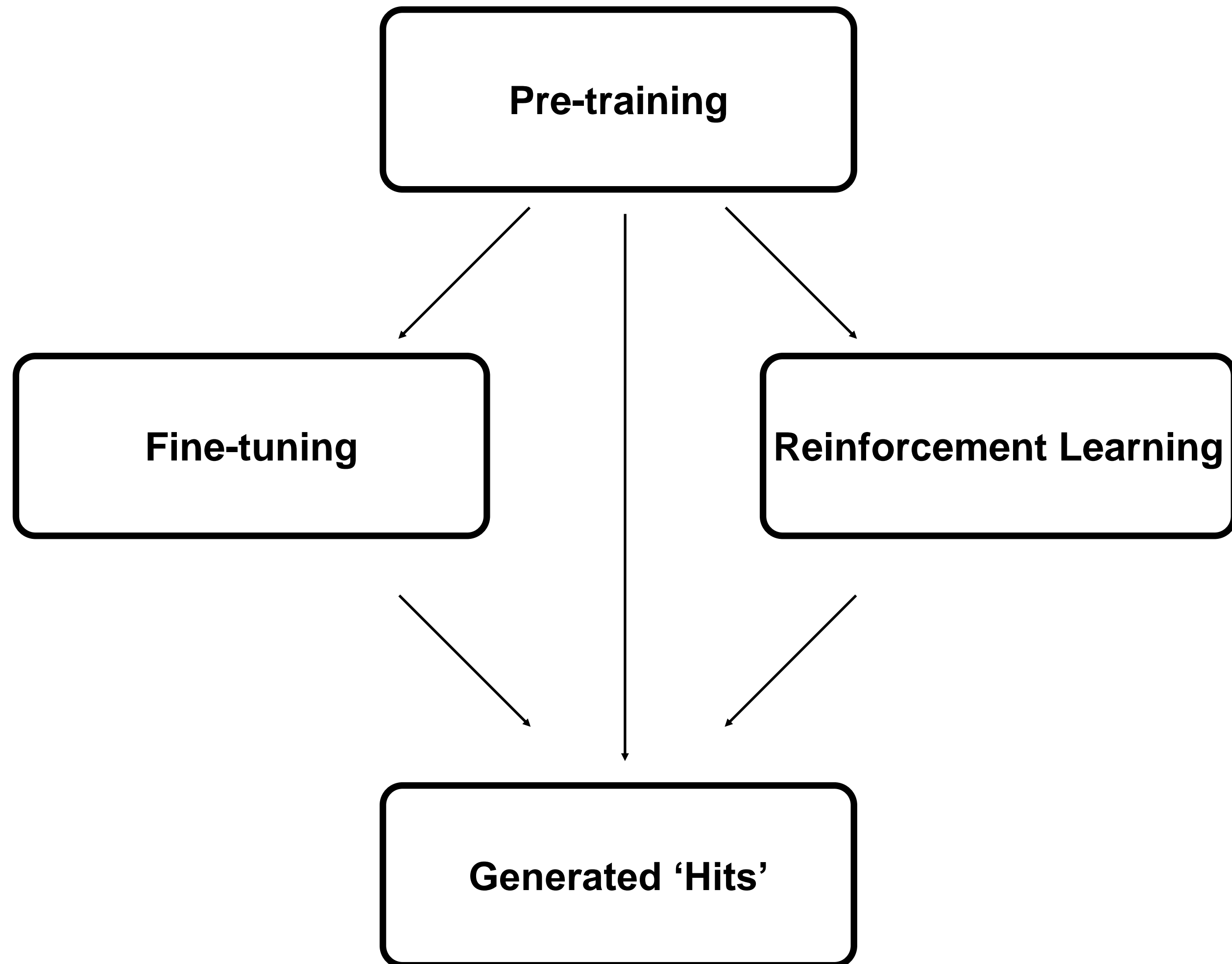


# Generative Modeling



Pre-training: Learn SMILES grammar and broad patterns, e.g. frequency of halogens, typical molecular weight, etc. Large, diverse dataset preferred

Fine-tuning: Bias toward patterns of interest by exposing model only to compounds satisfying certain properties, e.g. containing a desired substructure. Need a dataset of such compounds, possibly from virtual screening

Reinforcement Learning: Similar to fine-tuning except the dataset used for biasing is generated on-the-fly. Only need a scoring function specifying the properties of interest

# Sample Hit Libraries

## Pre-training datasets:

- ChEMBL (1.5M),
- Enamine diverse REAL drug-like (20M)

## Fine-tuning datasets:

- Virtually screened hits from Enamine 20M (100k)

### **Library A**

Pre-trained on ChEMBL

Hits from filtered,  
unbiased model using  
pharmacophore scoring  
function

~15k

### **Library B**

Virtual screening hits  
from Enamine 20M

~15k

### **Library C**

Pre-trained on ChEMBL

Use Library B to fine-  
tune generative model

Filter generated  
compounds for hits

~15k

### **Library D**

Pre-trained on Enamine  
20M

Hits from filtered,  
unbiased model using  
pharmacophore scoring  
function

~1.5k

### **Library E**

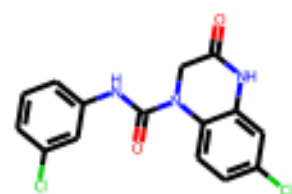
Pre-trained on ChEMBL

Reinforcement learning  
using pharmacophore  
scoring function

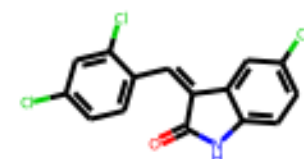
~15k

# Library A: ChEMBL pre-trained, filtered

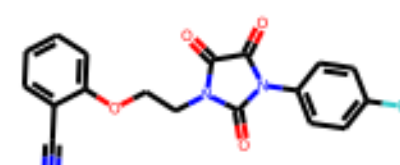
40 random samples



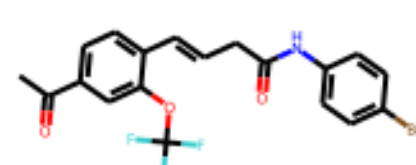
13707\_Library\_A



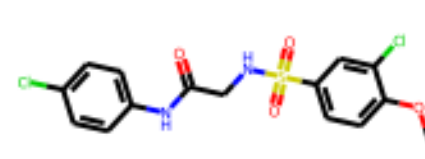
14207\_Library\_A



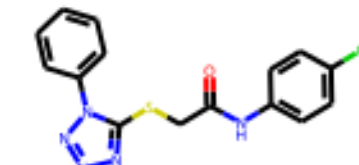
7929\_Library\_A



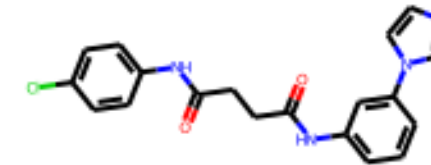
370\_Library\_A



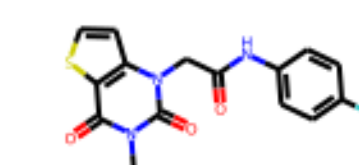
4838\_Library\_A



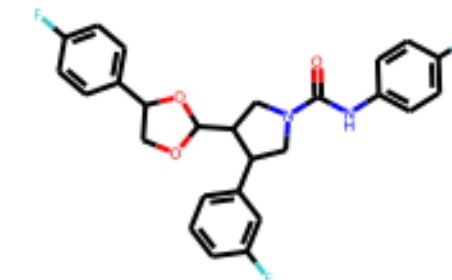
9620\_Library\_A



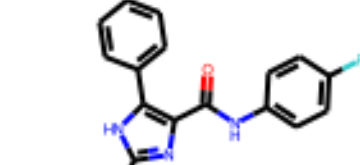
8691\_Library\_A



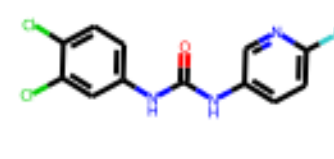
7519\_Library\_A



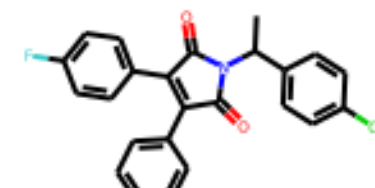
12009\_Library\_A



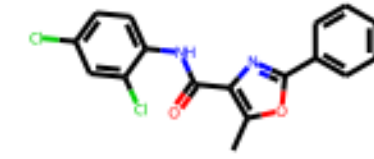
7206\_Library\_A



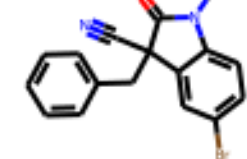
12583\_Library\_A



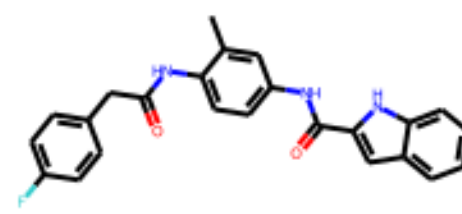
1113\_Library\_A



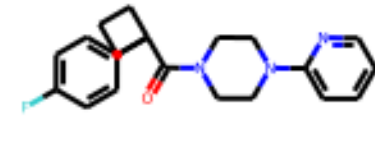
7424\_Library\_A



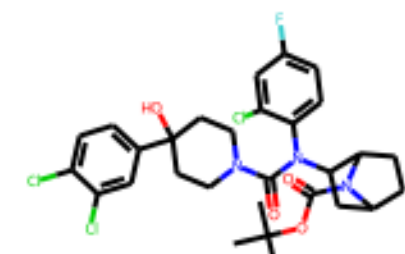
3011\_Library\_A



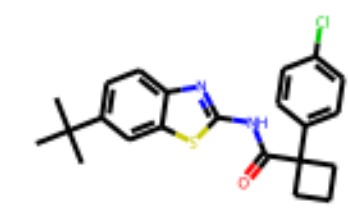
6085\_Library\_A



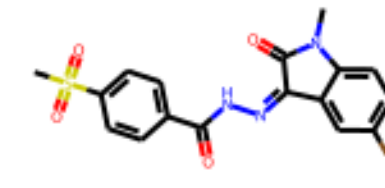
10391\_Library\_A



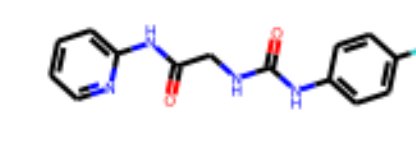
570\_Library\_A



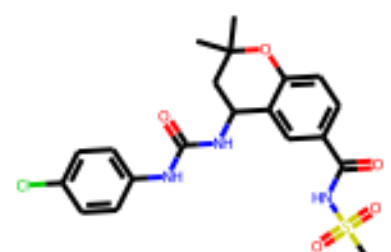
641\_Library\_A



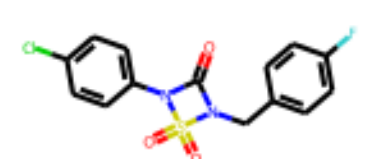
3004\_Library\_A



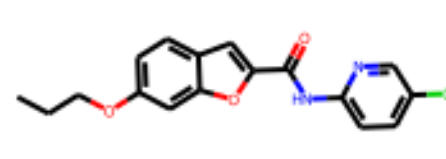
9012\_Library\_A



1191\_Library\_A



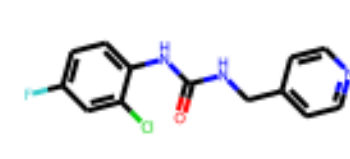
13837\_Library\_A



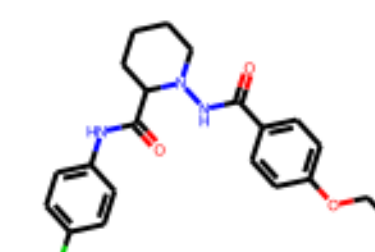
1824\_Library\_A



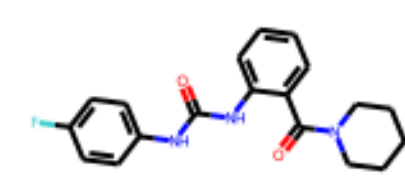
13783\_Library\_A



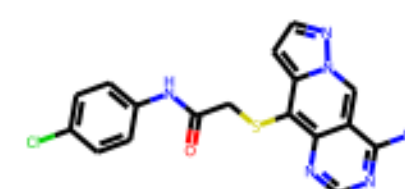
10672\_Library\_A



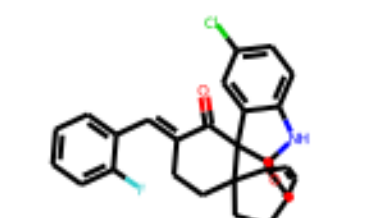
2227\_Library\_A



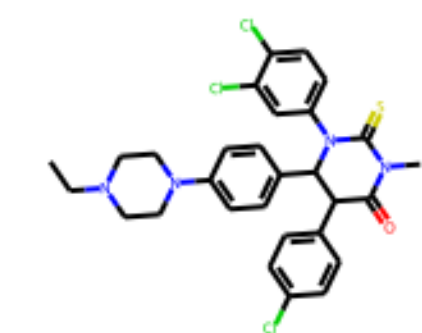
1524\_Library\_A



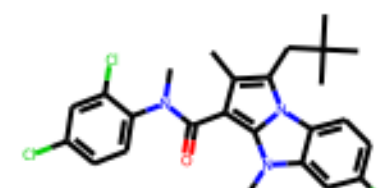
8372\_Library\_A



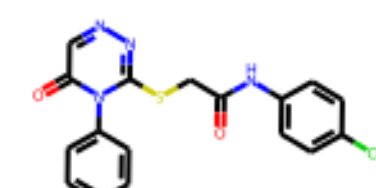
14165\_Library\_A



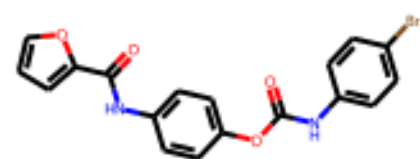
1993\_Library\_A



5521\_Library\_A



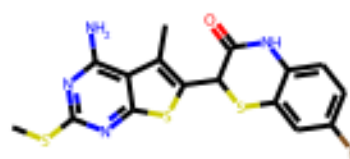
6221\_Library\_A



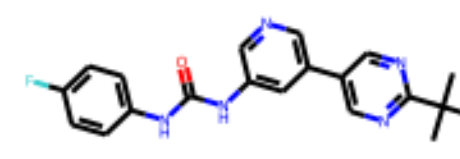
10941\_Library\_A



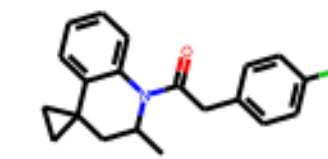
1354\_Library\_A



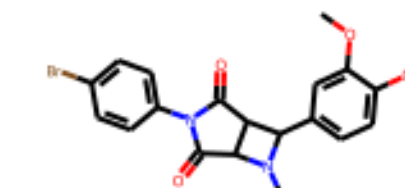
5459\_Library\_A



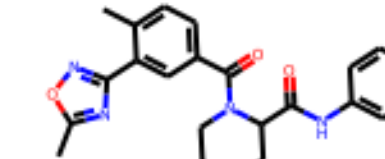
648\_Library\_A



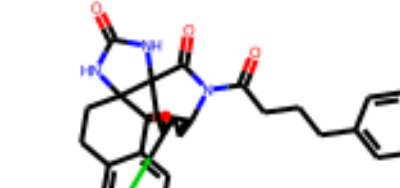
1496\_Library\_A



59\_Library\_A



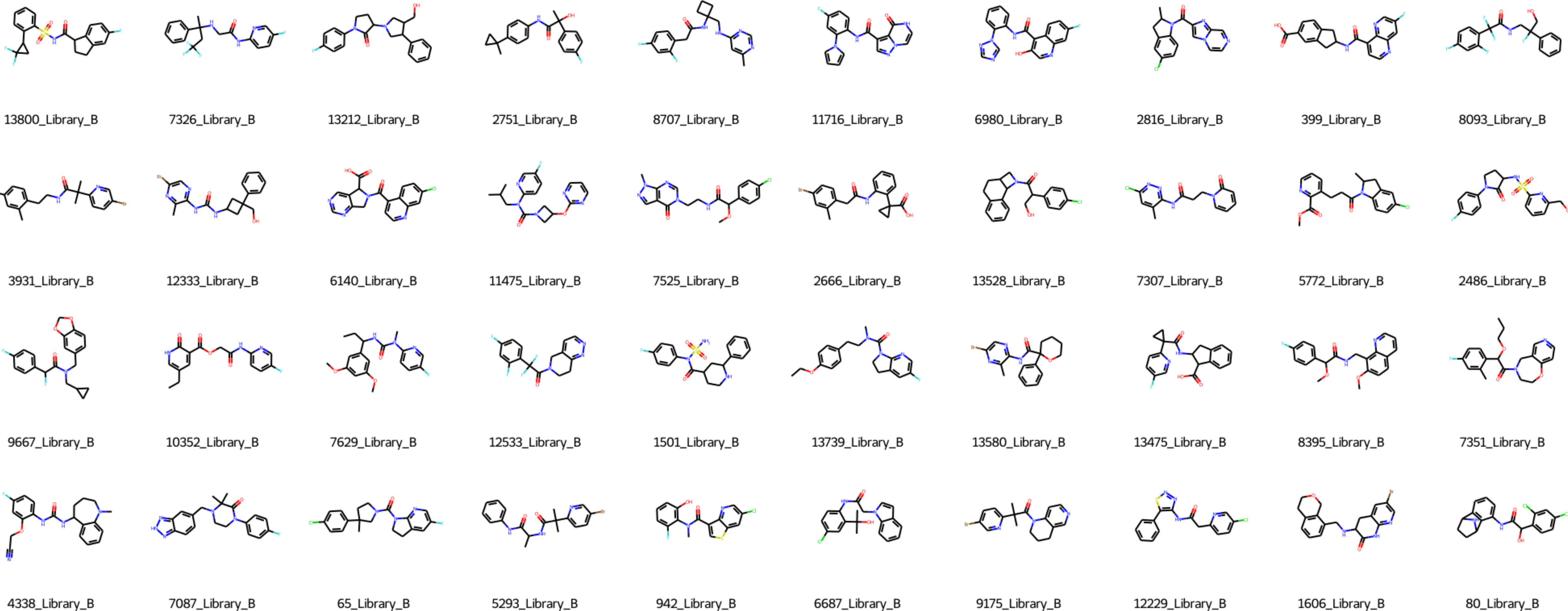
7176\_Library\_A



14011\_Library\_A

# Library B: Enamine filtered

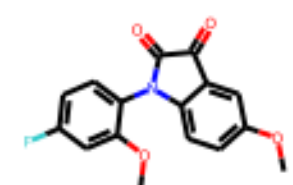
40 random samples





# Library C: ChEMBL pre-trained, library B fine-tuned

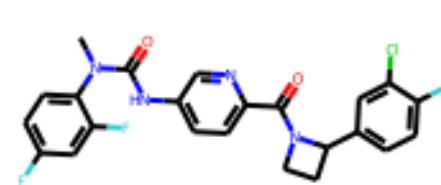
40 random samples



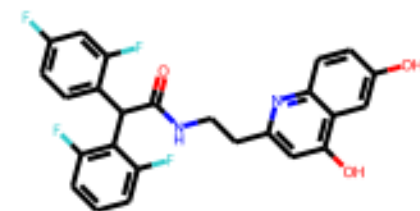
7213\_Library\_C



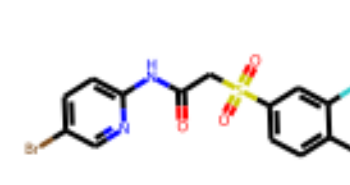
5226\_Library\_C



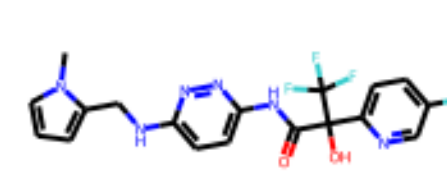
4540\_Library\_C



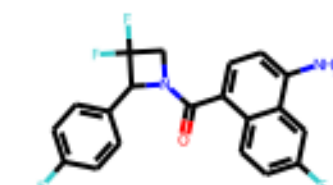
12238\_Library\_C



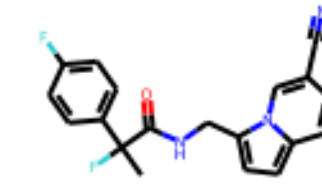
8924\_Library\_C



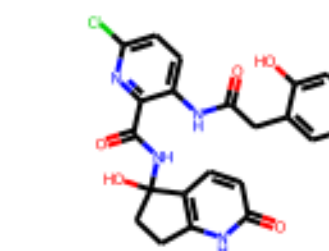
9878\_Library\_C



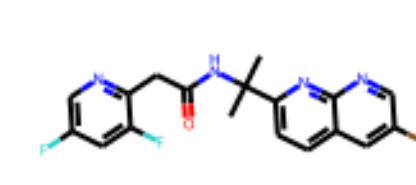
10449\_Library\_C



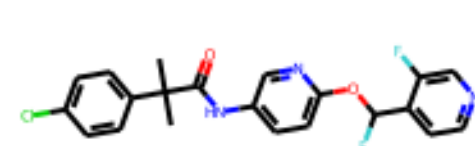
1377\_Library\_C



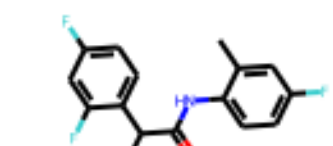
11801\_Library\_C



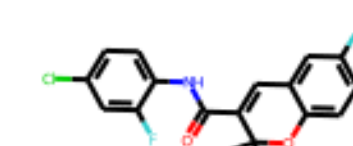
981\_Library\_C



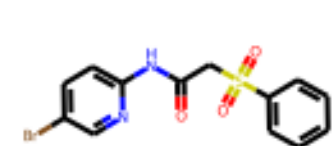
821\_Library\_C



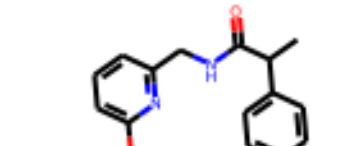
2337\_Library\_C



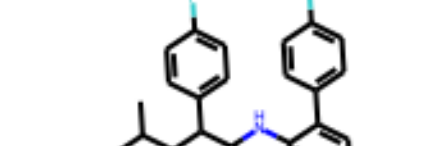
1812\_Library\_C



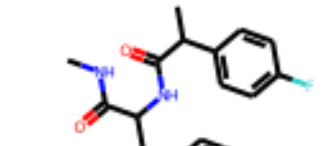
11284\_Library\_C



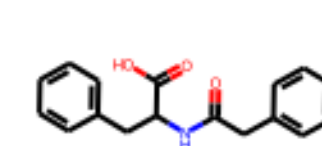
7390\_Library\_C



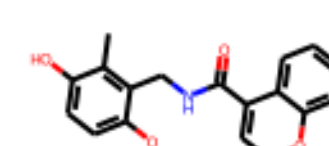
1094\_Library\_C



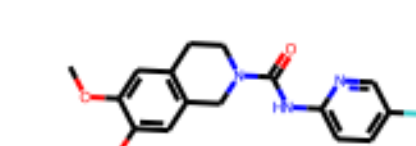
5342\_Library\_C



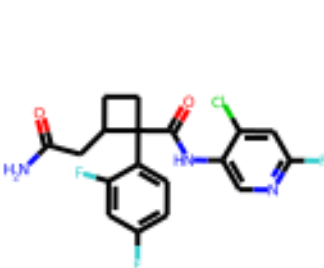
11600\_Library\_C



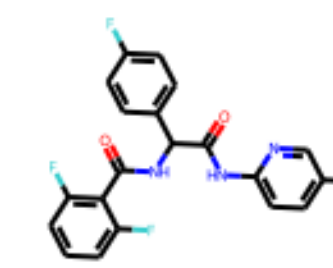
7121\_Library\_C



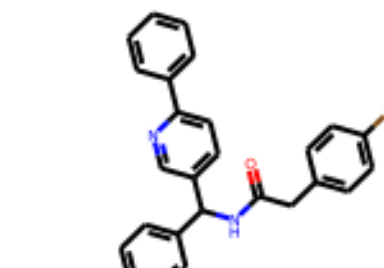
6472\_Library\_C



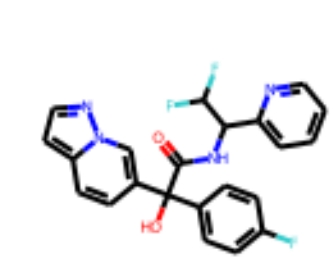
10209\_Library\_C



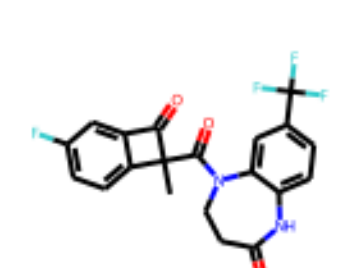
12108\_Library\_C



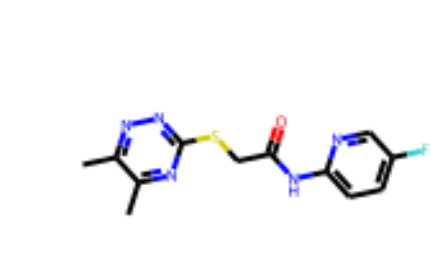
11429\_Library\_C



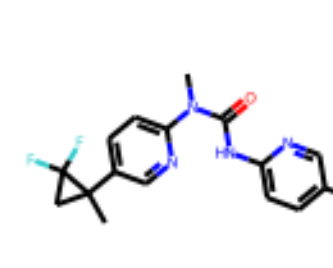
12133\_Library\_C



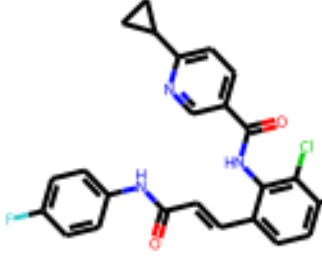
1723\_Library\_C



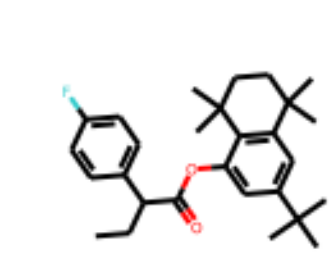
9713\_Library\_C



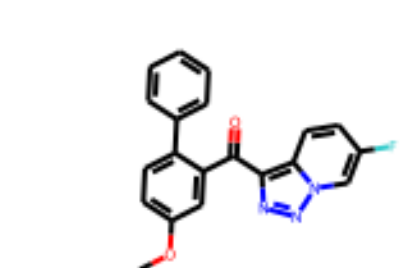
4552\_Library\_C



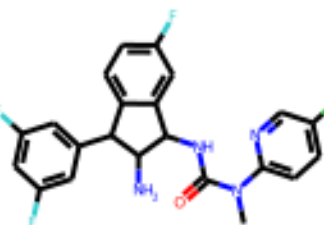
10683\_Library\_C



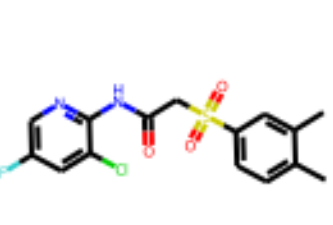
2388\_Library\_C



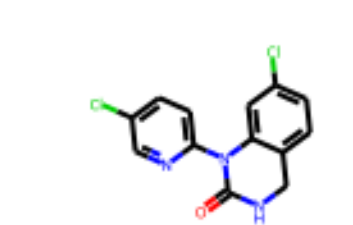
6519\_Library\_C



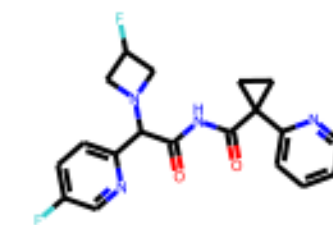
4530\_Library\_C



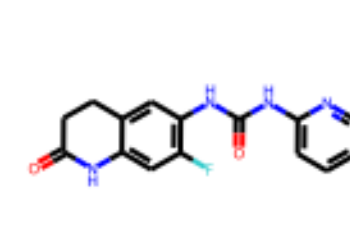
8926\_Library\_C



14699\_Library\_C



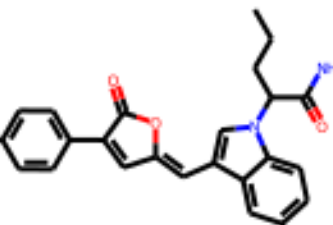
12044\_Library\_C



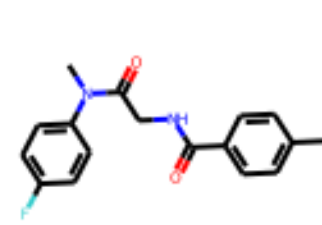
14500\_Library\_C



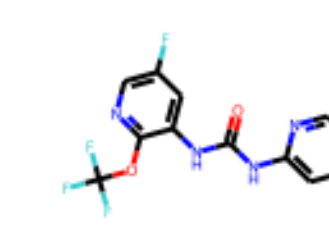
13111\_Library\_C



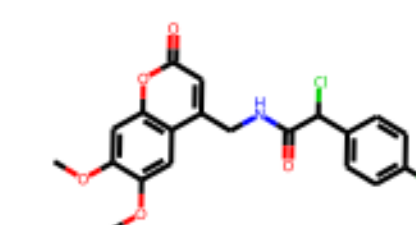
2644\_Library\_C



4372\_Library\_C



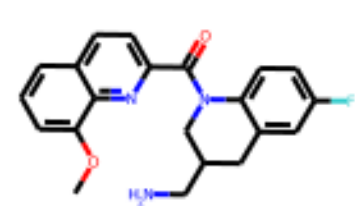
13479\_Library\_C



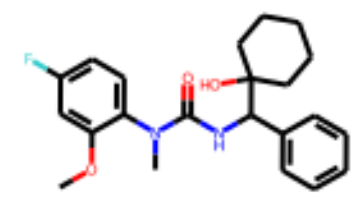
6499\_Library\_C

# Library D: Enamine pre-trained, filtered

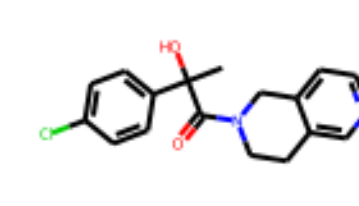
40 random samples



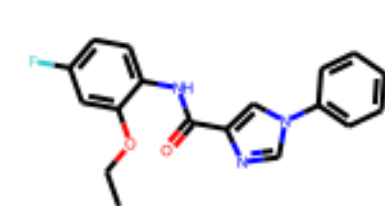
760\_Library\_D



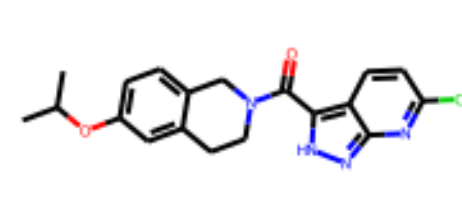
689\_Library\_D



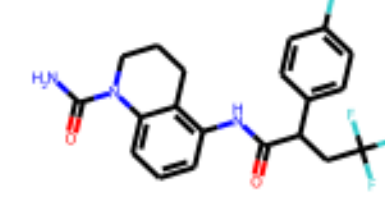
179\_Library\_D



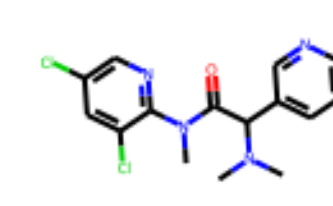
352\_Library\_D



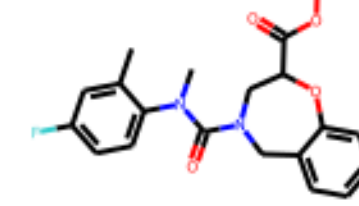
110\_Library\_D



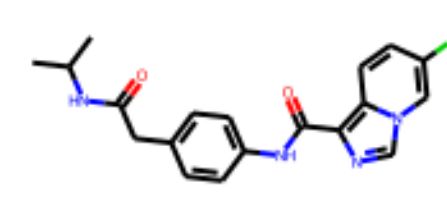
1171\_Library\_D



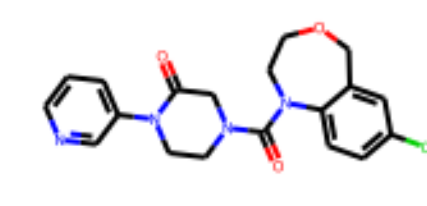
417\_Library\_D



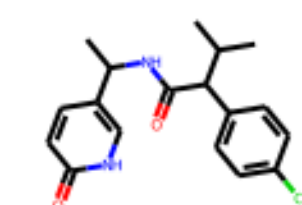
526\_Library\_D



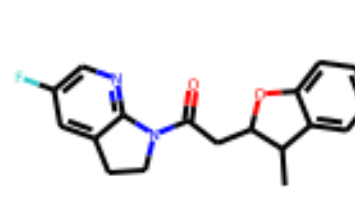
105\_Library\_D



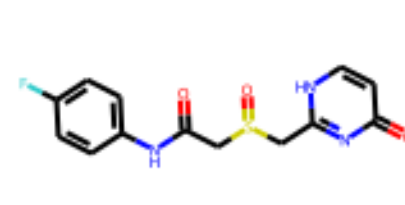
1492\_Library\_D



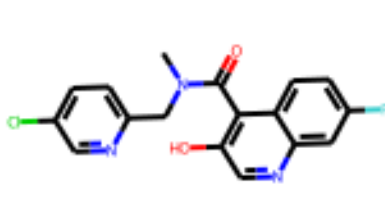
156\_Library\_D



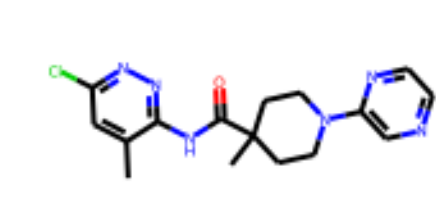
248\_Library\_D



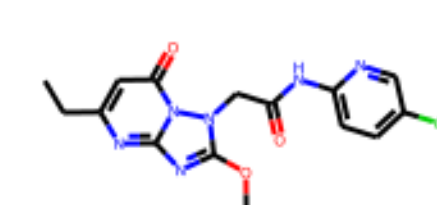
1245\_Library\_D



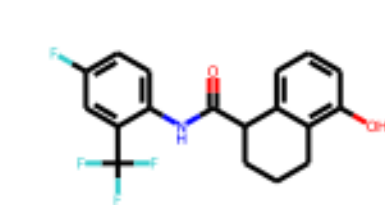
480\_Library\_D



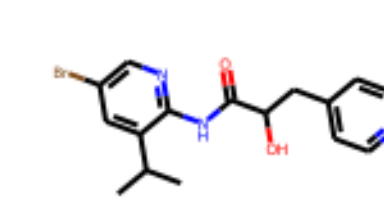
853\_Library\_D



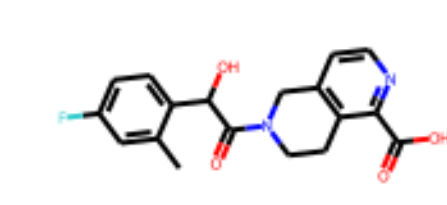
365\_Library\_D



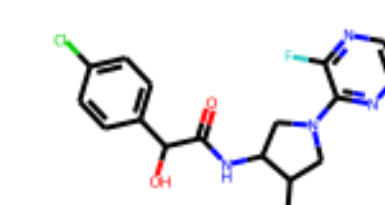
1376\_Library\_D



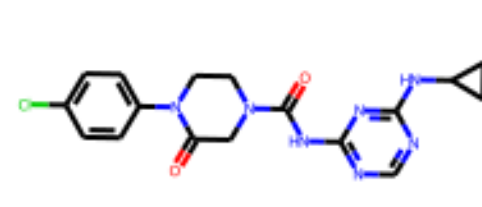
120\_Library\_D



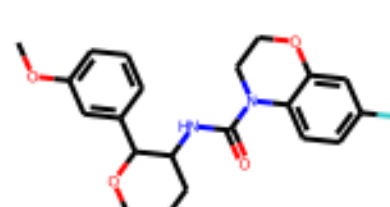
908\_Library\_D



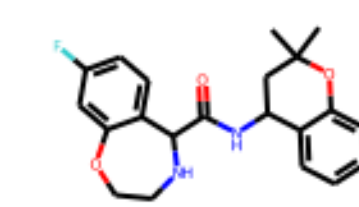
281\_Library\_D



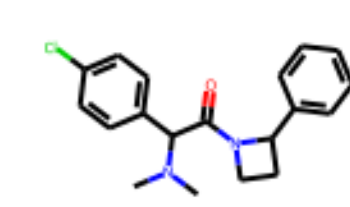
1436\_Library\_D



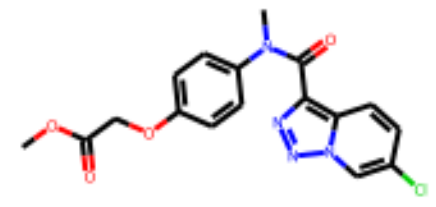
742\_Library\_D



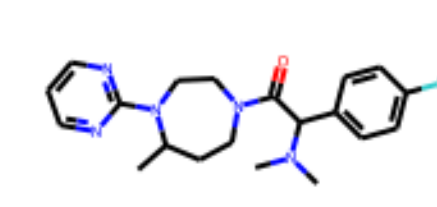
211\_Library\_D



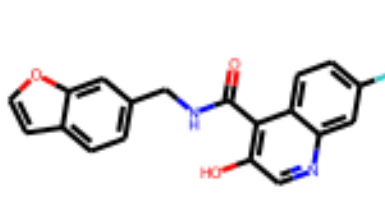
451\_Library\_D



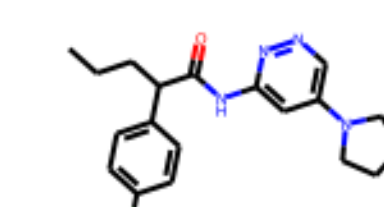
534\_Library\_D



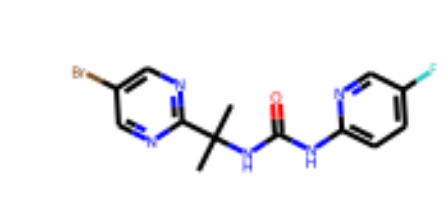
225\_Library\_D



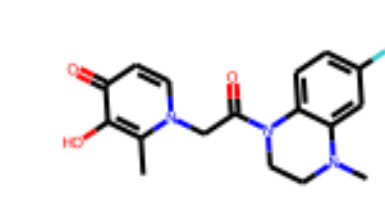
1337\_Library\_D



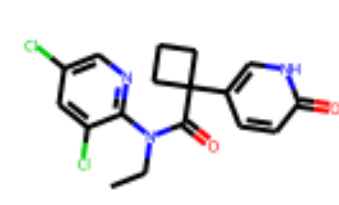
283\_Library\_D



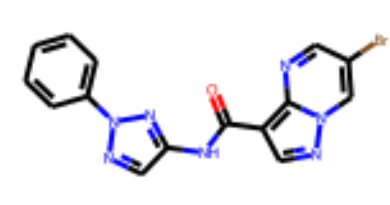
67\_Library\_D



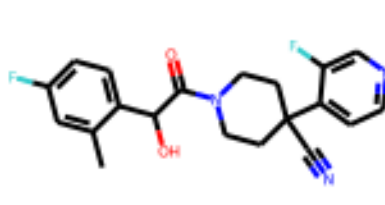
806\_Library\_D



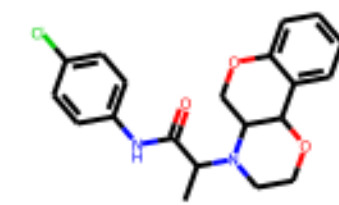
303\_Library\_D



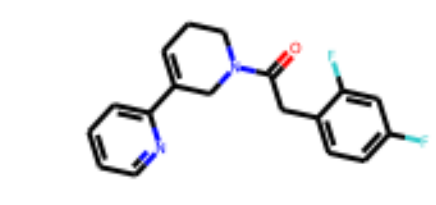
1422\_Library\_D



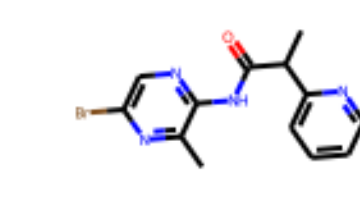
902\_Library\_D



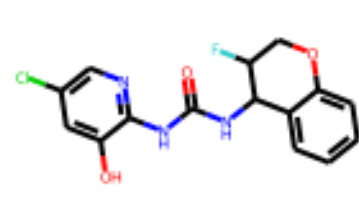
27\_Library\_D



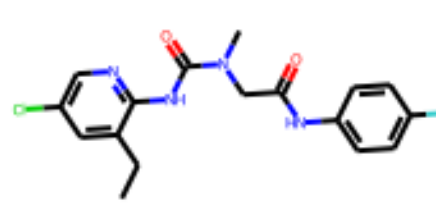
1269\_Library\_D



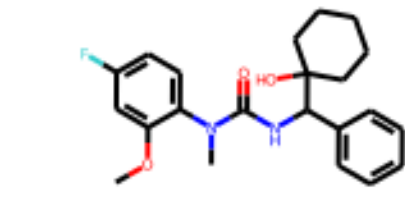
1061\_Library\_D



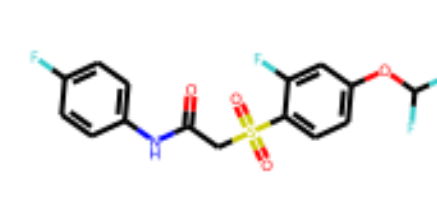
1429\_Library\_D



375\_Library\_D



689\_Library\_D

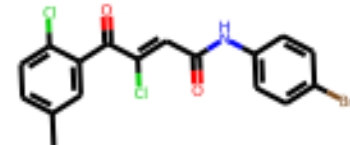


1244\_Library\_D

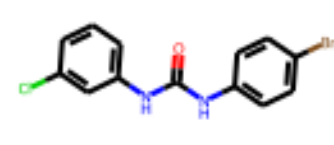


# Library E: ChEMBL pre-trained, reinforcement learning

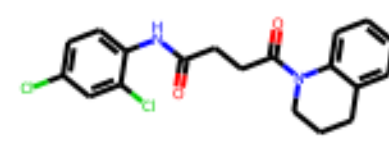
40 random samples



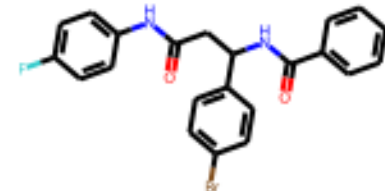
5088\_Library\_E



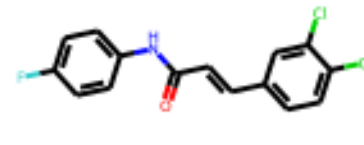
10085\_Library\_E



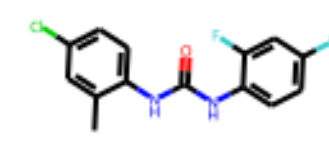
6813\_Library\_E



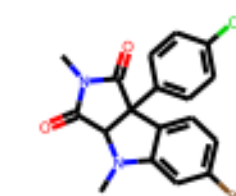
6725\_Library\_E



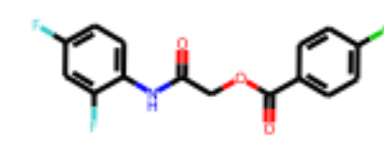
6569\_Library\_E



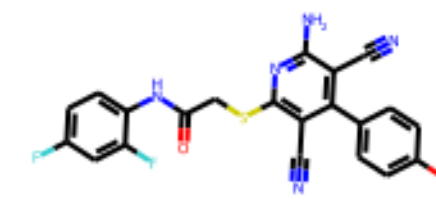
4768\_Library\_E



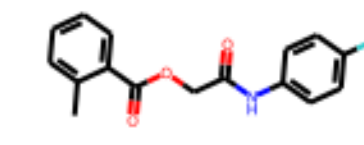
2710\_Library\_E



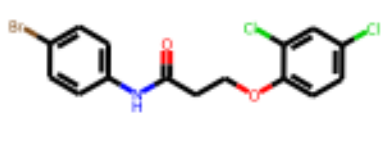
7711\_Library\_E



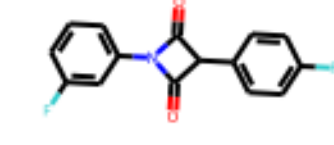
3348\_Library\_E



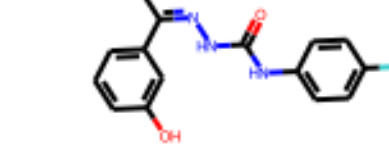
5482\_Library\_E



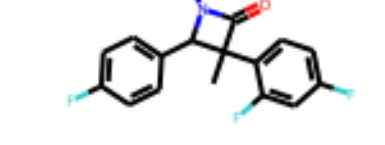
6999\_Library\_E



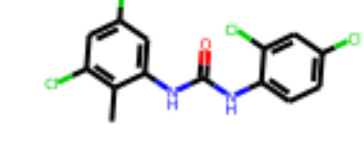
12702\_Library\_E



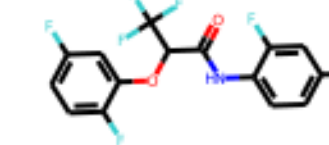
81\_Library\_E



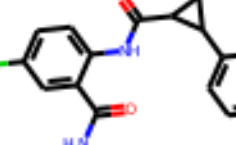
2677\_Library\_E



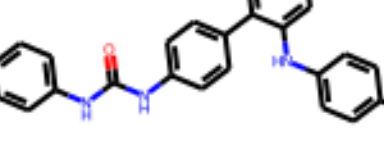
4407\_Library\_E



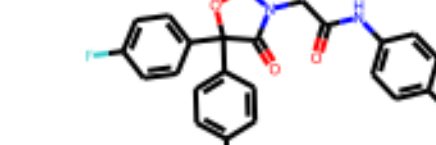
11648\_Library\_E



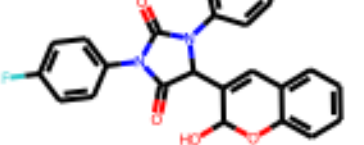
6248\_Library\_E



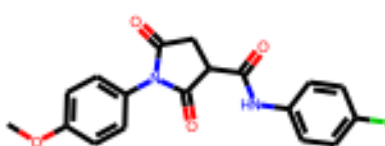
3864\_Library\_E



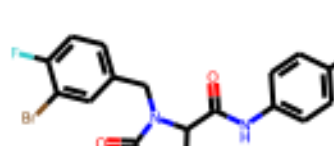
2058\_Library\_E



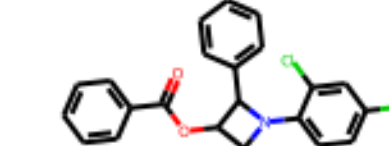
12590\_Library\_E



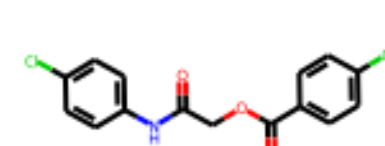
3769\_Library\_E



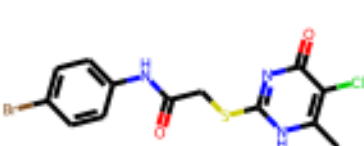
11300\_Library\_E



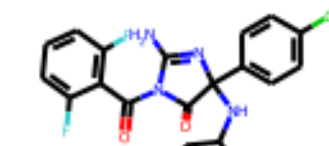
12172\_Library\_E



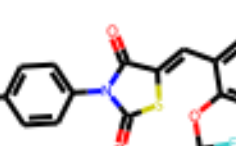
7708\_Library\_E



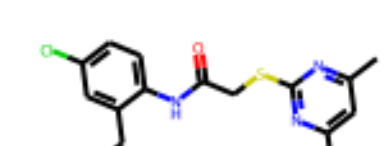
4277\_Library\_E



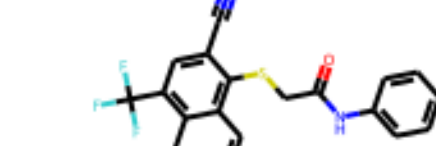
6306\_Library\_E



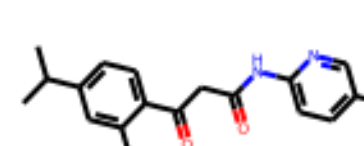
13879\_Library\_E



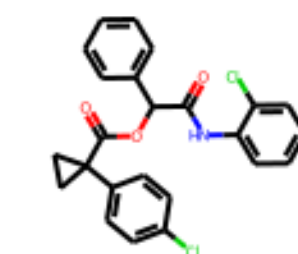
2923\_Library\_E



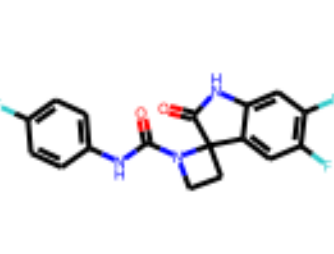
6133\_Library\_E



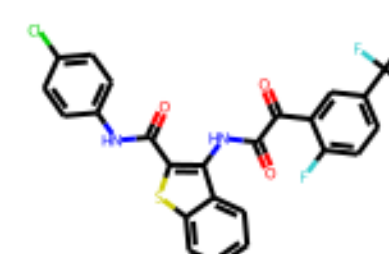
4558\_Library\_E



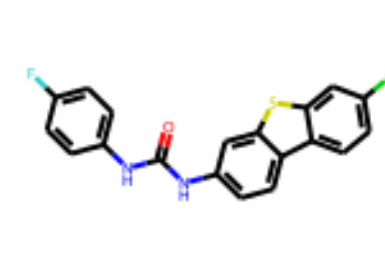
11997\_Library\_E



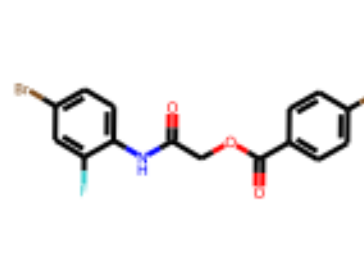
11358\_Library\_E



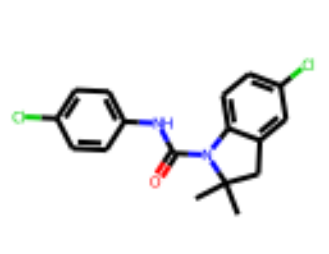
9918\_Library\_E



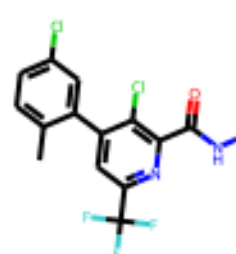
11428\_Library\_E



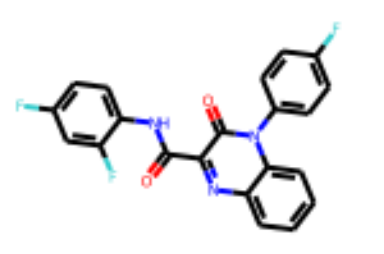
7684\_Library\_E



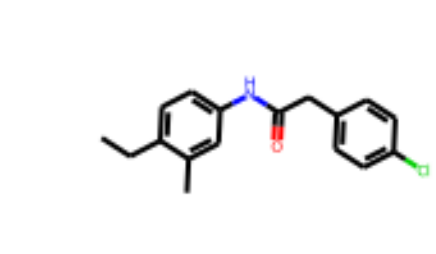
1139\_Library\_E



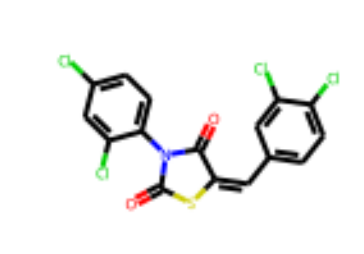
5101\_Library\_E



11839\_Library\_E



2073\_Library\_E

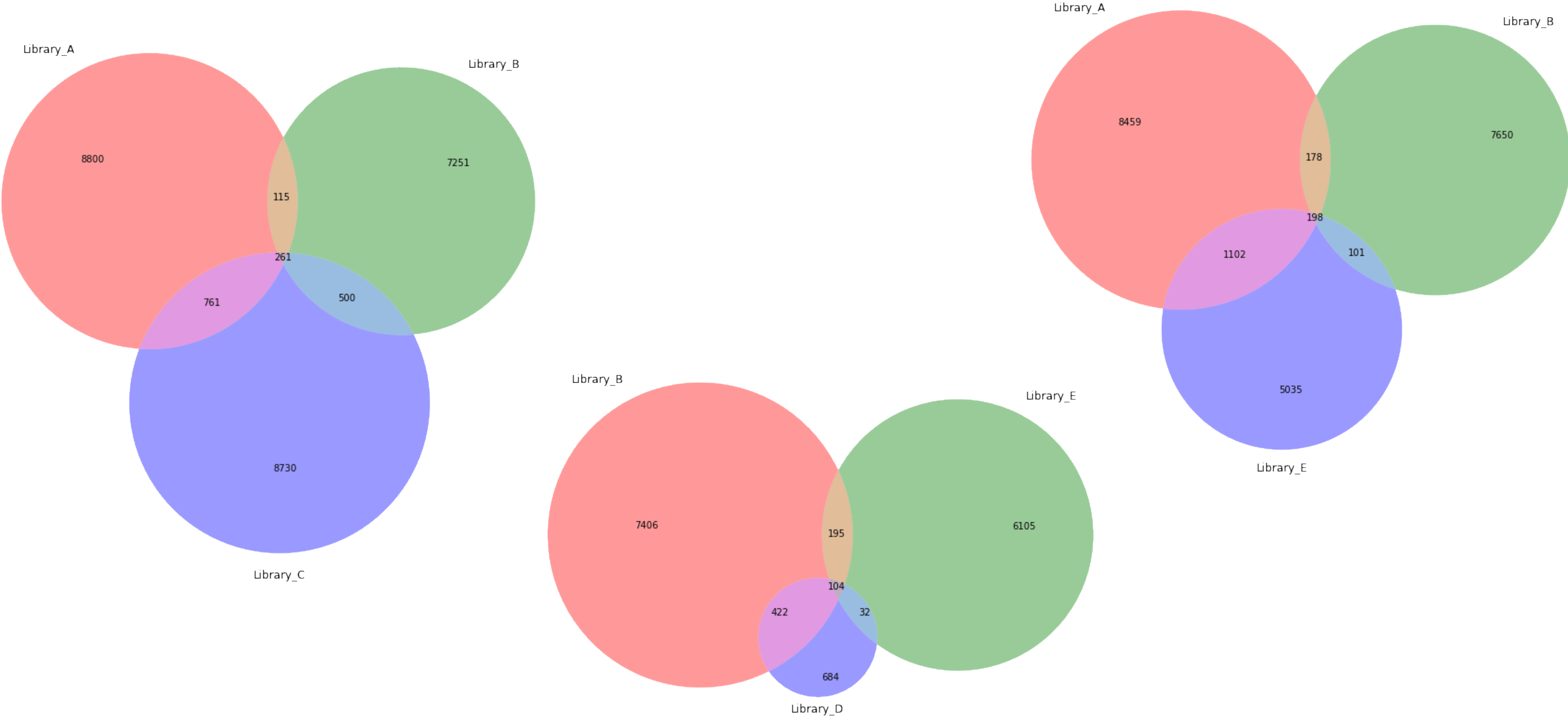


13883\_Library\_E

# Murcko Scaffold Overlap

Find all unique Murcko scaffolds present in each library, compare overlap between libraries

**Library A:** ChEMBL pre-trained, filtered  
**Library B:** Enamine pre-trained, filtered  
**Library C:** ChEMBL pre-trained, library B fine-tuned  
**Library D:** Enamine pre-trained, filtered  
**Library E:** ChEMBL pre-trained, reinforcement learning



Generally low overlap between scaffolds in each library