

---

## **Build a Molecular Search Engine in Minutes**

@Towhee.io

6/23/2022

---

## Build a Molecular Search Engine in Minutes

### Introduction

duration: 1

Drug discovery, as the source of medical innovation, is an important part of new medicine research and development. Drug discovery is implemented by target selection and confirmation. In order to discover available compounds in the fragment space from billion-scale compound libraries, chemical fingerprint is usually retrieved for substructure search and similarity search.

This example will show you how to find the similar, sub or super molecular formula. Moreover, we managed to make the core functionality as simple as 10 lines of code with Towhee, so that you can start hacking your own molecular search engine.

### Preparations

duration: 2

### Install Dependencies

First we need to install dependencies such as pymilvus, towhee, rdkit and gradio.

```
1 $ python -m pip install -q pymilvus towhee rdkit-pypi gradio
```

### Prepare the Data

There is a subset of the Pubchem dataset (10000 SMILES) used in this demo, everyone can download on Github.

```
1 $ curl -L https://github.com/towhee-io/examples/releases/download/data/pubchem_10000.smi -O
```

**pubchem\_10000.smi**: a file containing SMILES and corresponding ids.

Let's take a quick look:

```

1 import pandas as pd
2
3 df = pd.read_csv('pubchem_10000.smi')
4 df.head()

```

	smiles	id
0	<chem>COc1cc(C=NNC(=O)CC2(C)OCCO2)cc(c1OCc1cccc2c1cc...</chem>	2860827
1	<chem>Fc1ccc(cc1)C=C1C(=O)NC(=S)N(C1=O)c1ccc(cc1Cl)Cl</chem>	2860828
2	<chem>O=C1NC(=C(C(N1)c1cccc(c1)Br)C(=O)Nc1ccc(cc1)Cl)C</chem>	2860829
3	<chem>COc1ccc(c(c1)OC)C=C1C(=O)NC(=S)N(C1=O)c1cccc(c...</chem>	2860830
4	<chem>Nc1ccn(n1)c1ccccc1.Cl</chem>	2860831

To use the dataset for molecular search, let's first define the dictionary and helper function:

- `id_smiles`: a dictionary of id and corresponding smiles;
- `to_images(input)`: convert the input smiles or results to towhee.Image for display.

```

1 from rdkit.Chem import Draw
2 from rdkit import Chem
3 from towhee.types.image_utils import from_pil
4
5 id_smiles = df.set_index('id')['smiles'].to_dict()
6
7 def to_images(inputs):
8     if isinstance(inputs, str):
9         smiles = inputs
10        mol = Chem.MolFromSmiles(smiles)
11        return from_pil(Draw.MolToImage(mol))
12    imgs = []
13    results = inputs
14    for re in results:
15        smiles = id_smiles[re.id]
16        mol = Chem.MolFromSmiles(smiles)
17        imgs.append(from_pil(Draw.MolToImage(mol)))
18    return imgs

```

## Create Milvus Collection

Before getting started, please make sure you have installed milvus. Let's first create a `molecular_search` collection that uses the L2 distance metric and an IVF\_FLAT index.

```
1 from pymilvus import connections, FieldSchema, CollectionSchema,
  DataType, Collection, utility
2
3 connections.connect(host='127.0.0.1', port='19530')
4
5 def create_milvus_collection(collection_name, dim):
6     if utility.has_collection(collection_name):
7         utility.drop_collection(collection_name)
8
9     fields = [
10        FieldSchema(name='id', dtype=DataType.INT64, description='ids',
11                    is_primary=True, auto_id=False),
12        FieldSchema(name='embedding', dtype=DataType.BINARY_VECTOR,
13                    description='embedding vectors', dim=dim)
14    ]
15    schema = CollectionSchema(fields=fields, description='molecular
16                             similarity search')
17    collection = Collection(name=collection_name, schema=schema)
18
19    return collection
20
21 collection = create_milvus_collection('molecular_search', 2048)
```

## Load Molecular Fingerprint into Milvus

duration: 2

We first generate fingerprint from SMILES with daylight algorithm and insert the fingerprints into Milvus. Towhee provides a method-chaining style API so that users can assemble a data processing pipeline with operators.

```
1 import towhee
2
3 dc = (
4     towhee.read_csv('pubchem_10000.smi')
5     .runas_op['id', 'id'](func=lambda x: int(x))
6     .molecular_fingerprinting['smiles', 'fp'](algorithm='daylight')
7     .to_milvus['id', 'fp'](collection=collection, batch=100)
8 )
9
10 print('Total number of inserted data is {}'.format(collection.
11              num_entities))
```

Total number of inserted data is 10000.

**Explanation of Data Processing Pipeline** Here is detailed explanation for each line of the code:

`towhee.read_csv('pubchem_10000.smi')`: read tabular data from the file (smiles and id columns);

`.runas_op['id', 'id'](func=lambda x: int(x))`: for each row from the data, convert the data type of the column id from `str` to `int`;

`.molecular_fingerprinting['smiles', 'fp'](algorithm='daylight')`: use the daylight algorithm to generate fingerprint with the rdkit operator in towhee hub.

`.to_milvus['id', 'fp'](collection=collection, batch=100)`: insert molecular fingerprints in to Milvus;

## Query Molecular from Milvus with Towhee

duration: 3

Now that fingerprint for candidate SMILES have been inserted into Milvus, we can query across it. Again, we use Towhee to load the input SMILES, compute a fingerprint, and use it as a query in Milvus. Because Milvus only outputs IDs and distance values, we provide the `id_smiles` dictionary to get the original smiles based on IDs and display.

### similarity search

```
1 ( towhee.dc['smiles'](['Cn1ccc(=O)nc1', 'CN1C=NC2=C1C(=O)N(C(=O)N2C)C',
2   'CCOC(=O)C1CN1C(C(=O)OCC)CC'])
3   .molecular_fingerprinting['smiles', 'fp'](algorithm='daylight')
4   .milvus_search['fp', 'result'](collection=collection, metric_type
5     = 'JACCARD')
6   .runas_op['result', 'similar_smile'](func=lambda res: [id_smiles[
7     x.id] for x in res])
8   .select['smiles', 'similar_smile']()
9   .show()
10 )
```

smiles	similar_smile
Cn1ccc(=O)nc1	[Cn1ccc(=O)nc1,Cn1c[n+](C)ccc1=O,Cn1cnc1=O,Cn1ccc(=O)[nH]c1=O,...] len=10
CN1C=NC2=C1C(=O)N(C(=O)N2C)C	[Cn1cnc2c1c(=O)n(C)c(=O)n2C,Cn1c2nc[nH]c2c(=O)n(C)c1=O,Cn1cnc2c1c(=O)[nH]c(=O)n2C,Cn1cnc2[nH]c(=O)n(C)c(=O)c21,...] len=10
CCOC(=O)C1CN1C(C(=O)OCC)CC	[CCOC(=O)C1CN1C(C(=O)OCC)CC,CCCCOC(=O)C1CN1C(C(=O)OCC)CC,CCCCCOC(=O)C1CN1C(C(=O)OCC)CC,CCC(N1CC1C(=O)OC)C(=O)OC,...] len=10

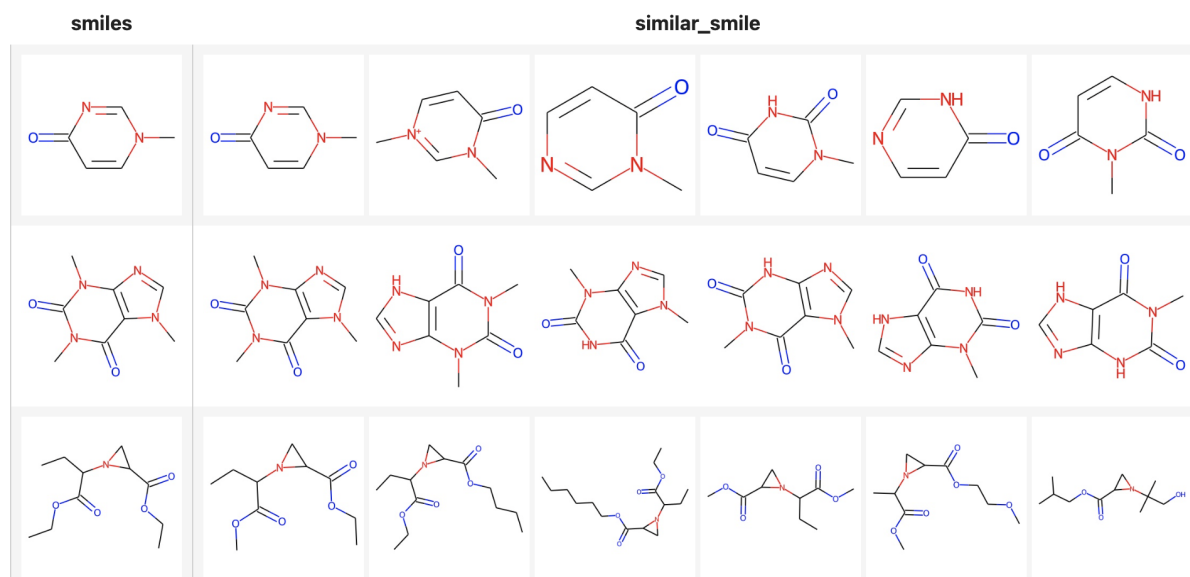
If you want to show the molecular structure with images, you can use the `to_images` function.

```
1 ( towhee.dc['smiles'](['Cn1ccc(=O)nc1', 'CN1C=NC2=C1C(=O)N(C(=O)N2C)C',
2   'CCOC(=O)C1CN1C(C(=O)OCC)CC'])
3   .molecular_fingerprinting['smiles', 'fp'](algorithm='daylight')
4   .milvus_search['fp', 'result'](collection=collection, metric_type
5     = 'JACCARD', limit=6)
6   .runas_op['result', 'similar_smile'](func=to_images)
7   .runas_op['smiles', 'smiles'](func=to_images)
```

```

6     .select['smiles', 'similar_smile']()
7     .show()
8 )

```



**Superstructure and Substructure search** Milvus not only supports searching similar structures of molecular formulas, but also superstructure and substructure searches, you only need to specify the metric types:

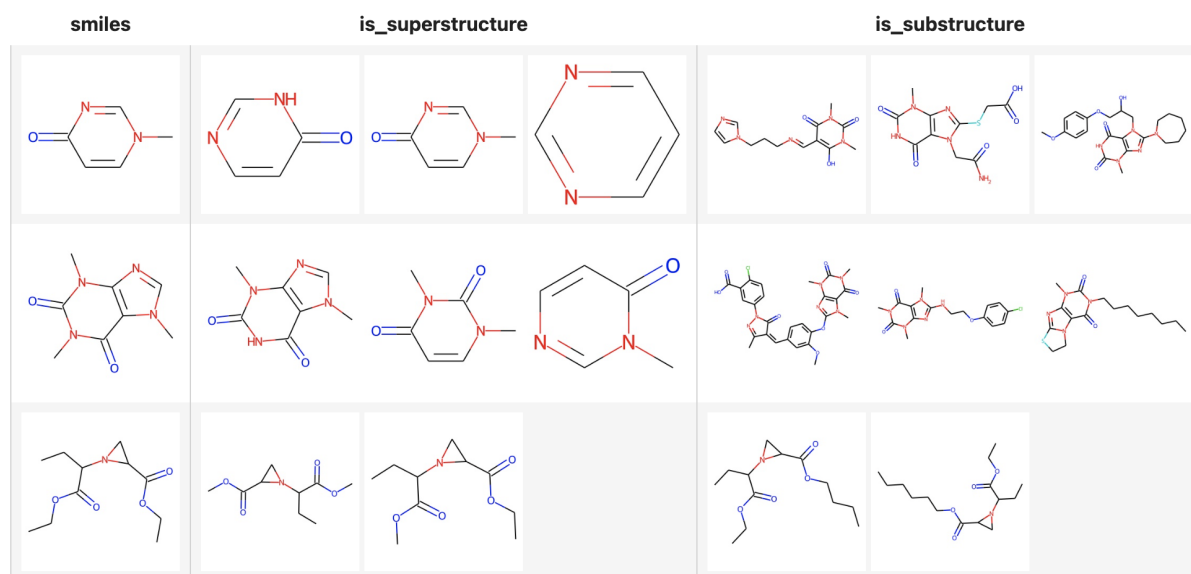
- Similarly search: “JACCARD”
- Superstructure search: “SUPERSTRUCTURE”
- Substructure search: “SUBSTRUCTURE”

In the following example, the limit is set to 3, but there are less than 3 substructures or superstructures of the query formula in the Milvus dataset.

```

1  ( towhee.dc['smiles'](['Cn1ccc(=O)nc1', 'CN1C=NC2=C1C(=O)N(C(=O)N2C)C',
2      'CCOC(=O)C1CN1C(C(=O)OCC)CC'])
3      .molecular_fingerprinting['smiles', 'fp'](algorithm='daylight')
4      .milvus_search['fp', 'result_super'](collection=collection,
5          metric_type='SUPERSTRUCTURE', limit=3)
6      .milvus_search['fp', 'result_sub'](collection=collection,
7          metric_type='SUBSTRUCTURE', limit=3)
8      .runas_op['result_super', 'is_superstructure'](func=to_images)
9      .runas_op['result_sub', 'is_substructure'](func=to_images)
10     .runas_op['smiles', 'smiles'](func=to_images)
11     .select['smiles', 'is_superstructure', 'is_substructure']()
12     .show()
13 )

```



## Release a Showcase

duration: 2

We've done an excellent job on the core functionality of our molecular search engine. Now it's time to build a showcase with interface. Gradio is a great tool for building demos. With Gradio, we simply need to wrap the data processing pipeline via a `search_smiles_with_metric` function:

```

1 def search_smiles_with_metric(smiles, metric_type):
2     def smiles_to_pil(smiles):
3         mol = Chem.MolFromSmiles(smiles)
4         return Draw.MolToImage(mol)
5
6     with towhee.api() as api:
7         milvus_search_function = (
8             api.molecular_fingerprinting(algorithm='daylight')
9             .milvus_search(collection='molecular_search',
10                           metric_type=metric_type, limit=5)
11             .runas_op(func=lambda res: [smiles_to_pil(id_smiles[x.
12                                     id]) for x in res])
13             .as_function()
14         )
15     return milvus_search_function(smiles)
16
17 import gradio
18 interface = gradio.Interface(search_smiles_with_metric,
                               [gradio.inputs.Textbox(lines=1, default='
                               CN1C=NC2=C1C(=O)N(C(=O)N2C)C'),

```

```
19         gradio.inputs.Radio(['JACCARD', '
20             substructure', 'superstructure'])),
21     [gradio.outputs.Image(type="pil", label=
22         None) for _ in range(5)]
23     )
24 interface.launch(inline=True, share=True)
```