

# **1 Probing the origins of human 2 acetylcholinesterase inhibition via QSAR 3 modeling and molecular docking**

**4 Saw Simeon<sup>1</sup>, Nuttapat Anuwongcharoen<sup>1</sup>, Watshara Shoombuatong<sup>1</sup>,  
5 Aijaz Ahmad Malik<sup>1</sup>, Virapong Prachayasittikul<sup>2</sup>, Jarl E. S. Wikberg<sup>3</sup>, and  
6 Chanin Nantasenamat\*<sup>1</sup>**

**7 <sup>1</sup>Center of Data Mining and Biomedical Informatics, Faculty of Medical Technology,  
8 Mahidol University, Bangkok 10700, Thailand**

**9 <sup>2</sup>Department of Clinical Microbiology and Applied Technology, Faculty of Medical  
10 Technology, Mahidol University, Bangkok 10700, Thailand**

**11 <sup>3</sup>Department of Pharmaceutical Biosciences, Uppsala University, Uppsala 751 24,  
12 Sweden**

## **13 ABSTRACT**

Alzheimer's disease (AD) is a chronic neurodegenerative disease which leads to the gradual loss of neuronal cells. Several hypotheses for AD exists (e.g. cholinergic, amyloid, tau hypotheses, etc.). As per the cholinergic hypothesis, the deficiency of choline is responsible for AD therefore the inhibition of AChE is a lucrative therapeutic strategy for the treatment of AD. Acetylcholinesterase (AChE) is an enzyme that catalyzes the breakdown of the neurotransmitter acetylcholine that is essential for cognition and memory. A large non-redundant data set of 2,570 compounds with reported IC<sub>50</sub> values against AChE was obtained from ChEMBL and employed in quantitative structure-activity relationship (QSAR) study so as to gain insights on their origin of bioactivity. AChE inhibitors were described by a set of 12 fingerprint descriptors and predictive models were constructed from 100 different data splits using random forest. Generated models afforded  $R^2$ ,  $Q_{CV}^2$  and  $Q_{Ext}^2$  values in ranges of 0.66–0.93, 0.55–0.79 and 0.56–0.81 for the training set, 10-fold cross-validated set and external set, respectively. The best model built using the substructure count was selected according to the OECD guidelines and it afforded  $R^2$ ,  $Q_{CV}^2$  and  $Q_{Ext}^2$  values of  $0.92 \pm 0.01$ ,  $0.78 \pm 0.06$  and  $0.78 \pm 0.05$ , respectively. Furthermore, Y-scrambling was applied to evaluate the possibility of chance correlation of the predictive model. Subsequently, a thorough analysis of the substructure fingerprint count was conducted to provide informative insights on the inhibitory activity of AChE inhibitors. Moreover, Kennard-Stone sampling of the actives were applied to select 30 diverse compounds for further molecular docking studies in order to gain structural insights on the origin of AChE inhibition. Site-moiety mapping of compounds from the diversity set revealed three binding anchors encompassing both hydrogen-bonding and van der Waals interaction. Molecular docking revealed that compound **13** exhibited the lowest binding energy of -12.2 kcal/mol against human AChE, which is modulated by hydrogen bonding,  $\pi$ - $\pi$  stacking and hydrophobic interaction inside the binding pocket. These information may be used as guidelines for the design of novel and robust AChE inhibitors.

**15 Keywords:** Acetylcholinesterase, Acetylcholinesterase inhibitor, Alzheimer's disease, Quantitative structure-activity relationship, QSAR, Data mining

## **16 INTRODUCTION**

**17 Neurodegenerative diseases is caused by** the progressive loss of neural cells **thereby** leading to nervous  
18 system dysfunction (Beal, 1995; Kuca et al., 2016). In particular, Alzheimer's disease (AD) is a debilitating  
19 illness that is expected to triple by the year 2050 (Brookmeyer et al., 2007). AD is characterized by gradual  
20 cognitive impairment, memory loss and decline in speech, behavioral abnormality and **eventually** death.  
21 The pathological changes in the AD are mainly because of the dramatic loss of neurons in many areas of

\*Corresponding author. E-mail: chanin.nan@mahidol.ac.th

22 the central nervous system accompanied by a great reduction in the levels of neurotransmitters. Thus, a  
23 promising therapeutic approach is to maintain the level of the acetylcholine (ACh) neurotransmitter by  
24 inhibiting the enzyme that is responsible for its breakdown.

25 Acetylcholinesterases (AChE) are mainly found in the joints of neuromuscular junction in which it  
26 not only serves as a synaptic transmission but also as an enzyme that catalyzes the breakdown of ACh to  
27 choline and acetic acid (Quinn, 1987). The structure of AChE showed that it is comprised of two active  
28 sites consisting of anionic and esteratic sites (Bourne et al., 1995). The anionic site is involved in the  
29 binding of the positive quaternary amine of ACh (Ordentlich et al., 1993). The substrate interacts with the  
30 14 aromatic residues that forms the active site. Of these 14 aromatic residues, Trp84 is important for the  
31 enzyme activity because when it is replaced by alanine, the activity of the enzyme decreased by 3000-fold  
32 (Tougu, 2001). In contrast, the esteratic site hydrolyzed ACh to acetate and choline in which the existence  
33 of a catalytic triad (i.e. Ser200, His440 and Glu327) resembled that of chymotrypsin and other serine  
34 proteases (Harel et al., 1993). The mechanism of the hydrolysis starts from the carboxyl ester leads to  
35 the formation of an acyl-enzyme and choline. Finally, the acyl-enzyme undergoes nucleophilic attack by  
36 water molecules thereby regenerating the enzyme (Tougu, 2001).

37 AChE inhibitors form one of the most actively investigated classes of compounds having been labeled  
38 as a potential agent for the treatment of AD by inhibiting AChE from hydrolyzing ACh as to increase the  
39 level of ACh and maintain homeostasis (Birks, 2006). Generally, they can be classified into reversible and  
40 irreversible inhibitors. Reversible inhibitor bind to the enzyme at allosteric sites as to reduce the activity  
41 of the enzyme whether or not the enzyme has already bind the substrate or not. For example, tacrine was  
42 synthesized nearly five decades ago and in 1993 it has become the first drug to be marketed as a form of  
43 treatment for AD with approval from the U.S. Food and Drug Administration (Racchi et al., 2004). On  
44 the other hand, irreversible inhibitors such as metrifonate (Morris et al., 1998) bind to the target enzymes  
45 and dissociates very slowly from the enzyme via either covalent or non-covalent interactions (Kitz and  
46 Wilson, 1962).

47 Quantitative structure-activity relationship (QSAR) is a paradigm that enables the prediction of  
48 biological activities for compounds of interest as a function of their descriptors through the use of  
49 statistical or machine learning methods (Nantaseamat et al., 2009). Aside from the ability to predict the  
50 activity, QSAR models have been instrumental in enabling understanding on the origin of these biological  
51 activities by means of interpreting the descriptors used in building such models.

52 Historically, the first QSAR investigation of AChE inhibitors was reported by Mundy et al. (1978)  
53 almost 40 years ago in which the log(1/LD<sub>50</sub>) for a series of twelve substituted 0,0-dimethyl 0-(*p*-  
54 nitrophenyl) phosphorothioates and 0-analogs was predicted as a function of the octanol/water partition  
55 coefficient. Analysis of the literature of QSAR studies of AChE revealed that much of the early studies  
56 are classical QSAR models (i.e. Hansch and Free-Wilson approach) that are based on small congeneric  
57 compound set and primarily aimed at predicting AChE inhibition as to investigate the toxic effect  
58 of pesticides of various chemotypes belonging to either organophosphates (Mager, 1983; Aaviksaar,  
59 1990) or carbamates (Su and Lien, 1980; Goldblum et al., 1981; Walters and Hopfinger, 1986). Recent  
60 QSAR studies are based on the use of large and heterogeneous data sets comprising of structurally  
61 diverse chemotypes. This include the study from Yan and Wang (2012) where they predicted AChE  
62 inhibition for a large set of 404 compounds using multiple linear regression and support vector machine.  
63 Furthermore, Lee and Barron (2016) performed a 3D-QSAR investigation on a large set of 341 compounds  
64 comprising of organophosphates and carbamates. Moreover, Veselinović et al. (2015) compiled a set of  
65 278 organophosphates for which they developed QSAR models for predicting AChE inhibition using  
66 SMILES-based descriptors.

67 Research in this field had experienced two distinct transitions when viewed from biological and  
68 computational viewpoints. Biologically, early QSAR studies treat AChE as a biomarker of toxicity from  
69 pesticides while investigations from later years had shifted the focus by viewing AChE as a therapeutic  
70 target for the treatment of AD. In regards to the latter point, viewpoint on targeting AChE as a single  
71 target for treating AD is starting to be replaced by the multi-target concept in which the treatment for AD  
72 can be approached by a panel of key targets (Fang et al., 2015; Huang et al., 2011). Computationally,  
73 early studies are predominantly based on simple 2D-QSAR (Mundy et al., 1978; Su and Lien, 1980) while  
74 later years started to use more sophisticated approach for understanding AChE inhibition encompassing  
75 3D-QSAR (Deb et al., 2012; Lee and Barron, 2016; Prado-Prado et al., 2012), molecular dynamics (Shen  
76 et al., 2002), molecular docking (Lu et al., 2011; Deb et al., 2012; Giacoppo et al., 2015), pharmacophore

77 modeling Lu et al. (2011); Gupta and Mohan (2014) and statistical molecular design (Andersson et al.,  
78 2014; J Prado-Prado et al., 2013).

79 Herein, we propose **the first** large-scale QSAR investigation for predicting AChE inhibition, which to  
80 the best of our knowledge represents the largest collection of **2,570 non-redundant** compounds. QSAR  
81 models were built using interpretable learning methods (e.g. random forest) and descriptors (i.e. molecular  
82 fingerprints) as to unravel the underlying AChE inhibitory activity, **which was performed in accordance**  
83 **with guidelines of the Organisation for Economic Cooperation and Development (OECD)**. Molecular  
84 docking was also performed on a chemically diverse set of compounds selected from active AChE  
85 inhibitors. Together, the ligand and structure-based approach employed in this study is anticipated to be  
86 useful in the design and development of robust AChE inhibitors.

## 87 MATERIALS AND METHODS

88 A summary of the workflow of **this study** is provided in Figure 1. Briefly, this included a large-scale  
89 **QSAR model for predicting and analyzing the AChE inhibition**, which was performed in accordance  
90 **with the OECD guidelines** as follows: (i) a data set with a defined endpoint; (ii) an unambiguous learning  
91 algorithm; (iii) a defined applicability domain of the QSAR model; (iv) using appropriate measures  
92 of goodness-of-fit, robustness and predictivity; (v) a mechanistic interpretation of the QSAR model.  
93 Furthermore, molecular docking was also performed on a chemically diverse data set as to elucidate the  
94 underlying binding mechanism. To facilitate the reproducibility of the research work performed herein,  
95 the data set and codes used to perform the multivariate analysis as well as the results tables and figures  
96 are made freely available on GitHub at <https://github.com/sawsimeon/AChE>.

### 97 Data set

98 A data set of inhibitors against human AChE (Target ID CHEMBL220) were compiled from the ChEMBL  
99 20 database (Gaulton et al., 2012) that is comprised of a total number of 9,242 bioactivity data points from  
100 5,049 compounds. **SMILES notations of the compounds were curated with ChemAxon's Standardizer** (?)  
101 **using the same parameter settings as described in our previous study** (Simeon et al., 2016). The initial  
102 data set was assembled from several bioactivity measurement units including (in order of decreasing data  
103 size) IC<sub>50</sub>, K<sub>i</sub>, % activity, % inhibition, MIC, EC<sub>50</sub>, etc. IC<sub>50</sub> was selected for further investigation as they  
104 constituted the largest subset with 4,910 compounds. **A closer look revealed that 1,301 compounds had no**  
105 **reported IC<sub>50</sub> values or had lesser/greater than signs which were subjected to removal to produce in 3,609**  
106 **compounds. Only data points having nM as the bioactivity unit were selected for further study, which**  
107 **produced 3,596 compounds. Furthermore, redundant compounds having different activity values were**  
108 **kept if the standard deviation of IC<sub>50</sub> was less than 2 and this resulted in 2,571 compounds. Moreover,**  
109 **some compounds were found to have no SMILES notation associated with it and were thus removed. A**  
110 **final data set comprising of **2,570** compounds was obtained.**

### 111 Description of inhibitors

112 AChE inhibitors were encoded by a vector of fingerprint descriptors accounting for its molecular constituents. Prior to calculating descriptors, salts were removed and tautomers were standardized using the built-in function of the PaDEL-Descriptor software (Yap, 2011).

113 Although fingerprint descriptors are able to capture the feature space of chemical compounds, however  
114 their ability to be used as descriptors for bioactivity modeling can vary. In fact, performance differences  
115 existing amongst the different fingerprint type has been the subject of several investigations into its utilization  
116 for bioactivity modeling. Riniker and Landrum (2013) benchmarked and assessed the performance of  
117 predictive models constructed from 2D fingerprint descriptors obtained from RDKit.

118 In this study, the suitability of 12 different fingerprint descriptors for predicting the bioactivity of AChE  
119 inhibitors was investigated. Table 1 summarizes the employed fingerprints along with their corresponding  
120 size, description and reference.

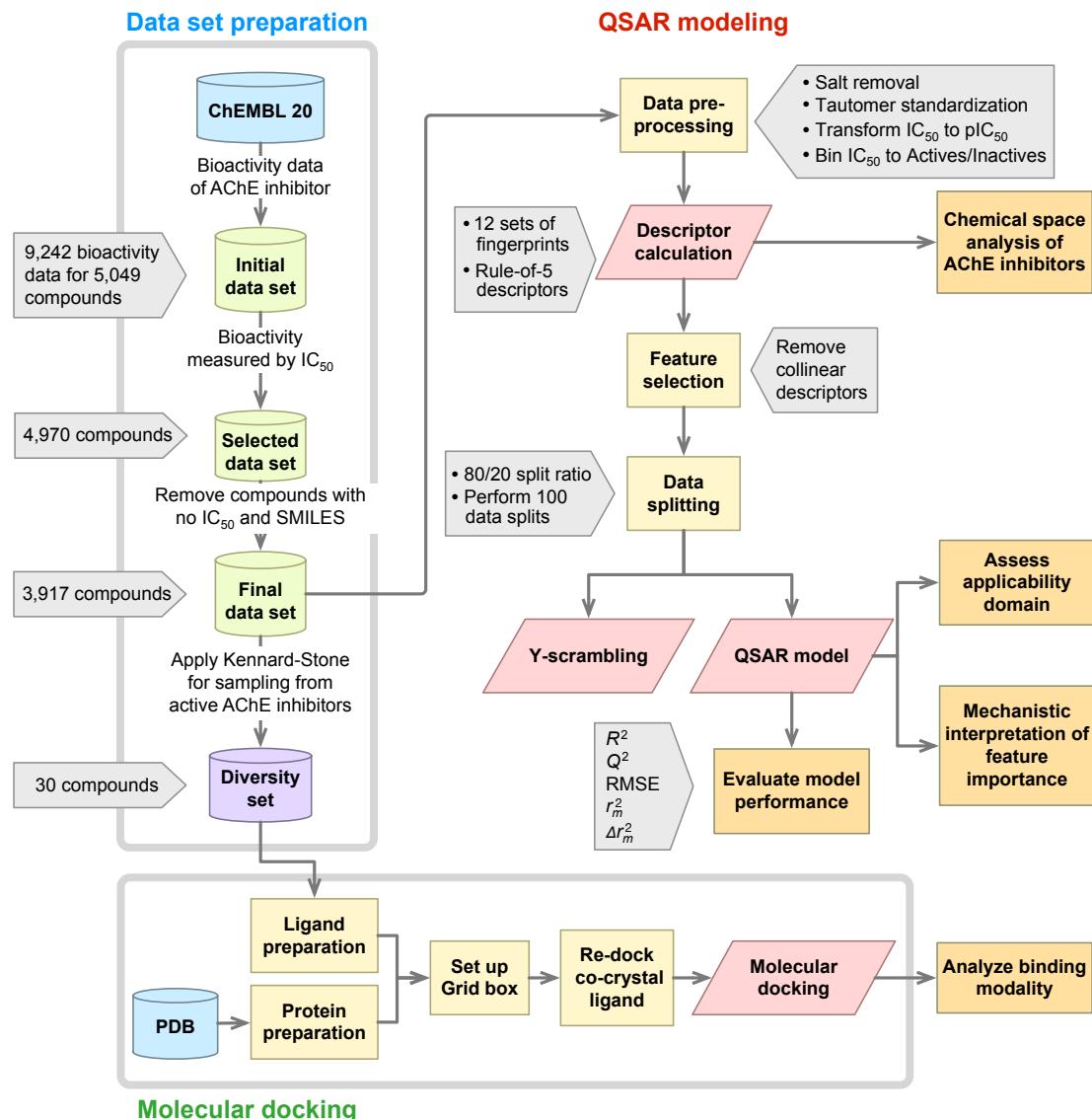
121 Additionally, the four molecular descriptors that are used to define the Lipinski's rule-of-five comprising  
122 of molecular weight (MW), logarithm of the octanol/water partition coefficient (ALogP), number of  
123 hydrogen bond donor (nHBDon) and number of hydrogen bond acceptor (nHBAcc) were also computed  
124 by the PaDEL-Descriptor software.

## 127 Feature selection

128 Collinearity is a condition where descriptor pairs are known to have intercorrelation, which not only add  
129 complexity to the model but could potentially give rise to bias. To remedy this, the *cor* function from  
130 the R package *stats* was used to find the pairwise correlation among descriptors, and descriptors in a  
131 pair with a Pearson's correlation coefficient greater than the threshold of 0.7 was filtered out using the  
132 *findCorrelation* function from the R package *caret* to obtain a smaller subset of descriptors (Kuhn, 2008).

## 133 Data splitting

134 To avoid the possibility of bias that may arise from a single data split when building predictive models  
135 (Puzyn et al., 2011), predictive models were constructed from 100 independent data splits and the mean  
136 and standard deviation values of statistical parameters were reported. The data set was split into internal  
137 and external sets in which the former comprises 80% whereas the latter constitutes 20% of the initial data  
138 set. The *sample* function from the R *base* package was used to split the data.



**Figure 1. Workflow of QSAR modeling and molecular docking for investigating AChE inhibitory activity.**

**Table 1.** Summary of 12 sets of fingerprint descriptors employed in this study.

No.	Fingerprint	Number	Description	Reference
1	CDK	1024	Fingerprint of length 1024 and search depth of 8	Steinbeck et al. (2003)
2	CDK extended	1024	Extends the fingerprint with additional bits describing ring features	Steinbeck et al. (2003)
3	CDK graph only	1024	A special version that considers only the connectivity and not bond order	Steinbeck et al. (2003)
4	E-state	79	Electrotopological state atom types	Hall and Kier (1995)
5	MACCS	166	Binary representation of chemical features defined by MACCS keys	Durant et al. (2002)
6	PubChem	881	Binary representation of substructures defined by PubChem	NCBI (2009)
7	Substructure	307	Presence of SMARTS patterns for functional groups	Laggner (2005)
8	Substructure count	307	Count of SMARTS patterns for functional groups	Laggner (2005)
9	Klekota-Roth	4860	Presence of chemical substructures	Klekota and Roth (2008)
10	Klekota-Roth count	4860	Count of chemical substructures	Klekota and Roth (2008)
11	2D atom pairs	780	Presence of atom pairs at various topological distances	Carhart et al. (1985)
12	2D atom pairs count	780	Count of atom pairs at various topological distances	Carhart et al. (1985)

### 139 Multivariate analysis

140 Supervised learning is to learn a model from labeled training data which can be used to make prediction  
141 about unseen or future data (James et al., 2013). This study constructs regression models, which affords  
142 the prediction of the continuous response variable (i.e. pIC<sub>50</sub>) as a function of predictors (i.e. fingerprint  
143 descriptors).

144 Random forest (RF) is an ensemble classifier that is composed of several decision trees (Breiman,  
145 2001). Briefly, the main idea behind RF is that instead of building a deep decision tree with an ever-  
146 growing number of nodes, which may be at risk for overfitting and overtraining of the data, rather multiple  
147 trees are generated as to minimize the variance instead of maximizing the accuracy. As such, the results  
148 will be more noisier when compared to a well-trained decision tree, yet, these results are usually reliable  
149 and robust. The *ranger* function from the R package *ranger*, which is a fast implementation of the RF  
150 algorithm that was used for constructing the models (Wright and Ziegler, 2015).

### 151 Validation of QSAR models

152 Model validation is an important process, which should be performed to ensure that a fitted model can  
153 accurately predict responses for future or unknown subjects. Two statistical parameters were used to  
154 evaluate the performance of the QSAR models consisting of Pearson's correlation coefficient (*r*) and  
155 root mean squared error (RMSE). The *r* value is a commonly used metric to represent the degree of  
156 relationship between two variables of interest. It can range from -1 to +1 in which negative values are  
157 indicative of negative correlation between two variables and vice versa. RMSE is a commonly used  
158 parameter to assess the relative error of the predictive model. The predictive performance of the QSAR  
159 models was verified by 10-fold cross-validation, external validation and Y-scrambling test.

160 The 10-fold cross-validation technique does not use the entire data set to build predictive model.  
161 Instead, it splits the data into training and testing data set by allowing model that is built with training  
162 data set us allow to assess the performance of the model on the testing data set. By performing repeats of  
163 the 10-fold validation, the average accuracies can be used to truly assess the performance of the predictive  
164 model.

Y-scrambling test was used to ensure the robustness of the predictive model not only to rule out the possibility of chance correlations but also to assess the statistical significance of  $R^2$  and  $Q^2$ , ensuring the generalizability of QSAR model. The true Y-dependent variable (i.e.  $\text{pIC}_{50}$ ) was randomly scrambled and the statistical assessment parameters are recalculated. Performance of the Y-scrambling test can be deduced from the regression line of the plot of  $R^2$  versus  $Q^2$ . Intercept values for  $R^2$  and  $Q^2$  as denoted by  $iR^2$  and  $iQ^2$ , respectively, were calculated. Negative  $iQ^2$  is indicative of an acceptable QSAR model and that there is no chance correlation from the real model (Eriksson et al., 2003). Furthermore,  $r_m^2$  and  $\Delta r_m^2$  metrics introduced as introduced by Roy et al. (2013) were used to verify the robustness of the proposed QSAR model in which an acceptable QSAR model should give  $r_m^2 > 0.5$  and  $\Delta r_m^2 < 0.2$ .

### 174 Applicability domain analysis

The applicability domain (AD) estimates the likelihood of reliable prediction for compounds on the basis of how similar they are to compounds used to build the model. Thus, compounds falling outside the AD may lead to unreliable predictions. The most common approach for determining AD is described by Gramatica (2007) and Tropsha et al. (2003), which is to compute the leverage values for each compound. The leverage value allows one to identify whether new compounds will lie within or outside the domain. Leverage values for all compounds are calculated via adjustment of  $X$  to give the hat matrix  $H$ :

$$H = X(X^T X)^{-1} X^T \quad (1)$$

where  $X$  is a two-dimensional matrix comprising of  $n$  compounds and  $m$  descriptors while  $X^T$  is the transpose of  $X$ . Meanwhile, the leverage value of the  $i^{\text{th}}$  compound ( $h_i$ ) is the  $i^{\text{th}}$  diagonal element of  $H$ :

$$h_i = x_i^T (X^T X)^{-1} x_i \quad (2)$$

where  $x_i$  is the descriptor row-vector of the  $i^{\text{th}}$  compound. The warning leverage  $h^*$  is calculated by:

$$h^* = 3(p + 1)/n \quad (3)$$

Practically, the leverage value along with the plot is often used to assess the AD of QSAR models. The plot is constructed by depicting the standardized residuals versus the leverage value for each compound's  $h_i$ . If the  $i^{\text{th}}$  compound has  $h_i > h^*$  then it means that the  $i^{\text{th}}$  compound exerts a great influence on the QSAR model and may be excluded from the AD. In spite of this, it does not appear to be an outlier because its standardized residual may be small.

### 180 Molecular docking

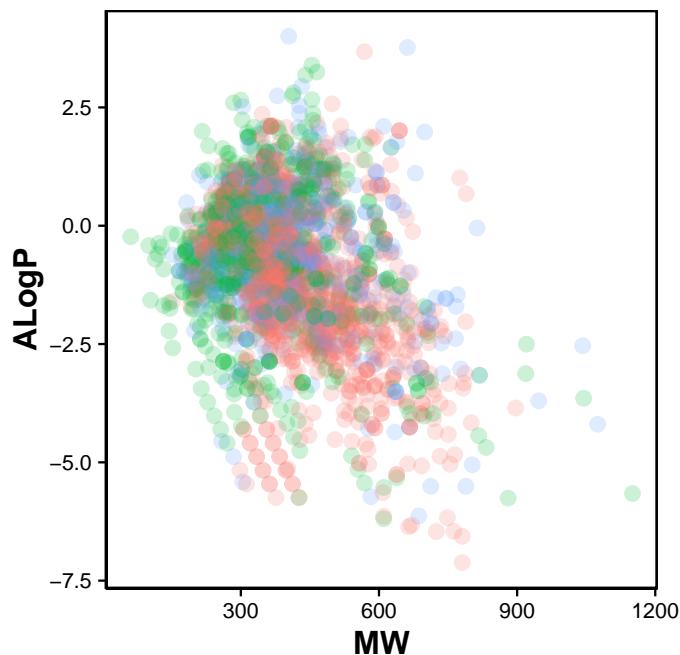
The co-crystal structure of human AChE with donepezil (PDB ID: 4EY7) was retrieved from the Protein Data Bank and initially prepared by removing alternative side chains and water molecules. The protein was prepared via the rebuilding of bonds and the addition of missing hydrogen atoms. Subsequently, the protein was cleaned by merging the atomic charge and removing lone pair atoms, non-polar hydrogen atoms and non-standard amino acid residues. Grid box was set up to provide coverage of the active site of human AChE with a dimension of  $40 \times 30 \times 40$  Å (X, Y and Z axes of -13.987, -41.668 and 27.109, respectively). Molecular docking was consequently performed with AutoDock Vina (Trott and Olson, 2010) using default parameters. The docking protocol was validated in order to ensure its reliability for subsequent analysis of the studied compounds. This was performed by extracting the co-crystal ligand, donepezil, from the PDB file and re-docked to the co-crystal human AChE protein. The root mean squared deviation (RMSD) of the atomic position between the original orientation of the co-crystal ligand and the re-docked ligand is computed and is deemed acceptable if the RMSD value is less than or equal to 2.0 Å.

A set of 30 representative and chemically diverse compounds, which will be referred hereafter as the diversity set, were extracted from the full set of active AChE inhibitors (i.e.  $\text{IC}_{50} < 1 \mu\text{M}$ ) using the Kennard-Stone algorithm (Kennard and Stone, 1969). These compounds were used as ligands for molecular docking against the human AChE. The binding energy (kcal/mol) of AChE inhibitors were calculated according to the built-in scoring function of Autodock Vina and conformers providing the lowest binding energy were selected for further analysis of the binding mode. Furthermore, key-interacting residues and their moiety preferences were analyzed using LigPlot+ (Wallace et al., 1995) and the SiMMAP web server (Bollback, 2006). Finally, three-dimensional structure of protein-ligand interaction was created and visualized using Pymol (Schrödinger, LLC, 2015a).

202 **RESULTS AND DISCUSSION**

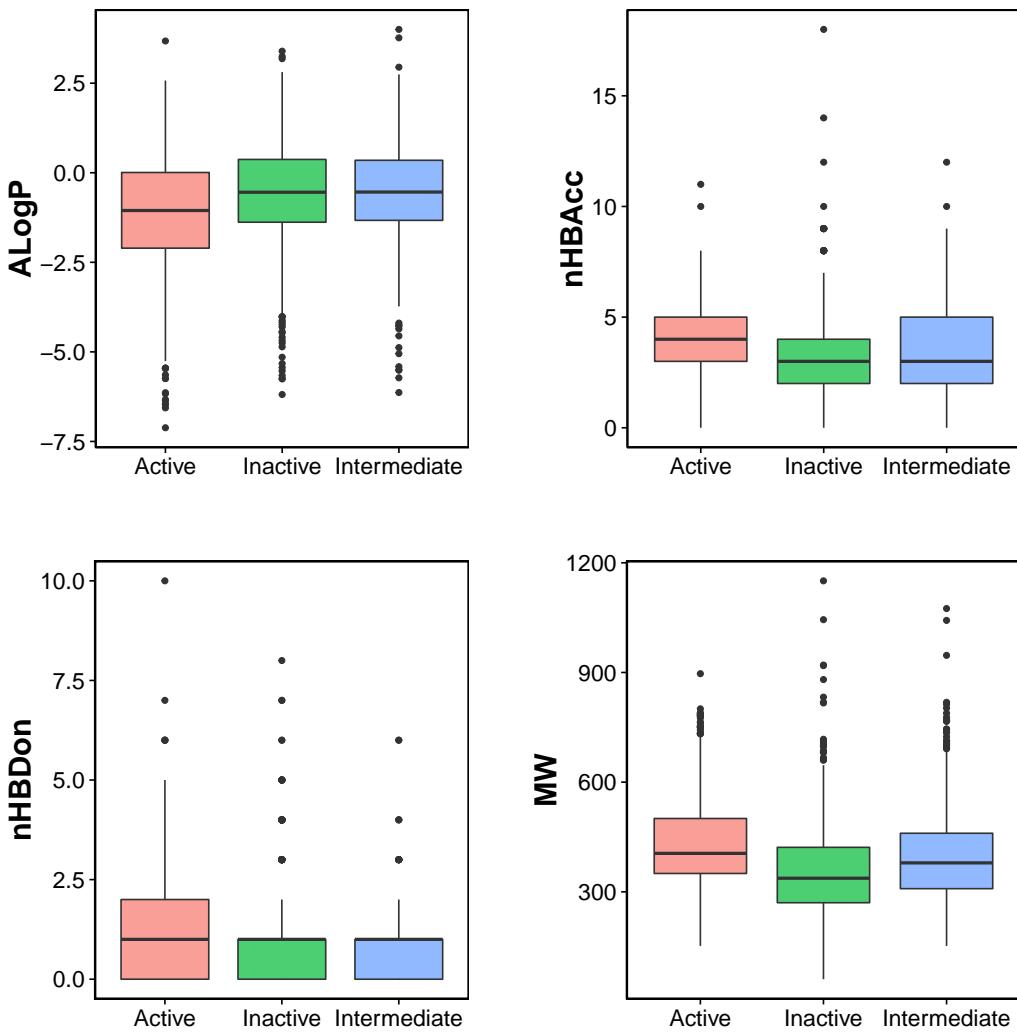
203 **Chemical space of AChE inhibitors**

204 Navigation of the chemical space of AChE inhibitors was performed to gain insights into the structure-  
205 activity relationship by analyzing the Lipinski's rule-of-five descriptors. Chemical space analysis may  
206 provide important knowledge on the general character of compounds governing inhibitory properties  
207 of compounds. Exploratory data analysis was performed using the Lipinski's rule-of-five descriptors  
208 comprising of MW, ALogP, nHBDon and nHBAcc. MW represents the molecular size of a compound  
209 that is commonly used because of it can be easily interpreted and calculated as well as appropriate size  
210 of a compound is important for its passage via lipid membrane. ALogP is a widely used parameter  
211 for determining the lipophilicity of a compound and used in calculating the membrane penetration and  
212 permeability of compounds. nHBDon and nHBAcc describe the number of hydrogen bond donors  
213 and hydrogen bond acceptors, respectively, which is used to measuring hydrogen bonding capacity.  
214 Visualization of the chemical space of ALogP as a function of MW is shown in Figure 2, as to investigate  
215 the chemical space of AChE inhibitors. A dense distribution of inhibitors was observed within the  
216 space of MW starting from approximately 300 to 600 Da and within the space of ALogP ranging from  
217 approximately -2.5 to 2.5. In addition, the box plot of the Lipinski's descriptors is shown in Fig. 3.  
218 Compounds with negative ALogP values approximately of closer to 0.0 can be found in inactive inhibitors  
219 whereas most of the active inhibitors tend to possess approximately lower values in average of ALogP  
220 values.



221 **Figure 2.** Chemical space of AChE inhibitors are shown as actives (green), inactives (red) and intermediates (blue)

222 Visual representation of the overall distribution of data values of Lipinski's descriptors is shown as  
223 box plots in Fig. 3 in which the ALogP, nHBAcc, nHBDon and MW are shown in the top-left, top-right,  
224 bottom-left and bottom-right corner, respectively. Analysis of the box plots revealed that there were  
225 no differences amongst the three bioactivity classes for nHBAcc and nHBAcc as deduced from the  
226 boundaries of the boxes (i.e. representing the first and third quartiles). ALogP and MW were found to  
227 display differences amongst the bioactivity classes. Particularly, ALogP values for actives were the lowest  
while negligible differences were observed for the other two classes. Furthermore, MW for actives were



**Figure 3.** Box plot of AChE inhibitors using Lipinski's rule-of-five descriptors (e.g. ALogP, nHBAcc, nHBDOn and MW)

the largest amongst the three bioactivity classes, which is followed by the intermediates while inactives were smallest.

#### QSAR model for predicting AChE inhibitory activity

A data set comprising of 2,570 compounds were used for construction of QSAR models. Particularly, twelve sets of fingerprint descriptors were benchmarked in order to find the best performing set. Prior to modeling, feature selection was applied to remove collinear descriptors. Each of the twelve models were then built using a data split ratio of 80/20 in which 80% of the data set was used as the internal set and 20% as the external set. This procedure was iteratively performed in which each of the 100 independent data splits were used for model construction and the performance results given in Table 2 are the mean and standard deviation values derived from these runs.

It can be observed that all twelve models are capable of capturing the inhibitory activity space of AChE inhibitors as they provided  $R^2$  and  $Q^2$  (i.e. both 10-fold CV and external sets) greater than the threshold values proposed by Golbraikh and Tropsha (2002) of 0.6 and 0.5, respectively, which is indicative of robust model performance. The possibility of chance correlation can be assessed from the  $R^2-Q^2$  margin as described by Eriksson and Johansson (1996) in which values < 0.2–0.3 are indicative of predictive and

243 reliable models while values > 0.2–0.3 suggests possible chance correlation or the presence of outliers in  
244 the data set. Furthermore, observation of the  $Q_{CV}^2$ – $Q_{Ext}^2$  margin revealed that the difference was negligible  
245 with values in the range of 0 and 0.01.

246 Generally, it can be seen that models with larger descriptor size, namely CDK and CDK extended,  
247 afforded the best performance with  $Q_{CV}^2$  of  $0.79 \pm 0.07$  and  $0.79 \pm 0.06$ , respectively, and  $Q_{Ext}^2$  of  $0.80 \pm 0.04$   
248 and  $0.81 \pm 0.04$ , respectively. The opposite also holds true as the model with the least number of descriptors  
249 were also found to perform the worst amongst the other fingerprints with  $Q_{CV}^2$  of  $0.55 \pm 0.09$  and  $Q_{Ext}^2$   
250 of  $0.56 \pm 0.05$ . In a nutshell, the model performance in order of decreasing value is as follows: CDK  
251 extended > CDK > MACCS ≈ Substructure count ≈ Klekota-Rota count > PubChem > Klekota-Roth  
252 ≈ 2D atom pairs count > CDK graph only > Substructure > 2D atom pairs > E-state.

253 The best performing model is not necessarily the best choice considering the fact that the descriptor  
254 size for the best models were quite high and is consequently prone to overfitting. It was found that the  
255 substructure count provided reasonably good predictive performance (i.e.  $Q_{CV}^2$  and  $Q_{Ext}^2$  of  $0.78 \pm 0.06$  and  
256  $0.78 \pm 0.05$ , respectively) with the advantage of making use of a small set of 26 descriptors. Therefore,  
257 this fingerprint was selected for further interpretation of the feature importance.

258 To further check the reliability and validity of the selected model, Y-scrambling test was performed  
259 for 100 iterations. Table demonstrates that QSAR models built using substructure count has a low  $Q^2$   
260 ( $-0.0013$ ), which rules out the possibility of chance correlation. Furthermore, model afforded an  $r_m^2$  value  
261 of  $0.61 \pm 0.06$  thereby revealing its robustness. It is observed in Table that the value of  $\Delta r_m^2$  is greater than  
262 0.2 but also close to 0.20.

263 As shown in Fig. 4, it can also be seen that scatter plots of experimental versus predicted  $pIC_{50}$  of  
264 panels A, C and F displayed narrower variance of the data points than the other methods as assessed via  
265 10-fold cross-validation and external set.

## 266 Applicability domain

267 The AD of the proposed QSAR model was defined as provided by the Williams plot shown in Fig. The  
268 employed data set consisting of 2,570 compounds was randomly split to two separate subset in which the  
269 first subset constituting 80% of the data set was used as an internal set while the second subset constituting  
270 the remaining 20% were used as an external set. Compounds representing the internal set (blue dots) and  
271 external set (red dots) are shown in the Williams plot and it can be clearly seen that almost all of the  
272 2,570 compounds were located within the boundaries of applicability domain, which indicated that our  
273 proposed QSAR model had a well-defined AD.

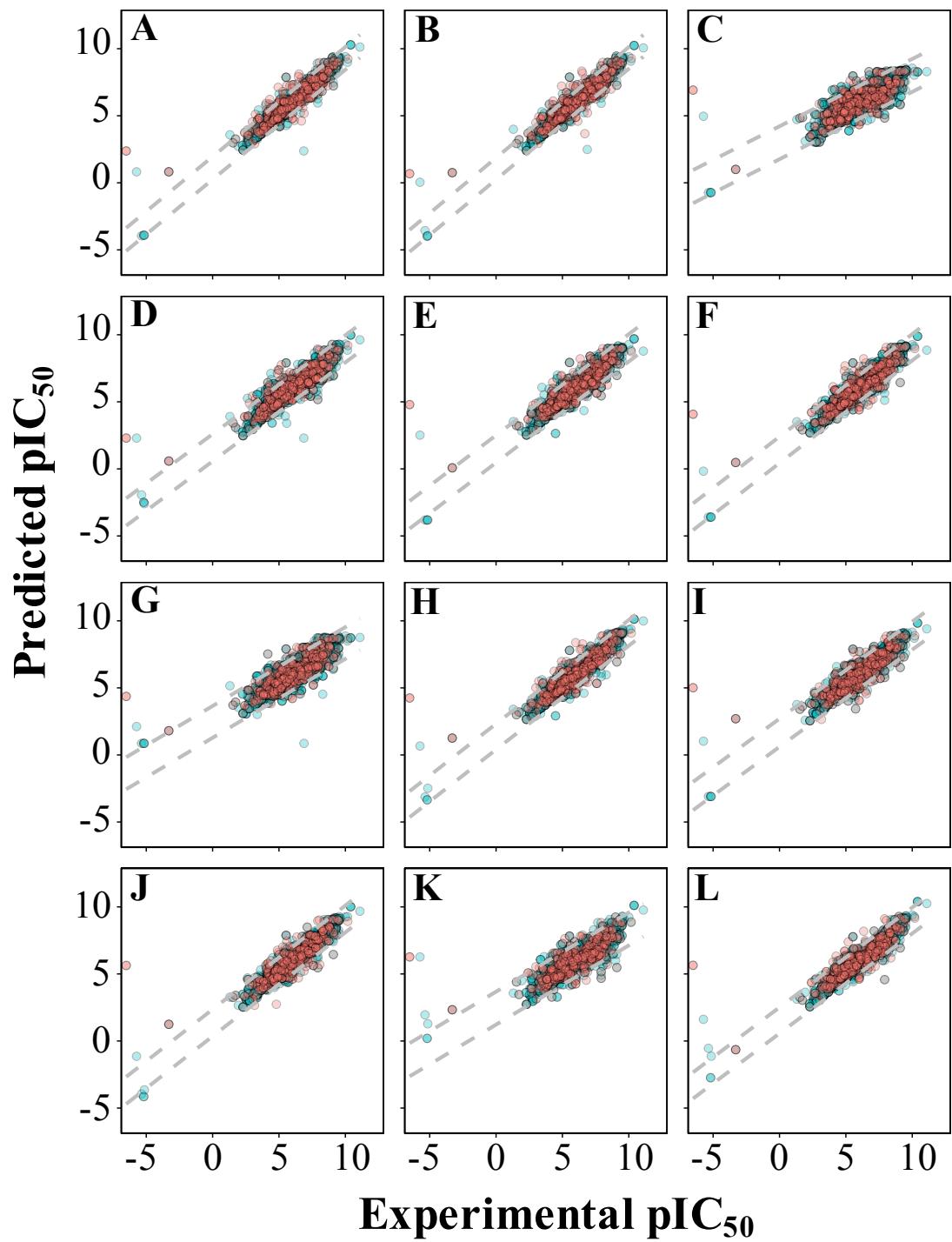
274 As can be seen in Fig. very few compounds indeed fall outside the  $\pm 3$  standardized residual range.  
275 This consisted of 6 compounds (997, 1829, 62, 1096, 13, 677) from the internal set and 11 compounds  
276 (2116, 2120, 2388, 2117, 2392, 2323, 2423, 2424, 2507, 2219, 2422) from the external set that had  
277 standardized residual higher than 3. On the other hand, 7 compounds (1567, 576, 1644, 1098, 2022,  
278 1447 and 322) from the internal set and 8 compounds (2486, 2353, 2130, 2553, 2389, 2072, 2103, 2125)  
279 from the external set had standardized residual lower than -3. The corresponding chemical structures are  
280 provided in Table S1.

## 281 Mechanistic interpretation of feature importance

282 Feature importance analysis help reveal features that are important toward bioactivity. There are essentially  
283 two parameters for evaluating the relative importance of features used in models using the RF algorithm:  
284 (i) accuracy and (ii) Gini index (i.e. variance of the responses). The latter was selected as a metric for  
285 ranking important features (i.e. mean decrease of the Gini index) for predicting the  $pIC_{50}$  of AChE  
286 inhibitors (Fig. 6). Table 4 lists the substructure fingerprints along with their respective descriptions.

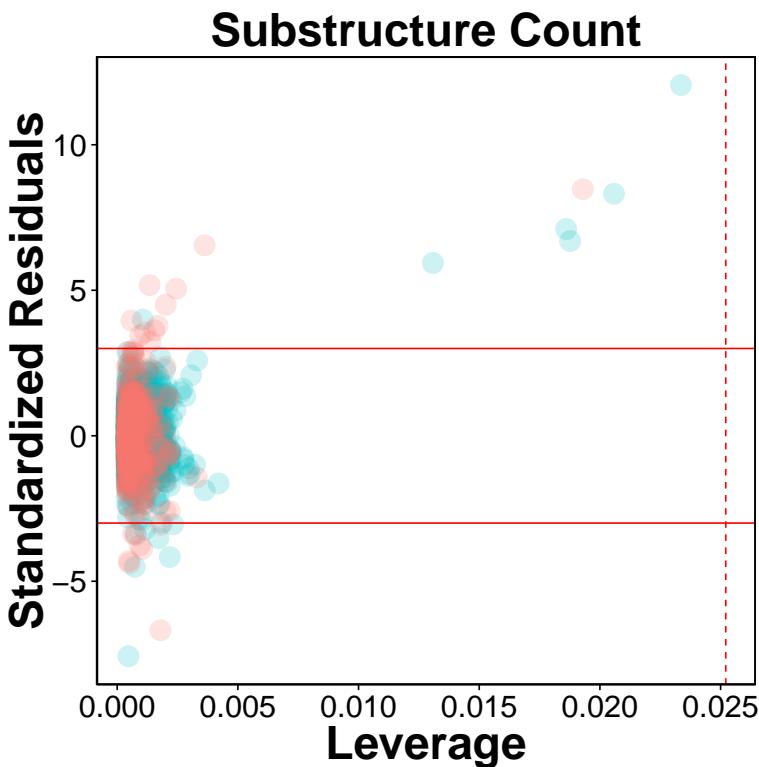
287 As can be seen in Fig. 6, the top ranking feature is secondary carbon (SubFPC2), which is a carbon  
288 atom with two carbon neighbors. In the context of drug design, such central carbon atom may be more  
289 difficult to be accessed and metabolized by cytochrome P450 (Utrecht and Trager, 2007) and therefore  
290 are more metabolically stable.

291 The second most important feature is the rotatable bond (SubFPC302). Based on the rule of three  
292 for defining lead-like compounds, a compound may have a lead-like characters if it does have rotatable  
293 bonds of no more than 3. On the other hand, Veber et al. (2002) noted that the upper limit of orally  
294 bioavailable drugs is of seven rotatable bonds. Nevertheless, it has been found that number of rotatable  
295 bonds provide better discrimination between compounds that are orally active and those that are not.  
296 Kryger et al. (1999) claimed that E2020 (i.e. also known as donepezil and marketed as Aricept) needs at



**Figure 4.** Plot of experimental versus predicted  $\text{pIC}_{50}$  values for models constructed with 12 different fingerprint descriptors. Shown are models built with CDK fingerprint (A), CDK extended fingerprint (B), E-State fingerprint (C), CDK graph only fingerprint (D), MACCS fingerprint (E), PubChem fingerprint (F), substructure fingerprint (G), substructure fingerprint count (H), Klekota-Roth fingerprint (I), Klekota-Roth fingerprint count (J), 2D atom pairs (K) and 2D atom pairs count (L).

297 least two rotatable bonds on each side of the piperidine in which two aromatic moieties of E2020 interact  
 298 with Trp86 and Trp286 (human AChE numbering), suggesting that links between aromatic systems of the  
 299 inhibitor against its AChE counterparts are essential to yield high affinity.



**Figure 5.** William plot for the internal (blue dots) and external (red dots) sets for QSAR model built using substructure fingerprint count. The solid and dashed lines correspond to the  $\pm 3$  standardized residual and the warning leverage value ( $h^* = 0.025$ ), respectively.

300        The third important substructure is the aromatic ring (SubFPC274). Findings from X-ray crystallographic study showed that in the binding site of the co-crystal structure of AChE with tacrine, the aromatic ring of acridine engages in a  $\pi$ - $\pi$  stacking interaction with the indole of Trp86 (human AChE numbering) thereby indicating the importance of the aromatic ring for AChE inhibition. Chen et al. (2012)

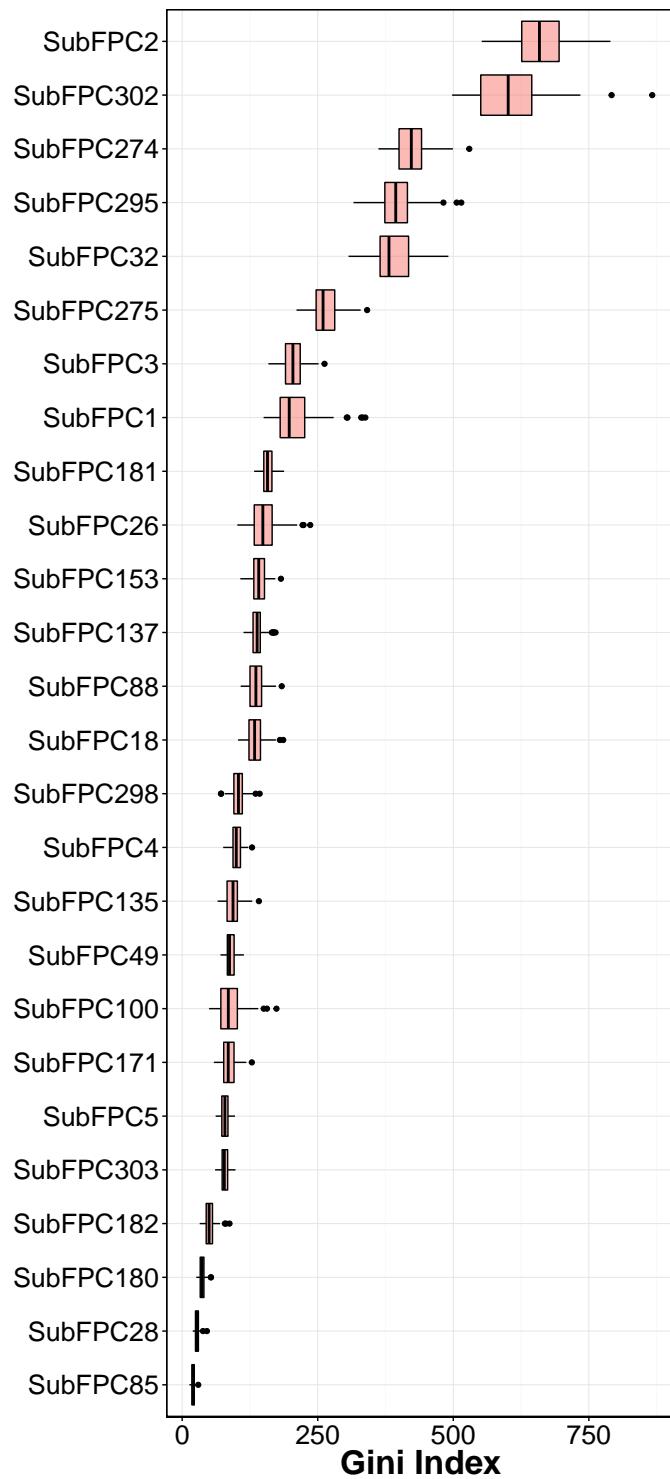
304        The fourth important feature is C ONS bond (SubFPC295), which is defined as the presence of any carbon connected with either oxygen, nitrogen or sulfur atom in a molecule. These atoms contribute 305 electrostatic charge, which deemed to increase the polarity

307        The fifth important substructure is secondary mixed amine (SubFPC32). The importance of the moiety 308 was demonstrated in the work of Bembeneck et al. (2008) in which a structure-based approach was used to 309 reveal that in the catalytic triad, Trp86 interacts with the quaternary amine of ACh through a cation- $\pi$  310 interaction. Furthermore, in the 'anionic' site Trp286 appears to attract the amine moiety via cation and/or 311 hydrophobic interactions.

312        The sixth, seventh and eighth important substructures are heterocyclic rings (SubFPC275), tertiary 313 carbon (SubFPC3) and primary carbon (SubFPC1). Heterocycles are of high relevance in the design of 314 AChE inhibitors as it allows  $\pi$ - $\pi$  stacking interaction with key amino acid residues in the binding site of 315 AChE. It is observed that the binding site of the AChE are highly hydrophobic in nature. Particularly, the 316 Trp286 (human AChE numbering) which is a part of the peripheral anionic site of the AChE is involved 317 in the  $\pi$ - $\pi$  interaction with heterocycles of AChE inhibitors (Lu et al., 2011). The aforementioned 318 explanation made for the secondary carbon is also applicable for the tertiary carbon in which the higher 319 number of carbon neighbors would also confer high stability against cytochrome P450 metabolism.

320        The ninth and tenth important substructures are Hetero N nonbasic (SubFPC181) and tertiary aliphatic 321 amine (SubFPC26) in which former is defined as aromatic nitrogen. Present of Nitrogen atom, which is 322 cation moiety deemed to interact with aromatic residues through  $\pi$ -cation interaction as observed in E2020 323 against *Torpedo californica* AChE (tcAChE). The charge of nitrogen atom located in piperidine ring 324 provide  $\pi$ -cation interaction with the side chain of Phe337 (Phe330 in tcAChE numbering) (Jianxin Guo 325 et al., 2004). Since the active site gorge of AChE comprised of several aromatic residues, also known as

<sup>326</sup> aromatic patch, adding cation moiety possibly increase the binding affinity when the ligand arranged in  
<sup>327</sup> suitable conformation against aromatic side chain of residue the active site.



**Figure 6.** Plot of feature importance as exemplified by the Gini index.

**Table 2.** Performance summary of QSAR models for predicting pIC<sub>50</sub>.

Descriptor class	N	Training set		10-fold CV set		External set	
		R <sup>2</sup>	RMSE <sub>Tr</sub>	Q <sup>2</sup> <sub>CV</sub>	RMSE <sub>CV</sub>	Q <sup>2</sup> <sub>Ext</sub>	RMSE <sub>Ext</sub>
CDK	960	0.93±0.01	0.44±0.04	0.79±0.07	0.76±0.15	0.80±0.04	0.73±0.09
CDK extended	948	0.94±0.01	0.42±0.03	0.79±0.06	0.76±0.12	0.81±0.04	0.72±0.08
CDK graph only	198	0.87±0.01	0.61±0.03	0.72±0.06	0.87±0.13	0.72±0.05	0.87±0.09
E-State	21	0.66±0.03	1.00±0.05	0.55±0.09	1.11±0.13	0.56±0.05	1.10±0.08
MACCS	77	0.89±0.01	0.56±0.03	0.77±0.07	0.81±0.15	0.77±0.04	0.80±0.09
PubChem	103	0.90±0.01	0.55±0.03	0.76±0.05	0.80±0.11	0.78±0.03	0.79±0.08
Substructure	30	0.75±0.01	0.85±0.03	0.64±0.06	1.00±0.13	0.65±0.05	0.98±0.08
Substructure count	26	0.92±0.01	0.50±0.02	0.78±0.06	0.77±0.14	0.78±0.05	0.77±0.10
Klekota-Roth	111	0.89±0.01	0.59±0.03	0.74±0.07	0.85±0.14	0.76±0.05	0.81±0.10
Klekota-Roth count	72	0.91±0.01	0.52±0.03	0.78±0.06	0.79±0.14	0.78±0.05	0.77±0.11
2D atom pairs	42	0.75±0.03	0.85±0.06	0.61±0.08	1.03±0.15	0.60±0.06	1.05±0.12
2D atom pairs count	38	0.92±0.01	0.51±0.02	0.74±0.07	0.84±0.15	0.76±0.05	0.82±0.10

**Table 3.** Performance summary of QSAR models assessed using  $\text{i}R^2$ ,  $\text{i}Q^2$ ,  
 $r_m^2$  and  $\Delta r_m^2$  metrics.

Descriptor class	N	Training set		10-fold CV set		External set		$\text{i}R^2$	$\text{i}Q^2$
		$r_m^2$	$\Delta r_m^2$	$r_m^2$	$\Delta r_m^2$	$r_m^2$	$\Delta r_m^2$		
CDK	960	0.82±0.02	0.07±0.01	0.62±0.09	0.20±0.06	0.64±0.05	0.19±0.03	0.0003	-0.0003
CDK extended	948	0.83±0.01	0.07±0.01	0.62±0.08	0.20±0.05	0.65±0.05	0.18±0.03	0.0006	-0.0005
CDK graph only	198	0.70±0.02	0.14±0.01	0.51±0.07	0.27±0.05	0.53±0.05	0.26±0.03	0.0007	-0.0006
E-State	21	0.35±0.03	0.38±0.02	0.27±0.07	0.40±0.04	0.28±0.05	0.41±0.03	0.0011	-0.0009
MACCS	77	0.73±0.01	0.12±0.01	0.57±0.09	0.23±0.05	0.58±0.05	0.23±0.03	0.0005	-0.0004
PubChem	103	0.74±0.02	0.12±0.01	0.57±0.07	0.23±0.04	0.59±0.05	0.22±0.03	0.0006	-0.0005
Substructure	30	0.50±0.02	0.28±0.01	0.39±0.07	0.34±0.05	0.41±0.05	0.33±0.03	0.0033	-0.0027
Substructure count	26	0.77±0.01	0.10±0.01	0.60±0.08	0.22±0.05	0.61±0.06	0.21±0.04	0.0015	-0.0013
Klekota-Roth	111	0.71±0.02	0.14±0.01	0.54±0.08	0.25±0.05	0.56±0.06	0.24±0.03	0.0006	-0.0004
Klekota-Roth count	72	0.76±0.02	0.11±0.01	0.60±0.08	0.22±0.05	0.61±0.07	0.21±0.04	0.0006	-0.0005
2D atom pairs	42	0.49±0.03	0.28±0.02	0.35±0.08	0.36±0.04	0.35±0.06	0.36±0.03	0.0010	-0.0008
2D atom pairs count	38	0.75±0.01	0.10±0.01	0.52±0.05	0.26±0.05	0.54±0.06	0.25±0.04	0.0006	-0.0005

**Table 4.** List of top substructure fingerprints and their corresponding description.

Fingerprints	Description
SubFPC1	Primary carbon
SubFPC2	Secondary carbon
SubFPC3	Tertiary carbon
SubFPC5	Alkene
SubFPC18	Alkylaryether
SubFPC23	Amine
SubFPC26	Tertiary aliphatic amine
SubFPC28	Primary aromatic amine
SubFPC32	Secondary mixed amine
SubFPC35	Ammonium
SubFPC49	Ketone
SubFPC88	Carboxylic acid derivative
SubFPC100	Secondary amide
SubFPC135	Carbonyl derivative
SubFPC137	Vinylogous ester
SubFPC143	Carbonic acid diester
SubFPC153	Urethan
SubFPC171	Arylchloride
SubFPC180	Hetero N basic H
SubFPC181	Hetero N nonbasic
SubFPC182	Hetero O
SubFPC184	Heteroaromatic ring
SubFPC274	Aromatic ring
SubFPC275	Heterocyclic ring
SubFPC276	Epoxide
SubFPC287	Conjugated double bond
SubFPC295	C ONS bond
SubFPC296	Charged
SubFPC298	Cation
SubFPC300	1,3-Tautomerizable
SubFPC301	1,5-Tautomerizable
SubFPC302	Rotatable bond
SubFPC303	Michael acceptor
SubFPC307	Chiral center specified

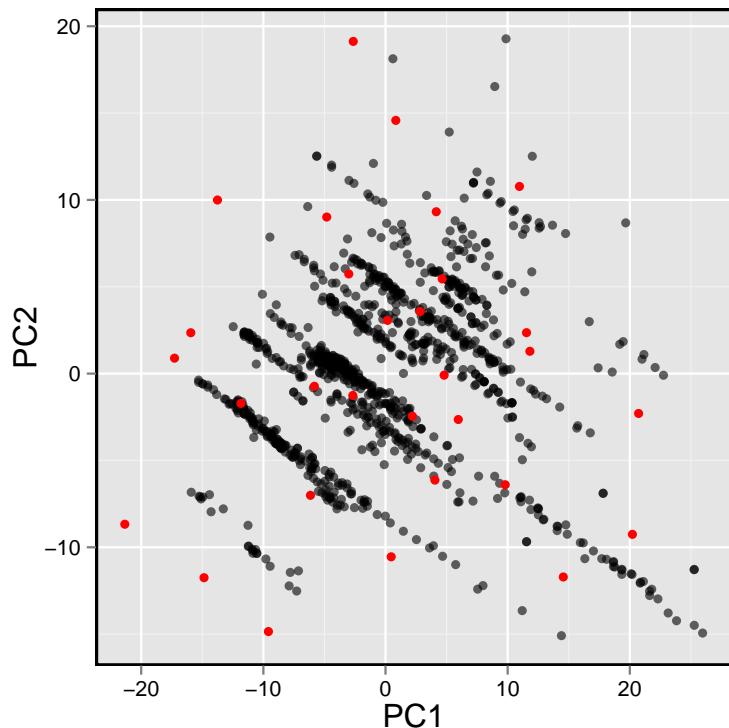
## 328 MOLECULAR DOCKING OF AChE INHIBITORS

329 To gain a further understanding on the non-covalent interaction between AChE and their inhibitors,  
330 a chemically diverse set of 30 representative compounds was extracted from active AChE inhibitors  
331 (i.e. having  $IC_{50} < 1 \mu M$ ) using the Kennard-Stone algorithm and subjected to an investigation on its  
332 binding modality against the active site of AChE. Figure 7 shows the distribution of the selected subset  
333 of compounds in the context of the full set of actives, which was found to provide a full coverage of the  
334 original chemical space. The chemical structures of these compounds are shown in Fig. 8.

335 The active site of this enzyme is buried inside a narrow gorge of 20 Å deep, which permits multiple  
336 enzyme-substrate interaction thereby facilitating the formation of the transition state of ACh (Silman and  
337 Sussman, 2008; Zhou et al., 2010; Cheung et al., 2012a). The entry of the active site gorge is lined up by  
338 peripheral anionic site (PAS), which is composed of Tyr72, Tyr124, Trp286 and Try341. The function  
339 of PAS is to trap the substrate via  $\pi$ -cation interaction and proceed through the constriction residues  
340 Tyr124 and Phe338 and onto the catalytic site (Silman and Sussman, 2008; Dvir et al., 2010). As a serine  
341 hydrolase, AChE contains Ser203, Glu334 and His447 in the catalytic triad that catalyzed the acylation  
342 and deacylation of ACh. The catalytic triad is surrounded by the catalytic anionic site (CAS) (i.e. contains  
343 Trp86, Glu202 and Tyr337, oxyanion hole (i.e. comprising of Glu121, Glu122 and Ala204) and the acyl  
344 pocket (i.e. comprised of Phe295 and Phe297). These sites help to tie up the ACh and make the substrate  
345 arranged in a suitable orientation for interacting with the catalytic triad as well as providing proton transfer  
346 that is essential for nucleophilic substitution during the catalytic reaction (Zhou et al., 2010). As a result  
347 of the acylation process, the nucleophilic attack from the O atom of Ser203 to the carbonyl group of  
348 ACh affords the proton transfer from Ser203 to His447, which consequently breaks down the choline  
349 moiety and forms a covalent acylenzyme complex between Ser203 and the acetyl group. This complex  
350 consequently proceeds with deacylation, which follows a similar mechanism with the acylation stage.  
351 The protonation of Glu202 provides a water molecule and the proton transfer from His447 to the water  
352 molecule leads to the nucleophilic attack against the acetyl group of the complex. Finally, this results in  
353 the breaking down of the complex thereby restoring wild-type AChE and causing the release of acetic  
354 acid from the active site (Zhou et al., 2010).

355 Prior to carrying out the molecular docking calculations, the docking protocol was validated by  
356 re-docking the co-crystal ligand and protein. It was found that the re-docked ligand exhibited negligible  
357 deviation from the co-crystal conformation with an RMSD value of 0.963 Å, which was deemed to be  
358 suitable for further molecular docking investigation and its subsequent interpretation. Consequently, the  
359 binding modality was analyzed in order to gain understanding on the contribution of key residues in  
360 interacting with the investigated set of 30 compounds. This was performed using the SiMMAP web server,  
361 which revealed three major binding anchors: Hbond1, vdW1 and vdW2 along with their site-moiety  
362 preferences. The first anchor site involves hydrogen bond interaction between Tyr124 (i.e. an important  
363 residue in the PAS that is spatially located as a bottleneck between the peripheral region and the catalytic  
364 site of AChE) and the following ligand moieties: secondary amide, secondary amine, nitrogen moiety in  
365 aromatic ring, ketone and ester. Such interaction can be observed in the co-crystal structure of huperzine  
366 A with human AChE (Cheung et al., 2012b). Interestingly, analysis of the important features from QSAR  
367 models also revealed the importance of "C ONS bond", "secondary mixed amine", "heterocyclic" and  
368 "hetero N non-basic" as they were found to be in the top ten important substructures and is therefore  
369 crucial for forming hydrogen bonds. Furthermore, the other anchor sites involve van der Waals interaction  
370 in which members of the first van der Waals interaction site (vdW1) are comprised of Tyr124, Phe338  
371 and Try341, which has a preference to interact with heterocyclic, aromatic, phenol and other non-polar  
372 moieties from representative inhibitors. The second van der Waals' interaction site (vdW2) is consisted of  
373 Trp86 and Gly121 with preference for the following ligand moieties: aromatic ring, heterocyclic ring,  
374 aliphatic moiety with alkene linkage and phenol moiety. These residues contain either bulky aromatic ring  
375 or non-polar moiety as their side chain to provide the van der Waal's surface contact against non-polar  
376 moiety from ligands. Notably, aromatic and heterocyclic substructures were also observed in the top ten  
377 important substructures for predicting the inhibitory activity of human AChE inhibitors.

378 Furthermore, analysis of the binding energy from the 30 representative compounds revealed that  
379 compounds **13**, **5** and **28** exhibited the lowest binding energy of  $-12.2$ ,  $-12.0$  and  $-12.0$  kcal/mol, respectively,  
380 when interacting with the human AChE binding site, which is comparable to donepezil ( $-12.2$   
381 kcal/mol) as indicated in Figure 8. Key interacting residues and their moiety preference was deduced  
382 from the protein-ligand interaction diagram generated by LigPlot+ (Wallace et al., 1995) in combination

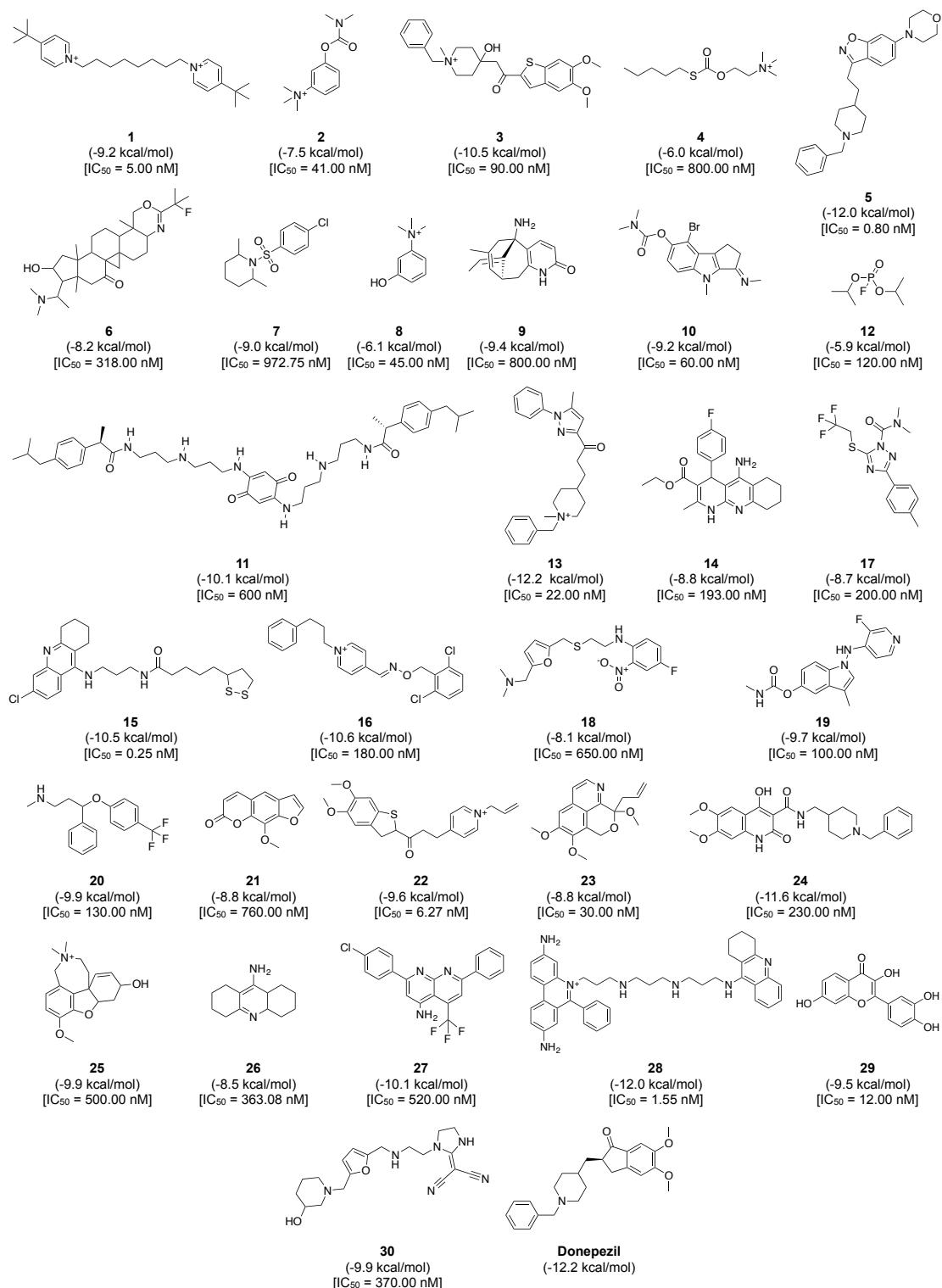


**Figure 7.** Plot showing the distribution of active AChE inhibitors (gray circles) and the diversity set (red circles) selected for molecular docking.

383 with Maestro (Schrödinger, LLC, 2015b) and their three-dimensional structure was visualized by PyMOL  
384 (Schrödinger, LLC, 2015a).

385 Figures 9A-B revealed three major interaction sites for **13** inside the binding pocket of human AChE.  
386 The first interaction site (Fig. 9C) is formed by residues from the PAS consisting of Trp286 and Tyr341  
387 both of which engages in  $\pi$ - $\pi$  stacking interaction where the terminal benzene and the attached pyrazole  
388 moiety of the ligand interacts with the former residue while the pyrazole moiety of **13** provide contact  
389 with the phenol moiety of the latter residue, which are deemed to increase the binding fitness against the  
390 active site of human AChE. The second interaction site (Fig. 9D) is dominated by Tyr124 and Phe338  
391 in which the former interacts with the ketone moiety of **13** by means of a hydrogen bond, which tends  
392 to increase the binding affinity of this compound. The side chain of Phe338 is involved in hydrophobic  
393 interaction with carbon atoms from the piperidine ring and the central aliphatic chain of **13**. It can be  
394 noted that these residues are members of constriction site, which arranged as bottleneck of active site.  
395 The third interaction site (Fig. 9E) is moderated by  $\pi$ - $\pi$  stacking between Trp86 and the terminal benzene  
396 with an attached piperidine moiety as well as hydrophobic interaction between Tyr337 and the piperidine  
397 moiety of **13** in which both residues belonged to the CAS. In addition, Gly121 of oxyanion hole also  
398 interact with terminal benzene of the ligand through hydrophobic contact increasing binding affinity  
399 against catalytic site of AChE.

400 Analysis of the binding modality of compound **5** revealed interactions with all subsites of the AChE  
401 active site gorge. PAS was the first subsite dominated by  $\pi$ - $\pi$  stacking between Trp286 and the ben-  
402 zisoxazole moiety of compound **5**, which is essential for stabilizing the binding affinity of the ligand  
403 against entry into the gorge. The second interaction site was observed at the constriction site in which the  
404 piperidine moiety makes contact with Phe338 via hydrophobic interaction thereby increasing the binding  
405 fitness against the bottleneck region of the active site. Similar hydrophobic interaction was also observed  
406 in the binding pocket in which Tyr337 from CAS interacts with the piperidine moiety and Gly121 of  
407 oxyanion pocket interacts with the terminal benzene of compound **5**. The  $\pi$ - $\pi$  interaction between the  
408 terminal benzene and Trp86 from the CAS was deemed to increase the binding fitness with the catalytic  
409 site of AChE. In addition, H-bond interaction is facilitated by the N atom from Phe295 at the acyl pocket



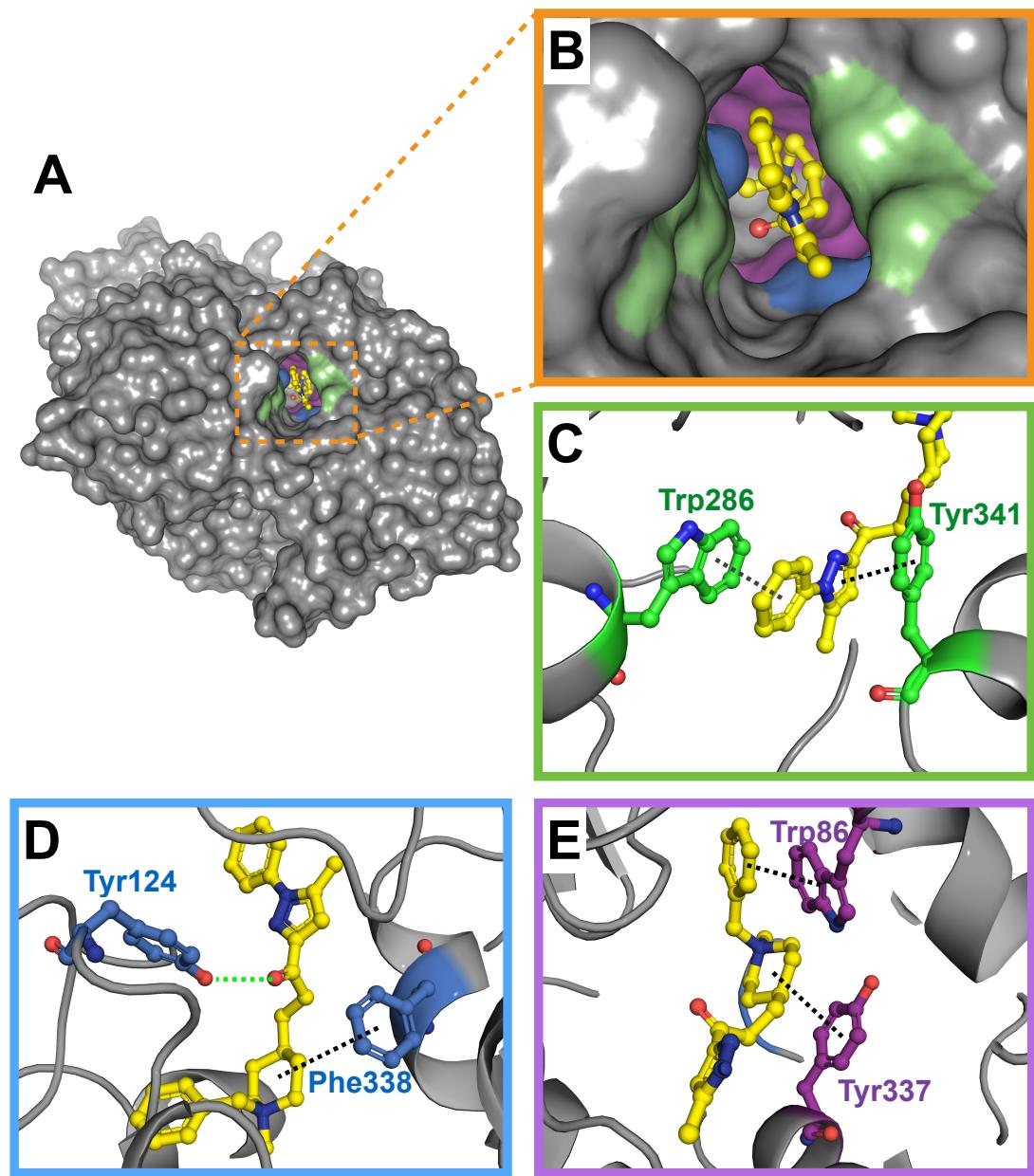
**Figure 8. Chemical structures, binding energy and bioactivity of the diversity set consisting of 30 representative compounds from active AChE inhibitors.**

to the O atom from the benzisoxazole moiety. Notably, all sites from the active site gorge are snugly bound by compound 5, which is deemed to exhibit strong intermolecular interaction with human AChE.

The binding energy of compound 5 was similar to that of compound 28. These compounds possessed

413 several aromatic rings at both terminal, which are favourable for interacting with aromatic residues lining  
414 up the surface of the gorge and these are known as the aromatic patch. The 5,7-dihydrophenanthridine  
415 moiety facilitates  $\pi$ - $\pi$  stacking with the side chain of Trp286 from PAS at the gorge opening. Meanwhile,  
416 this moiety also engages in hydrophobic interaction with Phe295, which tends to increase the binding  
417 fitness for the acyl pocket. Aside from the former moiety, 1,2,3,4-tetrahydroacridine at the opposite  
418 terminal provides  $\pi$ - $\pi$  interaction with Trp86 and hydrophobic contact with Tyr337 where both of which  
419 are members of CAS in the catalytic site. Furthermore, long aliphatic chain linking the two aromatic  
420 moieties provide hydrophobic contact with several aromatic residues in the aromatic patch consisting of  
421 Tyr72, Tyr124, Trp286, Tyr337, Phe338 and Tyr341 (i.e. these residues are the members of PAS, CAS  
422 and constriction site of the gorge). Moreover, this chain contain several N atoms, which can act as H-bond  
423 donor to Tyr124 and Tyr337 from PAS and CAS, respectively. This would tighten the binding between  
424 compound **28** and the active site gorge.

425 It should be noted that compounds exhibiting strong binding fitness against AChE are those that  
426 interact with residues from both PAS and CAS at the entry and inner pocket of the gorge, respectively, as  
427 dual-binding site inhibitor through either  $\pi$ - $\pi$  stacking or  $\pi$ -cation interaction together with hydrophobic  
428 contact. These compounds competes with the natural substrate in interacting with these residues. For non-  
429 covalent inhibitors, the aromatic moiety is preferred for occupying the interaction sites while hydrophobic  
430 moieties are preferred for making contact with the aromatic residues surrounding the catalytic site. H-bond  
431 donor moieties such as secondary amine and heterocyclic ring can be employed for interacting with the O  
432 atom on the side chain of Tyr residues. Interestingly, this finding is corroborated by the feature importance  
433 results obtained from the QSAR model as shown in Fig. 6 in which the aromatic moiety, C ONS bond,  
434 secondary mixed amine, heterocyclic ring and the hetero N non-basic moiety were found amongst the top  
435 ten important substructures that are essential for the bioactivity of AChE inhibitor.



**Figure 9. Molecular docking of the representative compound 13 against AChE.** Illustration of the binding between AChE (gray colored surface representation) and 13 (yellow colored ball-and-stick) (A) and a zoom-in view showing the interaction of 13 with three spatially juxtaposed areas (blue, green and purple) of the AChE binding pocket (B). The first interaction site (green shaded area) situates at the entrance of the binding pocket and primarily engages in  $\pi$ - $\pi$  stacking interaction (C), the second interaction site (blue shaded area) constitutes hydrogen bond interaction between Tyr124 and ketone moiety while hydrophobic interaction between Phe338 and piperidine ring plus central aliphatic chain of 13 (D) and the third interaction site (purple shaded area) constitutes  $\pi$ - $\pi$  stacking between Trp86 and terminal aromatic ring while also affording hydrophobic interaction between Tyr337 and the piperidine moiety of 13 (E).

## CONCLUSION

In conclusion, twelve sets of fingerprint descriptors were used for constructing QSAR models and their performances were comparatively evaluated. It was observed that several fingerprint descriptors afforded

439 good performance for the constructed models indicating that they could capture the feature space of AChE  
440 inhibitors. By taking advantage of the built-in feature importance estimator from RF known as the Gini  
441 index, the following important features that are critical for AChE inhibition were identified: secondary  
442 carbon (SubFPC2), rotatable bond (SubFPC302), aromatic (SubFPC274), C ONS bond (SubFPC295),  
443 secondary mixed amine (SubFPC32) and heterocyclic (SubFPC275). Results from molecular docking  
444 also support the aforementioned findings from the QSAR models in which the aromatic, heteroaromatic  
445 and heterocyclic rings were preferable moieties for interacting with the hydrophobic pocket of AChE.  
446 These can be used as a guide for designing novel inhibitors.

## 447 ACKNOWLEDGEMENTS

448 This research is supported by a Research Career Development Grant (No. RSA5780031) from the  
449 Thailand Research Fund and the Swedish Research Links program (No. C0610701) from the Swedish  
450 Research Council.

## 451 REFERENCES

- 452 Aaviksaar, A. (1990). QSAR in reactions of organophosphorus inhibitors with acetylcholinesterase.  
453 *Phosphorus, Sulfur, and Silicon and the Related Elements*, 51(1-4):47–50.
- 454 Andersson, C. D., Hillgren, J. M., Lindgren, C., Qian, W., Akfur, C., Berg, L., Ekström, F., and Linusson,  
455 A. (2014). Benefits of statistical molecular design, covariance analysis, and reference models in QSAR:  
456 a case study on acetylcholinesterase. *Journal of Computer-Aided Molecular Design*, 29(3):199–215.
- 457 Beal, M. F. (1995). Aging, energy, and oxidative stress in neurodegenerative diseases. *Annals of Neurology*,  
458 38(3):357–366.
- 459 Bembeneck, S. D., Keith, J. M., Letavic, M. A., Apodaca, R., Barbier, A. J., Dvorak, L., Aluisio, L., Miller,  
460 K. L., Lovenberg, T. W., and Carruthers, N. I. (2008). Lead identification of acetylcholinesterase  
461 inhibitors-histamine H3 receptor antagonists from molecular modeling. *Bioorganic & Medicinal  
462 Chemistry*, 16(6):2968–2973.
- 463 Birks, J. (2006). Cholinesterase inhibitors for Alzheimer's disease. *Cochrane Database of Systematic  
464 Reviews*, (1):CD005593.
- 465 Bollback, J. P. (2006). SIMMAP: stochastic character mapping of discrete traits on phylogenies. *BMC  
466 Bioinformatics*, 7:88.
- 467 Bourne, Y., Taylor, P., and Marchot, P. (1995). Acetylcholinesterase inhibition by fasciculin: crystal  
468 structure of the complex. *Cell*, 83(3):503–512.
- 469 Breiman, L. (2001). Random forests. *Machine Learning*, 45(1):5–32.
- 470 Brookmeyer, R., Johnson, E., Ziegler-Graham, K., and Arrighi, H. M. (2007). Forecasting the global  
471 burden of Alzheimer's disease. *Alzheimer's & Dementia*, 3(3):186–191.
- 472 Carhart, R. E., Smith, D. H., and Venkataraghavan, R. (1985). Atom pairs as molecular features in  
473 structure-activity studies: definition and applications. *Journal of Chemical Information and Computer  
474 Sciences*, 25(2):64–73.
- 475 Chen, Y., Fang, L., Peng, S., Liao, H., Lehmann, J., and Zhang, Y. (2012). Discovery of a novel  
476 acetylcholinesterase inhibitor by structure-based virtual screening techniques. *Bioorganic & Medicinal  
477 Chemistry Letters*, 22(9):3181–3187.
- 478 Cheung, J., Rudolph, M. J., Burshteyn, F., Cassidy, M. S., Gary, E. N., Love, J., Franklin, M. C., and  
479 Height, J. J. (2012a). Structures of human acetylcholinesterase in complex with pharmacologically  
480 important ligands. *Journal of Medicinal Chemistry*, 55(22):10282–10286. PMID: 23035744.
- 481 Cheung, J., Rudolph, M. J., Burshteyn, F., Cassidy, M. S., Gary, E. N., Love, J., Franklin, M. C., and  
482 Height, J. J. (2012b). Structures of human acetylcholinesterase in complex with pharmacologically  
483 important ligands. *Journal of Medicinal Chemistry*, 55(22):10282–10286.
- 484 Deb, P. K., Sharma, A., Piplani, P., and Akkinepally, R. R. (2012). Molecular docking and receptor-specific  
485 3D-QSAR studies of acetylcholinesterase inhibitors. *Molecular Diversity*, 16(4):803–823.
- 486 Durant, J. L., Leland, B. A., Henry, D. R., and Nourse, J. G. (2002). Reoptimization of MDL keys for use  
487 in drug discovery. *Journal of Chemical Information and Computer Sciences*, 42(6):1273–1280.
- 488 Dvir, H., Silman, I., Harel, M., Rosenberry, T. L., and Sussman, J. L. (2010). Acetylcholinesterase:  
489 From 3d structure to function. *Chimico-Biological Interactions*, 187(1–3):10 – 22. 10th International  
490 Meeting on Cholinesterases.

- 491 Eriksson, L., Jaworska, J., Worth, A. P., Cronin, M. T., McDowell, R. M., and Gramatica, P. (2003).  
492 Methods for reliability and uncertainty assessment and for applicability evaluations of classification-and  
493 regression-based qsars. *Environmental Health Perspectives*, 111(10):1361.
- 494 Eriksson, L. and Johansson, E. (1996). Multivariate design and modeling in QSAR. *Chemometrics and*  
495 *Intelligent Laboratory Systems*, 34(1):1–19.
- 496 Fang, J., Li, Y., Liu, R., Pang, X., Li, C., Yang, R., He, Y., Lian, W., Liu, A.-L., and Du, G.-H. (2015).  
497 Discovery of multitarget-directed ligands against Alzheimer's disease through systematic prediction of  
498 chemical–protein interactions. *Journal of Chemical Information and Modeling*, 55(1):149–164.
- 499 Gaulton, A., Bellis, L. J., Bento, A. P., Chambers, J., Davies, M., Hersey, A., Light, Y., McGlinchey, S.,  
500 Michalovich, D., Al-Lazikani, B., et al. (2012). ChEMBL: a large-scale bioactivity database for drug  
501 discovery. *Nucleic Acids Research*, 40(D1):D1100–D1107.
- 502 Giacoppo, J. O., CC França, T., Kuča, K., Cunha, E. F. d., Abagyan, R., Mancini, D. T., and Ramalho,  
503 T. C. (2015). Molecular modeling and in vitro reactivation study between the oxime bi-6 and acetyl-  
504 cholinesterase inhibited by different nerve agents. *Journal of Biomolecular Structure and Dynamics*,  
505 33(9):2048–2058.
- 506 Golbraikh, A. and Tropsha, A. (2002). Beware of  $q^2$ ! *Journal of Molecular Graphics and Modelling*,  
507 20(4):269–276.
- 508 Goldblum, A., Yoshimoto, M., and Hansch, C. (1981). Quantitative structure-activity relationship  
509 of phenyl N-methylcarbamate inhibition of acetylcholinesterase. *Journal of Agricultural and Food  
510 Chemistry*, 29(2):277–288.
- 511 Gramatica, P. (2007). Principles of qsar models validation: internal and external. *QSAR & Combinatorial  
512 Science*, 26(5):694–701.
- 513 Gupta, S. and Mohan, C. G. (2014). Dual binding site and selective acetylcholinesterase inhibitors  
514 derived from integrated pharmacophore models and sequential virtual screening. *BioMed Research  
515 International*, 2014:291214.
- 516 Hall, L. H. and Kier, L. B. (1995). Electrotopological state indices for atom types: a novel combination of  
517 electronic, topological, and valence state information. *Journal of Chemical Information and Computer  
518 Sciences*, 35(6):1039–1045.
- 519 Harel, M., Schalk, I., Ehret-Sabatier, L., Bouet, F., Goeldner, M., Hirth, C., Axelsen, P., Silman, I.,  
520 and Sussman, J. (1993). Quaternary ligand binding to aromatic residues in the active-site gorge of  
521 acetylcholinesterase. *Proceedings of the National Academy of Sciences of the United States of America*,  
522 90(19):9031–9035.
- 523 Huang, W., Tang, L., Shi, Y., Huang, S., Xu, L., Sheng, R., Wu, P., Li, J., Zhou, N., and Hu, Y. (2011).  
524 Searching for the multi-target-directed ligands against Alzheimer's disease: Discovery of quinoxaline-  
525 based hybrid compounds with AChE, H3R and BACE 1 inhibitory activities. *Bioorganic & Medicinal  
526 Chemistry*, 19(23):7158–7167.
- 527 J Prado-Prado, F., Escobar, M., and Garcia-Mera, X. (2013). Review of bioinformatics and theoretical  
528 studies of acetylcholinesterase inhibitors. *Current Bioinformatics*, 8(4):496–510.
- 529 James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An introduction to statistical learning: With  
530 Applications in R*. Springer, New York.
- 531 Jianxin Guo, Margaret M. Hurley, Jeffery B. Wright, , , and Gerald H. Lushington\*, 2004). A docking  
532 score function for estimating ligandprotein interactions: application to acetylcholinesterase inhibition.  
533 *Journal of Medicinal Chemistry*, 47(22):5492–5500. PMID: 15481986.
- 534 Kennard, R. W. and Stone, L. A. (1969). Computer aided design of experiments. *Technometrics*,  
535 11(1):137–148.
- 536 Kitz, R. and Wilson, I. B. (1962). Esters of methanesulfonic acid as irreversible inhibitors of acetyl-  
537 cholinesterase. *Journal of Biological Chemistry*, 237(10):3245–3249.
- 538 Klekota, J. and Roth, F. P. (2008). Chemical substructures that enrich for biological activity. *Bioinformatics*,  
539 24(21):2518–2525.
- 540 Kryger, G., Silman, I., and Sussman, J. L. (1999). Structure of acetylcholinesterase complexed with  
541 E2020 (aricept®): implications for the design of new anti-Alzheimer drugs. *Structure*, 7(3):297–307.
- 542 Kuca, K., Soukup, O., Maresova, P., Korabecny, J., Nepovimova, E., Klimova, B., Honegr, J., Ramalho,  
543 T. C., and França, T. C. (2016). Current approaches against alzheimer's disease in clinical trials. *Journal  
544 of the Brazilian Chemical Society*, 27(4):641–649.
- 545 Kuhn, M. (2008). Building predictive models in R using the caret package. *Journal of Statistical Software*,

- 546 28(5):1–26.
- 547 Laggner, C. (2005). Smarts patterns for functional group classification. *Inte: Ligand Software-*  
548 *Entwicklungs und Consulting GmbH, Maria Enzersdorf, Austria.*
- 549 Lee, S. and Barron, M. G. (2016). A mechanism-based 3D-QSAR approach for classification and  
550 prediction of acetylcholinesterase inhibitory potency of organophosphate and carbamate analogs.  
551 *Journal of Computer-Aided Molecular Design*, 30(4):347–363.
- 552 Lu, S.-H., Wu, J. W., Liu, H.-L., Zhao, J.-H., Liu, K.-T., Chuang, C.-K., Lin, H.-Y., Tsai, W.-B., and Ho,  
553 Y. (2011). The discovery of potential acetylcholinesterase inhibitors: a combination of pharmacophore  
554 modeling, virtual screening, and molecular docking studies. *Journal of Biomedical Science*, 18(8):b22.
- 555 Mager, P. (1983). QSAR applied to aging of phosphorylated acetylcholinesterase. *Pharmazie*, 38(4):271–  
556 272.
- 557 Morris, J. C., Cyrus, P. A., Orazem, J., Mas, J., Bieber, F., Ruzicka, B. B., and Gulanski, B. (1998).  
558 Metrifonate benefits cognitive, behavioral, and global function in patients with Alzheimer's disease.  
559 *Neurology*, 50(5):1222–1230.
- 560 Mundy, R., Bowman, M., Farmer, J., and Haley, T. (1978). Quantitative structure activity study of a series  
561 of substituted 0,0-dimethyl 0-(p-nitrophenyl) phosphorothioates and 0-analogs. *Archives of Toxicology*,  
562 41(2):111–123.
- 563 Nantasesamat, C., Isarankura-Na-Ayudhya, C., Naenna, T., and Prachayassitkul, V. (2009). A practical  
564 overview of quantitative structure-activity relationship. *EXCLI J*, 8(7):74–88.
- 565 NCBI (2009). *PubChem Substructure Fingerprint, version 1.3.* [ftp://ftp.ncbi.nlm.nih.gov/pubchem/specifications/pubchem\\_fingerprints.txt](ftp://ftp.ncbi.nlm.nih.gov/pubchem/specifications/pubchem_fingerprints.txt).
- 566 Ordentlich, A., Barak, D., Kronman, C., Flashner, Y., Leitner, M., Segall, Y., Ariel, N., Cohen, S., Velan,  
567 B., and Shafferman, A. (1993). Dissection of the human acetylcholinesterase active center determinants  
568 of substrate specificity. Identification of residues constituting the anionic site, the hydrophobic site, and  
569 the acyl pocket. *Journal of Biological Chemistry*, 268(23):17083–17095.
- 570 Prado-Prado, F., Garcia-Mera, X., Escobar, M., Alonso, N., Caamano, O., Yanez, M., and Gonzalez-  
571 Diaz, H. (2012). 3d mi-dragon: new model for the reconstruction of us fda drug-target network and  
572 theoretical-experimental studies of inhibitors of rasagiline derivatives for ache. *Current Topics in  
573 Medicinal Chemistry*, 12(16):1843–1865.
- 574 Puzyn, T., Mostrag-Szlichtyng, A., Gajewicz, A., Skrzynski, M., and Worth, A. P. (2011). Investigating  
575 the influence of data splitting on the predictive ability of QSAR/QSPR models. *Structural Chemistry*,  
576 22(4):795–804.
- 577 Quinn, D. M. (1987). Acetylcholinesterase: enzyme structure, reaction dynamics, and virtual transition  
578 states. *Chemical Reviews*, 87(5):955–979.
- 579 Racchi, M., Mazzucchelli, M., Porrello, E., Lanni, C., and Govoni, S. (2004). Acetylcholinesterase  
580 inhibitors: novel activities of old molecules. *Pharmacological Research*, 50(4):441–451.
- 581 Riniker, S. and Landrum, G. A. (2013). Open-source platform to benchmark fingerprints for ligand-based  
582 virtual screening. *Journal of Cheminformatics*, 5(1):1–17.
- 583 Roy, K., Chakraborty, P., Mitra, I., Ojha, P. K., Kar, S., and Das, R. N. (2013). Some case studies  
584 on application of “rm2” metrics for judging quality of quantitative structure–activity relationship  
585 predictions: emphasis on scaling of response data. *Journal of Computational Chemistry*, 34(12):1071–  
586 1082.
- 587 Schrödinger, LLC (2015a). PyMOL, version 1.7.6.3. New York, NY.
- 588 Schrödinger, LLC (2015b). Schrödinger Release 2015-4: Maestro, version 10.4. New York, NY.
- 589 Shen, T., Tai, K., Henchman, R. H., , and McCammon, J. A. (2002). Molecular dynamics of acetyl-  
590 cholinesterase. *Accounts of Chemical Research*, 35(6):332–340.
- 591 Silman, I. and Sussman, J. L. (2008). Acetylcholinesterase: How is structure related to function?  
592 *Chemico-Biological Interactions*, 175(1–3):3 – 10. Proceedings of the {IX} International Meeting on  
593 Cholinesterases.
- 594 Simeon, S., Möller, R., Almgren, D., Li, H., Phanus-umporn, C., Prachayassitkul, V., Bülow, L., and  
595 Nantasesamat, C. (2016). Unraveling the origin of splice switching activity of hemoglobin  $\beta$ -globin  
596 gene modulators via qsar modeling. *Chemometrics and Intelligent Laboratory Systems*, 151:51–60.
- 597 Steinbeck, C., Han, Y., Kuhn, S., Horlacher, O., Luttmann, E., and Willighagen, E. (2003). The Chemistry  
598 Development Kit (CDK): An open-source Java library for chemo- and bioinformatics. *Journal of  
599 Chemical Information and Computer Sciences*, 43(2):493–500.
- 600

- 601 Su, C. and Lien, E. (1980). QSAR of acetylcholinesterase inhibitors: A reexamination of the role of  
602 charge-transfer. *Research Communications in Chemical Pathology and Pharmacology*, 29(3):403–415.
- 603 Tougu, V. (2001). Acetylcholinesterase: mechanism of catalysis and inhibition. *Current Medicinal  
604 Chemistry-Central Nervous System Agents*, 1(2):155–170.
- 605 Tropsha, A., Gramatica, P., and Gombar, V. K. (2003). The importance of being earnest: validation is the  
606 absolute essential for successful application and interpretation of qspr models. *QSAR & Combinatorial  
607 Science*, 22(1):69–77.
- 608 Trott, O. and Olson, A. J. (2010). AutoDock Vina: improving the speed and accuracy of docking with a  
609 new scoring function, efficient optimization, and multithreading. *Journal of Computational Chemistry*,  
610 31(2):455–461.
- 611 Utrecht, J. and Trager, W. (2007). *Drug Metabolism: Chemical and Enzymatic Aspects*. CRC Press,  
612 Boca Raton, Florida.
- 613 Veber, D. F., Johnson, S. R., Cheng, H.-Y., Smith, B. R., Ward, K. W., and Kopple, K. D. (2002).  
614 Molecular properties that influence the oral bioavailability of drug candidates. *Journal of Medicinal  
615 Chemistry*, 45(12):2615–2623.
- 616 Veselinović, J., Nikolić, G., Trutić, N., Živković, J., and Veselinović, A. (2015). Monte Carlo QSAR mod-  
617 els for predicting organophosphate inhibition of acetylcholinesterase. *SAR and QSAR in Environmental  
618 Research*, 26(6):449–460.
- 619 Wallace, A. C., Laskowski, R. A., and Thornton, J. M. (1995). LIGPLOT: a program to generate schematic  
620 diagrams of protein-ligand interactions. *Protein Engineering*, 8(2):127–134.
- 621 Walters, D. and Hopfinger, A. (1986). Case studies of the application of molecular shape analysis to  
622 elucidate drug action. *Journal of Molecular Structure: THEOCHEM*, 134(3-4):317–323.
- 623 Wright, M. N. and Ziegler, A. (2015). ranger: A fast implementation of random forests for high  
624 dimensional data in C++ and R. *arXiv e-print*. <http://arxiv.org/abs/1508.04409>.
- 625 Yan, A. and Wang, K. (2012). Quantitative structure and bioactivity relationship study on human  
626 acetylcholinesterase inhibitors. *Bioorganic & Medicinal Chemistry Letters*, 22(9):3336–3342.
- 627 Yap, C. W. (2011). PaDEL-descriptor: An open source software to calculate molecular descriptors and  
628 fingerprints. *Journal of Computational Chemistry*, 32(7):1466–1474.
- 629 Zhou, Y., Wang, S., and Zhang, Y. (2010). Catalytic reaction mechanism of acetylcholinesterase deter-  
630 mined by bornoppenheimer ab initio qm/mm molecular dynamics simulations. *Journal of Physical  
631 Chemistry B*, 114(26):8817–8825. PMID: 20550161.