

Supporting Information - Chemical Space Analysis and Property Prediction for Carbon Capture Amine molecules

James L. McDonagh, Stamatia Zavitsanou, Alexander Harrison, Dimitry Zubarev,
Theodore van Kessel, Benjamin H. Wunsch and Flaviu Cipcigan
flaviu.cipcigan@ibm.com, james.mcdonagh@serna.bio

December 29, 2023

1 Unit Conversions

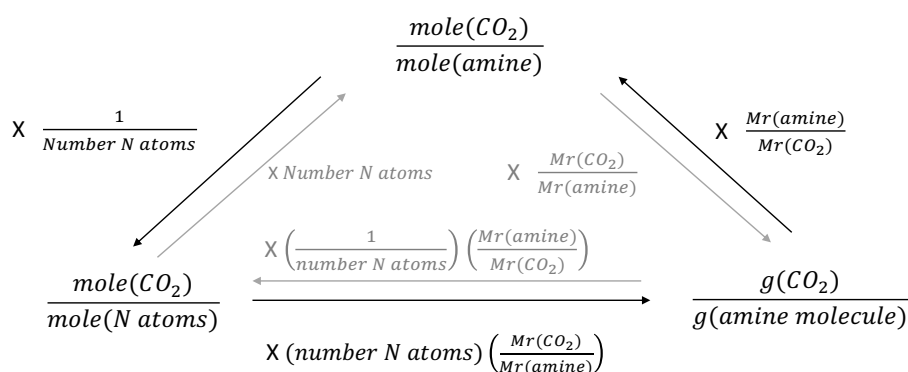


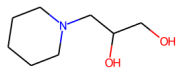
Figure 1: The conversion factors used in this work to move between common capacity measurement units.

Where units were given that need knowledge of the solution density for converting, such as $\frac{g(\text{CO}_2)}{L(\text{solution})}$, the following was used assuming the density (were not given) was that of water, and therefore has a density of 1 g/ml.

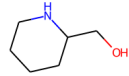
$$\text{capacity} \times \frac{1}{\text{amine weight fraction} \times 1000} \times \frac{\text{Mr}(\text{amine})}{\text{Mr}(\text{CO}_2)} = \frac{\text{mol}(\text{CO}_2)}{\text{mol}(\text{amine})} \quad (1)$$

Table 1: Experimental data and data set for 98 amine molecules tested for carbon capture

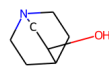
Index	IUPAC Name	InChIKey	Absorption $\left(\frac{\text{moles CO}_2}{\text{moles amine}}\right)$	capacity	Observed $\left(\frac{\text{moles CO}_2}{\text{moles amine} \cdot \text{sec}}\right)$	initial rate
68	piperidin-3-ol	BIWOSRSDCZIFM-UHFFFAOYSA-N	0.28		0.29	
69	2-(aminomethyl)aniline	GVQYKJPMUUIXBS-UHFFFAOYSA-N	0.28		0.17	
70	2-pyridin-4-ylethylamine	IDLHTECVNDEOY-UHFFFAOYSA-N	0.26		0.27	
71	1-amino-2-methylpropan-2-ol	LXQMHOKEXZETKB-UHFFFAOYSA-N	0.26		0.17	
72	2-[2-hydroxyethyl(methyl)amino]ethanol	CRVGTESFCXCCTH-UHFFFAOYSA-N	0.24		0.00	
73	1-[bis(2-hydroxyethyl)amino]propan-2-ol	ZFECCLNALLETDE-UHFFFAOYSA-N	0.22		0.00	
74	2-[bis(2-hydroxyethyl)amino]ethanol	GSEJCLTVZPLZKY-UHFFFAOYSA-N	0.20		0.05	
75	pyrazin-2-ylmethanamine	HQIBSDCOMQYSPF-UHFFFAOYSA-N	0.19		0.09	
76	2-(1H-imidazol-5-yl)ethanamine	NTYJJOPPIAHURM-UHFFFAOYSA-N	0.17		0.24	
77	2-[4-(2-hydroxyethyl)piperazin-1-yl]ethanol	VARKICWTYBUWNT-UHFFFAOYSA-N	0.11		0.03	
78	(2S)-2-amino-3-methyl-3-sulfanylbutanoic acid	VVNCNSJFMFPHL-VKHMVHEASA-N	0.09		0.00	
79	piperidin-1-amine	LWMPFIOTEAXAGV-UHFFFAOYSA-N	0.07		0.09	
80	1H-imidazol-5-ylmethanol	QDYTUZCWBVJRHKK-UHFFFAOYSA-N	0.07		0.00	
81	2-morpholin-4-ylethanol	KKPDCBRMNSAAW-UHFFFAOYSA-N	0.07		0.04	
82	2-[3-[[1,3-dihydroxy-2-(hydroxymethyl)propan-2-yl]amino]propylamino]-2-(hydroxymethyl)propane-1,3-diol	HHKZCCWKTZRCCCL-UHFFFAOYSA-N	0.06		0.15	
83	N-(2-hydroxyethyl)acetamide	PVCJKHHOXFKFRP-UHFFFAOYSA-N	0.05		0.23	
84	1-pyridin-3-yl-N-(pyridin-3-ylmethyl)methanamine	FEBQXMFOLRYSGC-UHFFFAOYSA-N	0.05		0.17	
85	tert-butyl N-(2-hydroxyethyl)carbamate	GPTXCAZYUMDUMN-UHFFFAOYSA-N	0.04		0.27	
86	(2S)-2,5-diamino-5-oxopentanoic acid	ZDXPRJPNDDTMRX-VKHMVHEASA-N	0.03		0.00	
87	pyridazin-3-amine	LETVJWLJLJADE-UHFFFAOYSA-N	0.02		0.38	
88	propanediamide	WRIRWRKPLXCTFD-UHFFFAOYSA-N	0.02		0.17	
89	N-(2-sulfanyylethyl)acetamide	AXFZADXLXIXTO-UHFFFAOYSA-N	0.02		0.00	
90	1,3-oxazolidin-2-one	IZXIZTKNFFYFOF-UHFFFAOYSA-N	0.02		0.07	
91	N,N-diethylhydroxylamine	FVCOIAYSJZGECG-UHFFFAOYSA-N	0.01		0.00	
92	1,4-diazabicyclo[2.2.2]octane	IMNIMPZHVRPE-UHFFFAOYSA-N	0.01		0.00	
93	pyridin-2-amine	ICSNLGPSSRYBMBD-UHFFFAOYSA-N	0.01		0.08	
94	1,3,5,7-tetraazatricyclo[3.3.1.1 ^{3,7}]decane	VKYKSIONXSXAKP-UHFFFAOYSA-N	0.01		0.08	
95	2-aminoethyl dihydrogen phosphate	SUHOOOKUPISOBE-UHFFFAOYSA-N	0.01		0.05	
96	pyrrol-1-amine	YNZAFFENDLJQG-UHFFFAOYSA-N	0.01		0.00	
97	pyrazin-2-amine	XFTQRUTUGRCSGO-UHFFFAOYSA-N	0.01		0.00	



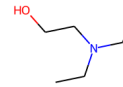
1) 3-piperidin-1-ylpropane-1,2-diol



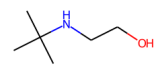
2) piperidin-2-ylmethanol



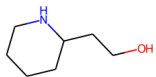
3) 1-azabicyclo[2.2.2]octan-3-ol



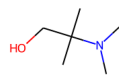
4) 2-(diethylamino)ethanol



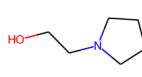
5) 2-(tert-butylamino)ethanol



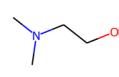
6) 2-piperidin-2-ylethanol



7) 2-(dimethylamino)-2-methylpropan-1-ol



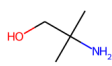
8) 2-pyrrolidin-1-ylethanol



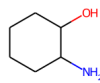
9) 2-(dimethylamino)ethanol



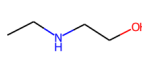
10) piperidine



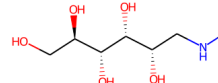
11) 2-amino-2-methylpropan-1-ol



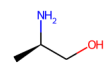
12) 2-aminocyclohexan-1-ol



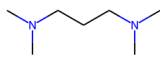
13) 2-(ethylamino)ethanol



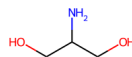
14) (2R,3R,4R,5S)-6-(methylamino)hexane-1,2,3,4,5-pentol



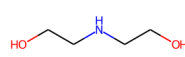
15) (2R)-2-aminopropan-1-ol



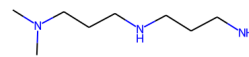
16) N,N,N,N-tetramethylpropane-1,3-diamine



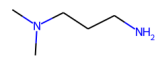
17) 2-aminopropane-1,3-diol



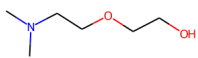
18) 2-(2-hydroxyethylamino)ethanol



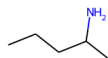
19) N-[3-(dimethylamino)propyl]propane-1,3-diamine



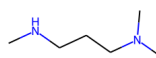
20) N,N-dimethylpropane-1,3-diamine



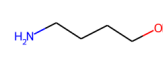
21) 2-[2-(dimethylamino)ethoxy]ethanol



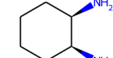
22) pentan-2-amine



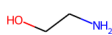
23) N,N,N-trimethylpropane-1,3-diamine



24) 4-aminobutan-1-ol



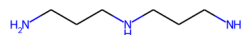
25) (1S,2R)-cyclohexane-1,2-diamine



26) 2-aminoethanol



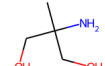
27) morpholine



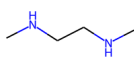
28) N-(3-aminopropyl)propane-1,3-diamine



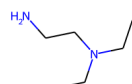
29) 1,3-diaminopropan-2-ol



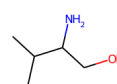
30) 2-amino-2-methylpropane-1,3-diol



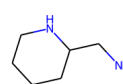
31) N,N-dimethylethane-1,2-diamine



32) N,N-diethylethane-1,2-diamine



33) 2-amino-3-methylbutan-1-ol



34) piperidin-2-ylmethanamine



35) azetidine



36) 1-aminobutan-2-ol



37) NN-diethylpiperidin-4-amine



38) butan-1-amine



39) pentan-1-amine



40) (2R)-3-aminopropane-1,2-diol



41) NN-diethylpropane-1,3-diamine



42) (2R)-2-amino-4-methylpentan-1-ol



43) NN-dimethylpropane-1,3-diamine



44) 1-methyl-1,4-diazepane



45) 2-(methylamino)ethanol



46) piperazine



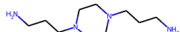
47) pyrrolidin-3-amine



48) 8-methyl-8-azabicyclo[3.2.1]octan-3-amine



49) NN-dimethylpyrrolidin-3-amine



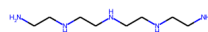
50) 3-[4-(3-aminopropyl)piperazin-1-yl]propan-1-amine



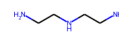
51) NN-dimethylpiperidin-4-amine



52) 1-(2-hydroxypropylamino)propan-2-ol



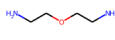
53) N-(2-[(2-(2-aminoethylamino)ethylamino)ethyl]ethane-1,2-diamine



54) N-(2-aminoethyl)ethane-1,2-diamine



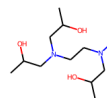
55) (1R,2R)-cyclohexane-1,2-diamine



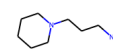
56) 2-(2-aminoethoxy)ethanamine



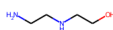
57) 2-aminobutan-1-ol



58) 1-(2-[(2-(2-hydroxypropylamino)ethyl-(2-hydroxypropylamino)propan-2-ol



59) 3-piperidin-1-ylpropan-1-amine



60) 2-(2-aminoethylamino)ethanol



61) (2R)-1-aminopropan-2-ol



62) 2,3,4,6,7,8-hexahydro-1,2-a]pyrimidine



63) 3-aminocyclohexan-1-ol



64) 2,3,4,6,7,8,9,10-octahydro-1,2-a]azepine



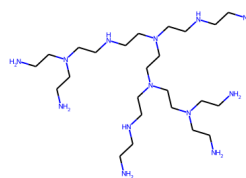
65) 3,4,6,7,8,9-hexahydro-2H-pyrimido[1,2-a]pyrimidine



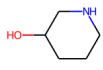
66) [(2S)-pyrrolidin-2-yl]methanol



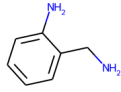
67) (2S)-1-aminopropan-2-ol



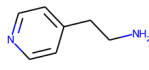
68) 1-(2-[(2-(2-aminoethylamino)ethylamino)ethyl]ethane-1,2-diamine



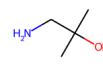
69) piperidin-3-ol



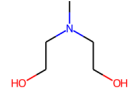
70) 2-(aminomethyl)aniline



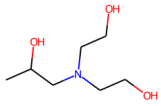
71) 2-pyridin-4-ylethanamine



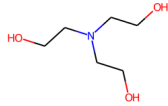
72) 1-amino-2-methylpropan-2-ol



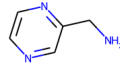
73) 2-[2-(2-hydroxyethyl)(methylamino)ethyl]ethanol



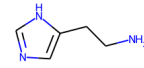
74) 1-[bis(2-hydroxyethyl)amino]propan-2-ol



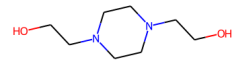
75) 2-[bis(2-hydroxyethyl)amino]ethanol



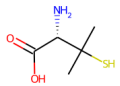
76) pyrazin-2-ylmethanamine



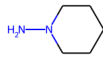
77) 2-(1H-imidazol-5-yl)ethanamine



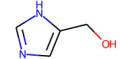
78) 2-[4-(2-hydroxyethyl)piperazin-1-yl]ethanol



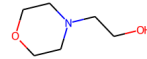
79) (2S)-2-amino-3-methyl-3-sulfanybutanoic acid



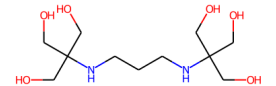
80) piperidin-1-amine



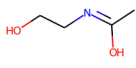
81) 1H-imidazol-5-ylmethanol



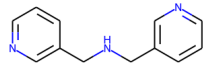
82) 2-morpholin-4-ylethanol



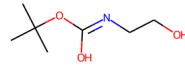
83) 2-[2-(1,3-dihydroxy-2-hydroxymethyl)propen-2-ylamino]propanoic acid



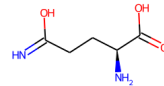
84) N-(2-hydroxyethyl)acetamide



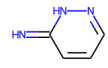
85) 1-pyridin-3-yl-N-(pyridin-3-ylmethyl)methanamine



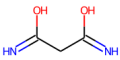
86) tert-butyl N-(2-hydroxyethyl)carbamate



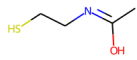
87) (2S)-2,5-diamino-5-oxopentanoic acid



88) pyridazin-3-amine



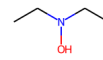
89) propanediamide



90) N-(2-sulfanyethyl)acetamide



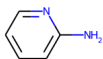
91) 1,3-oxazolidin-2-one



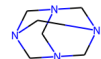
92) NN-diethylhydroxylamine



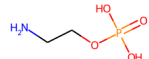
93) 1,4-diazabicyclo[2.2.2]octane



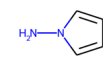
94) pyridin-2-amine



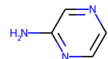
95) 1,3,5,7-tetraazatricyclo[3.3.1.1]decane



96) 2-aminoethyl hydrogen phosphate



97) pyrrol-1-amine



98) pyrazin-2-amine

3 Metrics for classification

The false positive rate (false alarm rate) is calculated as the number of false positives divided by the sum of the number of false positives and the number of true negatives, (Eq.2) and it summarizes how often a positive class is predicted when the actual outcome is negative.

$$\text{False Positive Rate} = \frac{\text{False Positives}}{(\text{False Positives} + \text{True Negatives})} \quad (2)$$

The ROC plot is a useful tool because the curves of different models can be compared directly or indirectly for different thresholds. Also, the area under the curve (AUC) can be used as a summary of the model skill. The shape of the curve contains information about the expected false positive rate, and the false negative rate. Smaller values on the x-axis of the plot indicate lower false positives and higher true negatives, and larger values on the y-axis of the plot indicate higher true positives and lower false negatives. Generally, a model is considered predictive when the curves are shown on the top left of the plot. A perfect model is represented by a line that travels from the bottom left of the plot to the top left and then across the top to the top right and has an AUC of 1. A less predictive classifier is one that cannot discriminate between the classes and would predict a random class or a constant class in all cases. Such a model is represented by a diagonal line from the bottom left of the plot to the top right and has an AUC of 0.5.

Sensitivity is calculated as the ratio of the number of true positives divided by the sum of the true positives and the false negatives (Eq.3), and specificity is calculated as the ratio of the number of true negatives divided by the sum of true negatives and false positives (Eq.4).

$$\text{Sensitivity} = \frac{\text{True Positives}}{(\text{True Positives} + \text{False Negatives})} \quad (3)$$

$$\text{Specificity} = \frac{\text{True Negatives}}{(\text{True Negatives} + \text{False Positives})} \quad (4)$$

Matthews Correlation Coefficient shows the correlation between observed and predicted values and is calculated from Eq.5. It returns a value between -1 and 1. A coefficient of 1 represents a perfect prediction and 0 no better than random prediction. Everything above 0.5 is considered a good correlation.

$$MCC = \frac{(TP \times TN - FP \times FN)}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (5)$$

Lastly, geometric mean (g-mean) is the geometric mean of the sensitivity and specificity,¹ this provides a measure of the combined performance of a model for both sensitivity and specificity (Eq.6).

$$\text{g-mean} = \sqrt{\text{sensitivity} \times \text{specificity}} \quad (6)$$

Table 2: Substructure in the ecs fingerprint

Indexes	Search
69	[#16]
70	[+]
71	[-]

5 Occurrence of fingerprint bits in datasets (fingerprint version 1)

5.1 ZINC dataset

Table 3: Fingerprint SMARTS bit number of occurrence over the ZINC dataset.

Bit	Occurrences
ammonia	0
primary amine	1666
secondary amine	3061
tertiary amine	1908
quaternary N	3
imine	128
nitrogen bonded to carbon	11725
aromatic N sp2	7713
carboxylic acid	3131
primary alcohol	1424
secondary alcohol	756
tertiary alcohol	58
t butyl	85
carbonyl	4465
halocarbon	5095
benezene ring	5827
6 member aromatic c and n ring	9410
6 member c and o ring	272
5 c ring	121
5 member aromatic c and n ring	2434
5 member c and o ring	272
Cyclohexane	242
Cyclohexylamine	472
Aniline	164
benzylamine	99
piperidine	839
pyridine	2880
pyrrole	339
primary amino alcohol two carbon separation	58
secondary amino alcohol two carbon separation	100
tertiary amino alcohol two carbon separation	486
primary amino alcohol three carbon separation	37
secondary amino alcohol three carbon separation	58
tertiary amino alcohol three carbon separation	295
aliphatic primary amino alcohol two carbon separation	53
aliphatic secondary amino alcohol two carbon separation	75
aliphatic tertiary amino alcohol two carbon separation	310
aliphatic primary amino alcohol three carbon separation	22
aliphatic secondary amino alcohol three carbon separation	28

Continued on next page

Table 3: Fingerprint SMARTS bit number of occurrence over the ZINC dataset.

Bit	Occurrences
aliphatic tertiary amino alcohol three carbon separation	83
primary amine one carbon aromatic group	429
primary amine two carbon aromatic group	263
primary amine three carbon aromatic group	33
secondary amine one carbon aromatic group	1051
secondary amine two carbon aromatic group	249
secondary amine three carbon aromatic group	88
tertiary amine one carbon aromatic group	1578
tertiary amine two carbon aromatic group	119
tertiary amine three carbon aromatic group	40
methyl branch one carbon from a N atom	121
methyl branch two carbon from a N atom	82
methyl branch three carbon from a N atom	24
methyl branch four carbon from a N atom	11
methyl branch five carbon from a N atom	0
methyl branch six carbon from a N atom	2
ethyl chain	5605
propyl chain	2251
butyl chain	680
pentyl chain	200
hexyl chain	66
poly primary and or secondary amine	362
poly primary and or secondary and or tertiary amine	1154
poly alcohol	166
pyrazine aliphatic C 2 and 5 substitution	5
pyridine aliphatic C 2 and 5 substitution	40
pyridine aliphatic C 2 substitution	478
Presence of Boron	136
Presence of Silicon	26
Presence of Phosphorus	10
Presence of Sulphur	1708
positive charge group	805
negative charge group	562

5.2 CCS literature datasets

Table 4: Fingerprint SMARTS bit number of occurrence over the ZINC dataset.

Bit	Occurrences
ammonia	0
primary amine	60

Continued on next page

Table 4: Fingerprint SMARTS bit number of occurrences over the ZINC dataset

Bit	Occurrences
secondary amine	45
tertiary amine	41
quaternary N	0
imine	0
nitrogen bonded to carbon	128
aromatic N sp2	18
carboxylic acid	5
primary alcohol	60
secondary alcohol	14
tertiary alcohol	0
t butyl	0
carbonyl	5
halocarbon	0
benezene ring	0
6 member aromatic c and n ring	17
6 member c and o ring	1
5 c ring	1
5 member aromatic c and n ring	1
5 member c and o ring	1
Cyclohexane	4
Cyclohexylamine	4
Aniline	0
benzylamine	0
piperidine	6
pyridine	13
pyrrole	0
primary amino alcohol two carbon separation	9
secondary amino alcohol two carbon separation	14
tertiary amino alcohol two carbon separation	24
primary amino alcohol three carbon separation	2
secondary amino alcohol three carbon separation	2
tertiary amino alcohol three carbon separation	10
aliphatic primary amino alcohol two carbon separation	9
aliphatic secondary amino alcohol two carbon separation	14
aliphatic tertiary amino alcohol two carbon separation	24
aliphatic primary amino alcohol three carbon separation	1
aliphatic secondary amino alcohol three carbon separation	2
aliphatic tertiary amino alcohol three carbon separation	6
primary amine one carbon aromatic group	3
primary amine two carbon aromatic group	0
primary amine three carbon aromatic group	1
secondary amine one carbon aromatic group	6
secondary amine two carbon aromatic group	0

Continued on next page

Table 4: Fingerprint SMARTS bit number of occurrences over the ZINC dataset

Bit	Occurrences
secondary amine three carbon aromatic group	0
tertiary amine one carbon aromatic group	0
tertiary amine two carbon aromatic group	0
tertiary amine three carbon aromatic group	0
methyl branch one carbon from a N atom	6
methyl branch two carbon from a N atom	6
methyl branch three carbon from a N atom	1
methyl branch four carbon from a N atom	1
methyl branch five carbon from a N atom	2
methyl branch six carbon from a N atom	0
ethyl chain	89
propyl chain	44
butyl chain	24
pentyl chain	11
hexyl chain	2
poly primary and or secondary amine	36
poly primary and or secondary and or tertiary amine	49
poly alcohol	23
pyrazine aliphatic C 2 and 5 substitution	1
pyridine aliphatic C 2 and 5 substitution	1
pyridine aliphatic C 2 substitution	4
Presence of Boron	0
Presence of Silicon	0
Presence of Phosphorus	0
Presence of Sulphur	0
positive charge group	0
negative charge group	0

6 Mordred features

We find 35 correlating features which have a significant p-value at 95%, shown in table 5. Following feature generation, we apply one-hot encoding for categorical features and min-max scaling for continuous features. There were 6 features (nBondsM, nBondsKD, C1SP2, HybRatio, FCSP3, ETA_beta_ns) considered as categorical out of the 35. Categorical in this case includes features with specific increments such as counts. Following one hot encoding the feature set extends to 84 as every unique value of the categorical features becomes a binary feature array.

Table 5: Features used as descriptors in the Mordred models based on Spearman correlation of ≥ 0.5 cutoff and significance at 95% on a two tail p-test of the Spearman correlation.

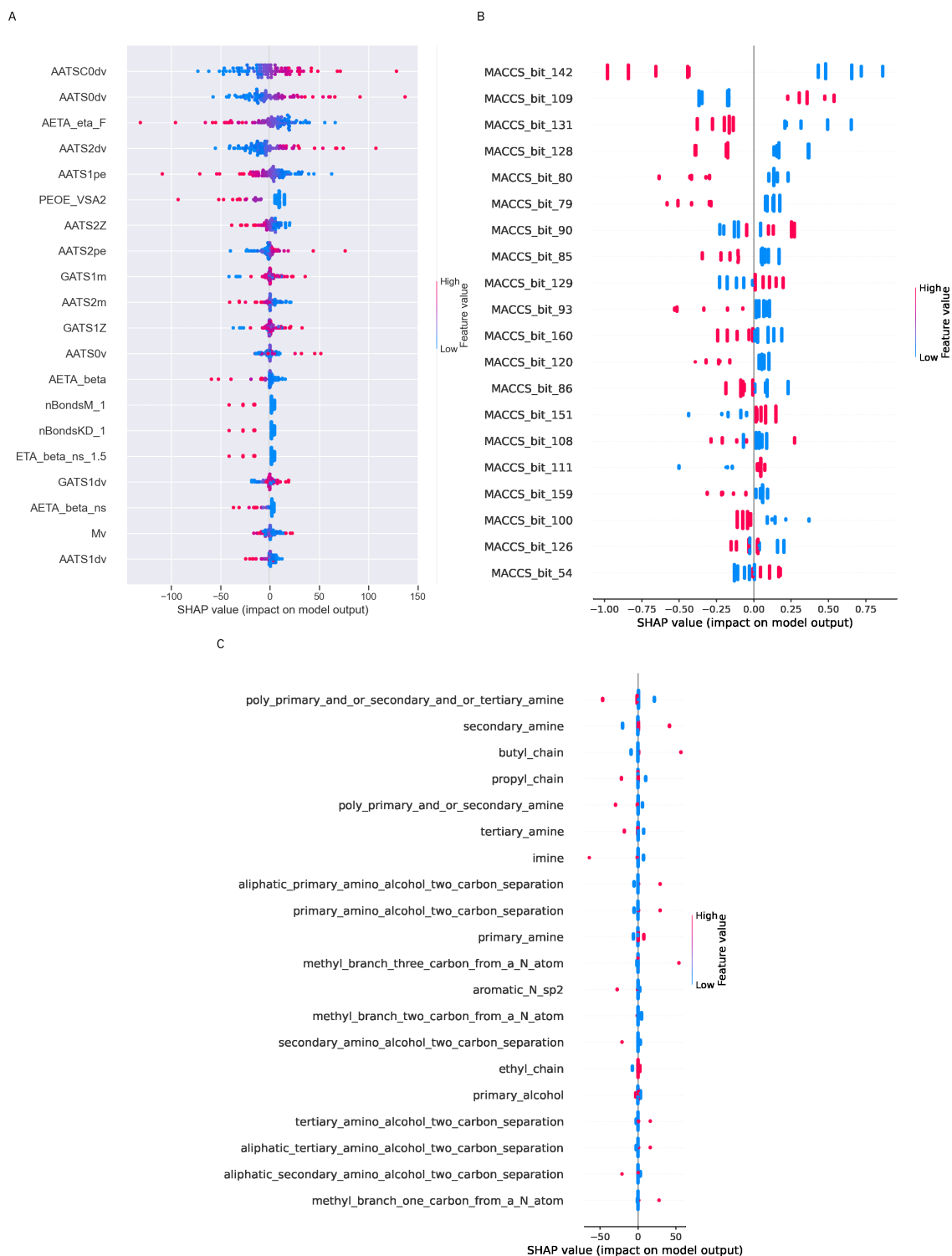
Index	Descriptor name
0	ATS2dv
1	AATS0dv
2	AATS1dv
3	AATS2dv
4	AATS0Z
5	AATS1Z
6	AATS2Z
7	AATS0m
8	AATS1m
9	AATS2m
10	AATS0v
11	AATS1pe
12	AATS2pe
13	AATSC0dv
14	GATS1dv
15	GATS1Z
16	GATS1m
17	BCUTdv-1h
18	nBondsM
19	nBondsKD
20	C1SP2
21	HybRatio
22	FCSP3
23	AXp-2dv
24	MZ
25	Mm
26	Mv
27	AETA_beta
28	ETA_beta_ns
29	AETA_beta_ns
30	AETA_eta_F
31	ETA_epsilon_1
32	PEOE_VSA2
33	SlogP_VSA6
34	AMW

7 Extra Trees Accuracy

Table 6: This table shows the headline accuracy results for using the Extra trees classifier. We see that this classifier is consistently the worst performer in 2 of three feature sets hence not included in the main analysis where we have focused on the best performing models.

Features	Accuracy
Mordred	0.72
MACCS	0.79
CCS	0.80

8 SHAP



References

- [1] Q. Zang, D. M. Rotroff and R. S. Judson, *Journal of chemical information and modeling*, 2013, **53**, 3244–3261.