

Advancing Energy Storage through Solubility Prediction: Leveraging the Potential of Deep Learning

Mesfin Diro Chaka ^{a,c,*}, Yedilfana Setarge Mekonnen ^b, Qin Wu ^d, Chernet Amente Geffe ^a

^aDepartment of Physics, College of Natural and Computational Sciences, Addis Ababa University, P. O. Box 1176, Addis Ababa, Ethiopia,


^bCenter for Environmental Science, College of Natural and Computational Sciences, Addis Ababa University, P. O. Box 1176, Addis Ababa, Ethiopia,

^cComputational Data Science Program, College of Natural and Computational Sciences, Addis Ababa University, P. O. Box 1176, Addis Ababa, Ethiopia,

^dCenter for Functional Nanomaterials, Brookhaven National Laboratory, Upton, NY 11973, USA,



*Corresponding author

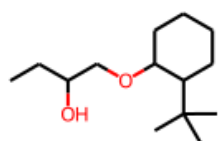
Email address: mesfin.diro@aaau.edu.et (Mesfin Diro Chaka )

1. AqSolDB Dataset

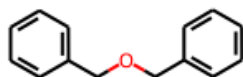
4

Table S 1: Randomly sampled smiles strings with their corresponding LogS in AqSolDB dataset

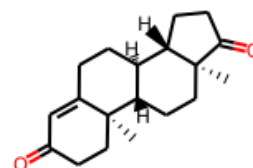
smiles	Experimental LogS(Mol/L)
<chem>O=C([O-])CO.[NH4+]</chem>	1.701
<chem>OCCCCCO</chem>	1.626
<chem>CC(N)=O</chem>	1.581
<chem>CNC</chem>	1.558
<chem>CN</chem>	1.369
<chem>C1CN1</chem>	1.366
<chem>CC=O</chem>	1.356
<chem>NC=O</chem>	1.346
<chem>CNN</chem>	1.337
<chem>c1cn[nH]c1</chem>	1.288
<chem>C#CCO</chem>	1.251
<chem>CC1CN1</chem>	1.243
<chem>C=CCN</chem>	1.243
<chem>NC1CC1</chem>	1.243
<chem>O=CC=O</chem>	1.236
<chem>C1COC1</chem>	1.236
<chem>CCO</chem>	1.234
<chem>CNC=O</chem>	1.229
<chem>CC(C)N</chem>	1.228
<chem>CCCN</chem>	1.228
<chem>NCCN</chem>	1.221
<chem>CN(C)N</chem>	1.221
<chem>NCCO</chem>	1.214
<chem>OCCO</chem>	1.207
<chem>OCCF</chem>	1.193
<chem>CN(C)C</chem>	1.178
<chem>O=CO</chem>	1.176
<chem>Clc1cc(Cl)c(Cl)c(-c2c(Cl)c(Cl)c(Cl)c(Cl)c2Cl)c1Cl</chem>	-10.412
<chem>Clc1cc(Oc2cc(Cl)c(Cl)c(Cl)c2Cl)c(Cl)c(Cl)c1Cl</chem>	-10.100
<chem>Clc1cc(Oc2c(Cl)c(Cl)cc(Cl)c2Cl)c(Cl)c(Cl)c1Cl</chem>	-10.100
<chem>Clc1cc(Cl)c(Oc2c(Cl)c(Cl)c(Cl)c2Cl)cc1Cl</chem>	-10.100
<chem>Clc1cc(-c2c(Cl)c(Cl)c(Cl)c2Cl)cc(Cl)c1Cl</chem>	-9.700
<chem>Clc1ccc(Oc2c(Cl)c(Cl)c(Cl)c2Cl)c(Cl)c1</chem>	-9.640
<chem>Clc1cc(Oc2cc(Cl)c(Cl)c(Cl)c2Cl)cc(Cl)c1Cl</chem>	-9.540
<chem>Clc1cc(Cl)c(-c2c(Cl)c(Cl)c(Cl)c2Cl)cc1Cl</chem>	-9.500



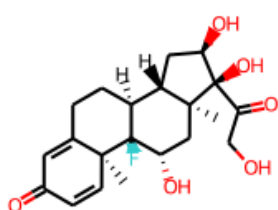
-3.6987



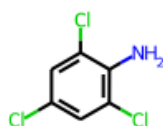
-3.695



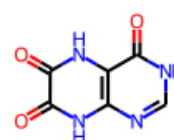
-3.6949



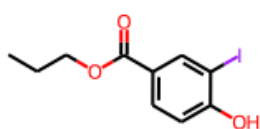
-3.6929



-3.6912



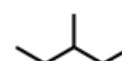
-3.69



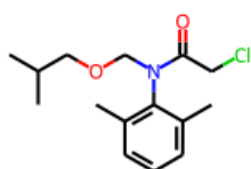
-3.6865



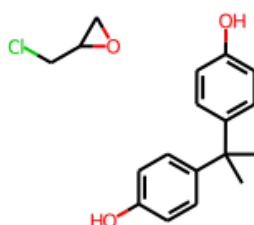
-3.684



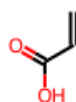
-3.6826



-3.6822



-3.6804



-3.68

Figure S 1: Some of randomly sampled molecules in AqSolDB dataset

2. Explanatory Analysis

5

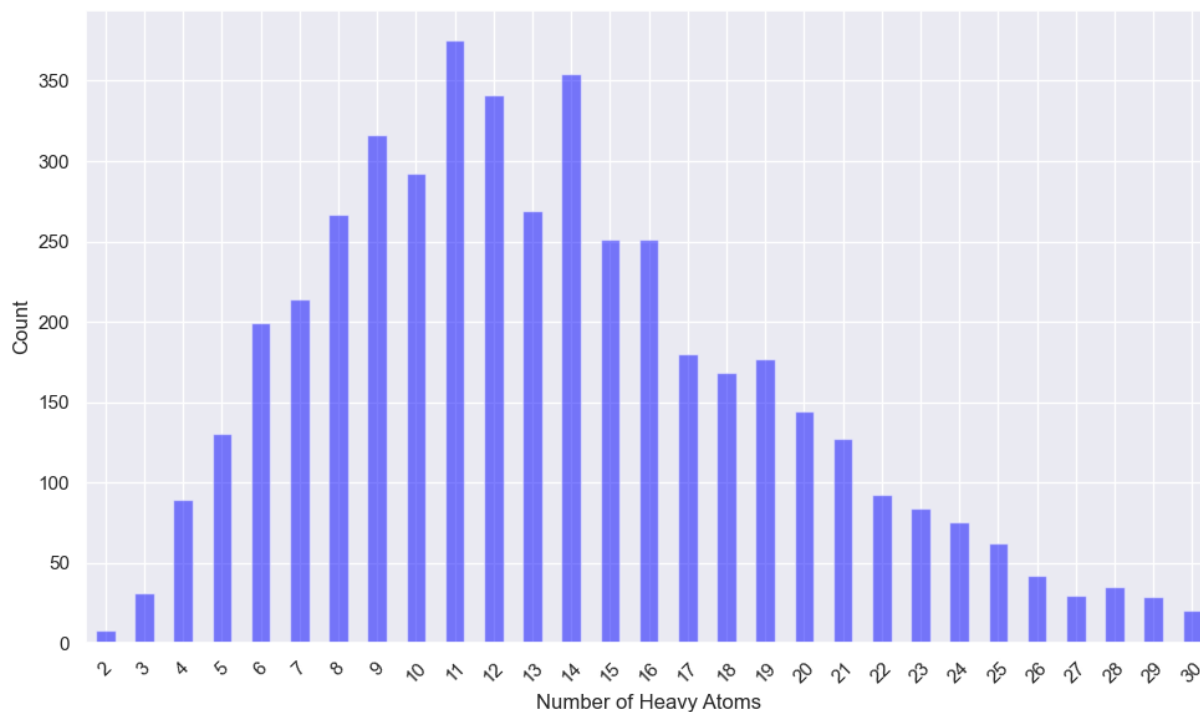


Figure S 2: The total number of heavy atoms in AqSolDB dataset used for the training

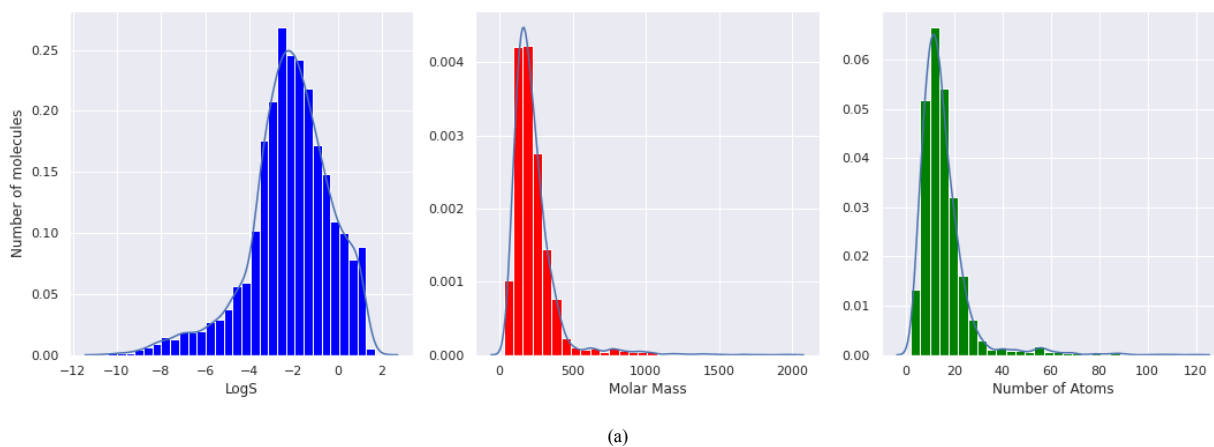


Figure S 3: Histogram of LogS, Molar mass and Number of atoms in AqSolDB dataset

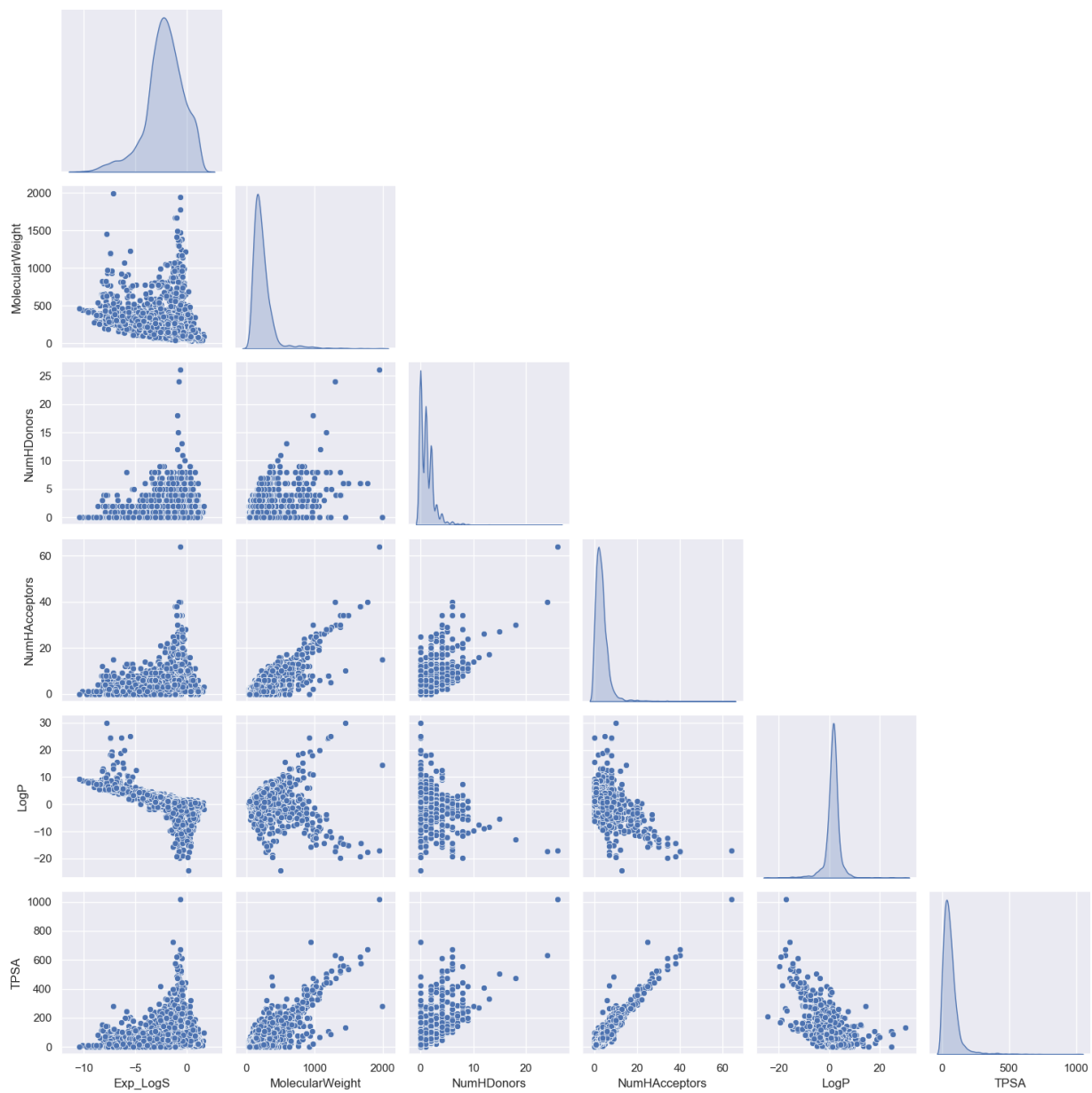


Figure S 4: Pair Plot of Selected Descriptors and Experimental LogS

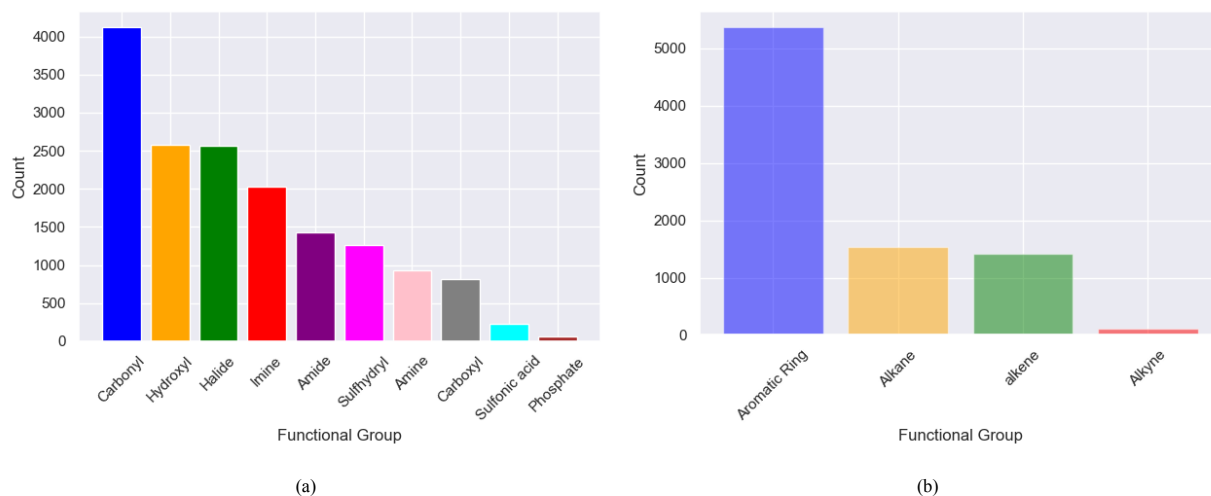


Figure S 5: The distribution of dominant functional groups in AqSolDB dataset a) polar functional groups b) non-polar functional groups

2.1. Pytorch_geometric(Pyg) graph representation of AqSolDB from the SMILE Strings

6

```
import os.path as osp
from custom_pygdata import AqSolDB
path = osp.join(osp.dirname(osp.realpath('__file__')), 'data', 'aqsolddb')
mol_sol = AqSolDB(root_dir=path,
                  name='AqSolDB2.csv',
                  smi_idx=-4,
                  target_idx=-3).shuffle()
```

7

Data(x=[290455, 30], edge_index=[2, 590090], edge_attr=[590090, 12], y=[8494, 1], smiles=[8494])

2.2. Target variable normalized to mean = 0 and std = 1

9

```
r_mean = mol_sol.data.y.mean()
r_std = mol_sol.data.y.std()
mol_sol.data.y = (mol_sol.data.y - r_mean) / r_std
print("Normalized LogS:\n",mol_sol.data.y)
```

10

```
Normalized LogS:
tensor([[ 0.8059],
        [-0.0870],
        [ 0.0549],
        ...,
        [-1.8041],
        [-0.1992],
        [-1.3752]])
```

11

12

13

14

15

16

17

18

2.3. Sample atomic attributes

19

Sample atomic features in Pyg graph format:

20

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

21

22

23

24

25

26

27

2.4. Sample edge attributes

28

Sample edge features in Pyg graph format:

29

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

30

31

32

33

34

35

36

3. MolGAT Model Implementation

37

MolGAT(

38

 (conv_list): ModuleList(

39

 (0): MolGATConv(30, 192, 12, heads=4)

40

 (1-2): 2 x MolGATConv(192, 192, 12, heads=4)

41

)

42

 (fc_list): ModuleList(

43

 (0): Linear(in_features=384, out_features=384, bias=True)

44

 (1): Linear(in_features=384, out_features=192, bias=True)

45

 (2): Linear(in_features=192, out_features=192, bias=True)

46

)

47

 (fc_out): Linear(in_features=192, out_features=1, bias=True)

48

)

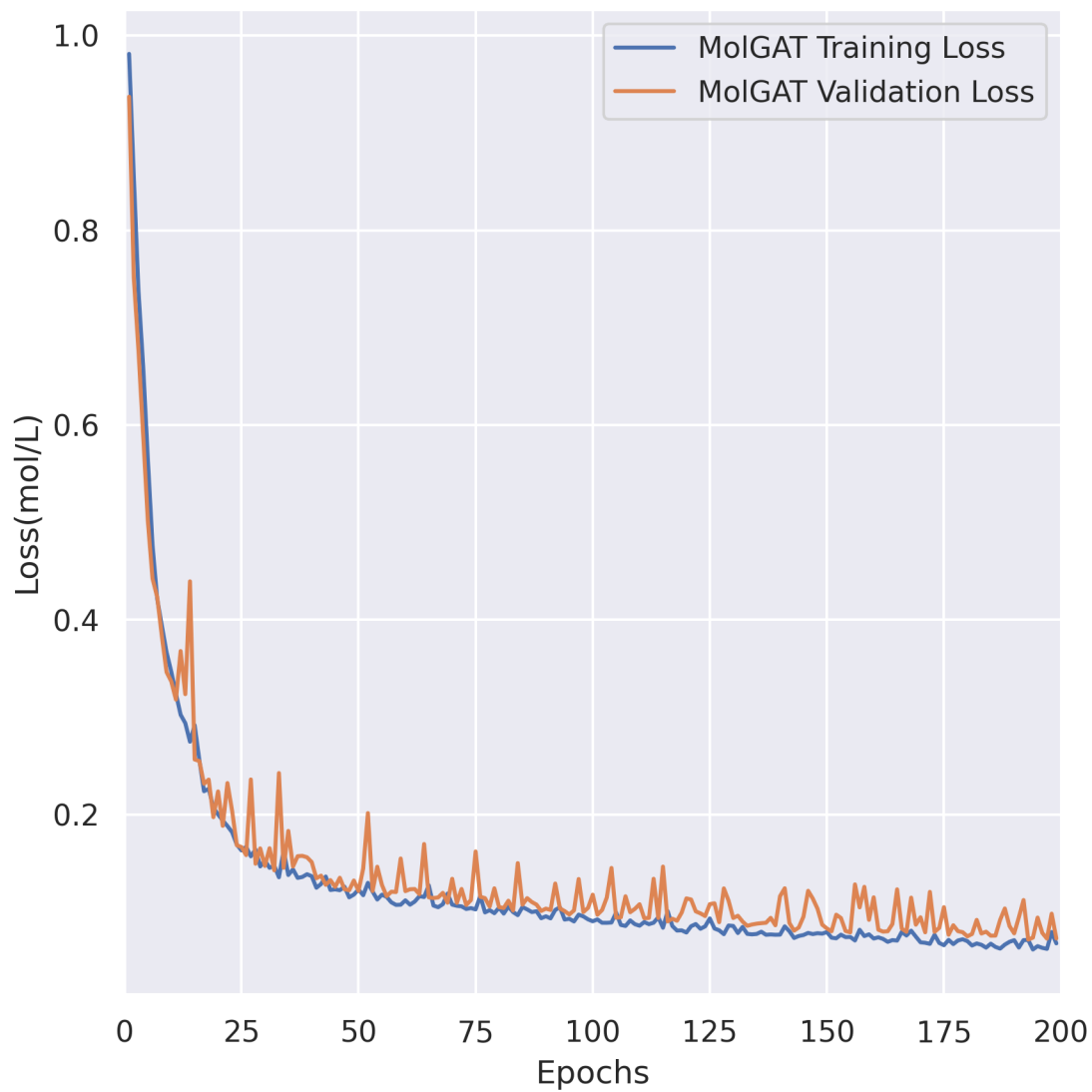
49

Number of parameters: 699793

50

4. Training MolGAT Model

51



(a)

Figure S 6: Error loss plot during MolGAT model training on AqSolDB

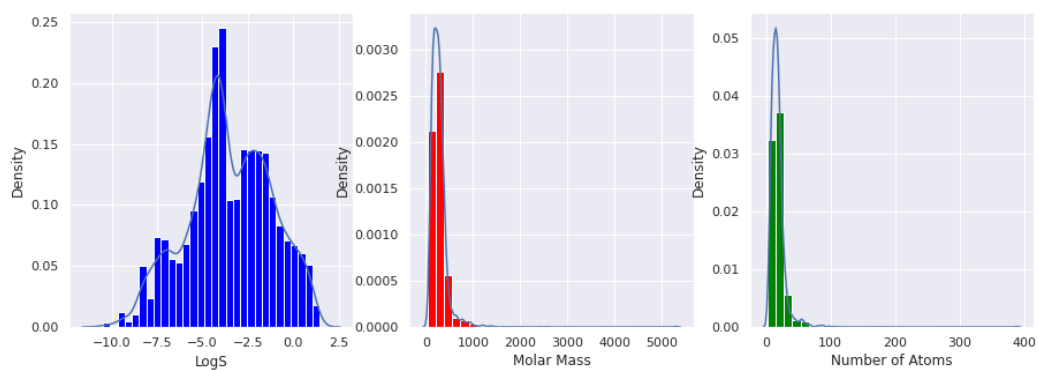


Figure S 7: Distributions of log solubility, molar mass (g/mol), and number of atoms for molecules in screened data set with MolGAT.

4.1. Benchmarking

52

4.1.1. Random Forest(RF) model training

53

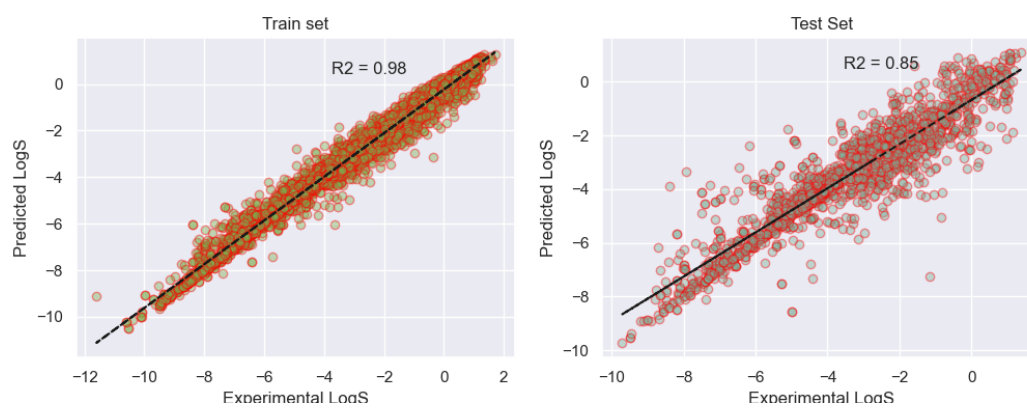


Figure S 8: Parity plot for Random Forest model with AqSolDB.

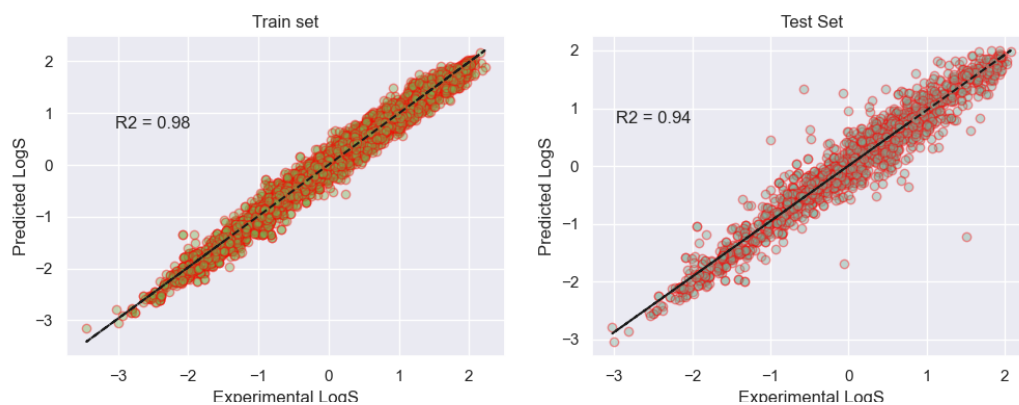


Figure S 9: Parity plot for MPNN model with AqSolDB.

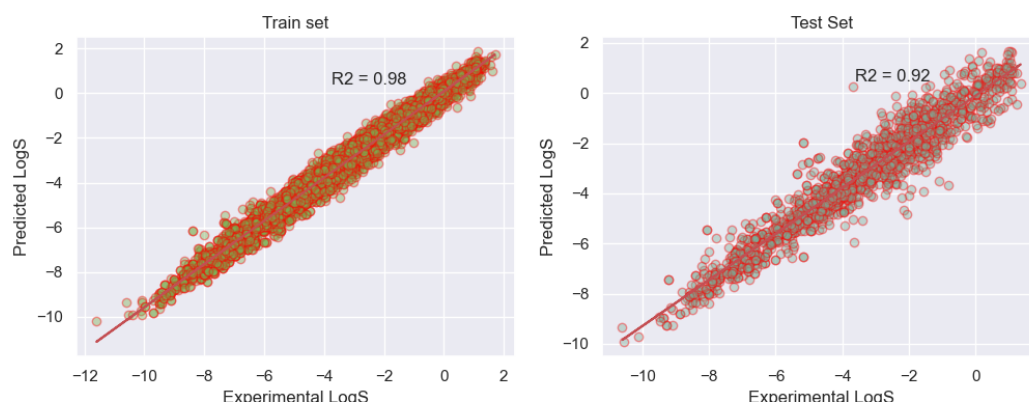


Figure S 10: Parity plot for GCM model with AqSolDB.

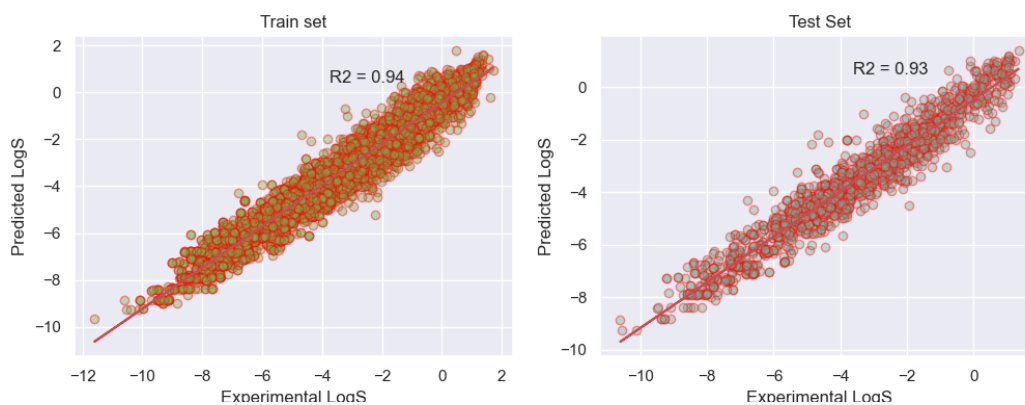


Figure S 11: Parity plot for AttentiveFP model with AqSolDB.

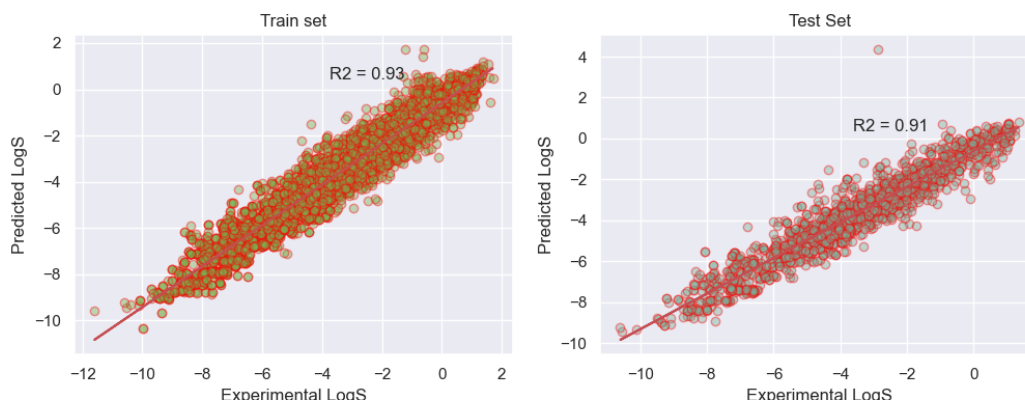


Figure S 12: Parity plot for GAT model with AqSolDB.

5. Screened Dataset

5.1. Sample molecules from screened dataset

Table S 2: Sample smiles, reaction energy and logS from a total screened dataset using the trained molgat model

SMILES	Reaction energy(eV)	Predicted LogS(Mol/L)
<chem>C1 = CC(= O)C(= CN = C2C = CC3 = NC(= NC3 = C2)C4 = C(C = C(C = C4)C)C)C = C1N(O)O</chem>	-2.35057	-3.63134
<chem>CONC1 = C2C = C(C = CC2 = NC1 = O)[N+](= O)[O-]</chem>	-2.31088	-2.31388
<chem>COC1 = CC(= CN = C2C = CC3 = NC(= NC3 = C2)C4 = CC = CC(= C4)C(= O)C(= C1)N(O)O</chem>	-2.30341	-2.66933
<chem>C1 = CC = C(C = C1)C2 = NC3 = CC(= NC = C4C = C(C = CC4 = O)N(O)O)C = CC3 = N2</chem>	-2.28478	-2.75551

SMILES	Reaction energy(eV)	Predicted LogS(Mol/L)
<chem>C1 = CC = C(C = C1)C2 = NC(= C3C = CC(= O)C(= C3)N(O)O)N = C2C4 = CC = CC = C4</chem>	-2.20292	-3.32659
<chem>C1 = CC2 = C(C(= C1)[N+](= O)[O-])C(= O)N = C2[O-]</chem>	-2.08322	-2.38020
<chem>C1 = CC2 = NC(= C3C = C(C = CC3 = O)N = CC4 = C(C(= CC = C4)N(O)O)O)N = C2C = C1</chem>	-2.06541	-3.33318
<chem>CC1 = CC(= CN = C2C = CC(= C3N = C4C = CC = CC4 = N3)C = C2)C(= O)C(= C1)N(O)O</chem>	-2.01272	-2.42499
<chem>C1 = CC2 = C3C(= CC(= N)C(= O)N3C4 = CC = C(C = C4)N(O)O)N = C2C = C1</chem>	-1.94489	-3.32238
<chem>C1 = CC2 = C3C(= C1)NC(= NC3 = CC = C2)C(= O)O</chem>	-1.78178	-3.04786
<chem>C1 = CC = C2C(= C1)C(= CC3 = C(N(C(= S)S3)CC(= O)[O-])O)C = N2</chem>	-1.76830	-3.55150
<chem>CC(= O)N = C1C = CC2 = NC(= NC(= S)NC(= O)C3 = CC = C(O3)C4 = CC = CC = C4N(O)O)SC2 = C1</chem>	-1.75821	-3.09294
<chem>CC1 = CC = CC2 = C(C(= O)N = C12)C3 = C(N(C(= S)S3)CC(= O)O)O</chem>	-1.71661	-3.48864
<chem>C1 = COC(= C1)C2 = NN3C(= N)C(= CC4 = CC = C(O4)[N+](= O)[O-])C(= O)N = C3S2</chem>	-1.71067	-3.51928
<chem>C1 = CC = C2C(= C1)C(= CC3 = C(N(C(= S)S3)C(CCC(= O)[O-])C(= O)[O-])O)C = N2</chem>	-1.70025	-3.29246
<chem>CCCN[C@H](e1enccc1C)C(C)(C)[NH+](C)C</chem>	2.29304	-1.58715
<chem>CCCC[NH2+][CO]e1cccnc1C</chem>	2.29404	-1.31283
<chem>COe1ccsc1CN[C@H]1CC[C@H]([NH+](C)C)C1</chem>	2.29995	-1.29218
<chem>Ce1enccc1CC[NH2+][C@@H](C(C)C)C1CC1</chem>	2.30251	-3.13770
<chem>CC[NH+](CC)CC[C@H](N)e1enccc1C</chem>	2.31325	-0.421087
<chem>CC[NH2+][C@@H](Oe1cccnc1C)C(C)C</chem>	2.31491	-1.53976
<chem>COe1ccsc1C[NH+](C)[C@H](C)C1(C)CC1</chem>	2.32320	-3.54068
<chem>Ce1nccc1O[C@H](C)C[NH2+][CC(C)C</chem>	2.33822	-1.31810
<chem>CCC(C)(C)[NH2+][CO]e1cccnc1C</chem>	2.36882	-1.49300
<chem>CN[C@@H](e1enccc1C)[C@]1([NH+](C)C)CCC[C@H](C)C1</chem>	2.79311	-2.82409
<chem>CC(C)(C)[NH2+][C]e1ccc(OCC[NH+][2CCCCC2])en1</chem>	2.79488	-1.70223
<chem>CCC(CO)(CCO)CNe1enccc(Br)c1</chem>	2.81312	-3.69122
<chem>CC[C@H](C)COe1ccc[C@H]([NH3+])CC)nc1</chem>	2.83353	-1.78672
<chem>Ce1ccc(O[C@H](C)[C@@H](C)O)c[C@H]([NH2+])CC(C)Cn1</chem>	2.85439	-1.72013
<chem>CC[C@@H](C)CN(CC)c1nc2c(s1)[C@H]([NH2+])CC(C)(C)C2</chem>	2.87374	-2.58608