

# High-throughput Screening of Promising Redox-Active Molecules with MolGAT

Mesfin Diro Chaka <sup>a,c,\*</sup>, Chernet Amente Geffe <sup>a</sup>, Alex Rodriguez <sup>d</sup>, Nicola Seriani <sup>d</sup>, Qin Wu <sup>e</sup>, Yedilfana Setarge Mekonnen <sup>b,\*</sup>

<sup>a</sup>Department of Physics, College of Natural and Computational Sciences, Addis Ababa University, P. O. Box 1176, Addis Ababa, Ethiopia,

<sup>b</sup>Center for Environmental Science, College of Natural and Computational Sciences, Addis Ababa University, P. O. Box 1176, Addis Ababa, Ethiopia,


<sup>c</sup>Computational Data Science, College of Natural and Computational Sciences, Addis Ababa University, P. O. Box 1176, Addis Ababa, Ethiopia,

<sup>d</sup>The Abdus Salam International Centre for Theoretical Physics (ICTP), Condensed Matter and Statistical Physics Section, 34100, Trieste, Italy,

<sup>e</sup>Brookhaven National Laboratory, Center for Functional Nanomaterials, Upton, NY, US,

---

\*Corresponding author

Email address: yedilfana.setarge@aau.edu.et (Yedilfana Setarge Mekonnen )

June 3, 2023

## 1. RedDB Dataset

4

Table S 1: Sample smiles and reaction energy in eV from RedDB dataset used for training

	Smiles	Reaction energy(eV)
0	<chem>O=C1CC(=O)C=C1C(=O)O</chem>	1.15667
1	<chem>O=C1CC(=O)C=C1</chem>	-0.496880
2	<chem>O=C1CC(=O)C(C(=O)O)C(=O)O</chem>	-0.710890
3	<chem>O=C1CC(=O)C=C1F</chem>	-0.446811
4	<chem>O=C1CC(=O)C(F)=C1F</chem>	-0.441369
5	<chem>O=C1CC(=O)C(N)=C1N</chem>	0.784777
6	<chem>O=C1CC(=O)C=C1N</chem>	0.344224
7	<chem>O=C1CC(=O)C=C1O</chem>	-0.073199
8	<chem>O=C1CC(=O)C(O)=C1O</chem>	0.201364
9	<chem>O=C1CC(=O)C(=C1)S(=O)(=O)O</chem>	-0.42741
10	<chem>O=C1CC(=O)C(S(=O)(=O)C(=O)O)C(=O)O</chem>	-0.869435
11	<chem>O=C1NC(=O)C(=C1)C(=O)O</chem>	-0.886275
12	<chem>O=C1NC(=O)C(C(=O)O)C(=O)O</chem>	-0.681160
13	<chem>O=C1NC(=O)C=C1</chem>	-0.198915
14	<chem>O=C1NC(=O)C(F)=C1F</chem>	-0.032109
15	<chem>O=C1NC(=O)C=C1F</chem>	-0.082451
16	<chem>O=C1NC(=O)C(=C1)N</chem>	0.576337
17	<chem>O=C1NC(=O)C(N)=C1N</chem>	0.655795
18	<chem>O=C1NC(=O)C(O)=C1O</chem>	0.465315
19	<chem>O=C1NC(=O)C(=C1)O</chem>	0.252794
20	<chem>O=C1NC(=O)C(S(=O)(=O)C(=O)O)C(=O)O</chem>	-0.693725
21	<chem>O=C1NC(=O)C(=C1)S(=O)(=O)O</chem>	-0.872942
22	<chem>O=C1SC(=O)C=C1</chem>	-0.506404
23	<chem>O=C1SC(=O)C(=C1)C(=O)O</chem>	-0.803691
24	<chem>O=C1SC(=O)C(C(=O)O)C(=O)O</chem>	-0.843670
25	<chem>O=C1SC(=O)C(F)=C1</chem>	-0.453070
26	<chem>O=C1SC(=O)C(F)=C1F</chem>	-0.454975
27	<chem>O=C1SC(=O)C(=C1)N</chem>	0.312659
28	<chem>O=C1SC(=O)C(N)=C1N</chem>	0.314564
29	<chem>O=C1SC(=O)C(O)=C1O</chem>	0.102315
30	<chem>O=C1SC(=O)C(=C1)O</chem>	-0.075376
31	<chem>O=C1SC(=O)C(S(=O)(=O)C(=O)O)C(=O)O</chem>	-0.821238
32	<chem>O=C1SC(=O)C(=C1)S(=O)(=O)O</chem>	-0.950068
33	<chem>O=C1OC(=O)C(C(=O)O)C(=O)O</chem>	-0.817820
34	<chem>O=C1OC(=O)C=C1</chem>	-0.192929
35	<chem>O=C1OC(=O)C(=C1)C(=O)O</chem>	-0.656067
36	<chem>O=C1OC(=O)C(F)=C1F</chem>	-0.067756

	Smiles	Reaction energy(eV)
37	<chem>O=C1OC(=O)C(F)=C1</chem>	-0.084355
38	<chem>O=C1OC(=O)C(N)=C1N</chem>	0.713211
39	<chem>O=C1OC(=O)C(=C1)N</chem>	0.681373
40	<chem>O=C1OC(=O)C(=C1)O</chem>	0.248712
41	<chem>O=C1OC(=O)C(O)=C1O</chem>	0.423137
42	<chem>O=C1OC(=O)C(S(=O)(=O)C1)S(=O)(=O)O</chem>	0.407379
43	<chem>O=C1OC(=O)C(=C1)S(=O)(=O)C1</chem>	0.824233
44	<chem>C1=NCN=C1</chem>	-0.751035
45	<chem>O=C(O)C1=NCN=C1C(=O)O</chem>	0.61596
46	<chem>C1=NCN=C1C(=O)O</chem>	-1.283562
47	<chem>FC1=NCN=C1F</chem>	-0.246263
48	<chem>C1=NCN=C1F</chem>	-0.507220
49	<chem>NC1=NCN=C1N</chem>	0.087893

## 2. Explanatory Analysis

5

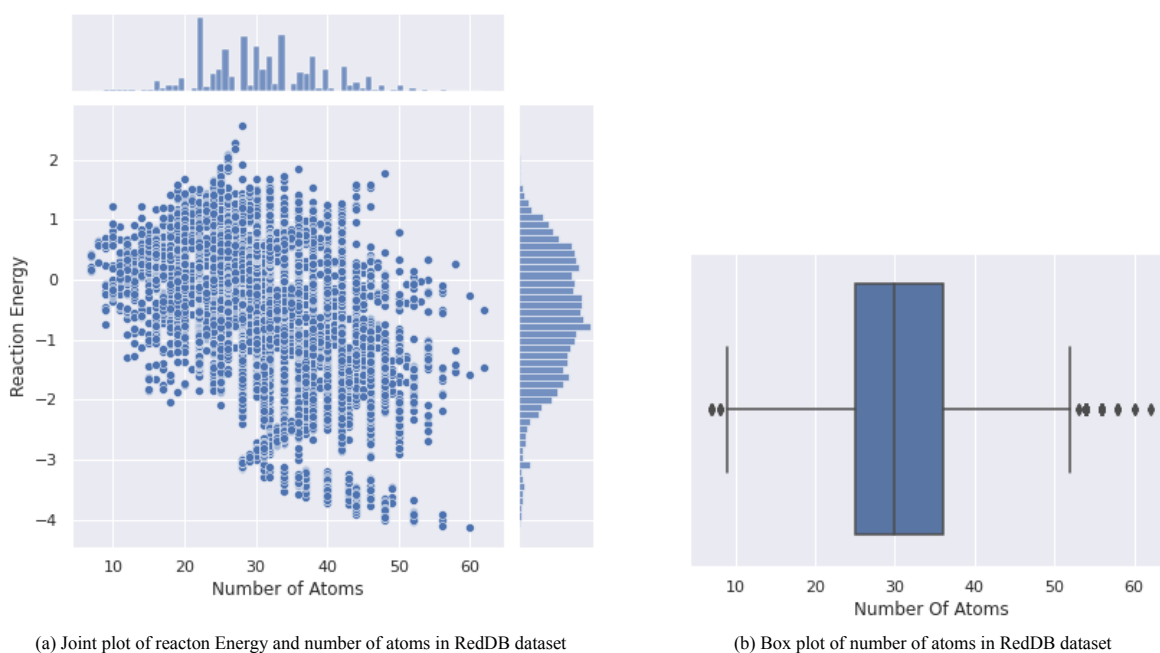


Figure S 1: Relationship of Number of atoms and Reaction energy in RedDB dataset a) Joint plot of number of atoms and reaction energy and b) box plot of number of atoms

2.1. Number of heavy atoms in RedDB

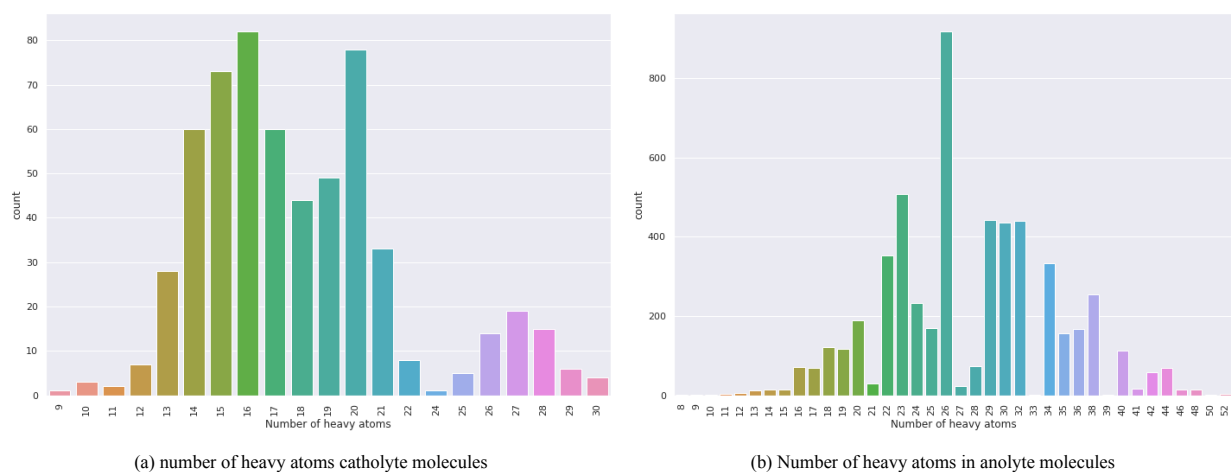
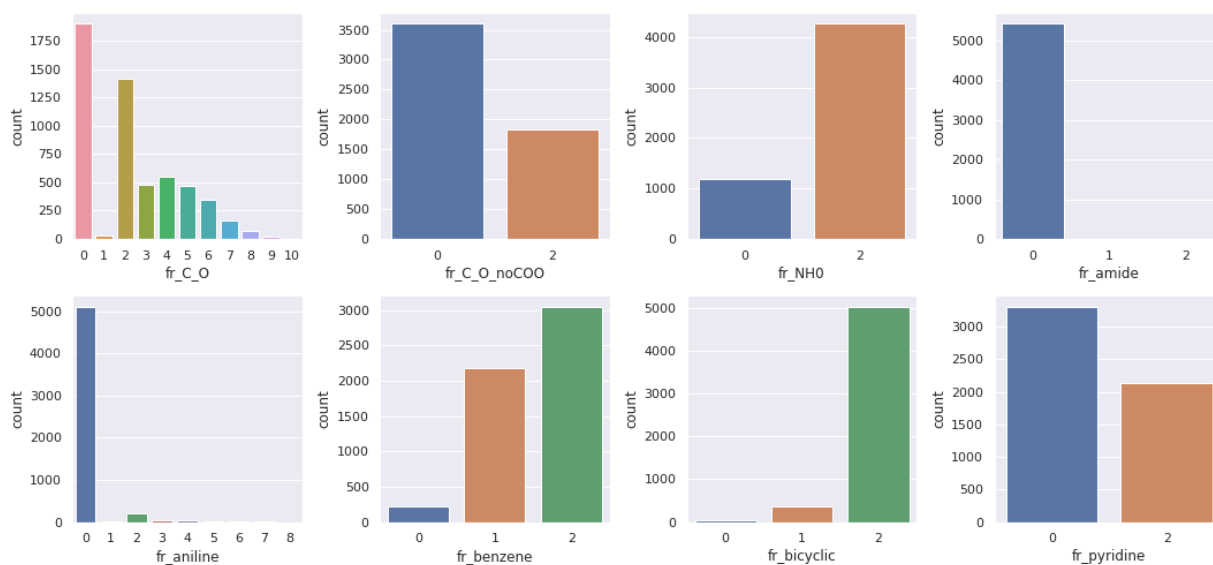


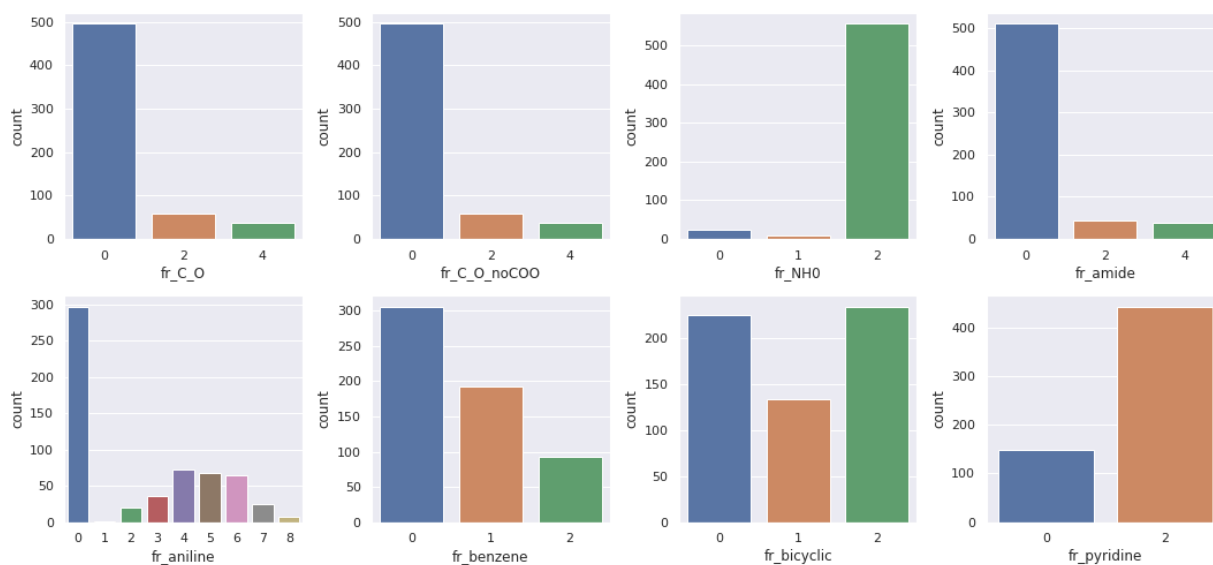
Figure S 2: Number of heavy atoms in RedDB dataset a) number of heavy atoms catholyte molecules and b) Number of heavy atoms in anolyte molecules

## 2.2. Types of functional groups in RedDB

7



(a) Functional groups in the analytes of RedDB dataset



(b) Functional groups in the catholytes of RedDB dataset

Figure S 3: The types and distributions of functional groups in RedDB dataset

## 2.3. Pytorch\_geometric(Pyg) graph representation of RedDB from the SMILE Strings

8

```
import os.path as osp
from pygdata import RedDB
from molfeatures import GenMolGraph, GenMolecules, GenMolFeatures
path = osp.join(osp.dirname(osp.realpath('__file__')), 'data', 'reddb')
```

9

---

```
mol_reddb = RedDB(root_dir=path,
                  name='reddb.csv',
                  smi_idx=-2,
                  target_idx=-1,pre_transform=GenMolFeatures()).shuffle()
```

#### 2.4. Target variable normalized to mean = 0 and std =1

```
mol_reddb.data.y = mol_reddb.data.y*27.2114 # 1 Hartree=27.2114eV
r_mean = mol_reddb.data.y.mean()
r_std = mol_reddb.data.y.std()
mol_reddb.data.y = (mol_reddb.data.y - r_mean) / r_std
print("Normalized redox potential:\n",mol_reddb.data.y)
```

Normalized redox potential:

```
tensor([[ -0.4768],
        [ 0.1099],
        [-1.0686],
        ...,
        [ 0.7942],
        [ 0.3313],
        [ 0.0276]])
```

#### 2.5. Sample atomic attributes

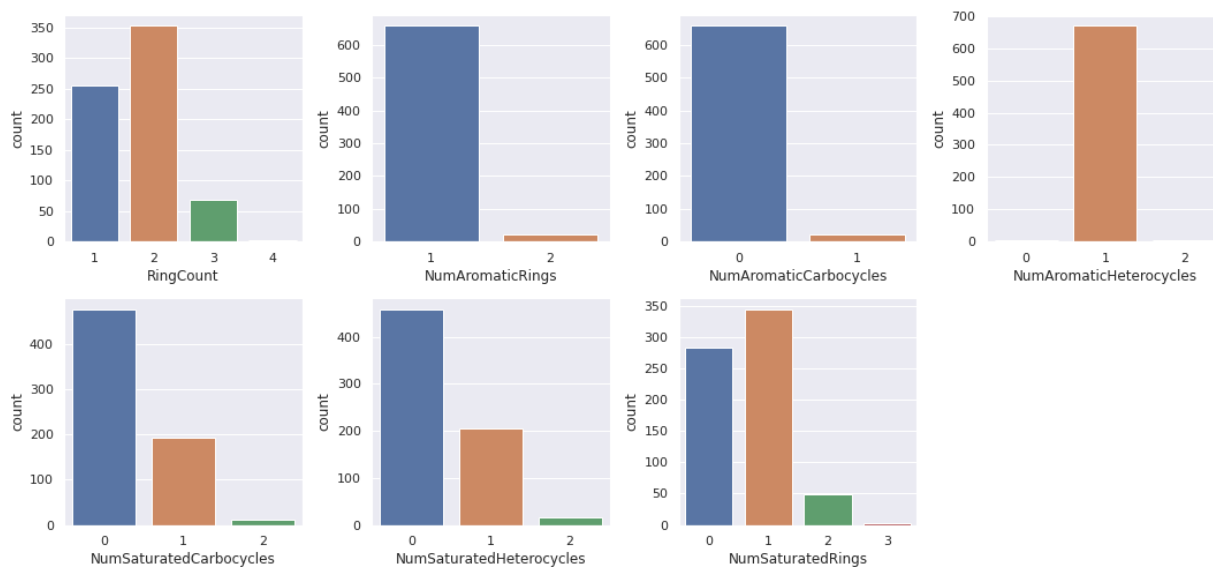
Sample atomic features in Pyg graph format:

```
tensor([[0.0769, 0.0000, 0.0000, ..., 0.0632, 0.2593, 0.2632],
        [0.0549, 0.0000, 1.0000, ..., 0.0464, 0.3704, 0.2632],
        [0.0549, 0.0000, 1.0000, ..., 0.0464, 0.3704, 0.2632],
        ...,
        [0.0000, 1.0000, 0.0000, ..., 0.0000, 0.0000, 0.0000],
        [0.0000, 1.0000, 0.0000, ..., 0.0000, 0.0000, 0.0000],
        [0.0000, 1.0000, 0.0000, ..., 0.0000, 0.0000, 0.0000]])
```

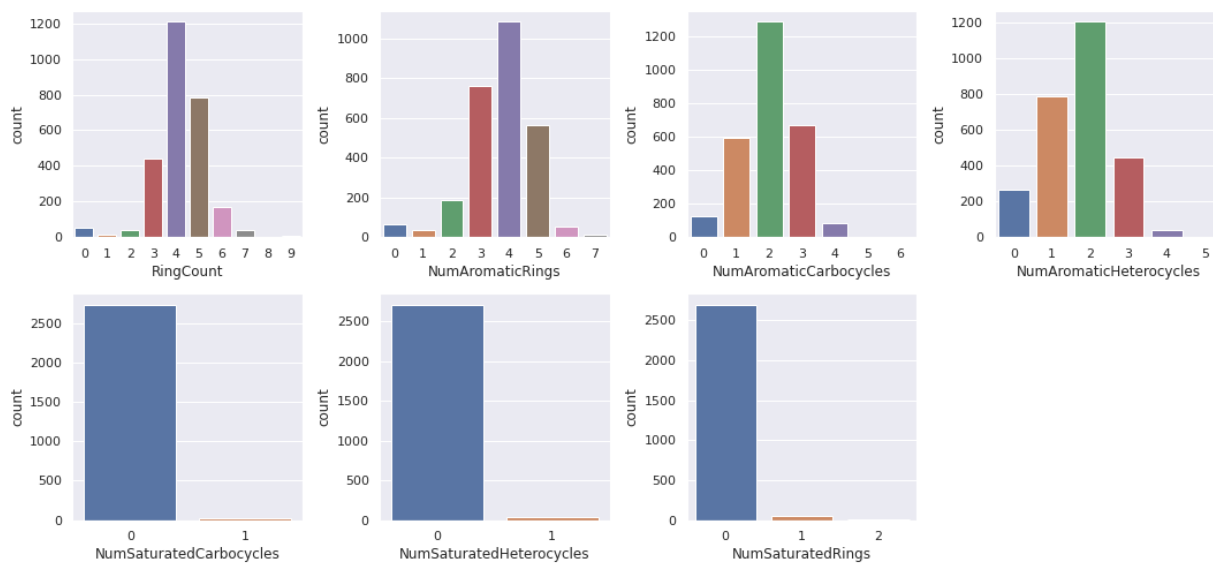
#### 2.6. Sample edge attributes

Sample edge features in Pyg graph format:

```
tensor([[0., 1., 0., ..., 0., 0., 0.],
        [0., 1., 0., ..., 0., 0., 0.],
        [1., 0., 0., ..., 0., 0., 0.],
        ...,
        [1., 0., 0., ..., 0., 0., 0.],
        [1., 0., 0., ..., 0., 0., 0.],
        [1., 0., 0., ..., 0., 0., 0.]])
```



(a) Structural descriptors in catholyte screened dataset



(b) Structural descriptors in anolyte screened dataset

Figure S 4: Structure Structural descriptors of screened dataset

### 3. Screened Dataset

40

#### 3.1. Sample molecules from screened dataset

41

	smiles	reaction_energy
0	<chem>O=C1C(c2ccccc2)=Nc2cc([N+](=O)[O-])c(O)cc21</chem>	-2.880693
1	<chem>CC1=C(C(=O)N(N1)C2=CC=CC=C2)C=C3C=NN=C3C4=CC(=...</chem>	-2.524888
2	<chem>C1=CC(=CC=C1C(=O)N=NC2=C3C=C(C=C(C3=NC2=O)C1)C...</chem>	-2.375739

42

43

44

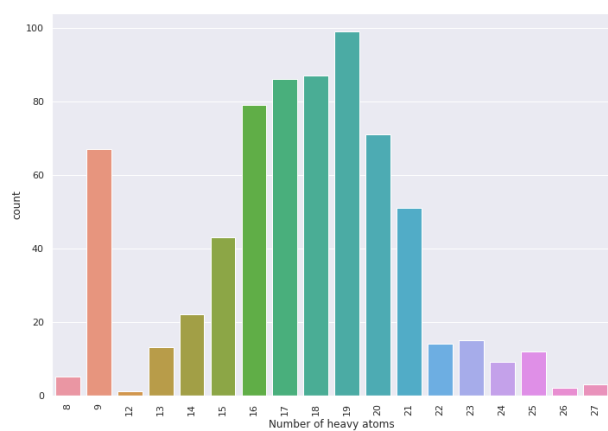
45

3	C1=CC(=O)C(=CN=C2C=CC3=NC(=NC3=C2)C4=C(C=C(C=C...	-2.350572	46
4	CONC1=C2C=C(C=CC2=NC1=O) [N+] (=O) [O-]	-2.310877	47
5	CONC1=C2C=C(C=CC2=NC1=O) [N+] (=O) [O-]	-2.310876	48
6	COC1=CC(=CN=C2C=CC3=NC(=NC3=C2)C4=CC=CC=C4)C(=...	-2.303406	49
7	C1=CC=C(C=C1)C2=NC3=CC(=NC=C4C=C(C=CC4=O)N(O)O...	-2.284782	50
8	C1=CC=C(C=C1)C2=NC(=C3C=CC(=O)C(=C3)N(O)O)N=C2...	-2.202916	51
9	C1=CC=C(C=C1)N2C(=O)C(=C(N2)C(=O) [O-])C=C3C=NC...	-2.197017	52
10	C1=CC=C(C=C1)N2C(=O)C(=C(N2)C(=O) [O-])C=C3C=NC...	-2.197017	53
11	C1=CC2=C(C(=C1) [N+] (=O) [O-])C(=O)N=C2 [O-]	-2.083223	54
12	C1=CC2=C(C(=C1) [N+] (=O) [O-])C(=O)N=C2 [O-]	-2.083223	55
13	COC1=CC=C(C=C1)N2C(=C(C(=O)NC2=O)C=C3C=NN=C3C4...	-2.070825	56
14	C1=CC2=NC(=C3C=C(C=CC3=O)N=CC4=C(C(=CC=C4)N(O)...	-2.065413	57
23452	CCCC [NH2+] CC0c1cccnc1C	2.294042	58
23453	COc1ccsc1CN [C@H] 1CC [C@H] ( [NH+] (C)C)C1	2.299948	59
23454	Cc1cnccc1CC [NH2+] [C@@H] (C(C)C)C1CC1	2.302506	60
23455	CC [NH+] (CC)CC [C@H] (N)c1cnccc1C	2.313252	61
23456	CC [NH2+] C [C@@H] (Oc1cccnc1C)C(C)C	2.314913	62
23457	COc1ccsc1C [NH+] (C) [C@H] (C)C1(C)CC1	2.323198	63
23458	CC [C@H] 1CCCC [C@@H] 10c1ccc(C)nc1C [NH2+] C(C)C	2.336868	64
23459	Cc1ncccc10 [C@H] (C)C [NH2+] CC(C)C	2.338219	65
23460	CCC(C) (C) [NH2+] CC0c1cccnc1C	2.368817	66
23461	CN [C@@H] (c1cnccc1C) [C@] 1 ( [NH+] (C)C)CCC [C@H] (C)C1	2.793114	67
23462	CC(C) (C) [NH2+] Cc1ccc(OCC [NH+] 2CCCCC2)cn1	2.794882	68
23463	CCC(CC) (CCO)CNc1cncc(Br)c1	2.813119	69
23464	CC [C@H] (C)COc1ccc(C [C@H] ( [NH3+] )CC)nc1	2.833529	70
23465	Cc1ccc(O [C@H] (C) [C@@H] (C)O)c(C [NH2+] CC(C)C)n1	2.854385	71
23466	CC [C@@H] (C)CN(CC)c1nc2c(s1) [C@H] ( [NH2+] )C)CC(C)...	2.873744	72

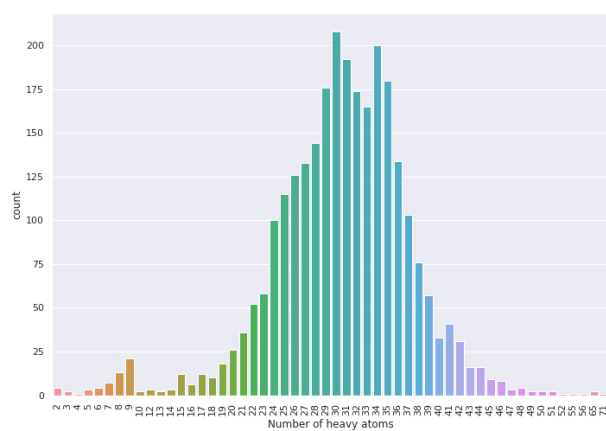


### 3.2. Number of heavy atoms in screened dataset

73



(a) number of heavy atoms catholyte molecules



(b) Number of heavy atoms in anolyte molecules

Figure S 5: Number of heavy atoms in screened dataset a) number of heavy atoms catholyte molecules and b) Number of heavy atoms in anolyte molecules

### 3.3. Types of functional groups in screened dataset

74

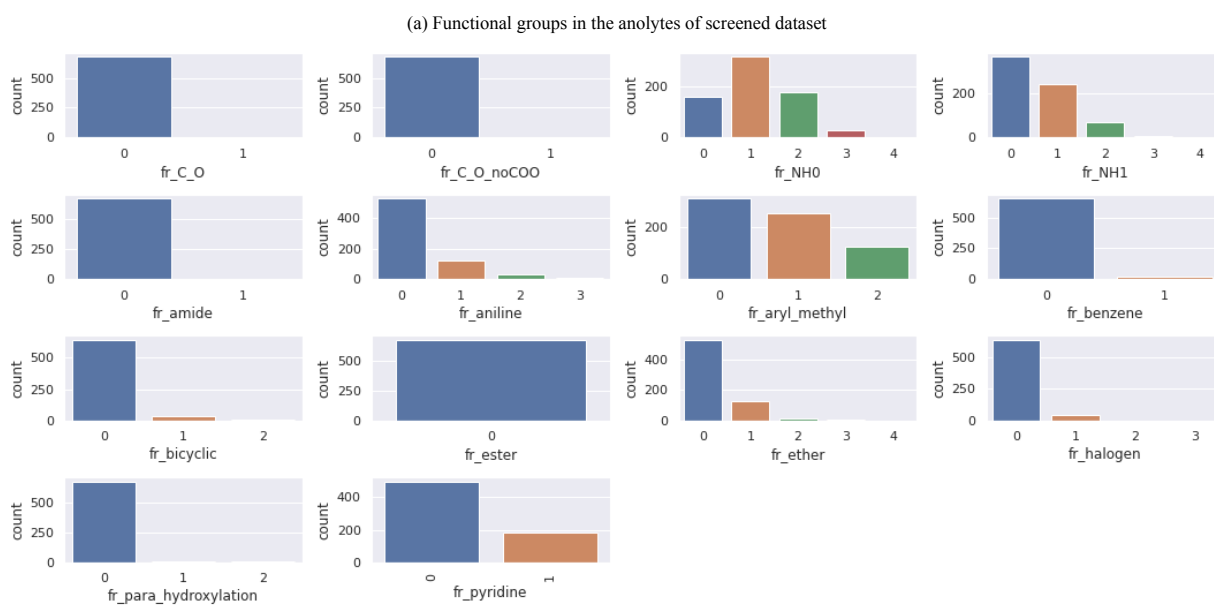
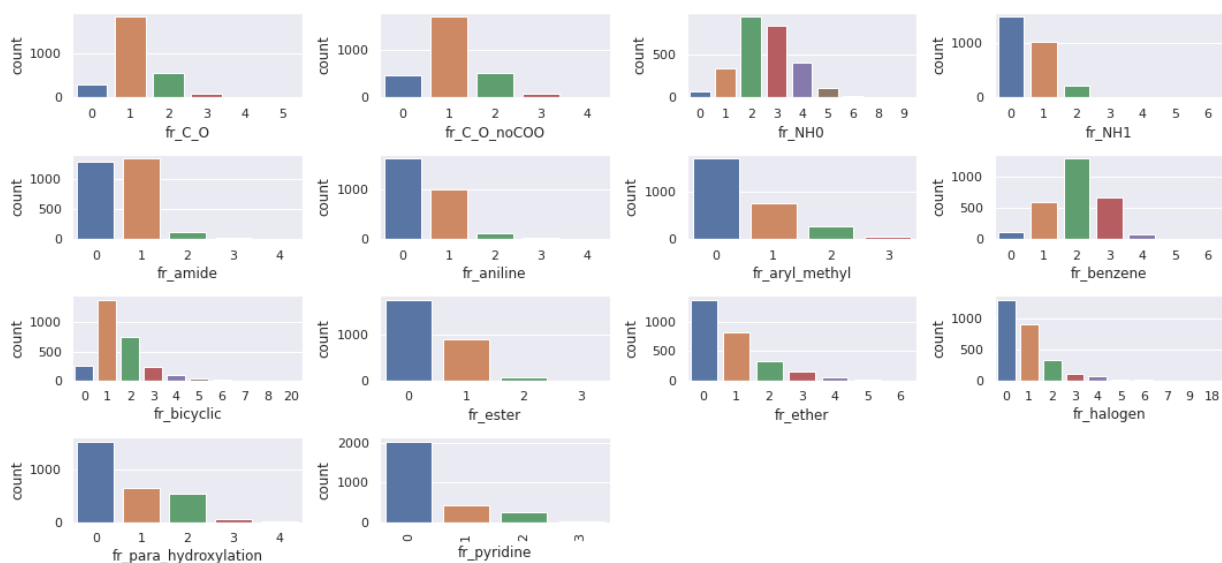


Figure S 6: The types and distributions of functional groups in screen dataset

## 4. MolGAT Model Implementation

75

```
class MolGAT(torch.nn.Module):  
    def __init__(self, node_features,  
                 hidden_dim,  
                 edge_features,
```

76

```

        num_heads,
        dropout,
        num_conv_layers,
        num_fc_layers):
super(MolGAT, self).__init__()
self.conv_list = torch.nn.ModuleList()
self.bn_list = torch.nn.ModuleList()
self.num_fc_layers = num_fc_layers

self.bn_list.append(BatchNorm(hidden_dim))
self.conv_list.append(MolGATConv(node_features, hidden_dim, edge_features, heads=num_heads))
for i in range(num_conv_layers-1):
    self.conv_list.append(MolGATConv(hidden_dim, hidden_dim, edge_features, heads=num_heads))

self.fc_list = torch.nn.ModuleList()
for i in range(num_fc_layers -1):
    if i == 0:
        self.fc_list.append(torch.nn.Linear(hidden_dim*2, hidden_dim*2))
    else:
        self.fc_list.append(torch.nn.Linear(hidden_dim*2, hidden_dim*2))

self.fc_out = torch.nn.Linear(hidden_dim*2, 1)

self.dropout = dropout

def forward(self, x, edge_index, batch_index, edge_attr):
    for i, (conv, bn) in enumerate(zip(self.conv_list, self.bn_list)):
        x = F.relu(conv(x, edge_index, edge_attr))
        x = F.dropout(x, p=self.dropout, training=self.training)
        if i != (self.num_fc_layers-1):
            x = bn(x)

    x = torch.cat([gmp(x, batch_index),
                    gap(x, batch_index)], dim=1)

    for i, fc in enumerate(self.fc_list ):
        x = F.relu(fc(x))
        if i != (self.num_fc_layers-1):
            x = F.dropout(x, p=self.dropout, training=self.training)

```

```

        x = self.fc_out(x)
        return x
# optimized training parameters
class TrainArgs:
    edge_features = mol_reddb.data.edge_attr.shape[1]
    num_features=mol_reddb.num_features
    dropout=0.1
    num_fc_layers=3
    num_conv_layers=3
    num_heads=4
    hidden_dim=512
    batch_size = 192
args = TrainArgs()

# define the device and the molgat model

device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
model = MolGAT(node_features=args.num_features,
               hidden_dim=args.hidden_dim,
               edge_features=args.edge_features,
               num_heads=args.num_heads,
               dropout=args.dropout,
               num_fc_layers=args.num_fc_layers,
               num_conv_layers=args.num_conv_layers).to(device)

## print model parameters
print(model)
print("Number of parameters: ", sum(p.numel() for p in model.parameters()))

```

```

MolGAT(
  (conv_list): ModuleList(
    (0): MolGATConv(99, 512, 12, heads=4)
    (1): MolGATConv(512, 512, 12, heads=4)
    (2): MolGATConv(512, 512, 12, heads=4)
  )
  (bn_list): ModuleList(
    (0): BatchNorm(512)
  )
  (fc_list): ModuleList(
    (0): Linear(in_features=1024, out_features=1024, bias=True)
    (1): Linear(in_features=1024, out_features=1024, bias=True)
  )

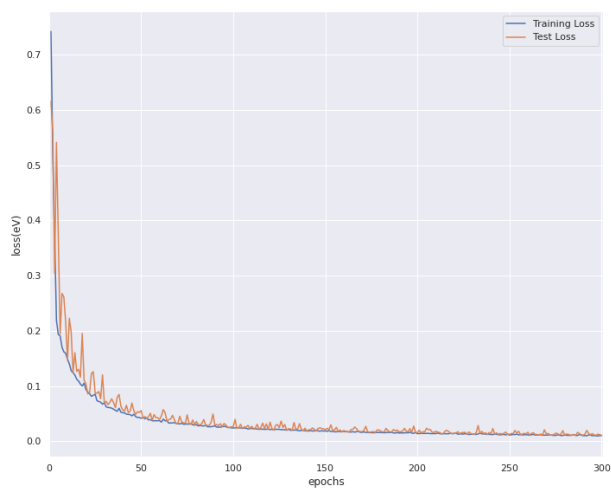
```

```
(fc_out): Linear(in_features=1024, out_features=1, bias=True)
)
Number of parameters: 5219985
```

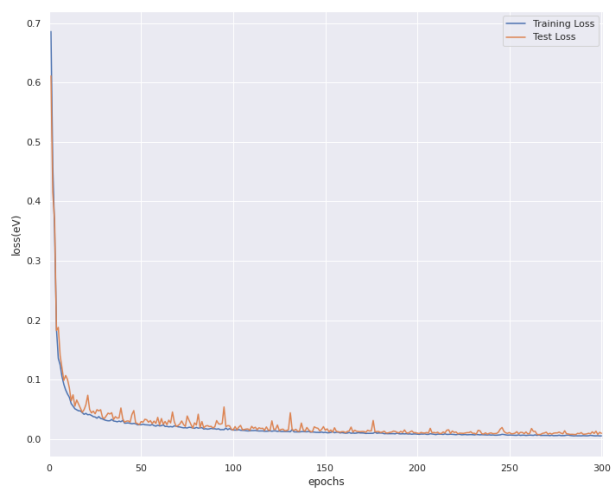
92  
93  
94

## 5. Training MolGAT Model

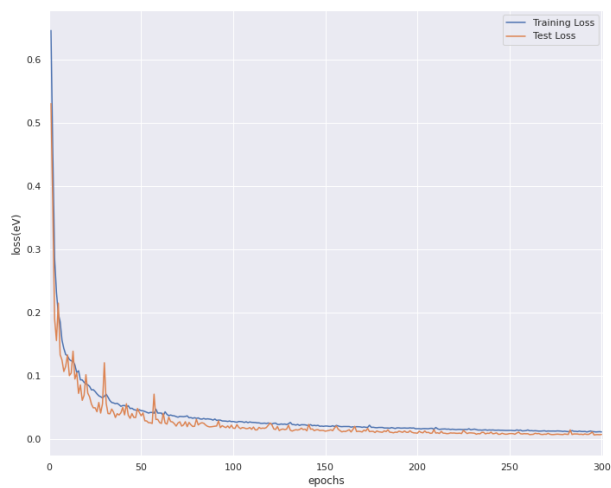
95



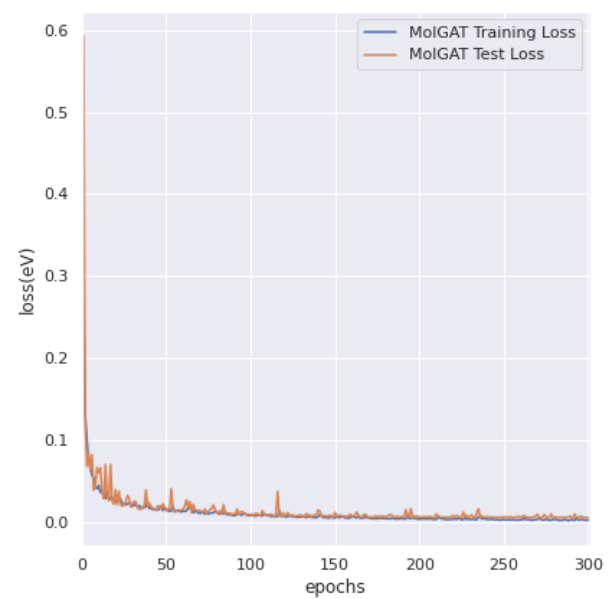
(a) Error loss plot during GCN model training



(b) Error loss plot during GAT model training

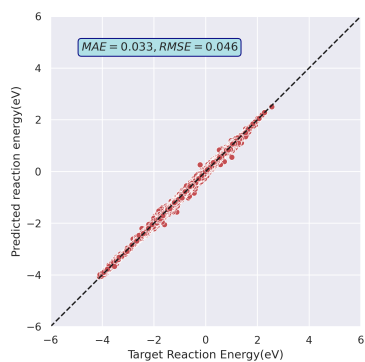


(c) Error loss plot during AttentiveFP model training

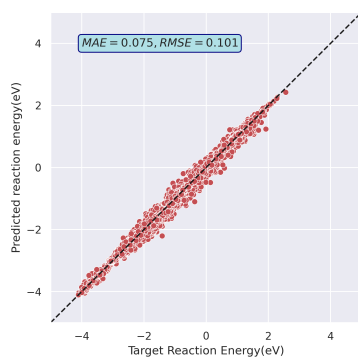


(d) Error loss plot for MolGAT model

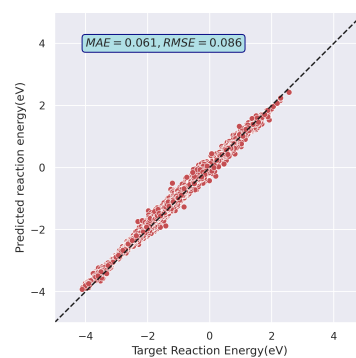
Figure S 7: Parity plots for Four different models a) Parity plot of MolGAT Model b) Parity plot of GCN Model c) Parity plot of GAT model d) Parity plot of AttentiveFP model



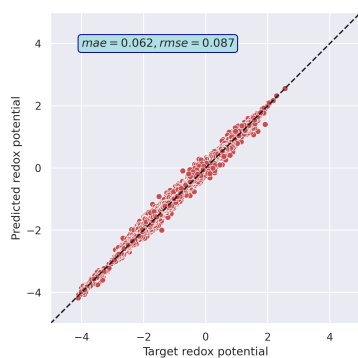
(a) Parity plot of MolGAT Model



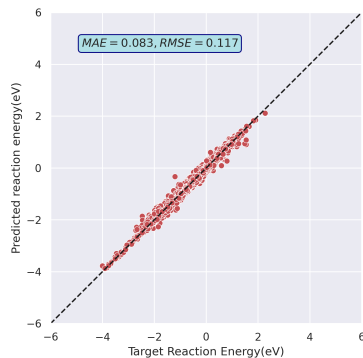
(b) Parity plot of GCN Model



(c) Parity plot of GAT model



(d) Parity plot of AttentiveFP model



(e) Parity plot of MPNN model

Figure S 8: Parity plots with their corresponding MAE and RMSE loss to benchmark four different models with a) Parity plot of MolGAT Model b) Parity plot of GCN Model c) Parity plot of GAT model d) Parity plot of AttentiveFP model e) Parity plot of MPNN Model