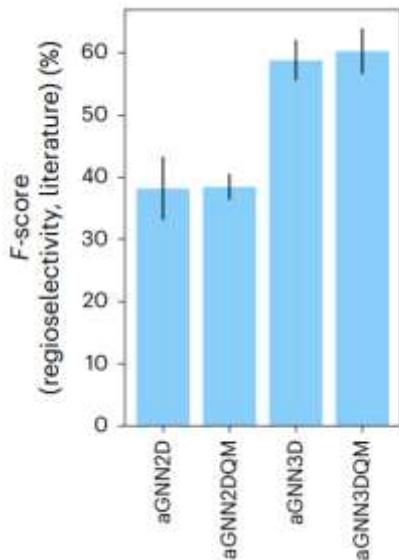


# Deep learning

## (B) Regioselectivity

	F-score (%)	PPV (%)	TPR (%)	Accuracy (%)
aGNN2D	38±5	56±1	30±6	88±1
aGNN2DQM	39±2	54±2	30±3	88±0.3
aGNN3D	59±3	<b>62±2</b>	56±4	<b>90±1</b>
aGNN3DQM	<b>60±4</b>	<b>62±2</b>	<b>59±6</b>	<b>90±1</b>

### Testing data from the experimental data



### Testing data from literature data

		预测	
		正样本	负样本
实际	正样本	True Positive (TP)	False Negative (FN)
	负样本	False Positive (FP)	True Negative (TN)

**a** : atomic features (No pooling)

**PPV** (positive predictive value): 阳性预测值。识别为正确的样品中真正正确的比值。

$$PPV = \frac{TP}{TP + FP}$$

**TPR** (true positive rate): 阳性预测正确值。真正正确的样品中被识别为正的比值。

$$TPR = \frac{TP}{TP + FN}$$

**F-score**: PPV和TPR的加权平均，权衡两个指标。

$$F - score = \frac{PPV * TPR * 2}{PPV + FPR}$$

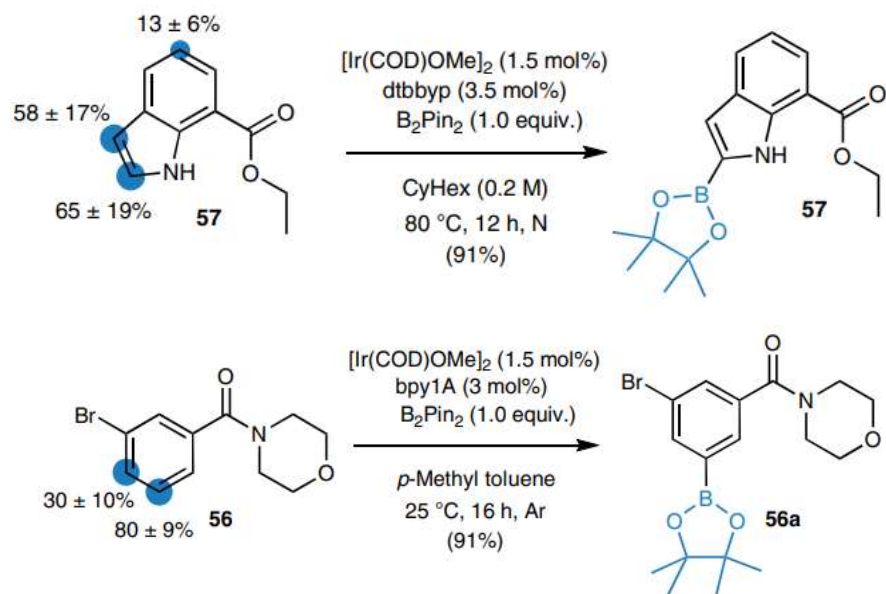
**Accuracy**: 被正确分类的样品占总体的比值。

$$Accuracy = \frac{TP + TN}{TP + FN + FP + TN}$$

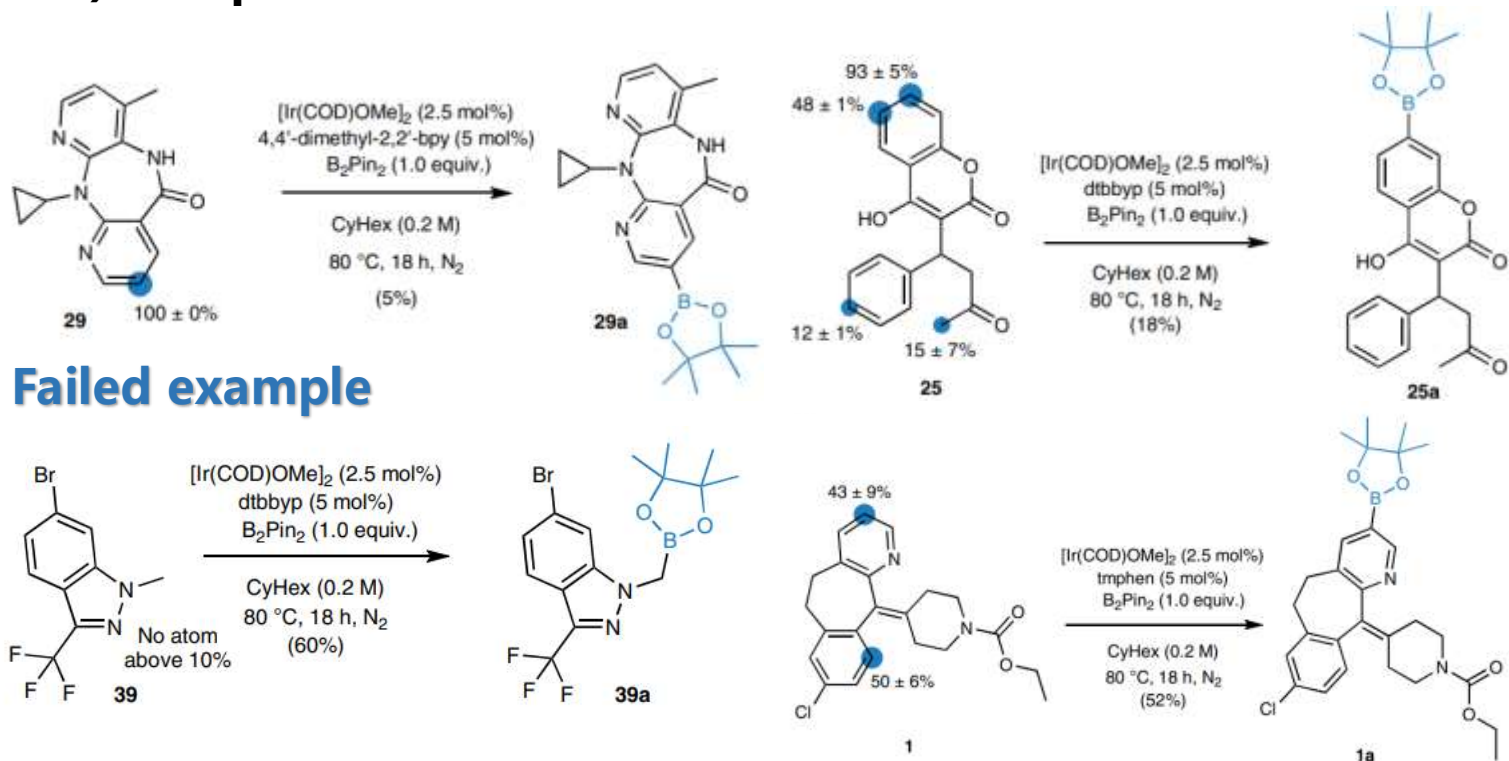
# Deep learning

## (B) Regioselectivity

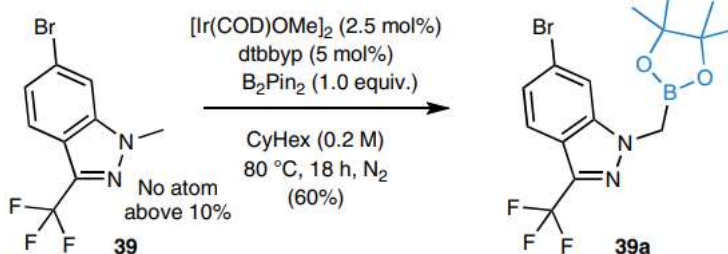
### 1) Retrospective from test data



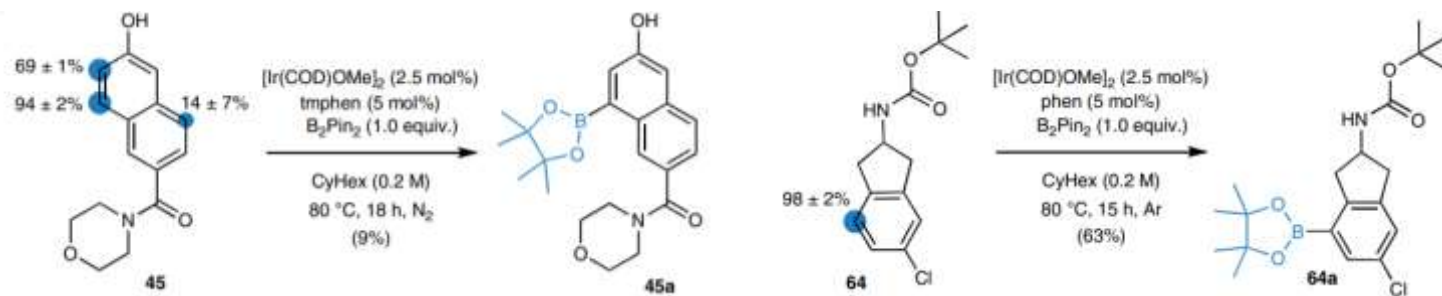
### 2) Prospective from literature data



### Failed example



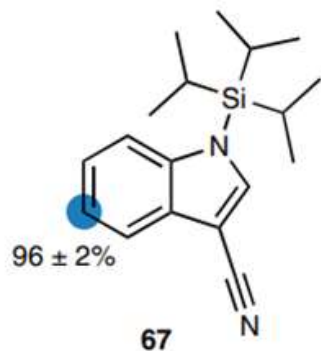
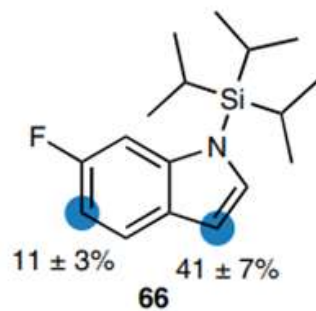
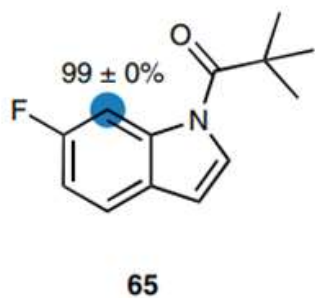
### 3) Retrospective from Roche dataset



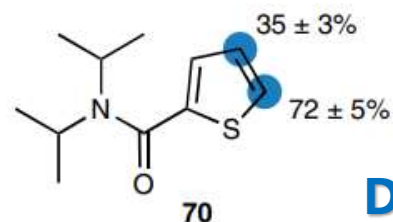
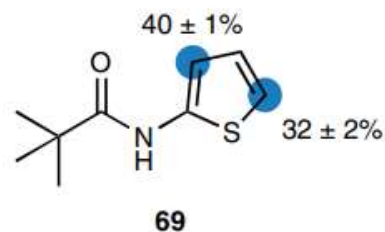
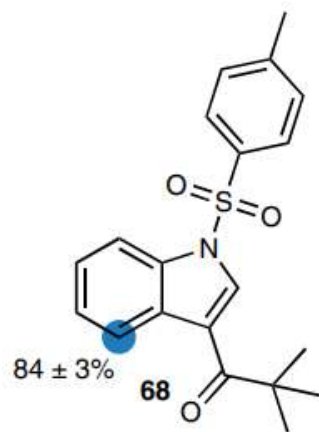
# Deep learning

## (B) Regioselectivity

### 4) Identification of 3D properties



**Steric hindrance**



**Directing group**

**Unseen molecules**

# Discussion

## (1) HTE VS Literature

### Prediction in outcome & yield

	Reaction yield <i>r</i> value	Reaction yield m.a.e. (%)
GTNN2D	0.896±0.006	4.53±0.09
GNN2D	0.866±0.005	5.61±0.06
GTNN3D	0.884±0.01	4.51±0.11
GNN3D	0.877±0.001	5.33±0.34
GTNN2DQM	<b>0.898±0.003</b>	4.41±0.17
GNN2DQM	0.876±0.01	5.41±0.10
GTNN3DQM	0.890±0.01	<b>4.23±0.08</b>
GNN3DQM	0.890±0.006	4.88±0.24
ECFP4NN	0.885±0.0006	4.55±0.14

### Validated by HTE results

Prediction error	Mean absolute error / %
GTNN2D	16.7 (±0.13)
GNN2D	16.4 (±0.2)
GTNN3D	16.4 (±0.24)
GNN3D	16.2 (±0.14)
GTNN2DQM	<b>16.1 (±0.02)</b>
GNN2DQM	16.3 (±0.04)
GTNN3DQM	16.2 (± 0.16)
GNN3DQM	16.2 (±0.14)
ECFP4NN	18.2 (±0.05)

### Validated by literature results

1. HTE: All in the same standard

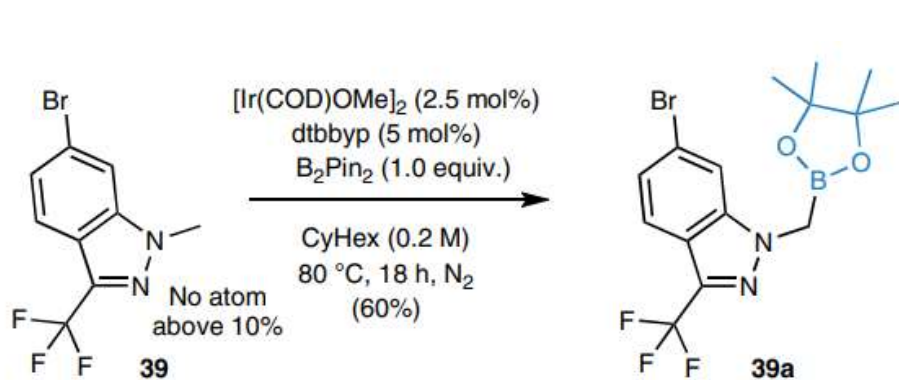
Lit.: standards variates

2. HTE covers a less diverse reaction parameter space

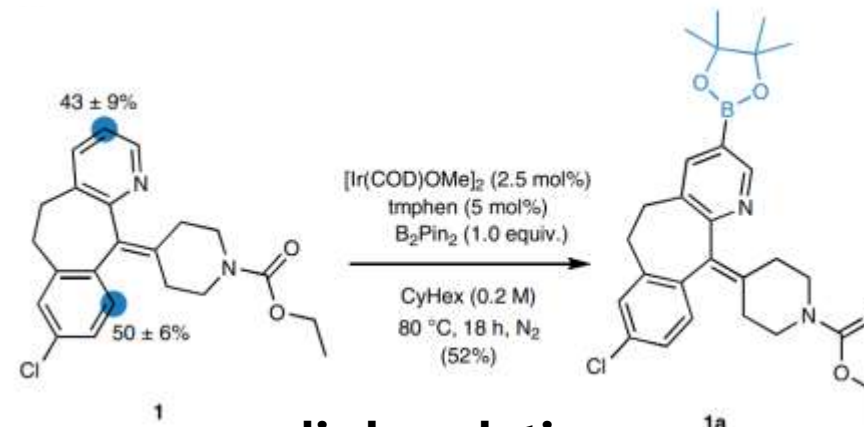
Why different ?

# Discussion

## (2) Deep learning greatly relies on the quality of input data!

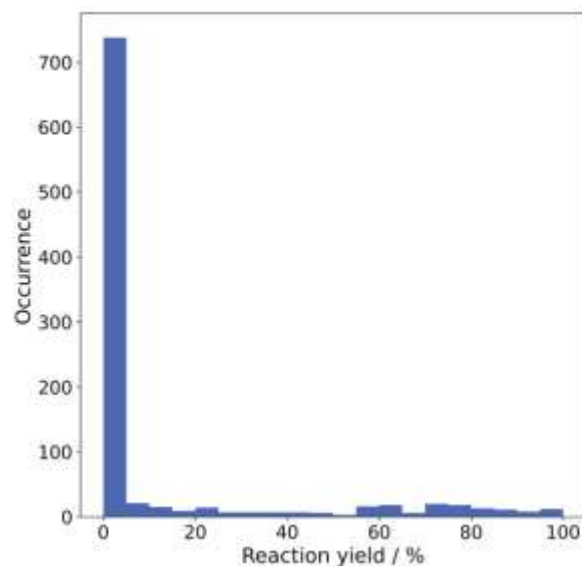


$sp^3$  borylation

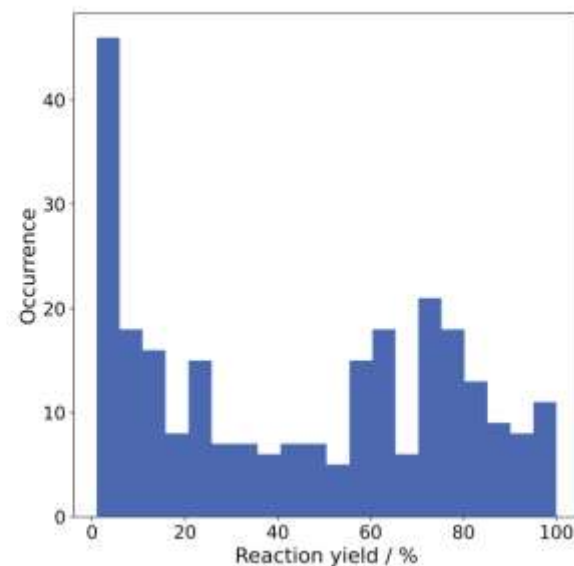


di-borylation

Whole dataset



Only positive reactions



# Discussion

## (3) Does this research deserve Nature Chemistry?

To be honest, I don't think so...

Combining **DL** and **HTS** in synthetic chemistry is not a fancy idea now...



**Highlights of this article: Application in C-H borylation, thus for LSF**

**However, no prospective application on out-of-dataset was reported**

**DL+HTS should have guidance significance to drug synthesis**

- **Tell what condition is the best?**
- **In combination with AI-promoted retro-synthesis analysis? (New route)**
- **In combination with virtual screening? (High-quality scaffold)**



**Thanks for your attention**

**Bin Yang**  
**2024.01.13**

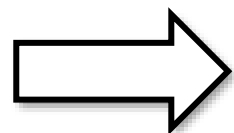




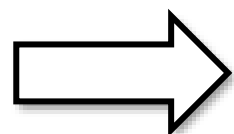
# Input for Deep Learning

## (A) Molecular fingerprint

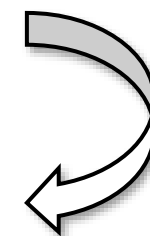
— Example: Retrosynthetic accessibility score



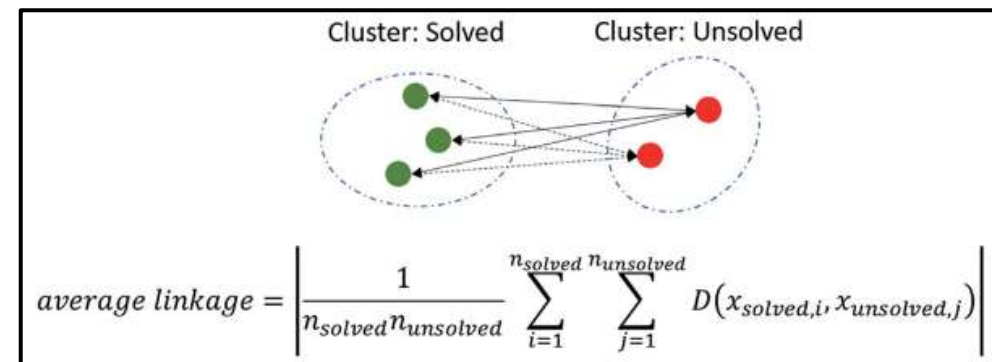
**Known retrosynthetic route  
(Solved)**



**Unknown retrosynthetic route  
(Unsolved)**



**RAscore**



# Input for Deep Learning

## (B) Density Functional Theory calculation

The **first principle** of Kohn-Hohenberg:

只要知道基态电荷数密度，那么就能确定基态的能量

The **second principle** of Kohn-Hohenberg:

使整体能量最小的电荷密度就是真实电荷密度



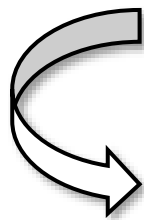
## DFT = Solving the Schrödinger equation

### 自洽求解:

- S1: 猜一个电荷数密度  $n_{trial}(r)$ , 代入方程, 求得波函数  $\Psi_i(r)$ ;
- S2: 将  $\Psi_i(r)$  代入方程, 求得  $n_{sol}(r)$ ;
- S3: 比较  $n_{trial}(r)$  和  $n_{sol}(r)$ , 若误差小于给定精度, 输出结果, 否则用  $n_{sol}(r)$  取代  $n_{trial}(r)$ , 循环往复。

### DFT in DL:

- 1).  $\Delta G$  in transition states
- 2). Atomic partial charges

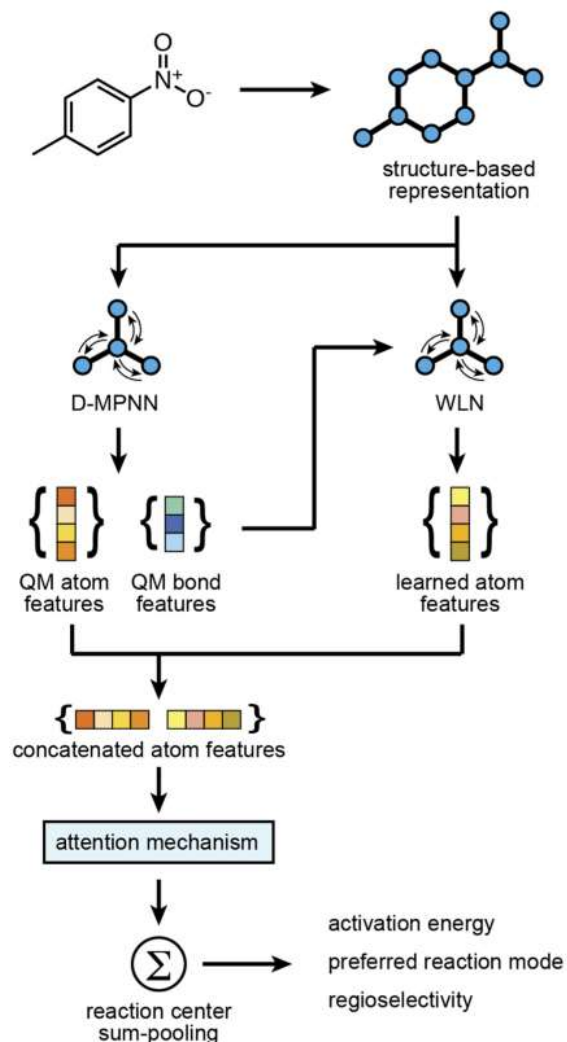


Regioselectivity, Diastereoselectivity

# Input for Deep Learning

## (B) Density Functional Theory calculation

### Example: Reactivity prediction



**Atom features:** Atomic number, ring status, etc.

**QM bond features:** Bond order, bond length from QM calculation, etc.

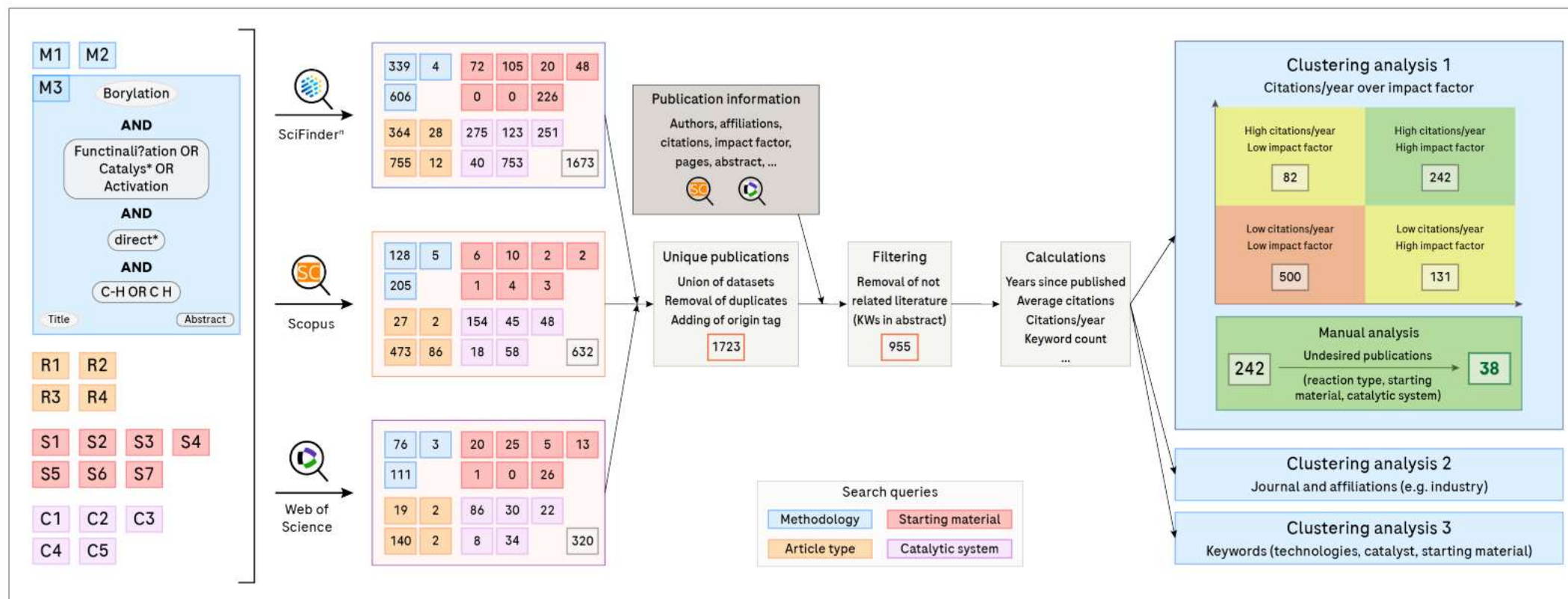
**QM atom features:** Hirshfeld atomic charges, nucleophilic and electrophilic Fukui functions, etc.

Table Average MAE (kcal/mol) when predicting SN2 barrier heights

Labeled points	Baseline GNN	ml-QM-GNN	Chemprop	BoB <sup>a</sup>	SLATM <sup>a</sup>
225 (180 + 45)	9.07 ± 0.04	3.61 ± 0.14	6.71 ± 0.08	4.89	4.44
450 (360 + 90)	8.89 ± 0.13	3.28 ± 0.04	4.01 ± 0.02	4.28	3.87
900 (720 + 180)	8.61 ± 0.07	2.97 ± 0.03	3.23 ± 0.07	3.78	3.21
1800 (1440 + 360)	8.49 ± 0.03	2.76 ± 0.01	2.85 ± 0.02	3.49	2.92

<sup>a</sup>Taken directly from the work of Heinen, von Rudorff, and von Lilienfeld.<sup>40</sup>

# Method—Systematic literature analysis



**M: Different methodologies (Direct/Indirect, etc.)**

**S: Different substrates (aromatic/aliphatic, etc.)**

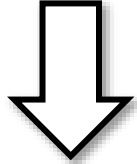
**R: Article types (review/article, etc.)**

**C: Catalytic systems (Ir/Rh, etc.)**

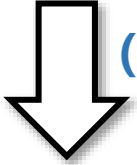
**Query keywords**

# Method—Substrate selection

Cortellis Drug Discovery Intelligence (CDDI database)

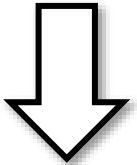


1174 drugs (200 < MW < 800)

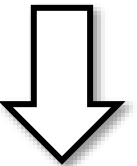


(Based on ECFP4)

8 clusters



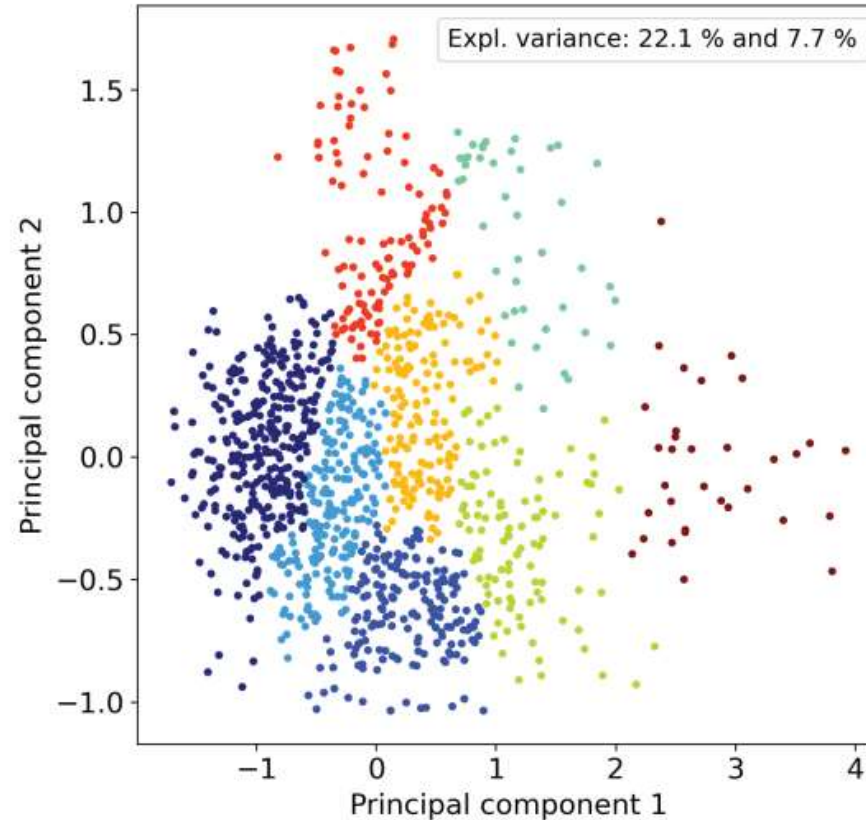
10 closest molecules  
to the cluster center



3/10 commercially available  
Total sample: 23

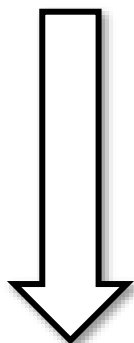
## DRUGS

Jaccard similarity:  $J = \frac{A \cap B}{A \cup B}$



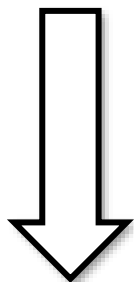
# Method—Substrate selection

Top 100 most popular ring assemblies in Roche's database



- MW < 300
- Heavy atoms < 20
- > 1 g stock available
- Not involved in other projects

268 fragments



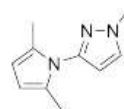
- Manually selection:**
- Halogen / OH on aromatic rings
  - Frequently used heterocycles

16 fragments

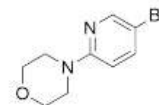
+ 5 frequently occurring literature substrates

## FRAGMENTS

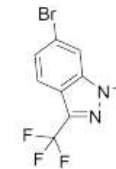
Fragments



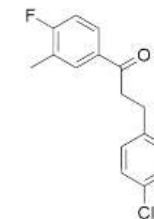
37



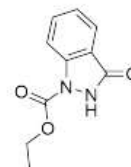
38



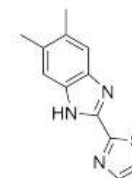
39



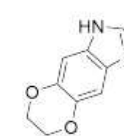
40



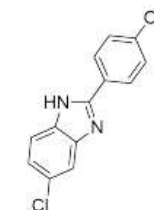
41



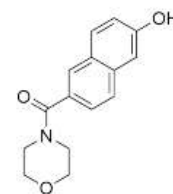
42



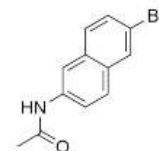
43



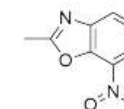
44



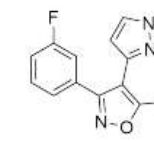
45



46



47



48

# Method—Substrate selection

# Input for Deep Learning

## (B) Density Functional Theory calculation

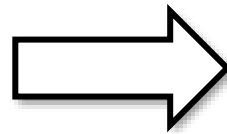
What does DFT do? An approximate solution for the Schrödinger equation

Schrödinger equation in describing electrons:

$$H\Psi(r_1, \dots, r_N) = E\Psi(r_1, \dots, r_N)$$

$$E = -\underbrace{\frac{\hbar^2}{2m_e} \sum_i^N \nabla^2}_{\text{电子动能}} + \underbrace{\sum_i^N V_{ext}(r_i)}_{\text{电子与相对固定的原子核的交互}} + \underbrace{\sum_{i=1}^N \sum_{j=i+1}^N U(r_i, r_j)}_{\text{电子-电子的交互}} \quad \text{Eq (1)}$$

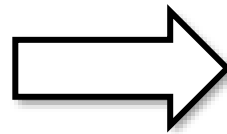
Considering a system with N electrons, there will be **3N** variables!



Almost impossible to accurately describe the situation of electrons!



All electrons as objects



**Electron density for ONE electron** as objects

Only **3** variables



# Input for Deep Learning

## (B) Density Functional Theory calculation

The first principle of Kohn-Hohenberg:

只要知道基态电荷数密度, 那么就能确定基态的能量

How can we get  $n(r)$  ?

$$E[\{\Psi_i\}] = E_{known}[\{\Psi_i\}] + E_{XC}[\{\Psi_i\}] \quad \text{Eq (2)} \quad ?$$

Where:

$$E_{known}[\{\Psi_i\}] = \underbrace{-\frac{\hbar^2}{m_e} \sum_i \Psi_i^* \nabla^2 \Psi_i d^3r}_{\text{电子动能}} + \underbrace{\int V(r) n(r) dr^3}_{\text{电子与相对固定的原子核的交互}} + \underbrace{\frac{e^2}{2} \iint \frac{n(r)n(r')}{|r-r'|} d^3r d^3r'}_{\text{电子-电子的交互}} \quad \text{Eq (3)}$$

$E_{XC}[\{\Psi_i\}]$  : 交换关联泛函, DFT的关键, 只能近似

Luo et al., *Science* 381, 1072–1079 (2023)

Fig. 4. Five proposed pathways for the oxidative addition of haloacetonitriles to ionic or neutral Cu(I) complexes and free energies of each species computed with DFT. The calculated activation free energies for the oxidative addition of BrCH<sub>2</sub>CN (**2-Br**) to the ionic Cu(I) complex [Ph<sub>4</sub>P]<sup>+</sup>[Cu(CF<sub>3</sub>)<sub>2</sub>]<sup>-</sup> (**1a**) in DMSO are given in blue, and those for oxidative

addition of ClCH<sub>2</sub>CN (**2-Cl**) to the neutral [(bpy)Cu(CF<sub>3</sub>)] (**1b**) in DMF are given in red. The energies are in kilocalories per mole and indicate the relative free energies calculated at the PBE0-D3(BJ)/Def2-TZVP(SMD, solvent)//PBE0-D3(BJ)/Def2-SVP(SMD, solvent) level. SMD, solvation model density.

# Input for Deep Learning

## (B) Density Functional Theory calculation

The second principle of Kohn-Hohenberg:

使整体能量最小的电荷密度就是真实电荷密度

$$n(\mathbf{r}) = 2 \times \sum_i \Psi_i^*(\mathbf{r}) \Psi_i(\mathbf{r}) \quad \text{Eq (4)}$$

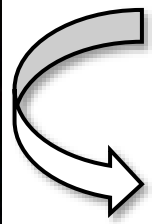


自洽求解:

- S1: 猜一个电荷数密度  $n_{trial}(\mathbf{r})$ , 代入KS方程与Eq.2, 求得波函数  $\Psi_i(\mathbf{r})$ ;
- S2: 将  $\Psi_i(\mathbf{r})$  代入Eq.4, 求得  $n_{sol}(\mathbf{r})$ ;
- S3: 比较  $n_{trial}(\mathbf{r})$  和  $n_{sol}(\mathbf{r})$ , 若误差小于给定精度, 输出结果, 否则用  $n_{sol}(\mathbf{r})$  取代  $n_{trial}(\mathbf{r})$ , 循环往复。

Calculation results in ML input:

- 1).  $\Delta G$  in transition states
- 2). Atomic partial charges



Regioselectivity, Diastereoselectivity

# Principal component analysis (PCA)

## Lipinski's Rule of Five:

- $MW < 500$
- $HBD < 5$
- $HBA < 10$
- $\log P < 5$

Adv Drug Deliver Rev, 1997, 23(1-3), 3  
J. Med. Chem. 2002, 45, 12, 2615  
J. Med. Chem. 2001, 44, 12, 1841  
J. Med. Chem. 2009, 52, 21, 6752

## Veber's Rule:

- Rotable bonds  $< 10$
- $PSA < 140$

## Muegge's Rule:

- Num Rings  $< 7$

## Fsp3:

- $0.3 < Fsp3 < 0.5$

