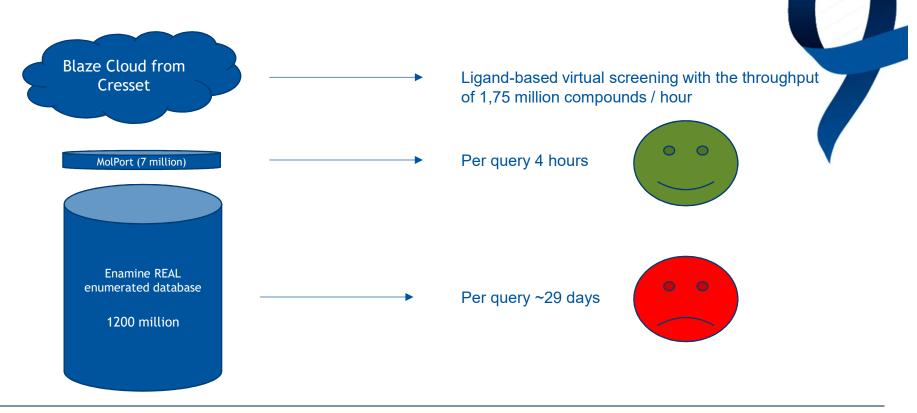
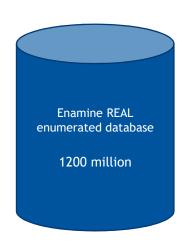


Background





Background





The current release of the *REAL* database comprises over 1.2 billion molecules which comply with "rule of 5" and Veber

criteria: MW≤500, SlogP≤5, HBA≤10, HBD≤5, rotatable bonds≤10, and TPSA≤140.

CNS-relevant chemical space for hits is much smaller!

Pick a CNS subset that can be screened in a day (max ~40 million or so)

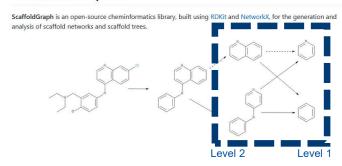


Workflow **GNU** awk: Using Enamine's pre-calculated values import multiprocessing from rdkit import Chem Additional from rdkit.Chem import Fragments 1200 million **Filtering** substructural filters (310 million compounds for line in gzip.open("c.smi.gz","rt"): (280 million l = line.strip().split() (SMILES) compounds) mol = Chem.MolFromSmiles(l[0]) compounds) if Fragments.fr_COO(mol)!=0: continue 34GB gzipped ScaffoldTree scaffold Diverse set of compounds from scaffold clusters clustering using (32 million compounds) ScaffoldGraph



ScaffoldTree clustering using ScaffoldGraph

O ScaffoldGraph O



```
import scaffoldgraph as sg
import networkx as nx

def extract_scaffolds(input_filename,output_filename):
    tree = sg.ScaffoldTree.from_smiles_file(input_filename,progress=False)
    w = gzip.open(output_filename,"wt")
    for molecule in tree.get_molecule_nodes():
        scaffolds = list(nx.bfs_tree(tree,molecule,reverse=True))
        molecule_smiles = tree.nodes[scaffolds[0]]["smiles"]
        scaffold="NoRings!"
        if len(scaffolds)==2:
            scaffold = scaffolds[-1]
        else:
            scaffold = scaffolds[-2]
        w.write(molecule_smiles + " " +molecule+" "+scaffold+"\n")
        w.close()
```

Scott and Chan. Bioinformatics 2020,. 36 3930 DOI: 10.1093/bioinformatics/btaa219

https://github.com/UCLCheminformatics/ScaffoldGraph

280 million compounds

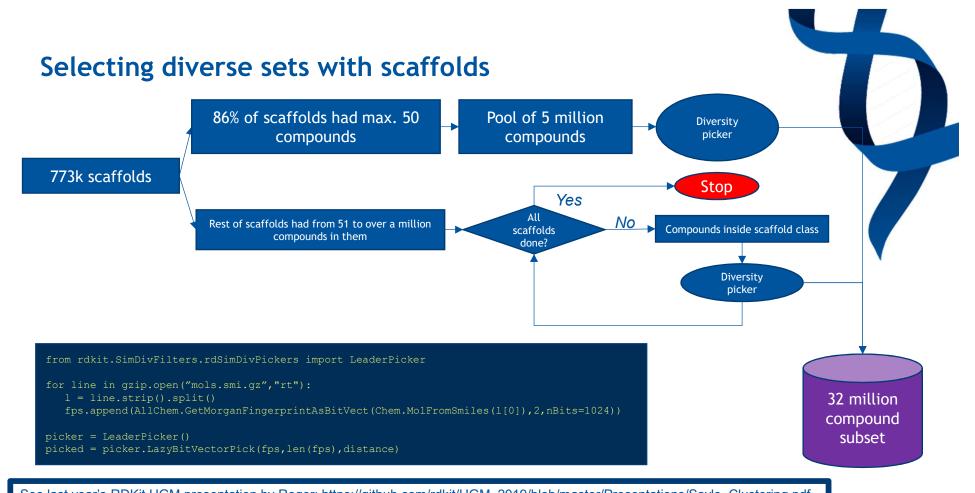
773k scaffolds

With 50 cores in a reasonable time

Painless to install via conda:

conda config --add channels conda-forge conda create --name scaffoldgraph python=3.7 conda activate scaffoldgraph conda install -c uclcheminformatics scaffoldgraph





See last year's RDKit UGM presentation by Roger: https://github.com/rdkit/UGM_2019/blob/master/Presentations/Sayle_Clustering.pdf

