# Answers

## 1.Explain the distance metric you utilized to calculate the similarity/dissimilarity between small molecules.

**I use the the Euclidean distance. It can be calculated from the Cartesian coordinates of the points using the Pythagorean theorem, therefore occasionally being called the Pythagorean distance. As below:**

$$d(p,q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \cdots + (p_i - q_i)^2 + \cdots + (p_n - q_n)^2}.$$

## 2.Use a dimensionality reduction algorithm (PCA, t-SNE, UMAP, etc) to generate a 2D visualization of the small molecule dataset. Each point should represent a single molecule. (Note: you may have to plot a subset of your data, depending on which dimensionality reduction algorithm you choose.)

**I took a subset of 2000 ligands for the analysis, and I also bin the scores to set up colors for the gerneral TSNE plot**

In [5]:

```python
import glob
import pandas as pd
import numpy as np
import os
import sys
from algs import Ligand, HierarchicalClustering, PartitionClustering,Si
import matplotlib.pyplot as plt
from __future__ import print_function
import time
import numpy as np
import pandas as pd
from sklearn.datasets import fetch_openml
from sklearn.decomposition import PCA
from sklearn.manifold import TSNE
%matplotlib inline
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
import seaborn as sns
import json

csv = pd.read_csv('../ligand_information.csv', index_col= False)
Data = list(np.random.randint(1,csv.shape[0],size=2000))

Ligands = []

for i in Data:

    a_ligand = Ligand(csv[csv.index==i])
    Ligands.append(a_ligand)

IDs = []
scores = []
dataset = []
for ligand in Ligands:
    IDs.append(ligand.get_ID())
    scores.append(ligand.get_scores())
    dataset.append(ligand.get_OnBits())
dataset = np.stack((dataset), axis=0)
onBits = [ 'pixel'+str(i) for i in range(dataset.shape[1]) ]
df = pd.DataFrame(dataset,columns=onBits)
_df = df[onBits].values


time_start = time.time()
tsne = TSNE(n_components=2, verbose=1, perplexity=40, n_iter=300)
tsne_results = tsne.fit_transform(_df)
print('t-SNE done! Time elapsed: {} seconds'.format(time.time()-time_st

df['scores'] = scores
#score_ir = pd.interval_range(start =-6.2, end = 534.3, freq = 100)

df['bins_scores'] = pd.cut(df['scores'], bins = [-7,37,81,450])
del df['scores']
df['tsne-2d-one'] = tsne_results[:,0]
df['tsne-2d-two'] = tsne_results[:,1]
plt.figure(figsize=(16,10))
```

```
57  sns.scatterplot(
58      x="tsne-2d-one", y="tsne-2d-two",
59      hue="bins_scores",
60      palette=sns.color_palette("hls", len(set(df['bins_scores']))),
61      data=df,
62      legend="full",
63      alpha=0.3
64  )
```
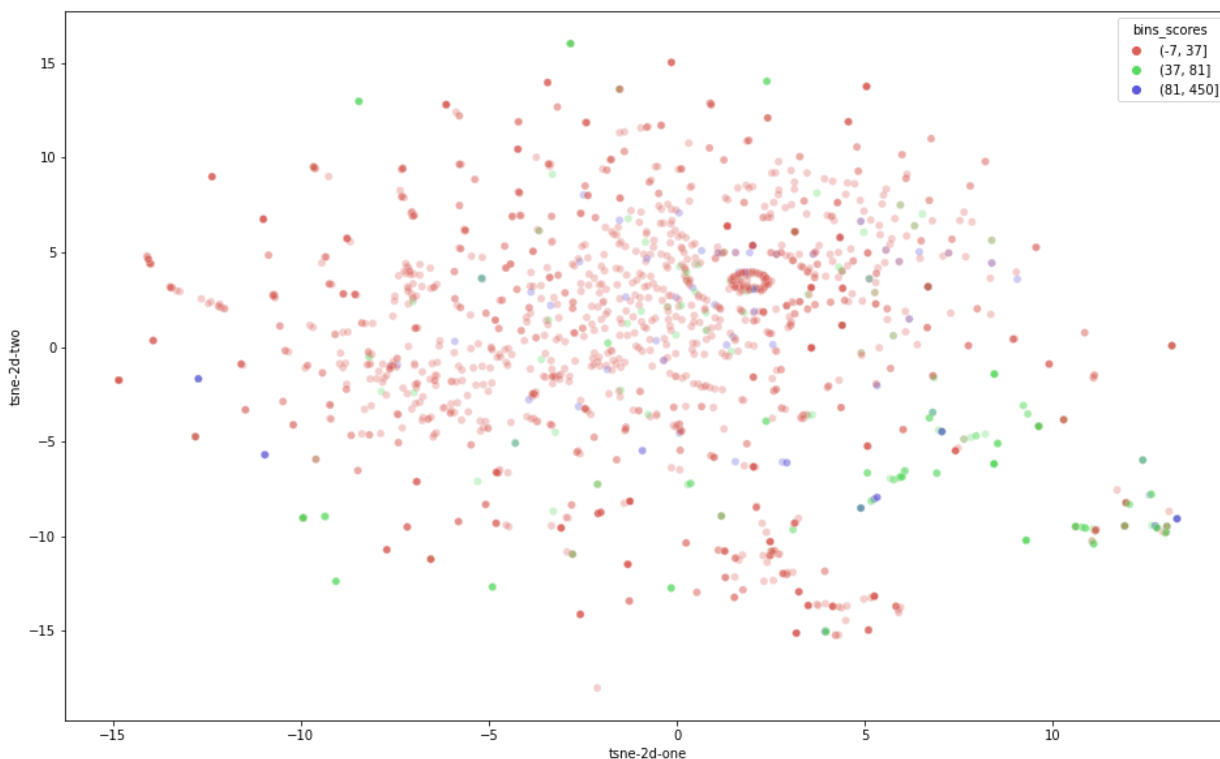
```
[t-SNE] Computing 121 nearest neighbors...
[t-SNE] Indexed 2000 samples in 0.172s...
[t-SNE] Computed neighbors for 2000 samples in 6.352s...
[t-SNE] Computed conditional probabilities for sample 1000 / 2000
[t-SNE] Computed conditional probabilities for sample 2000 / 2000
[t-SNE] Mean sigma: 0.000000
[t-SNE] KL divergence after 250 iterations with early exaggeration: 74.61
1977
[t-SNE] KL divergence after 300 iterations: 1.613455
t-SNE done! Time elapsed: 42.167349100112915 seconds
```

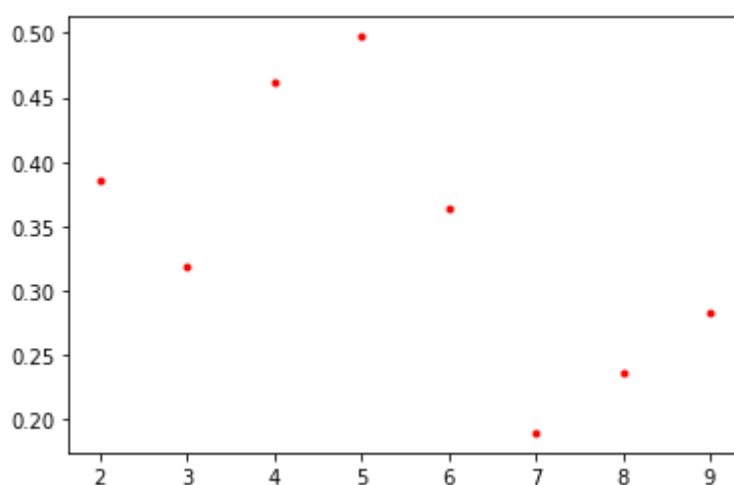Out[5]: <AxesSubplot:xlabel='tsne-2d-one', ylabel='tsne-2d-two'>



**3.Cluster the small molecules using your implementation of a partitioning clustering algorithm. Visualize this clustering by coloring clusters on the 2D visualization generated in question 2.**

In [9]:
```python
# Trying to find the most suitable cluster numbers
num_Cs = []
Silhouette_Coefficients = []

for num in range(2,10):
    print(num)
    kmeans = PartitionClustering(Ligands, num, 10)
    kmeans.implement()
    Sil_co = Silhouette_Coefficient(Ligands, kmeans.labels)
    num_Cs.append(num)
    Silhouette_Coefficients.append(Sil_co)
```
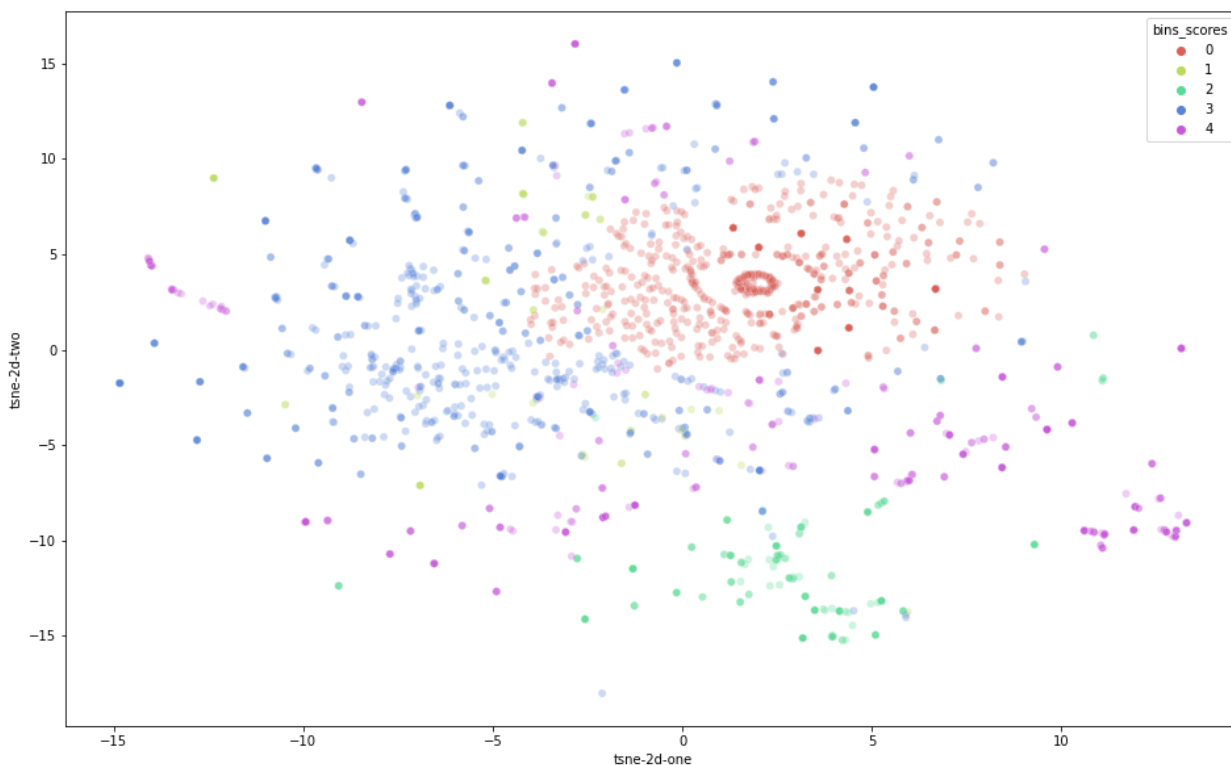
In [7]:
```python
plt.plot(num_Cs, Silhouette_Coefficients,'.r')
```

Out[7]: [<matplotlib.lines.Line2D at 0x7f55e91da0d0>]



In [10]:
```python
kmeans = PartitionClustering(Ligands, 5, 10)
kmeans.implement()
```

In [ ]:
```python
df['bins_scores'] = kmeans.labels
df['tsne-2d-one'] = tsne_results[:,0]
df['tsne-2d-two'] = tsne_results[:,1]
plt.figure(figsize=(16,10))
sns.scatterplot(
    x="tsne-2d-one", y="tsne-2d-two",
    hue="bins_scores",
    palette=sns.color_palette("hls", len(set(df['bins_scores']))),
    data=df,
    legend="full",
    alpha=0.3
)
```
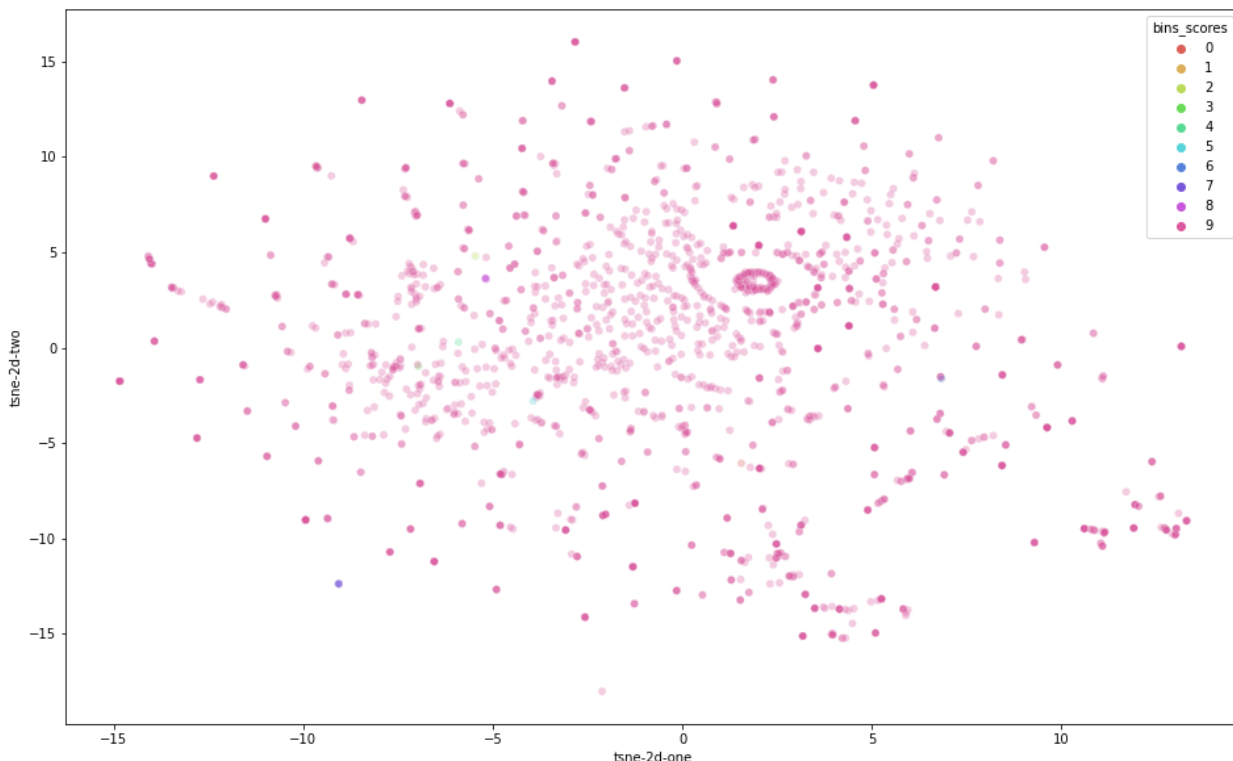
## 4.Explain your choice of partitioning clustering algorithm. Is it sensitive to initialization conditions? How do you select the number of clusters?

It is sensitive to the initialization conditions. It will be better after several times of iterations. For selecting the number of clusters. I plot the clustering quality scores with the cluster number and I choose 5 clusters based on this results.

## 5.Cluster the small molecules using your implementation of a hierarchical clustering algorithm. Visualize this clustering in the same way as question

```
1  Hiera_C = HierarchicalClustering(Ligands, 10)
2  Hiera_C.implement()
```

```
1  df['bins_scores'] = Hiera_C.labels
2  df['tsne-2d-one'] = tsne_results[:,0]
3  df['tsne-2d-two'] = tsne_results[:,1]
4  plt.figure(figsize=(16,10))
5  sns.scatterplot(
6      x="tsne-2d-one", y="tsne-2d-two",
7      hue="bins_scores",
8      palette=sns.color_palette("hls", len(set(df['bins_scores']))),
9      data=df,
10     legend="full",
11     alpha=0.3
12 )
```

## 6.Explain your choice of hierarchical clustering algorithm. Is it sensitive to initialization conditions? How do you select the number of clusters?

It is a determinate algorthims so it is not sensitive to the initialization conditions. For deciding the number of clusters, I do a plots of the number of clusters verse the clustering quality socres. And I realize that after a certain number, there will just be more single dots adding to the big cluster. So in the end I decide to do 10, in which case it would be ok to seprate the different ligands but also is not a huge cluster number.

## 7.Evaluate the quality of both clusterings using your implementation of a clustering quality metric. Explain your choice of quality metric. Which clustering performed 'best' according to your metric?

From the results (Silhouette_Coefficient) I got, the HierarchicalClustering work a little better than the k means, however, another thing to keep in mind is that the way that HierarchicalClustering is computed is very similar to the way that Silhouette_Coefficient is computed. Given that tsne and kmeans shows a more clear clustering result: although HierarchicalClustering has a higher Silhouette_Coefficient. I still believe that k means sever is a better model for this case.

Silhouette_Coefficient(Ligands, Hiera_C.labels)

Silhouette_Coefficient(Ligands, kmeans.labels)

## 8. Compare the two clusterings using your implementation of clustering similarity. How similar are the two clusterings using this function?

In [27]:
```
1  Rand_Index(kmeans.labels, Hiera_C.labels)
```

Out[27]:  0.2877113556778389

**Based on the rand index also the tsne results, I think this two cluster methods give very different results.**
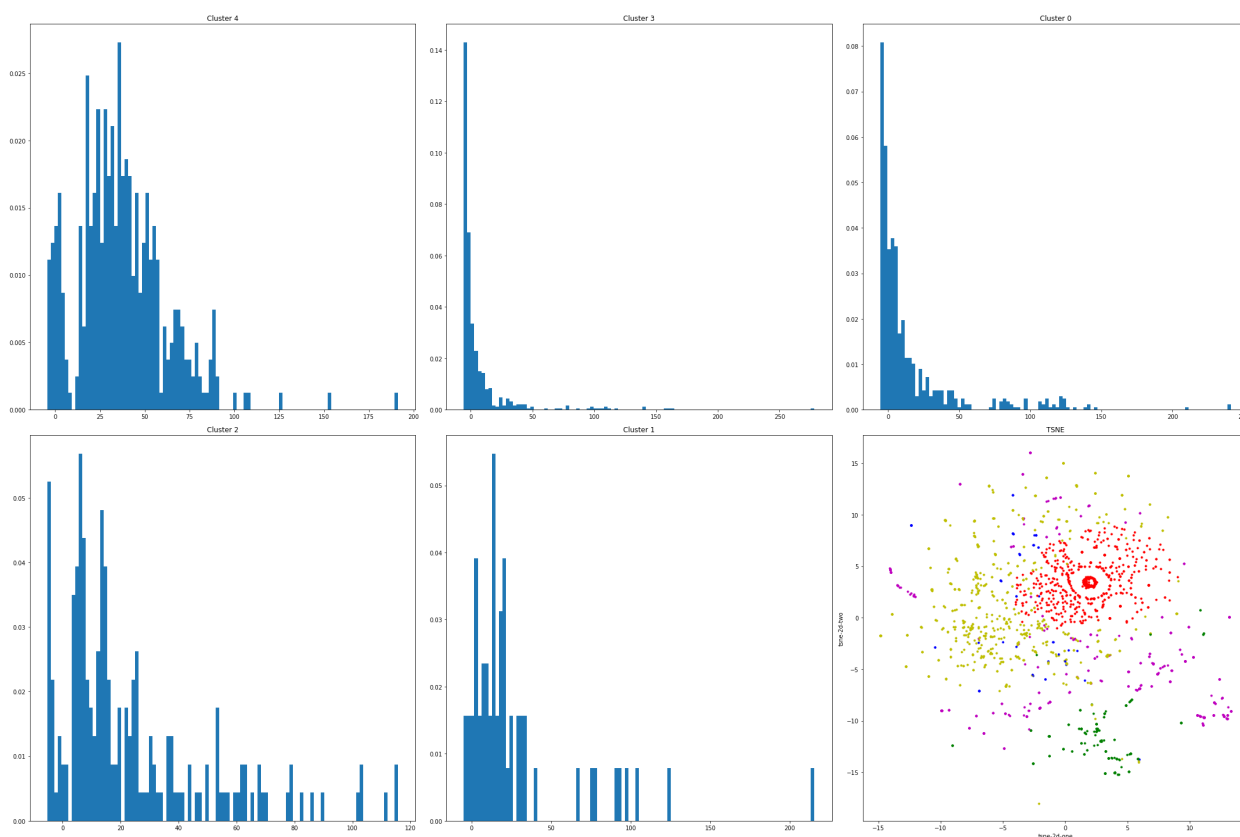
## 9. For the "best" clustering, as determined by your quality metric, visualize the distribution of Autodock Vina scores in each cluster. Do members of the same cluster have similar docking scores? Why or why not?

In [29]:
```python
1   # Kmeans
2   index = 2000
3   Cluster_scores = {}
4   for i in range(index):
5       c = kmeans.labels[i]
6       if c not in Cluster_scores:
7           Cluster_scores[c] = [kmeans.scores[i]]
8       else:
9           Cluster_scores[c].append(kmeans.scores[i])
10  num = 0
11  keys = []
12  for index, key in enumerate(Cluster_scores):
13      if len(Cluster_scores[key]) == 1:
14          pass
15      else:
16          keys.append(key)
17          num = num +1
```

```python
In [ ]:  1  plt.rcParams["figure.figsize"]=30,20
         2
         3  fig, axs = plt.subplots(nrows = 2, ncols = 3)
         4  p = 0
         5  for ax in axs.flatten():
         6      if p <5:
         7          values = Cluster_scores[keys[p]]
         8          ax.hist(values, 100, density = True)
         9          ax.set_title('Cluster '+str(keys[p]))
        10          p = p +1
        11      else:
        12          df['bins_scores'] = kmeans.labels
        13          cols = ['r','b','g','y','m']
        14          for i  in range(2000):
        15
        16              ax.plot(df["tsne-2d-one"][i], df["tsne-2d-two"][i], cols[km
        17          ax.set_title('TSNE')
        18          ax.set_xlabel('tsne-2d-one')
        19          ax.set_ylabel('tsne-2d-two')
        20
        21
        22  fig.tight_layout()
        23  plt.show()
```



**Based on the distribution, all of these clusters are having varied scores. However, cluster 2 which was colored in red 1) seems to form a good cluster 2) seems to be more enriched in higher scores. It is not surprised to me that the scores are varied becasue in the k means for calculating the distance, the onBites regrading different location are weighted uniformly, which might not be powerful enough to reflecting Autodock Vina scores.**

**10. Select the top scoring molecule from each cluster. This is your list of cluster heads. Visualize the top 5 by score in PyMOL and pick your favorite. Are they structurally diverse?**

```
In [46]:   1  num = 0
           2  for index, key in enumerate(Cluster_scores):
           3      if len(Cluster_scores[key]) == 1:
           4          continue
           5      else:
           6          num = num +1
           7
```

```
In [47]:   1
           2  top_scores = []
           3  for i in range(p):
           4      values = Cluster_scores[keys[i]]
           5      score = max(values)
           6      top_scores.append(score)
           7      top_scores.sort()
```

```
In [48]:   1  top_scores
```

Out[48]: [115.7, 191.3, 215.5, 241.0, 278.2]
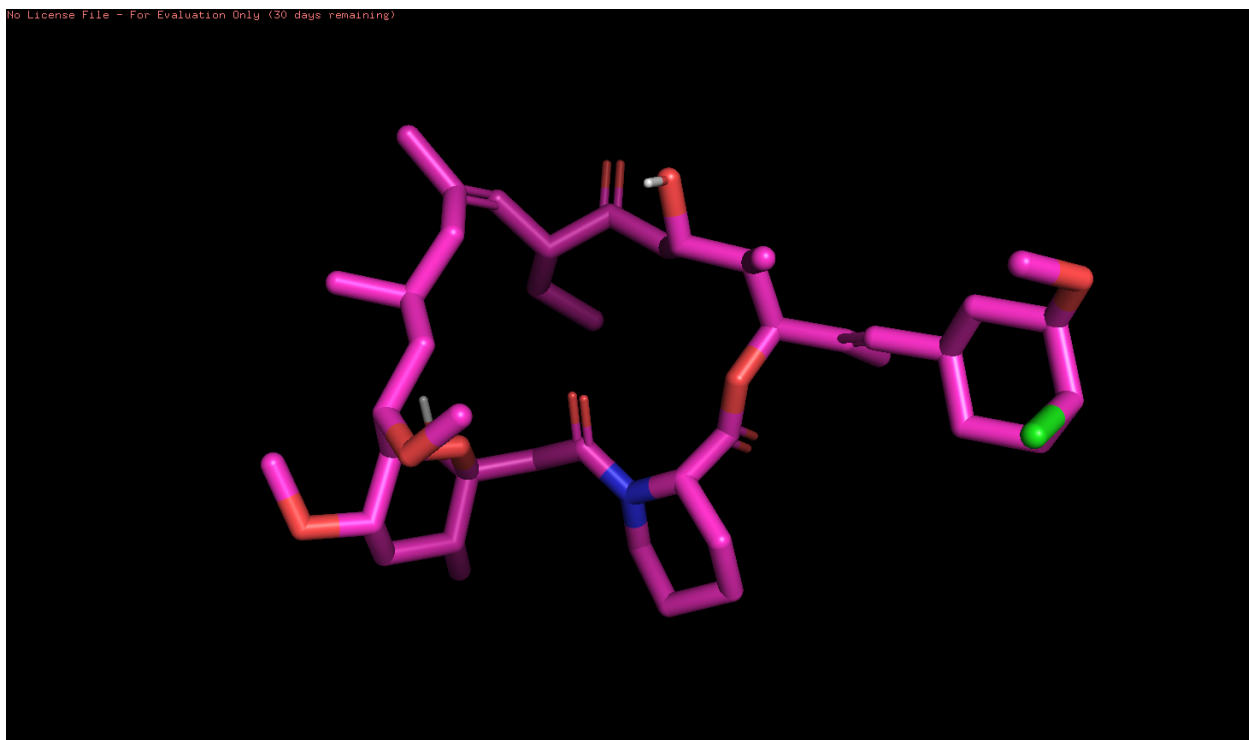
```
In [49]:   1  Leads = []
           2  for Ligand in Ligands:
           3      if Ligand.get_scores() in top_scores:
           4          Leads.append(Ligand.get_ID())
```
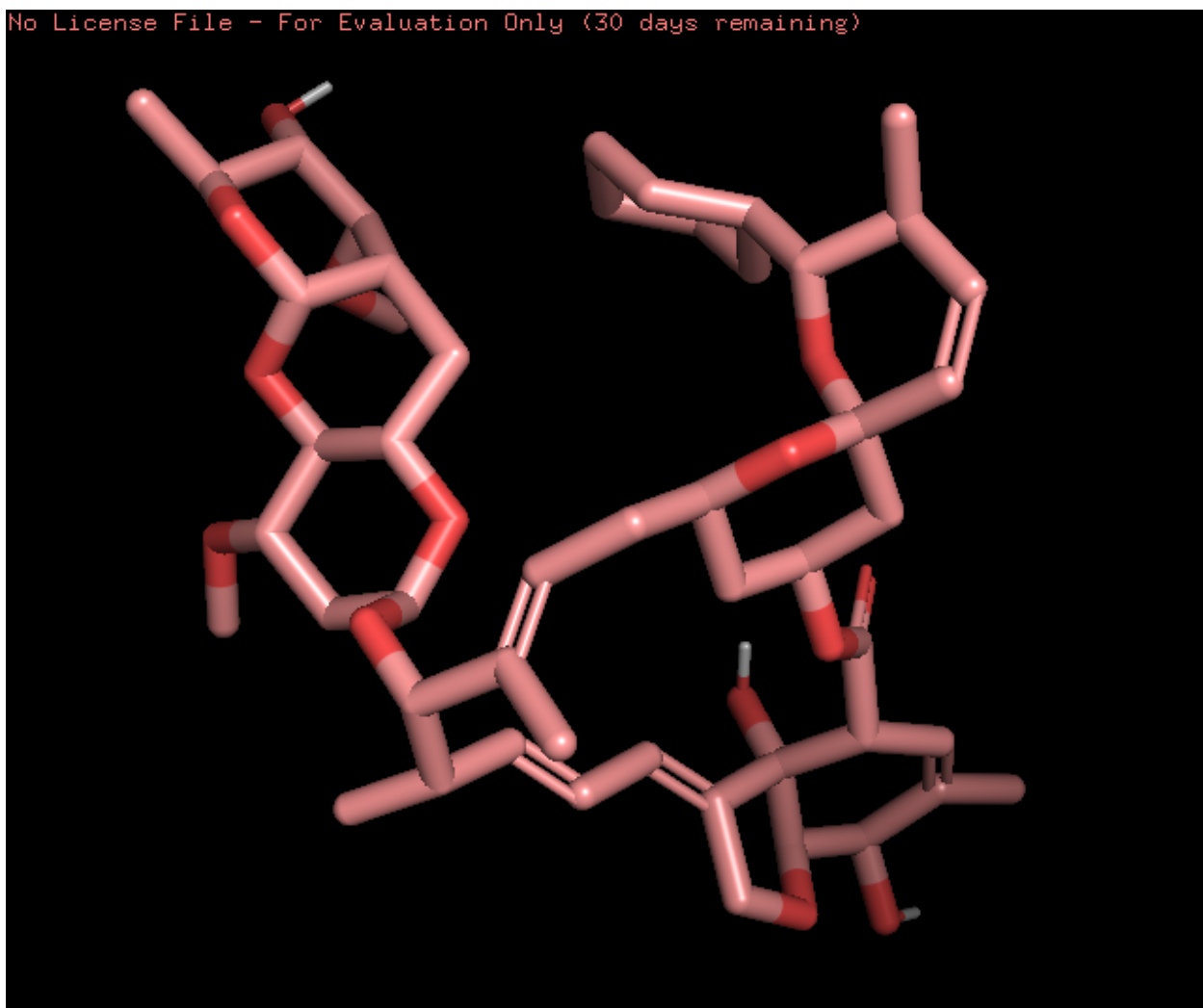
```
In [50]:   1  Leads
```

Out[50]: [8418, 8569, 7701, 8901, 6843]

**Answers: yes, they are all structurally diverse and have a complex structure.**

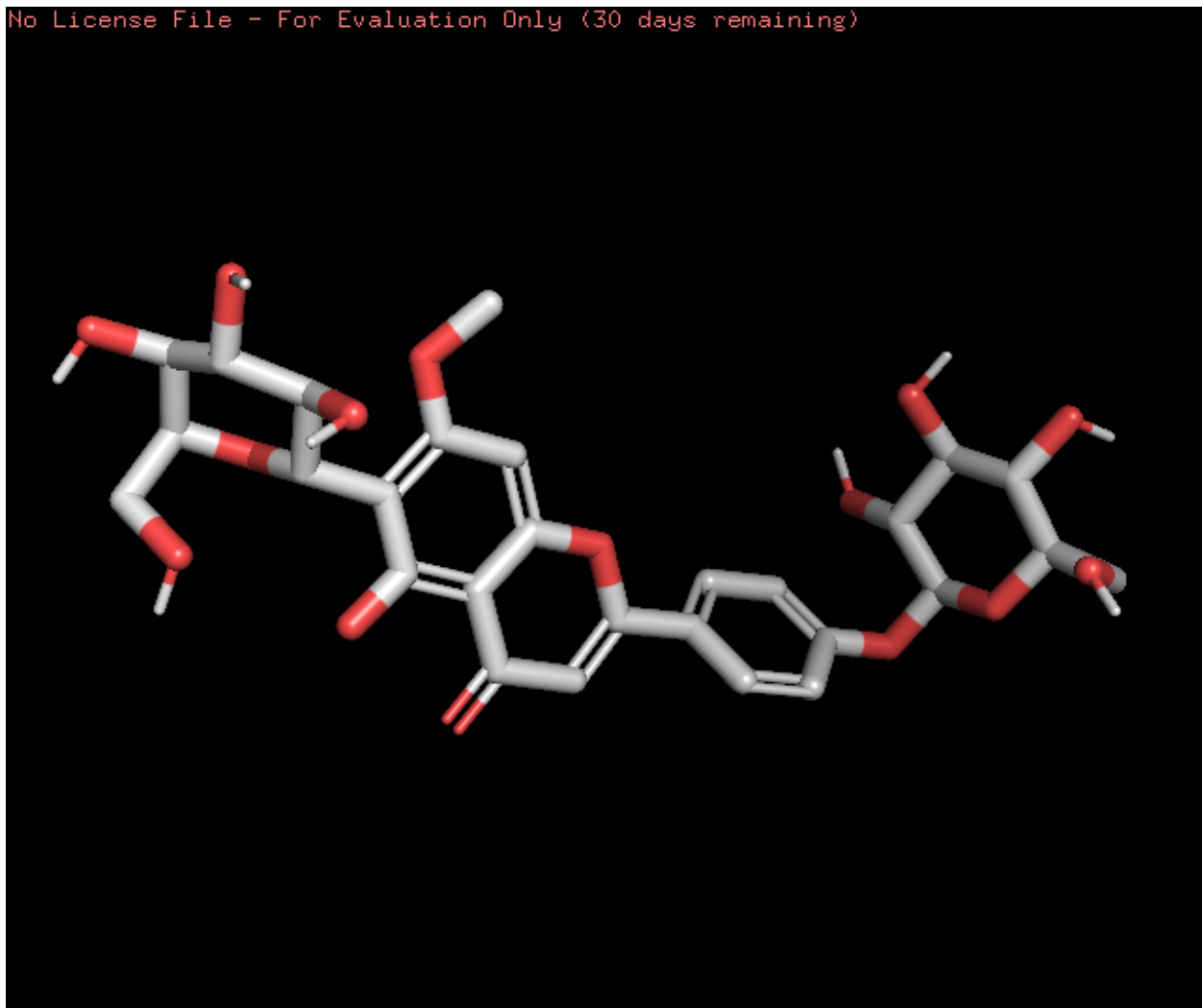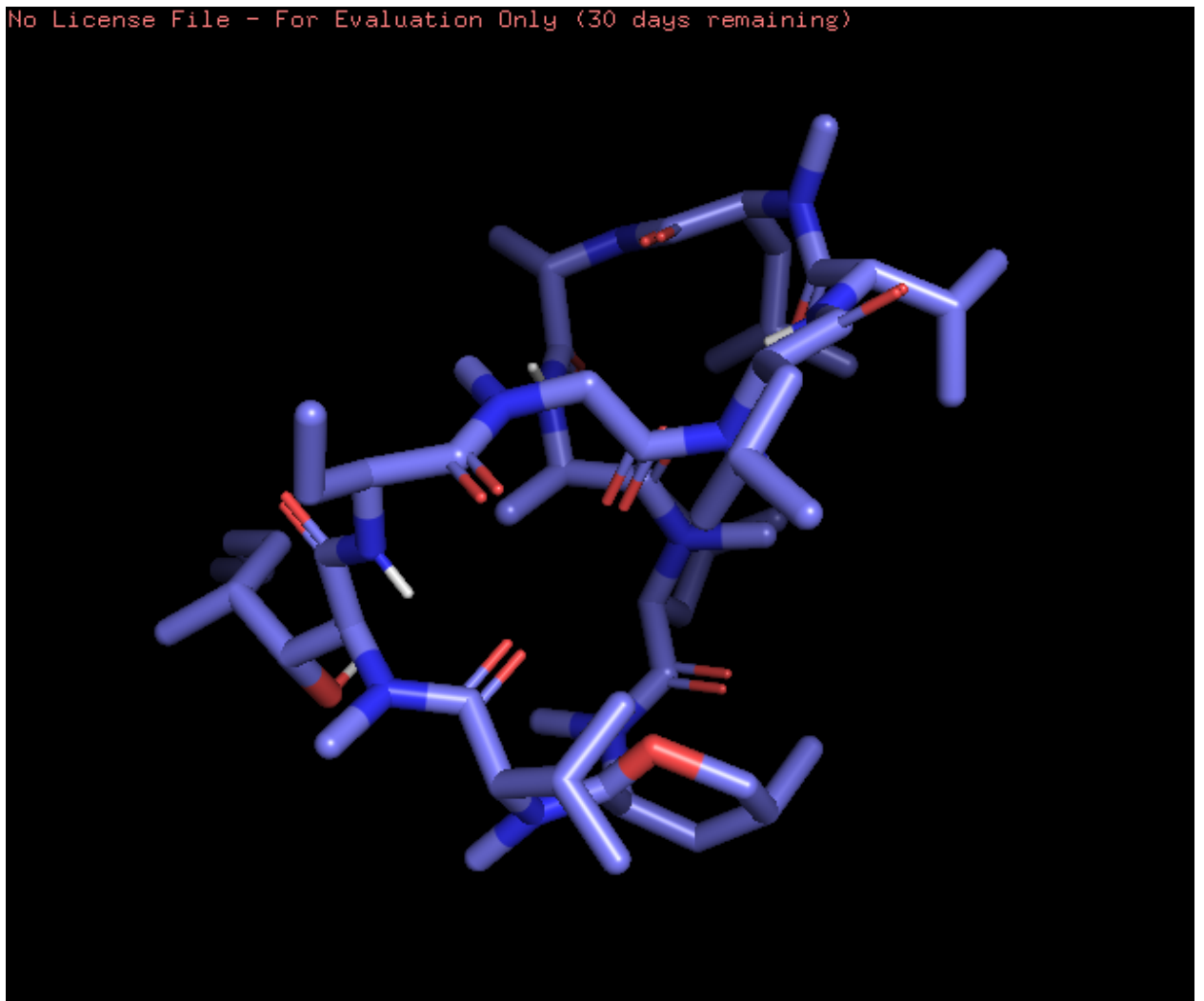**8418**

**8569**

**7701**



**8901**

No License File - For Evaluation Only (30 days remaining)

**6843**