

# Generating Molecules with Optimized Solubility using Iterative Graph Translation

*Camille Bilodeau<sup>1</sup>, Wengong Jin<sup>2</sup>, Hongyun Xu<sup>3</sup>, Jillian Emerson<sup>3</sup>, Sukrit Mukhopadhyay<sup>3</sup>, Tom Kalantar<sup>3</sup>,*

*Tommi Jaakkola<sup>2</sup>, Regina Barzilay<sup>2</sup>, Klavs F. Jensen (kfjensen@mit.edu)<sup>\*1</sup>*

<sup>1</sup>Department of Chemical Engineering, Massachusetts Institute of Technology, 77 Massachusetts Avenue, Cambridge, MA 02139, United States

<sup>2</sup>Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA 02139, United States

<sup>3</sup>Dow Chemical Company, Midland, MI 48674, United States

\*Corresponding Author

## Supplementary Information

### 1. Iterative Graph-to-Graph Translation

All framework, architecture, and training parameters for the iterative graph-to-graph translation procedure can be found in Table 1. The code used for training can be found here: [github.com/cbilodeau2/g2g\\_optimization](https://github.com/cbilodeau2/g2g_optimization).

**Table S1.** Framework, architecture, and training parameters for the iterative graph-to-graph translation.

Framework Parameters	
Cutoff Improvement, $\alpha$	0.8 LogS
Number of Iterations	3
Cutoff SA Score	3.5
Training Parameters	
Batch Size	32
Number of CPU	9
Dropout	0.0
Learning Rate	1e-3
Clip Norm	2.0
Beta	0.3
Number of Epochs	10 (per iteration)

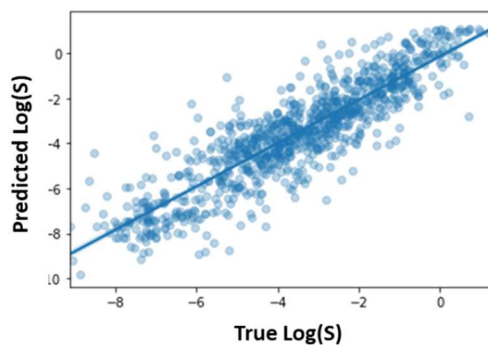
Anneal Rate	0.9
<hr/>	
Architecture Parameters	
<hr/>	
Activation	ReLU
Aggregation	Sum
RNN Type	LSTM
Hidden Size	270
Embedding Size	270
Latent Size	4
Tree Depth	20
Graph Depth	20
Tree Iterations	1
Graph Iterations	3
<hr/>	

## 2. Directed Message Passing Neural Network (DMPNN)

All architecture and training parameters for the DMPNN can be found in Table 2. The resulting model for predicting solubility had  $R = 0.790$ ,  $MAE = 0.789$ , and  $RMSE = 1.038$  (Figure 1).

**Table S2.** Architecture and training parameters for the DMPNN.

Training Parameters	
Batch Size	50
Dropout	0.0
Learning Rate	1e-3
Max Learning Rate	1e-2
Number of Epochs	100
Architecture Parameters	
Activation	ReLU
Aggregation	Mean
Ensemble Size	1
Depth	3
FFN Hidden Size	300
FFN Depth	2
Split Sizes	0.8, 0.1, 0.1
Split Type	Random



**Figure S1.** Comparison of the true and predicted Log solubilities of the test set molecules from AqSolDB.

**Table S3.** Source, catalog number, and physical properties for molecules used in solubility experiments.

<b>Selected Molecules</b>	<b>Supplier</b>	<b>Catalog number</b>	<b>Molar Mass (g/mol)</b>	<b>Density (g/mL)</b>
Dodecane	Fisher	117590250	170.33	0.75
Methyl butyrate	Sigma Aldrich	246093	102.13	0.90
4-Chlorobutyric acid	Sigma Aldrich	C29835	122.55	1.24
2-Heptanone	Sigma Aldrich	537683	114.19	0.82
<i>N</i> -Butoxymethyl acrylamide	Fisher	B106025ML	157.21	0.98
Ethyl Formate	Sigma Aldrich	112682	74.08	0.92
Octamethyltrisiloxane	Sigma Aldrich	O025725ML	236.53	0.82
Pinacolone	Sigma Aldrich	P45605	100.16	0.80
Dimethyl Glutarate	Fisher	G018525G	160.17	1.09
<i>N</i> -Methylformamide	Fisher	F005925G	59.07	1.01
Toluene	Fisher	T290-1	92.14	0.87
4-Fluoroaniline	Fisher	F3800	111.12	1.17

### 3. Aqueous solubility measurement by experiments

For each selected molecules, aqueous solutions of varies concentrations are prepared to determine the solubility of the molecule. Table S3 summarize the solubility behavior and Figure 2 & 3 display the image of each solution in Table 3.

**Table S4.** Solubility behavior of the aqueous solution with the selected molecules at each concentration.

Selected Molecules	Target & Actual Concentration															
	0.20%		0.50%		1%		3%		7%		15%		34%		95%	
<b>Dodecane*</b>	0.18%	X	0.44%	X	0.92%	X	3.11%	X	6.89%	X	15.77%	X	33.60%	X	94.93%	X
Methyl butyrate	0.17%	O	0.54%	O	1.14%	O	3.60%	X	7.14%	X	15.18%	X	33.52%	X	94.95%	X
4-Chlorobutyric acid	0.20%	O	0.70%	O	1.37%	O	4.07%	O	6.15%	O	13.78%	O	32.20%	X	94.89%	O
<b>2-Heptanone*</b>	0.17%	O	0.47%	O	0.92%	X	2.91%	X	7.27%	X	14.85%	X	33.25%	X	94.98%	X
<i>N</i> -Butoxymethyl acrylamide	0.16%	O	0.48%	O	0.96%	O	3.00%	X	7.27%	X	13.68%	X	34.68%	X	95.17%	X
Ethyl Formate	0.18%	O	0.30%	O	0.84%	O	2.73%	O	7.08%	O	14.94%	X	34.82%	X	95.18%	O
<b>Octamethyl trisiloxane*</b>	0.19%	X	0.41%	X	0.97%	X	2.93%	X	6.91%	X	14.82%	X	33.04%	X	94.96%	X
Pinacolone	0.14%	O	0.43%	O	0.90%	X	2.89%	X	6.80%	X	15.05%	X	34.24%	X	95.04%	X
Dimethyl Glutarate	0.21%	O	0.46%	O	0.98%	O	2.92%	O	6.49%	X	14.34%	X	34.05%	X	94.90%	X
<i>N</i> -Methyl formamide	0.20%	O	0.53%	O	0.96%	O	2.89%	O	6.79%	O	14.64%	O	33.93%	O	94.93%	O
<b>Toluene*</b>	0.13%	O	0.43%	X	0.88%	X	2.92%	X	6.99%	X	13.07%	X	34.12%	X	94.99%	X
4-Fluoroaniline	0.18%	O	0.49%	O	0.98%	O	2.84%	X	6.68%	X	14.47%	X	32.70%	X	94.92%	X

Note: "X" represents insoluble, where the solution is cloudy or phase separation. "O" represents soluble, where the solution is clear.

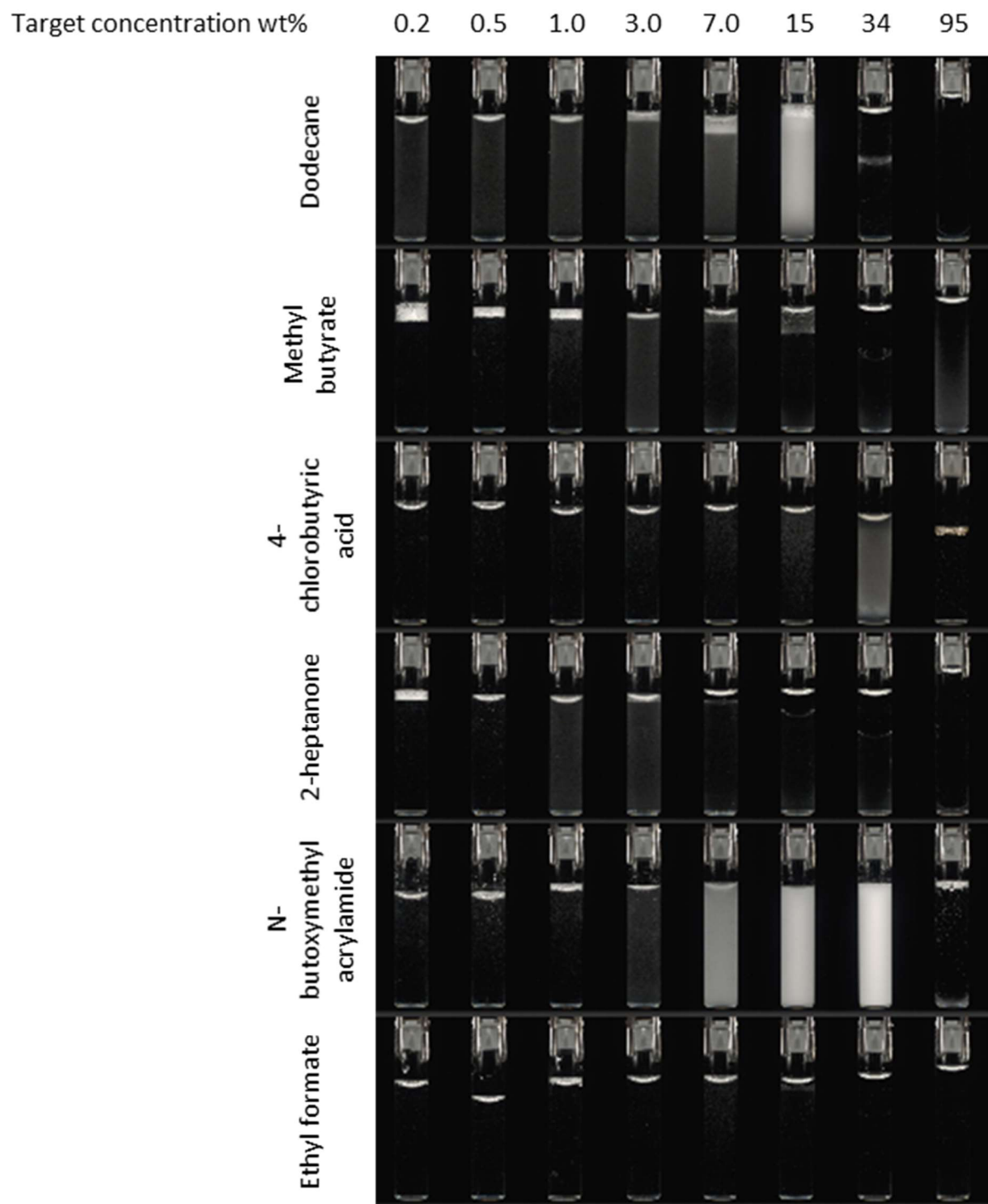
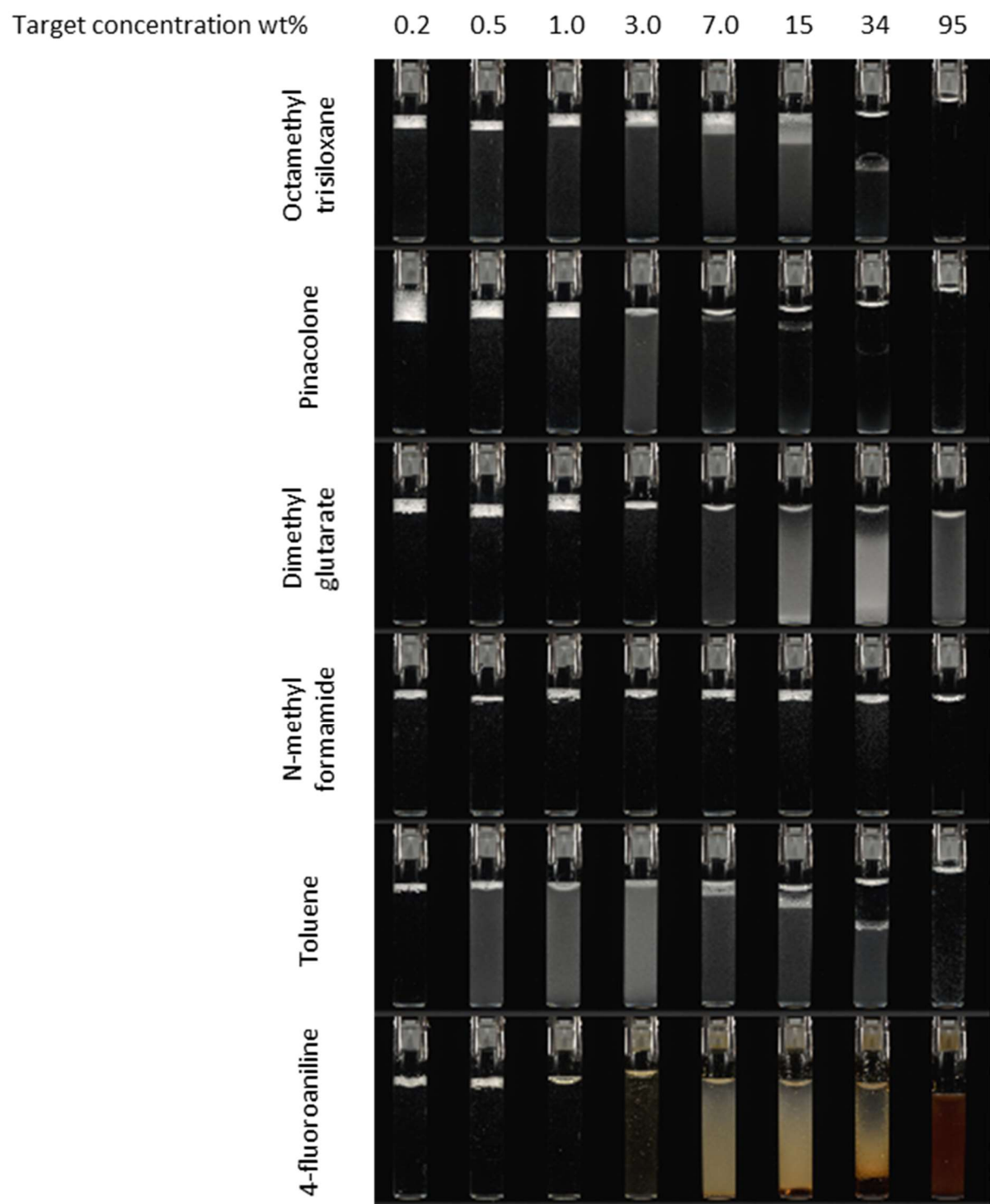


Figure S2. Images of first set of solvents in aqueous solutions.



**Figure S3.** Images of second set of solvents in aqueous solutions.





COc1ccc(NCP(=O)(O)O)cc1N(C)c1ccccc1	TRUE
O=C(NCCO)Nc1ccc(O)c2c1-c1ccccc1C2CO	FALSE
CC(=O)N(C)C1CC=C(NC(=O)CNC(=O)CN)CC1	FALSE
O=C(O)c1cc(O)nc(O)n1	TRUE
N=Cc1cc2cc(S(=O)(=O)O)ccc2c2ccccc12	FALSE
CNC(=O)Oc1ccccc1OC(C=O)CO	FALSE
NCC(O)COC(=O)c1cccnc1	TRUE
CNC(=O)C(C)(C(=O)NCC(=O)O)C1CCC(O)CC1	FALSE
CNCCC(O)(COC)c1ccc(O)c(O)c1	FALSE
CNCCNCCNCCNC(=O)Cn1c(=O)ccn(C)c1=O	TRUE
O=C(O)c1c(C(=O)N(CCO)c2ccc(O)c(O)c2)ccc(O)c1O	FALSE
Nc1ccc(S(=O)(=O)NC(N)N)cc1	FALSE
CNCC(=O)OCNC(=O)c1ccco1	TRUE
Cc1cnn(-c2cccc(S(=O)(=O)Nc3ccc(S(=O)(=O)O)nc3)c2)c1N	FALSE
CCNCC1CO1	TRUE
Cc1cc(O)c(O)c(Cc2cc(O)c(O)c(O)c2)c1	TRUE
CNCC(O)COC(=O)c1cnc2ccccc2n1	TRUE
CNCC(=O)Nc1ccc2c(c1)C(=O)NC2=O	TRUE
O=P(O)(O)CSCCO	FALSE
CCNC(=O)OCC(=O)NCC(=O)OCCOC(=O)CN	TRUE
O=S(=O)(O)NO	FALSE
COCNC(=O)c1ccc2cccnc2c1	TRUE
O=Cc1cc(O)n(CCO)n1	FALSE
O=C(Nc1ccc(S(=O)(=O)NC(CO)NCCO)cc1)Nc1ccccc(S(=O)(=O)O)c1	FALSE
CNC(=O)NCC(=O)OC	TRUE
CNC(=O)N(C)c1ccc(OC)cc1	TRUE
OCC(O)(COc1ccc(O)c(O)c1O)c1ccc(O)c(O)c1	FALSE
CCNc1cc(NCCO)n(CCO)n1	FALSE
NCCC(N)NC(=O)CCC(=O)O	FALSE
NC(=O)Sc1c(O)c(S(=O)(=O)O)cc2ncn(CO)c12	FALSE
O=C(NCCCCO)C(O)c1ccccc1	TRUE
CNCC(N)(C(=O)O)c1ccccc2cccnc12	FALSE
CCOCCN(CCO)CCOCCOCCOCCOCCOCCOCCOCCOCCOCCOCCOCCO	TRUE
Cc1c(S(=O)(=O)O)cc2ccc(S(=O)(=O)O)cc2c1N	FALSE
NC(NCC(=O)O)c1ccccc(S(=O)(=O)Nc2ccc(S(N)(=O)=O)nc2)c1	FALSE
CC(=O)OCN(CCO)C(=O)C1CCCCC(O)C1	FALSE
O=C(O)COC1CCCCC1OCCO	FALSE
CC(=O)CCc1ccn[nH]1	FALSE
O=S(=O)(O)c1nc2ccccc2s1	TRUE
CC(C)NCC(O)COCNC(=O)c1ccco1	FALSE
COc1cc([N+](=O)[O-])ccc1S(=O)(=O)O	TRUE
NNc1cc(O)c(NC(=O)c2ccc(O)cc2)cc1S(=O)(=O)O	FALSE
N=C(N)NCc1ccc(O)c2c1-c1ccc(S(=O)(=O)O)cc1CC2	FALSE









CNC(=O)C(O)Cn1cncc1OC	TRUE
CCC(C)(O)CC(N)C(N)=O	FALSE
NC(=O)C(N)(NCCCC(=O)O)c1ccnn1C(=O)[O-]	FALSE
NC(N)OC1=c2ccccc2=C(S(=O)(=O)O)CC1	FALSE
CCOP(=O)(OCC)OCOCN(C)CCOCCOCCOCCOCCOc1cnn(CO)c1	FALSE
O=[SH]CC(CO)(CC(=O)O)c1c(O)cc(C(=O)O)c2ccccc12	FALSE
OCNc1ncc2ncn(C(Cl)(CO)CO)c2n1	FALSE
O=S(=O)(O)C1=NC(S(=O)(=O)[O-])CN1CCO	FALSE
CN(CC(=O)O)C(N)Nc1ccc(O)c(NC(=O)CNC(=O)CN)c1N=Nc1ccc(O)cc1	FALSE
O=C(NCO)Nc1cc(N2CC2S(=O)(=O)O)cc(S(=O)(=O)O)c1	FALSE
O=[N+][[O-])C(O)CO	TRUE
CCNC1(N)CCC(S(=O)(=O)O)c2ccccc21	FALSE
NC(O)NCCO	FALSE
O=C(NO)NC1=CCC(S(=O)(=O)O)C1	FALSE
OCCNC(CO)NC1=CCC(CO)CC1	FALSE
N#CCNCN	FALSE
COC(=O)CN(C)C(=O)CNC(=O)N(C)CC(=O)OCCN(C)C(=O)CNC(=O)N(C)CC(=O)OCCNC(=O)C	
NC(=O)COC(=O)c1cc(O)cc2oc(=O)ccc12	FALSE
O=C(O)CC(O)(C(=O)O)C1C=CC(=O)c2ccccc21	FALSE
COC(=O)C(O)Cc1nc(O)nc(O)n1	FALSE
CNCCC(O)C1CCC(C(=O)OC)C1	TRUE
OCC(O)OCOC(O)(c1ccccc1)c1ccc(O)cc1	FALSE
CCC(O)C(N)C(C)(C)O	TRUE
NCCNC(N)NCc1ccc(S(=O)(=O)O)c2c1C(=O)OC2	FALSE
CC(C)NCC(O)COc1ccc(NC(=O)NCNC(=O)CN)c2nc(S(=O)(=O)O)sc12	FALSE
CP(=O)(O)c1cc2c(cc1O)C=C(S(=O)(=O)O)CC2	FALSE
O=[SH]Cc1nc2c(ncn2CO)n1CCO	FALSE
Cc1nc(O)nc(C(C)C(=O)N(C)C)n1	FALSE
CC(O)(CNC(=O)CNC(=O)CN)C(O)(c1cc(O)cc(O)c1)c1cc(O)cc(O)c1O	FALSE
O=C(NCO)N(CO)c1ncc[nH]1	FALSE
O=C(O)C(O)(C(=O)NCO)c1cc2scnc2cc1O	FALSE
COC(=O)C(NC(=O)C1=CCCC(O)C1)C(=O)O	FALSE
NC(=O)NC(N)(CO)NC(=O)C1CCC(O)CC1	FALSE
NCC(O)Cc1c[nH]ccc1=O	FALSE
O=C1CC(=O)N(O)CN1	FALSE
O=C(Br)N(O)C(CO)CO	FALSE
OCC(O)SC1CCC(O)CC1	FALSE
CNC(N)NN	FALSE
CNC(N)n1nccc1N	FALSE
CNCC(O)c1nc2cncnc2[nH]1	TRUE
CC(=N)C1CCCC(N)C1	FALSE
Cc1nc2ncn(CCO)c2n1CC(C)(O)CN	FALSE
NCn1nc(S(=O)(=O)O)cc1S(=O)(=O)O	FALSE







