

Supplementary Material for:

## Predicting $pK_a$ Using a Combination of Semi-Empirical Quantum Mechanics and Radial Basis Function Methods

Peter Hunt<sup>1</sup>, Layla Hosseini-Gerami<sup>2</sup>, Tomas Chrien<sup>1</sup>, Jeffrey Plante<sup>3</sup>, David J. Ponting<sup>3</sup>, Matthew Segall<sup>1</sup>

1 Optibrium Ltd, F5-6 Blenheim House, Cambridge Innovation Park, Denny End Road, Cambridge, CB25 9PB, UK  
Tel. +44 (0) 1223 815 900

2 Department of Chemistry, Lensfield Road, Cambridge, CB2 1EW, UK Tel. +44 (0) 1223 336 300

3 Lhasa Limited, Granary Wharf House, 2 Canal Wharf, Leeds, LS11 5PS, UK Tel. +44 (0)113 394 6020

Figure S1

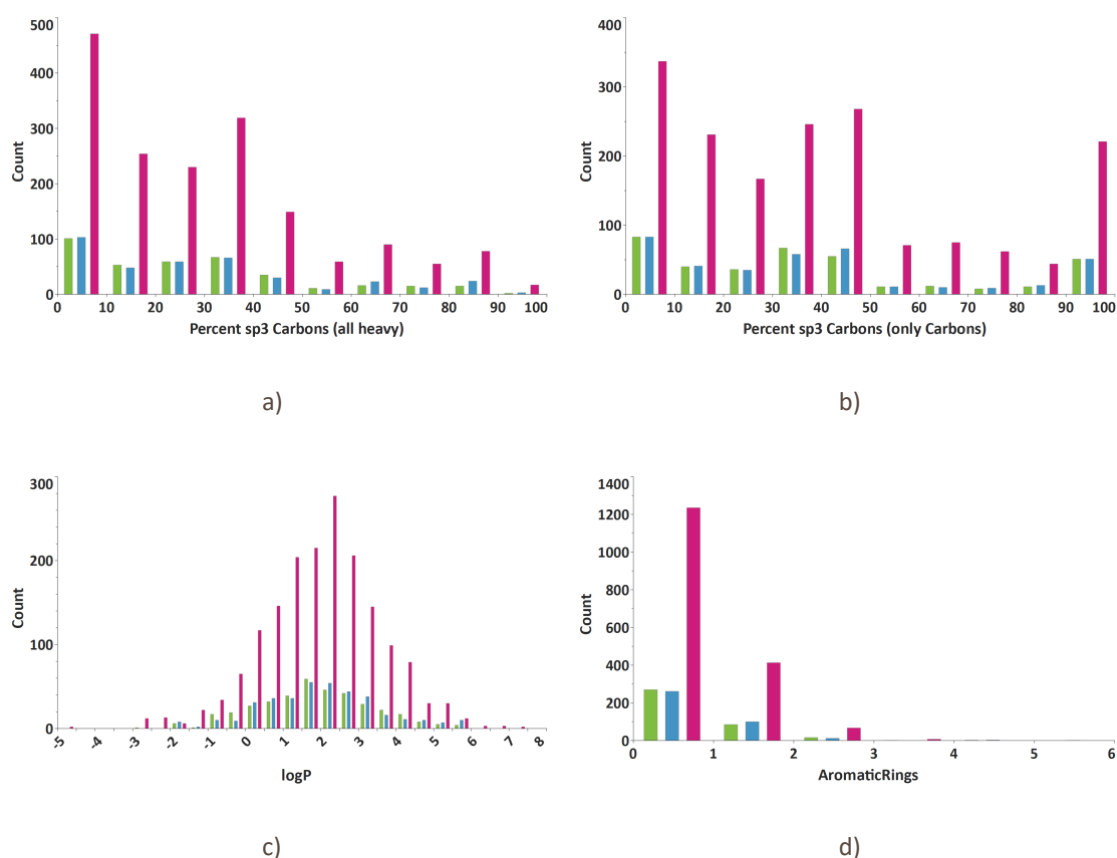


Figure S1. Plots of the distributions of simple properties across the training (pink bars), validation (blue bars) and test (green bars) sets. a) percentage of  $sp^3$  carbon atoms taking all heavy atoms into account, b) percentage of  $sp^3$  carbon atoms taking only carbon atoms into account, c) calculated  $\log P$ , d) Number of aromatic rings.

Figure S2

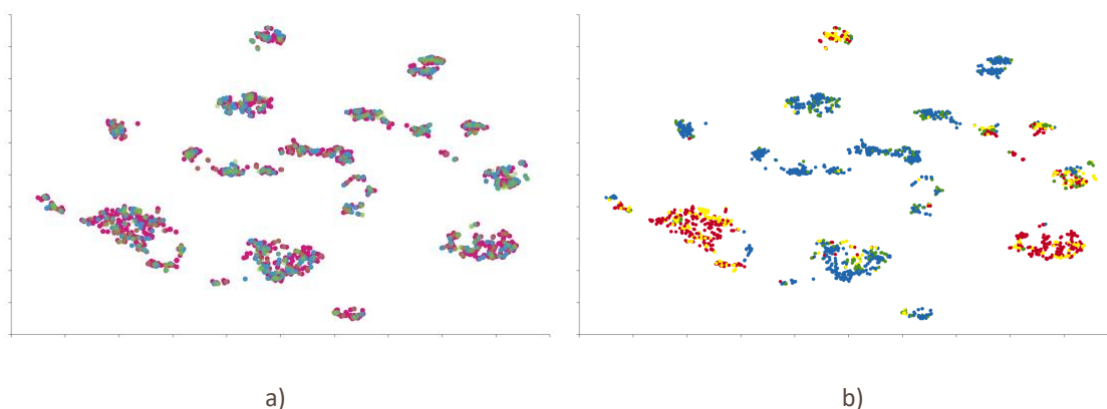


Figure S2. A pair of property space representations derived from the descriptor matrix used in the model generation. Plots a) and b) show the whole  $pK_a$  data set and the points are coloured in a) by their set membership either training (pink), validation (blue) and test (green) and in b) by the experimental  $pK_a$  values associated with the ionisable site(s) within the molecule  $pK_a \leq 4$  (red),  $pK_a > 4$  and  $\leq 6$  (yellow),  $pK_a > 6$  and  $\leq 8$  (green), and  $pK_a > 8$  (blue).

Table T1

Torsion of central bond should be 155 degrees:	[#7,#8,#16,#15,F,Cl,Br,I]-=[#6,#7,#8,#16,#15]-!@[#6,#7,#8,#16,#15]-,[#7,#8,#16,#15,F,Cl,Br,I]
Torsion should be 25 degrees:	[!#1]~[!#1;!#6]-[#6,#7](:[!#1]):[!#1]
	[#1]-[#7H2]-[#6,#7](:[!#1]):[!#1]
	[#1]-[#8H]-[#6,#7](:[!#1]):[!#1]
	[#1]-[#16H]-[#6,#7](:[!#1]):[!#1]
	[#7,#6,#8,#16]=[#6]-[#6,#7](:[!#1]):[!#1]

Table T1. SMARTS patterns used to assign the torsion angles shown prior to geometry optimisation with AM1 in order to minimise failed descriptor calculations due to proton transfer during optimisation.

Figure S3

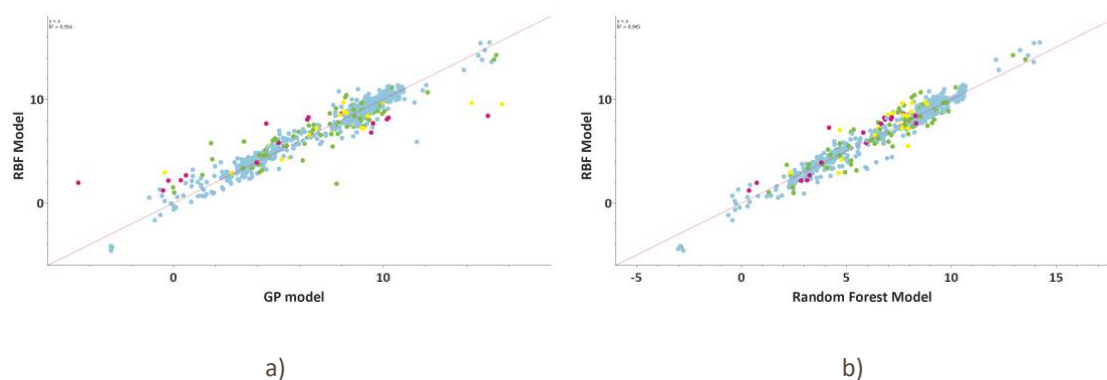


Figure S3. Plots of the correlation between predictions for the test and validation sets by different methods. The points are coloured by the absolute deviation from the experimental value for the RBF model prediction

(deviation  $\leq 1$  log unit blue,  $1 < \text{deviation} \leq 2$  log units green,  $2 < \text{deviation} \leq 3$  log units yellow,  $>3$  log units pink). The models are generally well correlated but the poorly predicted sites in the RBF model are not always poorly predicted by the models from the other methods.

Figure S4

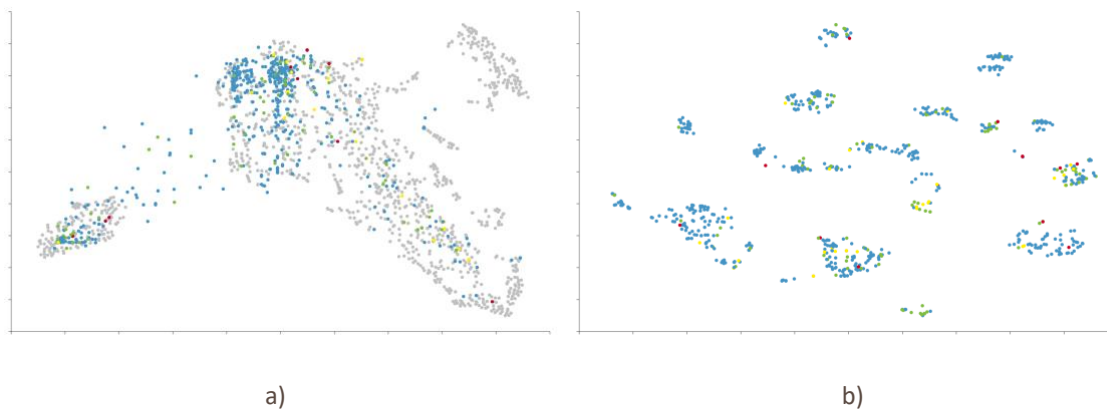


Figure S4 A set of property space representations of the data set. In a) the test and validation set compounds are displayed in a launched drug chemical space along with the approximately 1,300 launched small molecule drugs that defined the space, whilst in b) only the test set and validation set compounds are displayed in a property space derived from the descriptor matrix used in the model generation. The colour scheme is consistent across both spaces with the grey points being the launched drugs whilst the  $pK_a$  data set points are coloured by the absolute error of the RBF model prediction compared to the experimental  $pK_a$  values associated with the ionisable site(s) within the molecule (deviation  $\leq 1$  log unit blue,  $1 < \text{deviation} \leq 2$  log units green,  $2 < \text{deviation} \leq 3$  log units yellow,  $>3$  log units pink). The plots show that there are no consistently poorly predicted areas in either chemical or descriptor property space.