



**Nano-QSAR modeling for predicting biological activity of diverse nanomaterials**

Journal:	<i>RSC Advances</i>
Manuscript ID:	RA-ART-02-2014-001274
Article Type:	Paper
Date Submitted by the Author:	04-Jan-2014
Complete List of Authors:	Singh, Kunwar; CSIR-Indian Institute of Toxicology Research, Environmental Chemistry Gupta, Shikha; CSIR-Indian Institute of Toxicology Research, Environmental Chemistry

1 **Nano-QSAR modeling for predicting biological activity of diverse**  
2 **nanomaterials**

3  
4 **Kunwar P. Singh\*<sup>1,2</sup>, Shikha Gupta<sup>1,2</sup>**

5  
6 **<sup>1</sup>Academy of Scientific and Innovative Research,**  
7 **Anusandhan Bhawan, Rafi Marg, New Delhi-110 001, India**

8 **<sup>2</sup>Environmental Chemistry Division,**  
9 **CSIR-Indian Institute of Toxicology Research**  
10 **Post Box 80, Mahatma Gandhi Marg, Lucknow-226 001, India**

11 **\*Author for all correspondence**

12 Tel: 0091-522-2476091; Fax: 0091-522-2628227

13 E-mail: [kpsingh\\_52@yahoo.com](mailto:kpsingh_52@yahoo.com); [kunwarpsingh@gmail.com](mailto:kunwarpsingh@gmail.com)

14

15

**Abstract**

16 This study reports robust reliable ensemble learning (EL) approach based nano-QSAR models  
17 for predicting the biological effects of diverse nanomaterials (NMs) using simple molecular  
18 descriptors. EL based nano-QSAR models implementing stochastic gradient boosting and  
19 bagging algorithms were constructed and used to establish statistically significant relationships  
20 between measured biological activity profiles of nanoparticles (NPs) and their simple structural  
21 properties. To demonstrate the predictive ability of the developed nano-QSAR models, five  
22 different representative data sets (case studies) of NMs (NPs with diverse metal cores, NPs with  
23 similar core but diverse surface modifiers, metal oxide NPs, surface modified multi-walled  
24 carbon nanotubes, and fullerene derivatives) studied recently using *in vitro* cell based assays  
25 were employed. Rigorous validation of the constructed classification and regression nano-QSAR  
26 models performed using various statistical parameters suggested robustness of the EL based  
27 models for their future use. Proposed nano-QSAR models showed high prediction accuracy  
28 (binary classification) of more than 93.18 % (case study 1), 97.25 % (case study 2), and yielded  
29 correlation ( $R^2$ ) of more than 0.851 between experimental and model predicted values of  
30 biological activity in complete data of different diverse sets of NPs. Results for all the five case  
31 studies demonstrated better predictive performance of the proposed nano-QSAR models  
32 compared to the previous studies. The proposed models reliably predicted the biological activity  
33 of all considered NPs, and the methodology is expected to provide guidance for the future design  
34 and manufacturing of NMs ensuring better and safer products.

35

36 **Keywords:** Nanomaterials, ensemble learning, nano-QSAR model, biological activity,  
37 nanostructure, modeling

## 38 1. Introduction

39 In recent years, nanomaterials (NMs) have gained much importance due to their  
40 widespread applications in different areas. These materials are used in a variety of fields due to  
41 their unique physical and chemical properties, such as shape, size distribution, surface area and  
42 structure, overall charge, porosity, agglomeration rate, and surface chemistry.<sup>1</sup> Currently, these  
43 materials can be designed to achieve desired properties and are used in electronics, opto-  
44 electronics, biomedical, environmental, material and energy related areas, cosmetics,  
45 pharmaceuticals and catalysts.<sup>2</sup> Moreover, the use of NMs in various industries is projected to  
46 increase dramatically in the future and as a consequence, contamination of environment by these  
47 materials is expected, or at least such possibility cannot be disregarded. Manufactured NMs  
48 intended for industrial applications may cause toxic effects in humans and public concern about  
49 the safety of these materials is increasing.<sup>3</sup> Recently, there appeared some reports in literature on  
50 the adverse effects of nanoparticles (NPs) on humans and environment.<sup>4,5</sup> Acute or repeated  
51 exposure to NPs present in commercial products may potentially cause systemic, cellular and/or  
52 genomic toxicities. There remain scientific gaps in understanding of toxicology of NMs that  
53 these are already contained in commercial products not intended for human exposure, could  
54 contaminate the environment while also not intended for human exposure, and are intended for  
55 biomedical applications such as drug delivery, imaging and sensing.<sup>6</sup> Thus, understanding the  
56 biological effects of exposure to NPs is of paramount importance. There is still limited  
57 information about experimentally measured toxic effects of NPs and some isolated toxicity  
58 studies considering single or a few NPs are published in last few years.<sup>7-20</sup> Experimental studies,  
59 especially toxicological, are time-consuming, costly, unethical, and often impractical, calling for  
60 the development of efficient computational approaches capable of predicting biological effects of

61 NMs. Thus, it is imperative to develop a comprehensive, and ideally, predictive knowledge of  
62 the effects of NPs on the environment as well as animals and humans. Modeling the biological  
63 effects of NPs is difficult task due to their structural heterogeneity, complexity and diversity and  
64 reports on computational modeling of nano toxicology are scarce.<sup>21</sup> These computational  
65 methods are based on the assumption that the variation in the properties or biological activities of  
66 a NP can be correlated with changes in its molecular structure and can be used to predict the  
67 activity/property of newly synthesized NPs without resorting to experimentation.<sup>1</sup> In most cases,  
68 the exact composition of a given NP is not known and three-dimensional nano-structures may  
69 include large number of highly complex atoms leading to stoichiometric variations between the  
70 NPs, rendering the classical molecular descriptors inappropriate for modeling. However, a few  
71 attempts have been made to develop quantitative-structure activity/toxicity relationships  
72 (QSARs/QSTRs) to correlate molecular nano-structures with activities of NPs using the limited  
73 data available in literature. QSARs for predicting nano-toxicity/biological activity of 48 fullerene  
74 derivatives,<sup>22-25</sup> 17 metal oxides NPs,<sup>2,26</sup> 51 NMs possessing varying core metal compositions,  
75 coatings and surface attachments,<sup>6,20,21,27</sup> 109 NPs with similar metal core with diverse surface  
76 modifiers,<sup>1,6,27-29</sup> and 80 surface modified multi-wall carbon nanotubes<sup>30,31</sup> have recently been  
77 reported. Modeling techniques, such as means of balance of correlation (MBC),<sup>22,25</sup> multiple  
78 linear regression (MLR),<sup>2,27,31</sup> logistic regression (LR), Naïve bayes (NB),<sup>29</sup> k-nearest neighbor  
79 (k-NN),<sup>6,29</sup> partial least square regression (PLSR),<sup>23,24</sup> multi-layered perceptron neural network  
80 (MLPN),<sup>1,27</sup> support vector machines (SVM)<sup>6,29</sup> have been found useful for the establishment of  
81 the relationships between the molecular structures and biological activities of NPs. Although, the  
82 predictive responses achieved using these modeling techniques have been within acceptable  
83 range, these methods have certain limitation. Linear regression methods do not fit the data with

84 nonlinear structure, a common feature of experimental toxicity data.<sup>32</sup> SVM uses only a limited  
85 data during model building phase.<sup>32</sup> MLPN, although, a universal nonlinear method, it suffers  
86 from over-fitting problem in training.<sup>32</sup> Therefore, keeping in view the rapidly emerging scope  
87 and applications of the NMs, there is a need to develop more précised and robust methods  
88 capable of predicting the nano-toxicities of various types of materials, which could help in  
89 designing safer materials.

90 Ensemble learning (EL) methods have emerged as powerful tools for mapping the  
91 relationship between the response and predictors, and have not yet been used for predicting the  
92 biological activity of NPs so far. EL-based techniques are applicable to both classification and  
93 regression problems.<sup>33</sup> These methods overcome problems with weak predictors<sup>34</sup> and have the  
94 advantage to alleviate the small sample size problem by averaging and incorporating over  
95 multiple models to reduce the potential for over-fitting the training data.<sup>35</sup> EL methods with  
96 bagging and stochastic gradient boosting techniques improve the prediction accuracy of weak  
97 learners.<sup>36</sup> The bagging minimizes prediction variance by generating bootstrapped replica data  
98 sets, whereas, boosting creates a linear combination out of many models, where each new model  
99 is dependent on the preceding model.<sup>37</sup> Decision tree forest (DTF) and decision treeboost (DTB)  
100 implementing bagging and boosting techniques, respectively are inherently non-parametric  
101 statistical methods and make no assumption regarding the underlying distribution of the values  
102 of predictor variables and can handle numerical data that are highly skewed or multi-model in  
103 nature.<sup>38</sup>

104 Selection of appropriate descriptors in toxicity prediction is yet another important issue.  
105 A large number and variety of such descriptors have been used in earlier studies, generally  
106 derived through highly complicated semi-empirical and empirical methods based on quantum

107 mechanical calculations.<sup>2,6</sup> Hence, it would be desirable to develop toxicologically relevant EL-  
108 based nano-QSARs to relate set of simple structural descriptors characterizing NPs with their  
109 measured biological effects (smooth muscle apoptosis, cellular uptake, cytotoxicity, cell  
110 viability).

111 In this study, the basic objectives were to construct the EL-based classification and  
112 regression nano-QSAR models (DTB and DTF) for predicting the biological effects of diverse  
113 NPs using simple structural descriptors. Accordingly, five different datasets (a) fifty one various  
114 NMs with diverse metal cores,<sup>20</sup> (b) one hundred nine NPs with similar core but diverse surface  
115 modifiers,<sup>6</sup> (c) seventeen diverse metal oxide NPs,<sup>2</sup> (d) eighty surface modified multi-walled  
116 carbon nanotubes,<sup>39</sup> and (e) forty eight different fullerene derivatives<sup>24</sup> available in literature for  
117 QSARs analysis were considered. QSAR calculations led to statistically validated and externally  
118 predictive models; these models quantitatively relate the structural properties of NPs with their  
119 experimentally measured biological effects in different cell based assays. To the best of the  
120 knowledge of the authors, this report is the first example of EL-based nano-QSARs analysis of  
121 different sets of NPs successfully demonstrates the high potential of proposed modeling  
122 approaches for improving the experimental design and prioritizing the biological testing of novel  
123 NPs.

124

## 125 **2. Methods**

### 126 **2.1 Datasets**

127 In this study, data from multiple sources were considered for the analysis. For developing  
128 predictive EL-based nano-QSAR models, five different datasets on biological activities of  
129 diverse NPs were used.

130 **Case study 1: NMs with diverse metal cores-** Shaw et al.<sup>20</sup> reported a study on the effect of 51  
131 different NMs (with diverse metal cores) in four cell lines (endothelial and smooth muscle cells,  
132 monocytes, and hepatocytes), using four biological assays (ATP content, reducing equivalents,  
133 caspase-mediated apoptosis, and mitochondrial membrane potential) in each cell line, at four  
134 concentrations per assay. These experiments generated potentially 64 biological response  
135 variables for each of the NMs. Of the possible combinations of biological assays and cell types,  
136 only the apoptosis assays (smooth muscle cell apoptosis, SMA) exhibited dose-response  
137 relationship. Similar to Epa et al.<sup>27</sup>, the slope of the dose-response curve (SMA) was considered  
138 as a dependent variable for predictive regression modeling. Moreover, for 44 of these NMs, four  
139 structural descriptors (size, relaxivities, R1, R2, and zeta potential) were available. Fourches et  
140 al.<sup>6</sup> calculated arithmetic mean of biological activity profile (64 features) designating as Z-mean,  
141 which was taken as basis for binary classification ( $Z\text{-mean} > -0.40$ ; class 1, and  $Z\text{-mean} < -$   
142  $0.40$ , class 0), rendering 22 NMs in each class.

143

144 **Case study 2: NPs with common metal core-** The dataset comprised of 109 NPs in which a  
145 supermagnetic NP (cross-linked iron oxide with amine group) was decorated with different  
146 synthetic small molecules.<sup>7</sup> NPs were made magnetofluorescent with the addition of fluorescence  
147 isothiocyanate molecules on their surfaces to enable measurement of cell uptake. All the NPs in  
148 the dataset have exactly the same metal core decorated with different synthetic small molecules.  
149 Each NP is represented by the structure of organic surface modifier, which in turn is  
150 characterized by conventional molecular descriptors. Then, NPs were screened against human  
151 pancreatic cancer cells (PaCa2). Cellular uptake is expressed as decadic logarithm of the  
152 concentration (pM) of NPs per cell, which varied from 2.23 to 4.44. For binary classification, a



153 criterion of Chau and Yap<sup>29</sup> was considered. According to this criteria, the NPs having cellular  
154 uptake of more than 5000 NPs per cell were considered to have good/moderate cellular uptake  
155 (positive class), while NPs with cellular uptake of less than 5000 particles per cell were  
156 considered to have poor cellular uptake (negative class). Thus, 59 NPs belong to positive and  
157 remaining 50 NPs were in negative class.

158

159 **Case study 3: Diverse metal oxides NPs** - The dataset contains 17 different diverse metal  
160 oxides based NPs<sup>2</sup> with sizes ranging from 15 to 90 nm reporting their cytotoxicity in  
161 *Escherichia coli* bacteria and expressed in terms of the logarithmic values of molar  $1/EC_{50}$   
162 (effective concentration of a given oxide that reduces bacterial viability by 50%), which varied  
163 from 1.74 to 3.51 mol L<sup>-1</sup>.

164

165 **Case study 4: Surface modified multi-walled carbon nanotubes CNTs** - The dataset contains  
166 80 distinct surface modified multi-walled carbon nanotubes,<sup>39</sup> where the surface decorators were  
167 made from a combination of eight amines and nine acylators with a common linking group to the  
168 nanotube. Zhou et al.<sup>39</sup> tested these 80 decorator-nanotube complexes (DNC) for their six  
169 different end-points and evaluated acute cytotoxicity (cell) of the DNC library in macrophages  
170 using WST-1 assay.<sup>40</sup> Cell viability was measured by determining the mitochondrial  
171 dehydrogenase activity. As described above, only the 29 most nanotoxic DNC based upon the  
172 cumulative index over all six end-points<sup>31,39</sup> were retained here for modeling. The experimental  
173 cell viability varied from 1 to 80.

174

175 **Case study 5: Fullerene derivatives NPs** – The data on measured binding affinity (minus  
176 decimal logarithm of the 50 % effective concentration,  $pEC_{50}$ ,  $\mu M$ ) for 48 different fullerene  
177 derivatives with the HIV-1 PR (human immunodeficiency virus type 1 aspartic protease) were  
178 taken from Durdagi al.<sup>24</sup>. The binding affinities were assessed by quantitative assay, based on the  
179 estimated binding energies of fullerene analogous with HIV-1 PR which were determined by  
180 molecular docking. The measured  $pEC_{50}$  values for the considered fullerene derivatives varied  
181 between 2.25 and 8.70.

182 Histograms of the experimental values of the biological activity end-points under  
183 different case studies are plotted in Fig. 1a-e. A histogram consists of tabular frequencies, shown  
184 as adjacent rectangles, erected over discrete intervals, with an area equal to the frequency of the  
185 observations in the interval. The height of a rectangle is also equal to the frequency density of the  
186 interval. The total area of the histogram is equal to the number of data. From the plotted  
187 histograms, it may be noted that the end-point values in Fig. 1a,b,e show nearly normal  
188 distribution pattern, whereas, those in Fig. 1c,d show multi-model distribution. In multi-model  
189 distribution several processes with normal distributions are combined, because there are many  
190 peaks close together, the top of the distribution resembles a plateau.

### 191 **Figure 1**

## 192 **2.2 Molecular descriptors, feature selection and data processing**

193 Molecular descriptors are the simple mathematical representation of a molecule and are  
194 used to encode significant features of molecules. In case of first dataset containing 44 NMs, four  
195 descriptors available in literature were used here for classification and regression modeling. In  
196 case of second (109 NPs) and fourth (29 DNC) datasets, 174 molecular descriptors (topological,  
197 electronic, geometrical, and constitutional) were calculated for each NP using Chemistry

198 Development Kit (CDK v 1.0.3).<sup>41</sup> For the third (17 NPs) and fifth (48 fullerene derivatives)  
199 dataset, 32 molecular descriptors (topological, geometrical, and constitutional) were calculated  
200 using Chemspider.<sup>42</sup> The electronic, constitutional, geometrical and topological descriptors were  
201 calculated by 2D structures of the molecules, which were taken in the form of SMILES  
202 (simplified molecular input line entry system). For case studies 2, 3 and 5, the SMILES were  
203 taken from the literature,<sup>2,6,25</sup> whereas for case study 4, molecular structures of the decorator  
204 portions of DNCs were taken from Shao et. al.<sup>31</sup> and SMILES were obtained using the  
205 ChemDoodle program (trial version). Chemical structures of the NPs were drawn in  
206 ChemDoodle program using the SMILES (Table SI1-SI4, in the Supplementary Information).

207 Since, all the molecular descriptors may not be relevant to the nano-QSARs analysis;  
208 elimination of less significant descriptors can improve the accuracy of prediction, and facilitate  
209 the interpretation of the model through focusing on the most relevant variables. For selection of  
210 initial features, model-fitting approaches were considered. In all the case studies, calculated  
211 descriptors were analyzed for the existence of a constant or near constant values and the  
212 descriptors with low variation were excluded from the original pool of descriptors. EL-based  
213 QSAR modeling was then performed. For optimal values of the model parameters, the EL  
214 models were trained by using the set of remaining features computing the respective scoring  
215 functions to rank the contribution of features in the current set. The lowest ranked features were  
216 then removed.<sup>32</sup> The EL-QSAR models were retained by using the remaining set of features, and  
217 the corresponding prediction accuracies (classification accuracy, and root mean squared error of  
218 prediction) were computed by means of 5-fold cross validation. Distribution of selected  
219 descriptors for different case studies considered for nano-QSARs modeling in this study are  
220 shown in the radar charts (Fig.2). A radar chart is a graphical method that displays multivariate

221 data in the form of a two-dimensional chart with several quantitative variables represented on  
 222 axis starting from the same point. The radar charts analysis shows that the NMs used in our case  
 223 studies covered a sufficiently large structural space.

224 **Figure 2**

225 Since the aim of present study is to develop robust models capable of making accurate and  
 226 reliable predictions of biological activities of new NMs, the QSAR model derived from a  
 227 training set should be validated/tested using new moieties for checking its predictive ability. The  
 228 validation strategies check the reliabilities of models for their possible application on a new  
 229 dataset, and confidence in the prediction can thus be judged. In this study, for classification and  
 230 regression modeling, data were split into training (80 %) and test (20 %) subsets using random  
 231 distribution approach. Such test sets (when defined prior to analysis) come close to external  
 232 validation set, which are commonly accepted as the gold standard to assess real predictivity.<sup>43</sup>

233

### 234 2.3 Structural diversity

235 The diversity of a dataset is very important for global model development.<sup>44</sup> Structural  
 236 diversity of the NPs can be measured by using the Tanimoto similarity index (TSI), which is an  
 237 appropriate distance metric for topology-based chemical similarity studies. In this method the  
 238 structure of the chemical compounds to be compared are decomposed into fragments. Each  
 239 chemical structure thus characterized by a vector  $\mathbf{y}$  with components  $y(j)$  being binary  
 240 substructure descriptors. The similarity of two structures, represented by vectors  $\mathbf{y}_A$  and  $\mathbf{y}_B$  can  
 241 be characterized<sup>45</sup> by the TSI,  $t_{A,B}$  as;  $t_{A,B} = \frac{\mathbf{y}_A^T \cdot \mathbf{y}_B}{\mathbf{y}_A^T \cdot \mathbf{1} + \mathbf{y}_B^T \cdot \mathbf{1} - \mathbf{y}_A^T \cdot \mathbf{y}_B}$ . The TSI ranges from 0 (no  
 242 similarity) to 1 (pair-wise similarity). Smaller TSI means compounds have good diversity.<sup>46</sup> A  
 243 good cut-off for biologically similar molecules is 0.7 or 0.8. TSI values for second, third, fourth,

244 and fifth datasets are 0.12, 0.21, 0.17, 0.10, respectively. These values suggest that the datasets  
245 used in this work represent NPs with sufficiently high structural diversity. It warrants model  
246 stability and that the external test set is suitable to assess the predictive performance of the  
247 developed model.

248

## 249 **2.4 Ensemble learning based nano-QSAR modeling**

250 An ensemble contains a number of base learners<sup>47</sup> and their generalization ability is  
251 usually much stronger. Stochastic gradient boosting and bagging algorithms are implemented  
252 here for constructing the classification and regression nano-QSAR models (DTB, DTF). Brief  
253 description of these methods is given below.

254

### 255 **2.4.1 DTB-nano-QSAR model**

256 In DTB, stochastic gradient boosting improves the accuracy of a predictive function by  
257 applying it repeatedly in a series and combining the output of each function with weighting, so  
258 that the total error of prediction is minimized.<sup>37</sup> The DTB algorithm creates a tree ensemble and  
259 it uses randomization during the tree creations (Fig. 3a). The goal is to minimize the loss  
260 function in the training set,  $\{\mathbf{x}, \mathbf{y}\}$ . After each iteration,  $F$  represents sum of all trees built so far:  
261  $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \text{Tree}_m(\mathbf{x})$ , where  $m$  is the number of trees in the model. Regardless of the loss-  
262 function, the trees fitting the gradient on pseudo residuals are regression trees trained to  
263 minimize mean squared error (MSE). Optimal size of the tree was decided using the criteria of  
264 minimal cross-validation error. The DTB model for classification is essentially the same as for  
265 regression except logit (probability) values are fitted rather than raw target values.

266

### **Figure 3a**

#### 267 **2.4.2 DTF-nano-QSAR model**

268 In DTF, a large number of independent trees are grown in parallel, and they do not  
269 interact until after all of them have been built (Fig. 3b). Different training sub-sets are drawn at  
270 random with replacement from the training dataset. Separate models are produced and used to  
271 predict the entire data from aforesaid sub-sets. Then various estimated models are aggregated. In  
272 bagging, a bootstrapped sample is constructed.<sup>48</sup> The DTFs use the out of bag data rows for  
273 model validation. This provides an independent test set without requiring a separate data set or  
274 holding back rows from the tree construction. The DTF algorithm makes it highly resistant to  
275 over-fitting.

#### 276 **Figure 3b**

#### 277 **2.5 Model validation and prediction verification**

278 The optimal architectures and parameters of the EL-based classification and regression  
279 nano-QSAR models constructed here were determined following both the internal and external  
280 validation procedures. For internal validation, a V-fold cross validation (CV) method was  
281 adopted. The V-fold CV is the most common procedure recommended to check the  
282 generalization ability of the model.<sup>43</sup> The advantage of this method is that it performs reliable  
283 and unbiased testing on dataset. For external validation, a separate validation (test) sub-set of the  
284 data was used which was kept out during the training process.<sup>32</sup> In case of the predictive models,  
285 validation step using external data set provides information about the predictive ability of the  
286 trained model for the unknown data.<sup>46</sup> Benigni et al.<sup>49</sup> pointed out that the prediction reliability  
287 should be checked by means of an external test set with new moieties not used in model building.  
288 Optimal models were selected on the basis of the classification accuracy (classification) and root  
289 mean squared error (regression) in the training and validation data.<sup>50</sup> Predictive performance of

290 the regression models constructed here for external sets was evaluated using various OECD  
 291 recommended validation criteria parameters proposed in QSAR literature,<sup>51</sup> such as  $Q_{F1}^2$ ,<sup>52</sup>  $Q_{F2}^2$ ,<sup>53</sup>  
 292  $Q_{F3}^2$ <sup>54</sup> and concordance correlation coefficient (CCC).<sup>55</sup>  $Q_{F1}^2$  uses average of the training data,  
 293 instead of that of the prediction set ( $Q_{F1}^2 = 1 - \frac{\sum_{i=1}^{n_{\text{ext}}}(y_i - \bar{y}_{\text{tr}})^2}{\sum_{i=1}^{n_{\text{ext}}}(y_i - \bar{y}_{\text{tr}})^2}$ , where  $n_{\text{ext}}$  is the number of compounds  
 294 in external (test) set,  $y_i$  and  $\hat{y}_i$  are the observed and model calculated value of the dependent  
 295 variable in external set, and  $\bar{y}_{\text{tr}}$  is the mean value of the dependent variable in training set),  
 296 whereas,  $Q_{F2}^2$  takes no account of the distance from the average of the training values  
 297 ( $Q_{F2}^2 = 1 - \frac{\sum_{i=1}^{n_{\text{ext}}}(y_i - \bar{y}_{\text{ext}})^2}{\sum_{i=1}^{n_{\text{ext}}}(y_i - \bar{y}_{\text{ext}})^2}$ , where,  $\bar{y}_{\text{ext}}$  represents the mean value of the dependent variable in  
 298 external (test) set). In  $Q_{F3}^2$  the denominator is calculated on the training set, and both numerator  
 299 and denominator are divided by the number of corresponding elements  
 300 ( $Q_{F3}^2 = 1 - \frac{[\sum_{i=1}^{n_{\text{ext}}}(y_i - \hat{y}_i)^2]/n_{\text{ext}}}{[\sum_{i=1}^{n_{\text{tr}}}(y_i - \bar{y}_{\text{tr}})^2]/n_{\text{tr}}}$ , where,  $n_{\text{tr}}$  is the number of compounds in training set). Consonni et  
 301 al.<sup>54</sup> demonstrated that results obtained by  $Q_{F3}^2$  are independent of the prediction set distribution  
 302 and sample size. CCC ( $\text{CCC} = \frac{2 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n (y_i - \bar{y})^2 + n(\bar{x} - \bar{y})^2}$ , where,  $x$  and  $y$  correspond to the  
 303 abscissa and ordinate value of the graph plotting the prediction experimental data values vs. the  
 304 ones calculated using the model,  $n$  is the number of chemicals, and  $\bar{x}$  and  $\bar{y}$  correspond to the  
 305 averages of abscissa and ordinate values) measures both precision and accuracy and involves no  
 306 training set information, so it can be considered a true external validation measure, independent  
 307 of the samples chemical space. In all simulations, the validation measures are calculated only if  
 308  $R^2 > 0.7$ .

309 Different statistical parameters were used to evaluate the performance of constructed  
310 nano-QSAR models. For binary classification, statistical parameters, such as sensitivity,  
311 specificity, accuracy and Matthews's correlation coefficient (MCC) are considered.<sup>32</sup> Sensitivity  
312 denotes the percentage of correctly classified active NPs among the total number of active NPs,  
313 whereas specificity is the percentage of correctly classified inactive NPs among the total number  
314 of inactive NPs. Accuracy represents the total number of active and inactive NPs correctly  
315 predicted among the total number of tested NPs. Performance of the regression models  
316 constructed here was evaluated using the mean absolute error (MAE), root mean squared error  
317 (RMSE), squared correlation coefficient ( $R^2$ ) between the measured and predicted values of the  
318 response.<sup>56</sup> The RMSE represents the error associated with the model. It is a measure of the  
319 goodness-of-fit, best describes an average measure of the error in predicting the dependent  
320 variable.

321

## 322 **2.6 Applicability domain of the EL nano-QSAR model**

323 The applicability domain (AD) of a predictive model defines the boundaries whereby the  
324 predicted values can be trusted with confidence. The AD was taken into account in order to  
325 consider the scope and limitations of the proposed models, i.e. the range of chemical structures  
326 for which the models are considered to be applicable.<sup>57</sup> This approach is based on the ranges of  
327 individual descriptors used for the model building. According to this method, a NP with  
328 descriptor values within the range of those of the training set NPs is considered as being inside  
329 the AD of the model.<sup>58</sup>

330

331



### 332 3. Results and discussion

333 Optimal architecture of the proposed EL-based nano-QSAR models in different case  
334 studies considered here were determined using 5-fold CV procedure. The optimal parameters of  
335 proposed classification and regression models (DTB, DTF) for different case studies considered  
336 here are presented in Table 1. The internal (CV-RMSE) and external ( $R^2$ , RMSE,  $Q_{F1}^2$ ,  $Q_{F2}^2$ ,  $Q_{F3}^2$   
337 and CCC) validation results of the developed EL-based regression models in different case  
338 studies are provided in Table 2.

#### 339 Table 1

#### 340 Table 2

341 These results indicate that both the nano-QSAR models (DTB, DTF) herein investigated are  
342 robust and showed no over-fitting of data in any of the five case studies. Model validation using  
343 the external data yielded criteria parameter values were above (except CCC values) their  
344 respective thresholds. The validation criteria threshold for  $Q_{F1}^2$ ,  $Q_{F2}^2$ , and  $Q_{F3}^2 > 0.6$  and an  
345 arbitrary cut-off value of 0.85 for CCC have been considered.<sup>51</sup> Moreover, criteria proposed by  
346 Eriksson et al.<sup>59</sup>, the difference between  $R^2$  (training) and  $R^2$  (validation) should not exceed 0.3.  
347 Model yielding  $R^2 \geq 0.81$  for in vitro and  $\geq 0.64$  for in vivo data can be regarded as acceptable.<sup>60</sup>  
348 As the proposed models fulfill these criteria and also positively pass internal and external  
349 validation, these were applied to predict the biological activity of new, untested NPs of diverse  
350 NMs. The performance parameters for the regression nano-QSARs both in training, test and  
351 complete data for each of the five case studies are summarized in Table 3. Plots of the measured  
352 and model predicted biological activities in different cases are shown in Fig. 4. The results  
353 obtained for various case-studies are discussed here.

#### 354 Table 3

355 **Figure 4**

356 **Case study 1: Nano-QSAR modeling of biological activity induced by diverse NMs-**

357 In this study, both the classification and regression nano-QSAR models were constructed using

358 four experimentally measured physical descriptors of NMs reported in literature.<sup>20</sup> Such

359 structural descriptors, namely NM size, relaxivities (R1 and R2), and zeta-potential were

360 available for 44 of the NMs. NMs size ranged between 22 to 74 nm. Relaxivities of NMs

361 represent their magnetic properties, and zeta-potential represents the intensity of charge on their

362 surface. The performance parameters of the classification models in training, test, and complete

363 data are summarized in Table 4. Both the models yielded considerably high accuracy, sensitivity,

364 specificity of binary classification for the considered NPs.

365 **Table 4**

366 The selected optimal binary classification models applied to complete data yielded accuracy of

367 95.45 % (DTB) and 93.18 % (DTF). The sensitivity, specificity and MCC values yielded by two

368 models in complete data were 100 %, 91.67 %, 0.91 (DTB), 100 %, 88.0 %, 0.87 (DTF),

369 respectively. Fourches et al.<sup>6</sup> developed SVM based model for binary classification of these NMs

370 using same set of descriptors and reported average values of accuracy, sensitivity and specificity:

371 73 %, 60 %, and 86 %, respectively. The performance parameters of the constructed nano-QSAR

372 models for predicting SMA induced by metal oxide NPs in model building and test phases are

373 presented in Table 3. The results showed high correlations (Fig. 4a) and low prediction errors,

374 suggesting for their adequacy for predicting SMA induced by new NPs. These models yielded

375 high correlation ( $R^2$ ) and low RMSE values of 0.939, 2.0 (DTB) and 0.851, 2.93 (DTF),

376 respectively in complete data, which suggests for relatively better performance of the proposed

377 QSAR models compared to the earlier (MLR) approach.<sup>27</sup> From the results obtained in present

378 study, it is evident that the proposed EL-based DTB and DTF nano-QSAR models performed  
379 relatively better.

380

### 381 **Case study 2: QSAR modeling of cellular uptake of NPs with similar core-**

382 Here, 109 NPs with the same core structure but diverse organic molecules attached to their  
383 surfaces that were tested for cellular uptake against PaCa2 cell were investigated. Each  
384 individual NP was represented by the structure of the organic molecule attached to its surface  
385 which in turn is characterized by molecular descriptors. In this case, relevant descriptors were  
386 selected using the minimum variance followed by model-fitting approach. The descriptors  
387 selected by this procedure are weighted partial negative surface area-3 (WNSA-3), weighted  
388 partial positive area-2 (WPSA-2), chi simple path descriptor of order 5 (SP-5), chi valance path  
389 descriptor of order 4 (VP-4), moment of inertia along X/Z-axis (MOMI-XZ), logarithmic form of  
390 octanol-water partition coefficient predicted by atomic method (XlogP), number of rotatable  
391 bonds (nRotB), number of hydrogen bond donors (nHBDon) for classification and VP-4, chi  
392 valance path cluster of order 6 (VPC-6), ionization potential (IP), nRotB, and number of  
393 hydrogen acceptors (nHBAcc) for regression modeling. IP is electronic, SP-5, VP-4 and VPC-6  
394 are chi-path and chi-path cluster descriptors belonging to topological descriptors, MOMI-XZ,  
395 WNSA-3 and WPSA-2 are geometrical and XlogP, nRotB, nHBDon, nHBAcc are constitutional  
396 descriptors. SP-5, VP-4, VPC-6 signify the total number of fragments of nth order (nth bond  
397 path) in NPs. IP is a measure of the energy needed for the removal of an electron from the  
398 cluster, yields valuable information on the electronic structure. nRotB is a measure of molecular  
399 flexibility. It is obtained simply by counting the non-terminal, non-cyclic, single bonds except C-  
400 N amide bond. nHBAcc represents the number of H-bond acceptors.

401 EL-based nano-QSAR models (DTB, DTF) were developed for binary classification  
402 (good/moderate cellular uptake and poor cellular uptake) of the NPs and to predict their cellular  
403 uptake in PaCa2 cells, expressed in terms of decadic logarithm of concentration (pM) of NP per  
404 cell. The performance parameters of the classification models in training, test, and complete data  
405 are summarized in Table 4. Both the models yielded considerably high accuracy, sensitivity,  
406 specificity and MCC of binary classification for the considered NPs. Both the classification  
407 models (DTB, DTF) applied to complete data array yielded accuracy, sensitivity, specificity and  
408 MCC values of 97.25 %, 96.67 %, 97.96 %, and 0.94. MCC value equal to 1 is regarded as a  
409 perfect prediction, whereas, 0 is for a completely random prediction. Chau and Yap<sup>29</sup> developed  
410 QSAR models (LR, k-NN, and SVM) for binary classification of 105 NPs considering 1D and  
411 2D PaDEL descriptors and reported average sensitivity, specificity, and MCC values of 86.7 %,  
412 67.3 %, and 0.559 achieved by the best consensus model.

413 The optimal DTB and DTF regression models were applied to the test and complete  
414 datasets, which explained 91.60 %, 88.94 % variance in training, 78.52 %, 72.14 % variance in  
415 test, and 89.23 %, 85.89 % variance in complete data. Proportion of variance explained by the  
416 model variables is the best single measure of how well the predicted values match the actual  
417 values. The two models yielded RMSE and  $R^2$  values of 0.14, 0.932 (DTB) and 0.16, 0.923  
418 (DTF) in complete data. From the values of the performance criteria parameters yielded by the  
419 QSARs in training, test and complete data (Table 3), it is evident that both the models yielded  
420 considerably low RMSE and MAE values in all the three phases. RMSE is a quadratic scoring  
421 rule which measures the average magnitude of the error. It gives a relatively high weight to large  
422 errors, hence most useful when large errors are particularly undesirable. MAE measures the  
423 average magnitude of the error in a set of predictions, without considering their direction. It is a

424 linear score which means that all the individual differences between predictions and  
425 corresponding measured values are weighted equally in the average.<sup>33</sup> Further, a closely followed  
426 pattern of variation by the measured and model predicted cellular uptake of NPs by the  
427 constructed QSAR models in the training and test phases (Fig. 4b) suggest that both the models  
428 performed reasonably well. The results suggest that the EL-based nano-QSAR models are  
429 superior and potentially useful for predicting the cellular uptakes of NPs.

430 Fourches et al.<sup>6</sup> developed QSAR model based on k-NN approach for the cellular uptake  
431 of NPs using same dataset but different set of descriptors and reported  $R^2$  and MAE values of  
432 0.72 and 0.18, respectively for the complete data. Ghorbanzadeh et al.<sup>1</sup> considering same dataset  
433 performed regression modeling based on MLR and MLPN approaches using 2D and 3D  
434 descriptors and reported  $R^2$  and RMSE values of 0.591, 0.364 (MLR) and 0.872, 0.150 (MLPN)  
435 for complete data. Recently, Toropov et al.<sup>28</sup> developed a QSAR model with SMILES based  
436 descriptor using CORAL software (which calculates descriptors as well as correlates them with  
437 end-points) and reported  $R^2$  and MAE values in sub-training, calibration, test and validation  
438 performed in five different splits (sub-sets) of the data. The overall corresponding  $R^2$  and MAE  
439 values ranged between 0.638-0.934 and 0.097 - 0.30, respectively. The statistical results of these  
440 studies suggest that although the performance of these QSAR models based on different  
441 modeling approaches for predicting the cellular uptake of NPs are within acceptability range, the  
442 developed QSAR models in the present work yielded better prediction of the end-point, and the  
443 proposed EL approach towards building nano-QSARs for NPs is more robust.

444

445

446

**447 Case study 3: QSAR modeling of cytotoxicity of diverse metal oxide NPs -**

448 In this case, three descriptors (oxygen percent, molar refractivity and polar surface area) were  
449 selected by the initial feature selection method. The oxygen percent (OP) is a constitutional  
450 descriptor and represents the elemental composition of the molecule. Molar refractivity (MolRef)  
451 is calculated based on the atomic method.<sup>61</sup> It is strongly related to the volume of the molecules  
452 and their polarizability. Therefore, this measure is also related to London dispersion forces,  
453 which have important effect in NP-receptor interaction processes. Polar surface area (PSA) is a  
454 geometrical descriptor and is known to show good correlation with passive molecular transport  
455 through membranes. This descriptor is formed by polar atoms of the molecule.

456 EL-based nano-QSAR models were developed to predict the pEC<sub>50</sub> in *E. coli* bacteria.  
457 Optimal DTB and DTF models using three descriptors (OP, PSA and MolRef) captured 96.98 %,   
458 86.50 % of the data variance in training, 91.33 %, 86.38 % in test, and 95.26 %, 86.51 % in  
459 complete data, respectively. The respective models yielded RMSE and R<sup>2</sup> values of 0.11, 0.955  
460 (DTB) and 0.19, 0.896 (DTF) in complete data. Values of the model performance criteria  
461 parameters in training, test and complete data are presented in Table 3. It may be noted that the  
462 present QSAR model yielded low RMSE and MAE values in all three phases. Further, a closely  
463 followed pattern of variation by the measured and model predicted pEC<sub>50</sub> values by the  
464 constructed QSAR models in the training and test phases (Fig. 4c) suggest that it performed  
465 reasonably well.

466 Puzyn et al.<sup>2</sup> developed and validated a MLR model to describe the relationship between  
467 the structures of 17 metal oxide NPs and their cytotoxicity to bacteria *E. coli* using single  
468 descriptor. The authors reported the R<sup>2</sup> and RMSE values of 0.85, 0.20 in training, and 0.83, 0.19  
469 in test set. In an another study, Toropov et al.<sup>26</sup> developed a QSAR model with SMILES based

470 descriptor using CORAL software and reported  $R^2$  values in the range of 0.83-0.96 for different  
471 test sets in six random splits. A direct comparison of our results with these studies is  
472 inappropriate, because the nature and number of descriptors, and modeling approaches  
473 considered differ to a large extent. Nevertheless, a simple comparison of the model statistics  
474 could provide some basic information about the accuracy of various prediction methodologies. It  
475 may be noted that both these studies considered complex descriptors (including SMILES derived  
476 and quantum mechanical) and in most of the data split folds prediction accuracies were not  
477 satisfactory, thus limiting the applicability of these models for prediction of end-point in new  
478 unknown NPs. Among these, the present study proposed EL-based nano-QSAR models  
479 considering the structurally diverse NPs and using simple structural descriptors yielded better  
480 prediction accuracy for the training, test, and complete data arrays.

481

#### 482 **Case study 4: QSAR modeling of cell viability modified multi-walled CNTs -**

483 In this study, six descriptors (Kier 3, MDEC-22, SP-5, XlogP, WTunity, MOMI-Y) were  
484 selected. The third Kier and Hall kappa molecular shape indices (Kier 3), molecular distance  
485 edge between all secondary carbons (MDEC-22), simple path descriptor of order 5 (SP-5), and  
486 Weighted holistic invariant molecular descriptor (WTunity) are topological, XlogP constitutional  
487 and moment of inertia along y/z-axis (MOMI-YZ) is geometrical descriptor.

488 Here, nano-QSAR (DTB and DTF) models were developed to predict the cell viability of  
489 CNTs. The optimal DTB and DTF models using six descriptors captured 89.91 %, 85.55 % of  
490 the data variance in training, 77.73 %, 92.21 % in test, and 88.56 %, 84.97 % in complete data.  
491 The respective models yielded  $R^2$  values of 0.903 and 0.922 in complete data. Values of the  
492 QSAR models performance criteria parameters in training, test and complete data are presented

493 in Table 3. It may be noted that the present QSAR models yielded low RMSE and MAE values  
494 in all the three data arrays. Further, a closely followed pattern of variation by the measured and  
495 model predicted end-point values by the constructed QSAR models in the training and test  
496 phases (Fig. 4d) suggest that it performed reasonably well.

497 QSARs developed earlier<sup>31</sup> using 4-D FP descriptors reported  $R^2$  values of 0.857 and  
498 0.759 in training and test phases for predicting the cell viability of CNTs. The performance  
499 criteria values (Table 3) suggest that the EL-based nano-QSAR models developed in present  
500 study performed relatively better.

501

#### 502 **Case study 5: QSAR modeling of cytotoxicity of Fullerene derivatives NPs –**

503 In this case study, ten descriptors belonging to the constitutional (aliphatic atom counts, chain  
504 bond count, hetero ring count, atom count, bond count), topological (Balaban index, Platt index,  
505 Wiener index), and geometrical (minimal projection area, molecular Polarizability) classes were  
506 considered for predictive modeling.

507 Nano-QSARs models (DTB and DTF) were developed to predict the  $pEC_{50}$  of fullerenes  
508 derivatives in *E. coli* bacteria. The optimal QSAR models (DTB, DTF) using ten descriptors  
509 captured 92.82 %, 88.89 % of data variance in training, 84.91 %, 75.20 % in test, and 92.16 %,  
510 87.74 % in complete data array. The performance parameters of the constructed nano-QSARs for  
511 predicting  $pEC_{50}$  of fullerenes in model training, test, and complete data are presented in Table 3.  
512 The results showed high correlations (Fig. 4e) and low prediction errors, suggesting for adequacy  
513 of these models. The respective models yielded  $R^2$  values of 0.958 and 0.943 in complete data.  
514 Both the models yielded low RMSE and MAE values in training, test and complete data (Table  
515 3), and a closely followed pattern of variation by the measured and model predicted  $pEC_{50}$  values



516 by the constructed nano-QSAR models in the training and test phases (Fig. 4e) suggest that these  
517 performed reasonably well.

518 Durdagi et al.<sup>23</sup> and Toropov et al.<sup>22</sup> earlier developed QSARs for predicting  $pEC_{50}$  of  
519 fullerenes using complex descriptors and reported high  $R^2$  values of 0.997, 0.906 in training and  
520 0.835, 0.992 in test phase, respectively. However, both these studies considered only selected  
521 (20) fullerenes data. In their later studies,<sup>24,25</sup> regression models were developed using PLSR and  
522 CORAL approaches, respectively considering all the 48 fullerene derivatives and reported  $R^2$  of  
523 0.993, 0.844 in training and 0.744, 0.792 in test sets. It is noticeable that in present study, EL-  
524 based nano-QSAR models derived using simple descriptors yielded comparable correlation ( $R^2$ )  
525 values, while considering dataset of all 48 fullerenes.

526

### 527 **3.1 Applicability domain of the proposed EL- nano-QSAR models**

528 The AD of a QSAR model is defined as the response and chemical structure spaces in  
529 which the model makes prediction with a given reliability.<sup>57</sup> To validate the predictive ability of  
530 the proposed EL-based nano-QSAR models for screening new NPs, analysis of the AD was  
531 performed following the methods based on the range of the descriptors in training sets for each  
532 of the NPs. According to this approach, the ranges of descriptors calculated for the NPs of  
533 training and test sets pertaining to all the case studies are shown in Table 5. The results depict  
534 that all NPs in training and test sets under all the five case studies are inside the AD of the  
535 proposed models, except 4 NPs in classification (case study 2), 1 CNT (case study 4), and 2  
536 fullerenes (case study 5). These results show that the proposed QSAR models under all the case  
537 studies here have wide applicability for predicting biological activities in new NMs.

538

#### **Table 5**

### 539 3.2 Mechanistic interpretation of the selected descriptors

540 Selection of the descriptors in QSAR modeling is an important aspect. The selected  
541 descriptors must not only contribute to the end-point in a quantitative manner, but also must have  
542 interpretability from mechanistic point of view. Contributions of the selected descriptors in  
543 constructed nano-QSAR models are shown in Fig. 5.

#### 544 Figure 5

545 For the first case study, four descriptors were taken to construct EL-based QSAR models  
546 as selected in earlier study<sup>6</sup> for the purpose of a comparison. The relevance of these descriptors is  
547 discussed elsewhere.<sup>6</sup> The contributions of these descriptors in classification and regression  
548 models constructed in present study are shown in Fig. 5a,b. It may be noted that the contribution  
549 of zeta potential (ZP) was highest in classification and that of NPs size in regression nano-  
550 QSARs. The mechanistic interpretation of the descriptors identified in other four case studies is  
551 discussed here. The toxicity induced by the NPs may be investigated by considering various  
552 possible mechanisms<sup>2,62</sup> such as (a) the release of chemical constituents from NPs, (b) the size  
553 and shape of the particle, which produces steric hindrances or interferences with the important  
554 binding sites of macromolecules, (c) the surface properties of the material, such as  
555 photochemical and redox properties, and (d) the capacity of NPs to act as vectors for the  
556 transport of other toxic chemicals to sensitive tissues. Once a NP enters a cell, toxicity could  
557 occur through one or combination of these mechanisms.

558 In the second case study, the classification QSARs were developed using WNSA-3,  
559 WPSA-2, SP-5, VP-4, MOMI-XZ, XlogP, nRotB, nHBDOn, whereas in regression QSARs, VP-  
560 4, VPC-6, IP, nRotB, and nHBAcc were considered (Fig. 5c,d). The topological descriptors (SP-  
561 5, VP-4, VPC-6) help to differentiate the molecules according mostly to their size, degree of

562 branching, flexibility and overall shape.<sup>63</sup> Chi cluster descriptor (VPC-6) is an indicator of the  
563 degree of  $n^{\text{th}}$  order branching, and thus implicates the effect of substitution in a molecule. A  
564 molecule that is relatively compact at the same point(s) will have a high value of this  
565 descriptor.<sup>64</sup> IP, an electronic descriptor is of critical importance in determining the type of inter-  
566 molecular forces which underlie in molecule-receptor interactions. Extensive studies using  
567 electronic parameters reveal that electronic attributes of molecules are intimately related to their  
568 chemical reactivities and biological activities.<sup>65</sup> WNSA-3 is defined as the total sum of partial  
569 areas of the NPs which possesses negative partial charges times the total solvation area of the NP  
570 divided by 1000. This descriptor is related to the stability of the chemical bond and surface area  
571 of the NP molecule.<sup>66</sup> WPSA-2 is the surface weighted charge partial positive surface area and  
572 related to the charge distribution describing the positively charged surface area of the NP. It is  
573 directly dependent on the H-bonding donor or acceptor ability of the molecule.<sup>67</sup> MOMI-XZ is a  
574 geometric parameter and its value depends on the total mass of the molecule, the distribution  
575 within the molecule and position of axis rotation of the molecule.<sup>68</sup> XlogP calculated using atom  
576 type prediction method denotes an important property in describing the affinity of the  
577 compounds in terms of their partitioning in the biological membranes.<sup>29</sup> nRotB refers to number  
578 of rotatable bonds in the molecules. The positive term associated with the descriptor in QNAR  
579 model indicate that fractional increase in the rotatable bonds in the molecule is beneficial for  
580 biological activity.<sup>69</sup> nHBAcc describes capability of moiety in participating in H-bonding. The  
581 highest H-bonding acceptor potential is defined as the maximum ion-pair electro negativity on an  
582 atom considering all N, O, and F atoms in a compound. The H-bond donor plays an important  
583 role in NP-receptor interaction, aqueous solubility and partitioning. Properties such as oral  
584 bioavailability or membrane permeability have often been correlated to the number of H-bond

585 donor and log P in a molecule.<sup>70</sup> Moreover, the topological descriptors (VPC-6 and VP-4) were  
586 found to be the most important factors, which could be considered to synthesize a new organic  
587 modifier to control PaCa2 cellular uptake of NPs. Also, molecular shape and size, and amount of  
588 branching in organic coatings, can be effective factors in cellular uptake of studied NPs in  
589 pancreatic cancer cells.

590 In third case study, among the identified descriptors, OP has highest contribution in the  
591 QSAR model followed by PSA and MolRef (Fig. 5e). Auffan et al.<sup>71</sup> suggested that the most  
592 important parameter controlling the *in vitro* cytotoxicity of metallic NPs (zero-valent metals,  
593 metal oxides) is their chemical stability, which is related to the dissolution of the particles  
594 (release of cations) and the catalytic properties and redox modifications of the surface.  
595 Moreover, the release of cations can occur by simple breaking of chemical bonds in the crystal  
596 lattice (without changing the oxidation state of the metal) or by redox reactions with the  
597 molecules in the biological media. In the later case, the release of ions is often accompanied by  
598 the generation of reactive oxygen species (ROS), such as superoxides and hydroxyl radicals. The  
599 generation of ROS may be increased by intimate contact of NP with a cell membrane.<sup>72</sup> The  
600 observed toxicity can be induced by the released cations themselves, ROS or both.<sup>71,73</sup> PSA is  
601 defined as the part of the surface area of the molecule associated with N, O, S, and the H-bonded  
602 to any of these atoms.<sup>74</sup> This descriptor correlates well the passive molecular transport through  
603 membranes and allows the prediction of the transport properties of molecules to the target cell.<sup>75</sup>  
604 MolRef represents the molar volume corrected by the refractive index. It is a measure of size and  
605 Polarizability of a fragment or molecule and can be used for a substituent or for the whole  
606 molecule.<sup>76</sup> This property is an atomic contribution model that assumes the correct protonation  
607 state. Its positive contribution suggests that the increment of the polarity of molecule lead to

608 increased activity, moreover, it can also be said that more the number of polar groups in the  
609 molecule, more will be the affinity of the molecules towards the biological activity.<sup>77</sup>

610 In fourth case study, regression QSAR models are developed using six descriptors (Kier  
611 3, MDEC-22, SP-5, XlogP, WTunity, MOMI-YZ). Contribution of these descriptors in  
612 developed nano-QSARs is shown in Fig. 5f. Kier 3 is the most sensitive to the molecular  
613 topology and in particular to the branching of the molecule. It describes the valance connectivity  
614 of the molecule of the coordination sphere and also reflects the molecular composition.<sup>78</sup>  
615 WTunity (WTU) is the weighted holistic invariant molecular descriptor (WHIM). These are 3-D  
616 descriptor based on the calculation of principal component axis computed from a weighted  
617 covariance matrix obtained by the molecule geometric coordinate. It contains chemical  
618 information concerning size, symmetry, shape and distribution of molecular atoms.<sup>79</sup> MDEC-22  
619 represents molecular distance edge between all secondary carbons, larger molecular size  
620 increases the toxicity while greater degree of H-bonding in a molecule reduces toxicity by  
621 increasing its polarity.<sup>80</sup>

622 In case study 5, the constitutional descriptors (aliphatic atom counts, AAC; chain bond  
623 count, CBC; hetero ring count, HRC; atom count, AC; bond count, BC) capture properties of the  
624 molecule that are related to elements constituting its structure. These descriptors depend  
625 fundamentally on the composition of the molecule. Topological descriptors (Balaban index, BI;  
626 Platt index, PI; Wiener index, WI) treat the structure of the compound as a graph, with atoms as  
627 vertices and co-valent bonds as edge. Wiener index counts the total number of bonds in shortest  
628 paths between all pairs of non-H atoms.<sup>81</sup> Analysis of the Balaban index shows that it will  
629 increase with the size of the molecule, degree of branching and unsaturation. The more branched  
630 molecules are less toxic, probably due to their lower membrane penetration abilities. The

631 geometrical descriptors (minimal projection area, Mpa; molecular polarizability, Mpol) rely on  
632 spatial arrangement of the atoms constituting the molecule. These descriptors include  
633 information of molecular surface obtained from atomic vander Waals areas and their overlap.<sup>82</sup>  
634 Molecular polarizability (Mpol) measures the ability of the electrons in a molecule to move  
635 easily as a result of stimulus. Because the electrons in the molecules of the compounds with high  
636 polarizability can relatively move easily, both excited singlet and triplet states of the molecules  
637 of such compounds may not be stable.<sup>83</sup> Hence, suggesting that chemicals with large Mpol  
638 values will have higher toxicity. Contribution of the selected descriptors in constructed nano-  
639 QSARs is shown in Fig. 5g.

640 In view of the above facts, it is clear that the descriptors selected in these case studies  
641 have high relevance in the developed nano-QSARs and quantitative contributions to the end-  
642 points along with the mechanistic interpretability towards the biological activity exhibited by the  
643 diverse NMs. The results obtained under all the five case studies here suggest that the proposed  
644 EL-based nano-QSARs performed relatively better than those considered in previous studies and  
645 can be used as reliable tools for predicting the biological activities of diverse NPs using simple  
646 molecular descriptors.

647

#### 648 4. Conclusions

649 In conclusion, EL-approach based robust and reliable nano-QSAR models have been  
650 proposed for predicting biological activities of NPs derived from diverse NMs using simple  
651 structural descriptors and demonstrated potential benefits of the EL approaches to obtain  
652 predictive knowledge for NPs that affect human cells and utilize this knowledge to improve the  
653 experimental design of new NPs enabling their prioritization for *in vivo* screening. This work has

654 demonstrated the suitability and superiority of the EL approach in developing nano-QSAR  
655 models for predicting biological activity of diverse NPs through their applications to five  
656 different datasets of diverse NMs. The quality of the nano-QSAR models derived in this study  
657 was rigorously estimated according to their external prediction abilities assessed by a five-fold  
658 cross-validation and external data validation procedures. Present study on five diverse datasets  
659 clearly indicated that the proposed approaches have successfully provided promising nano-  
660 QSAR modeling tools in this challenging area. The case studies considered here successfully  
661 introduced a new approach to construct robust QSAR models both for classification and  
662 regression problems in the area of computational nano-toxicology. The superiority of the  
663 proposed EL approach over the earlier ones may be attributed to the fact that the DTB and DTF  
664 models incorporate stochastic gradient boosting and bagging algorithms, respectively, which  
665 improves generalization ability of weak learners. The proposed EL approach may be considered  
666 as a potential method for predictive modeling in the area of nano-technology.

667

### 668 **Supporting Information**

669 Additional tables are available in the online version of this article.

670

671 **Acknowledgement:** The authors thank the Director, CSIR-Indian Institute of Toxicology  
672 Research, Lucknow (India) for his keen interest in this work and providing all necessary  
673 facilities. The SMA data provided by Dr. Vidana C. Epa, CSIRO, Australia and SMILES of the  
674 fullerene derivatives provided by Dr. Andrey A. Toropov, IRCCS, Italy are thankfully  
675 acknowledged.

676

677 **References**

678

679 [1] M. Ghorbanzadeh, M. H. Fatemi, M. Karimpour, *Ind. Eng. Chem. Res.* 2012, **51**,  
680 10712–10718.

681

682 [2] T. Puzyn, B. Rasulev, A. Gajewicz, X. Hu, T. P. Dasari, A. Michalkova, H. Hwang, A.  
683 Toropov, D. Leszczynska, J. Leszczynski, *Nat. Nanotechnol.* 2011, **6**, 175–178.

684

685 [3] M. Heinemann, H.G. Schafer, *Hum. Exp. Toxicol.* 2009, **28**, 407–411.

686

687 [4] S. J. Klaine, P. J. J. Alvarez, G. E. Batley, T. F. Fernandes, R. D. Handy, D. Y. Lyon, S.  
688 Mahendra, M. J. McLaughlin, J. A. Lead, *Environ. Toxicol. Chem.* 2008, **27**, 1825–1851.

689

690 [5] H.L. Karlsson, P. Cronholm, J. Gustafsson, L. Moller, *Chem. Res. Toxicol.* 2008, **21**,  
691 1726–1732.

692

693 [6] D. Fourches, D. Pu, C. Tassa, R. Weissleder, S. Y. Shaw, R. J. Mumper, A. Tropsha, *ACS*  
694 *Nano* 2010, **4**, 5703–5712.

695

696 [7] R. Weissleder, K. Kelly, E. Y. Sun, T. Shtatland, L. Josephson, *Nat. Biotechnol.* 2005, **23**,  
697 1418–1423.

698

699 [8] L. Braydich-Stolle, S. Hussain, J. J. Schlager, M. C. Hofmann, *Toxicol. Sci.* 2005, **88**,  
700 412–419.

701

702 [9] H. W. Chen, S. F. Su, C. T. Chien, W. H. Lin, S. L. Yu, C. C. Chou, *FASEB J.* 2006, **20**,  
703 2393–2395.

704

705 [10] M. N. Moore, *Environ. Int.* 2006, **32**, 967–976.

706

707 [11] A. S. Karakoti, L. L. Hench, S. Seal, *JOM* 2006, **58**, 77–82.



- 708 [12] O. R. Moss, V. A. Wong, *Inhal. Toxicol.* 2006, **18**, 711–716.  
709
- 710 [13] J. K. Lee, M. H. Cho, *Toxicol. Sci.* 2006, **89**, 338–347.  
711
- 712 [14] J. S. Tsuji, A. D. Maynard, P. C. Howard, J. T. James, C. W. Lam, D. B. Warheit, A. B.  
713 Santamaria, *Toxicol. Sci.* 2006, **89**, 42–50.  
714
- 715 [15] R. Duffin, L. Tran, D. Brown, V. Stone, K. Donaldson, *Inhal. Toxicol.* 2007, **19**, 849–856.  
716
- 717 [16] S. Fiorito, Carbon nanoparticles: benefits and risks for human health. In *Nanotoxicology:  
718 Interactions of Nanomaterials with Biological Systems*; Zhao, Y., Nalwa, H. S., Eds.;  
719 American Scientific Publishers, New York, 2007; pp. 167–180.  
720
- 721 [17] S. Gill, R. Lobenberg, T. Ku, S. Azarmi, W. Roa, E. J. Prenner, *J. Biomed. Nanotechnol.*  
722 2007, **3**, 107–119.  
723
- 724 [18] K. W. Powers, M. Palazuelos, B. M. Moudgil, S. M. Roberts, *Nanotoxicology* 2007, **1**, 42–  
725 51.  
726
- 727 [19] C. Medina, M. J. Santos-Martinez, A. Radomski, O. I. Corrigan, M. W. Radomski, *Br. J.*  
728 *Pharmacol.* 2007, **150**, 552–558.  
729
- 730 [20] S. Y. Shaw, E. C. Westly, M. J. Pittet, A. Subramanian, S. L. Schreiber, R. Weissleder,  
731 *Proc. Natl. Acad. Sci. USA.* 2008, **105**, 7387–7392.  
732
- 733 [21] H. Meng, T. Xia, S. George, A. Nel, *ACS Nano.* 2009, **3**, 1620–1627.  
734
- 735 [22] A. A. Toropov, A. P. Toropova, E. Benfenati, D. Leszczynska, J. Leszczynski, *J. Comput.*  
736 *Chem.* 2010, **31**, 381–392.  
737

- 738 [23] S. Durdagi, T. Mavromoustakos, M. G. Papadopoulos, *Biorg. Med. Chem. Lett.* 2008, **18**,  
739 6283-6289.  
740
- 741 [24] S. Durdagi, T. Mavromoustakos, N. Chronakis, M. G. Papadopoulos, *Biorg. Med. Chem.*  
742 2008, **16**, 9957-9974.  
743
- 744 [25] A. P. Toropova, A. A. Toropov, E. Benfenati, D. Leszczynska, J. Leszczynski, *J. Math.*  
745 *Chem.* 2010, **48**, 959-987.  
746
- 747 [26] A. A. Toropov, A. P. Toropova, E. Benfenati, G. Gini, T. Puzyn. D. Leszczynska, J.  
748 Leszczynski, *Chemosphere* 2012, **89**, 1098–1102.  
749
- 750 [27] V. C. Epa, F. R. Burden, C. tassa, R. Weissleder, S. Shaw, D. A. Winkler, *Nano lett.* 2012,  
751 **12**, 5808-5812.  
752
- 753 [28] A. A. Toropov, A. P. Toropova, T. Puzyn. E. Benfenati, G. Gini, D. Leszczynska, J.  
754 Leszczynski, *Chemosphere* 2013, **92**, 31–37.  
755
- 756 [29] Y. T. Chau, C. W. Yap, *RSC Adv.* 2012, **2**, 8489-8496.  
757
- 758 [30] D. Fourches, D. Pu, A. Tropsha, *Comb. Chem. High Throughput Screening* 2011, **14**, 217-  
759 225.  
760
- 761 [31] C-Y. Shao, S-Z. Chen, B-H. Su, Y. J. Tseng, E. X. Esposito, A. J. Hopfinger, *J. Chem. Inf.*  
762 *Model.* 2013, **53**, 142–158.  
763
- 764 [32] K. P. Singh, S. Gupta, P. Rai, *Ecotox. Environ. Safe.* 2013, **95**, 221–233.  
765
- 766 [33] K. P. Singh, S. Gupta, P. Rai, *Atmos. Environ.* 2013, **80**, 426-437.  
767

- 768 [34] T. Hancock, R. Put, D. Coomans, Y. Vander Heyden, Y. Everingham, *Chemometr. Intell.*  
769 *Lab. Syst.* 2005, **76**, 185–196.  
770
- 771 [35] T. G. Dietterich, *Lect. Notes Comput. Sci.* 2000, **1857**, 1–15.  
772
- 773 [36] L. Breiman, *Mach. Learn.* 1996, **24**, 123–140.  
774
- 775 [37] J. H. Friedman, *Comput. Stat. Data An.* 2002, **38**, 367-378.  
776
- 777 [38] J. Mahjoobi, A. Etemad-Shahidi, *Appl. Ocean Res.* 2008, **30**, 172-177.  
778
- 779 [39] H. Zhou, Q. Mu, N. Gao, A. Liu, Y. Xing, S. Gao, Q. Zhang, G. Qu, Y. Chen, G. Liu, B.  
780 Zhang, B. Yan, *Nano Lett.* 2008, **8**, 859–865.  
781
- 782 [40] J. M. Worle-Knirsch, K. Pulskamp, H. F. Krug, *Nano Lett.* 2006, **6**, 1261–1268.  
783
- 784 [41] C. Steinbeck, C. Hoppe, S. Kuhn, M. Floris, R. Guha, E. L. Willighagen, *Curr. Pharm.*  
785 *Des.* 2006, **12**, 2110-2120.  
786
- 787 [42] [www.chemspider.com](http://www.chemspider.com)  
788
- 789 [43] R. Benigni, C. Bossa, T. Netzeva, E. Benfenati, R. Franke, C. Helma, E. Hulzebos, C.  
790 Marchant, A. M. Richard, Y. Woo, C. Yang, *J. Environ. Sci. Health. C. Environ. Carcinog.*  
791 *Revs.* 2007, **25**, 53-97.  
792
- 793 [44] C. Y. Zhao, H. X. Zhang, X. Y. Zhang, M. C. Liu, Z. D. Hu, B. T. Fan, *Toxicology* 2006,  
794 **217**, 105-119.  
795
- 796 [45] K. Varmuza, M. Karlovits, W. Demuth. *Anal. Chim. Acta* 2003, **490**, 313–324.  
797
- 798 [46] K. P. Singh, S. Gupta, P. Rai, *Toxicol. Appl. Pharmacol.* 2013, **272**, 465-475.

- 799 [47] H. Ishwaran, U. B. Kogalur, *Stat. Probabil. Lett.* 2010, **80**, 1056-1064.  
800
- 801 [48] H. I. Erdal, O. Karakurt, *J. Hydrol.* 2013, **477**, 119-128.  
802
- 803 [49] R. Benigni, C. Bossa, N. Jeliaskova, T. Netzevac, A. Worth, The Benigni/Bossa rulebase  
804 for mutagenicity and carcinogenicity – a module of toxtree. Technical Report EUR  
805 23241EN, European Commission–Joint Research Centre, 2008.  
806
- 807 [50] K. P. Singh, N. Basant, S. Gupta, *Anal. Chim. Acta* 2011, **703**, 152–162.  
808
- 809 [51] N. Chirico, P. Gramatica, *J. Chem. Inf. Model.* 2011, **51**, 2320–2335.  
810
- 811 [52] L. M. Shi, H. Fang, W. Tong, J. Wu, R. Perkins, R. M. Blair, W. S. Branham, S. L. Dial, C.  
812 L. Moland, D. M. Sheehan, *J. Chem. Inf. Comput. Sci.* 2001, **41**, 186–195.  
813
- 814 [53] G. Schuurmann, R. Ebert, J. Chen, B. Wang, R. Kuhne, *J. Chem. Inf. Model.* 2008, **48**,  
815 2140–2145.  
816
- 817 [54] V. Consonni, D. Ballabio, R. Todeschini, *J. Chem. Inf. Model.* 2009, **49**, 1669–1678.  
818
- 819 [55] L. I. Lin, *Biometrics* 1992, **48**, 599–604.  
820
- 821 [56] K. P. Singh, N. Basant, A. Malik, G. Jain, *Anal. Chim. Acta* 2010, **658**, 1-11.  
822
- 823 [57] T. I. Netzeva, A. Worth, T. Aldenberg, R. Benigni, M. T. Cronin, P. Gramatica, J.S.  
824 Jaworska, S. Kahn, G. Klopman, C. A. Marchant, G. Myatt, N. Nikolova-Jeliaskova, G. Y.  
825 Patlewicz, R. Perkins, D. Roberts, T. Schultz, D. W. Stanton, J. J. van de Sandt, W.  
826 Tong, G. Veith, C. Yang, *ATLA, Altern. Lab. Anim.* 2005, **33**, 155–173.  
827
- 828 [58] S. Kovarich, E. Papa, P. Gramatica, *J. Hazard. Mater.* 2011, **190**, 106–112.  
829

- 830 [59] L. Eriksson, J. Jaworska, A. P. Worth, M. T. D. Cronin, R. M. McDowell, P. Gramatica,  
831 *Environ. Health Perspect.* 2003, **111**, 1361–1375.
- 832
- 833 [60] H. Kubinyi, *Methods and Principles in Medicinal Chemistry Vol. 1* (eds Mannhold, R.,  
834 Kroogsgard-Larsen, P. & Timmerman, H.) VCH, 1993.
- 835
- 836 [61] V. N. Viswanadhan, A. K. Ghose, G. R. Revankar, R. K. Robins, *J. Chem. Inf. Comput.*  
837 *Sci.*, 1989, **29**, 163–172.
- 838
- 839 [62] Linkov, J. Steevens, G. Adlakha-Hutcheon, E. Bennett, M. Chappell, V. Colvin, J. M.  
840 Davis, T. Davis, A. Elder, S. Foss Hansen, P. B. Hakkinen, S. M. Hussain, D. Karkan, R.  
841 Korenstein, I. Lynch, C. Metcalfe, A. B. Ramadan, F. K. Satterstrom, *J. Nanopart. Res.*  
842 2009, **11**, 513–527.
- 843
- 844 [63] L. Gupta, A. Patel, C. Karthikeyan, P. Trivedi, *J. Curr. Pharm. Res.* 2010, **01**, 19-25.
- 845
- 846 [64] Y. Yuan, P. D. Mosier, Y. Zhang, *J. Biophys. Chem.* 2012, **3**, 49-57.
- 847
- 848 [65] D. M. Patel, N. M. Patel, *J. Sci. Res.* 2009, **1**, 594-605.
- 849
- 850 [66] M. Karelson, G. Karelson, T. Tamm, I. Tulp, J. Janes, K. Tamm, A. Lomaka, *ARKIVOC*  
851 2009, **ii**, 218-238.
- 852
- 853 [67] F. M. Abu-Awwad, *Der. Pharma Chemica* 2010, **2**, 1-13.
- 854
- 855 [68] A. P. G. Nikalje, M. Pathan, A. S. Narute, M. S. Ghodke, D. Rajani, *Der. Pharma. Sinica*  
856 2012, **3**, 229-238.
- 857
- 858 [69] M. C. Sharma, D. V. Kohli, *World Appl. Sci. J.* 2011, **12**, 2111–2119.
- 859

- 860 [70] S. K. Paliwal, S. Singh, S. Kumari, A. A. Siddiqui, S. K. Paliwal, *Ind. J. Chem.* 2010, **49B**,  
861 554-560.
- 862
- 863 [71] M. Auffan, J. Rose, M.R. Wiesner, J. Y. Bottero, *Environ. Pollut.* 2009, **157**, 1127–1133.
- 864
- 865 [72] M. Heinlaan, A. Ivsak, I. Bilnova, H. C. Dubourguier, A. Kahru, *Chemosphere* 2008, **71**,  
866 1308-1316.
- 867
- 868 [73] A. L. Neal, *Ecotoxicology* 2008, **17**, 362–371.
- 869
- 870 [74] P. Ertl, B. Rohde, P. Selzer, *J. Med. Chem.* 2000, **43**, 3714–3717.
- 871
- 872 [75] S. Winiwarter, F. Ax, H. Lennernas, A. Hallberg, C. Pettersson, A. Karlen, *J. Mol. Graph.*  
873 *Model.* 2003, **21**, 273–287.
- 874
- 875 [76] B. Niu, Y. Jin, W. Lu, G. Li, *Chemom. Intell. Lab. Syst.* 2009, **96**, 43–48.
- 876
- 877 [77] M. Bhaisare, C. Karthikeyan, O. Tanwar, S. Waghulde, S. Laddha, A QSAR analysis of  
878 some amino substituted Pyrido[3,2-b]pyrazinones as potent and selective PDE-5 inhibitors.  
879 *In Proceedings of the 14<sup>th</sup> Int. Electron. Conf. Synth. Org. Chem. Sciforum Electronic*  
880 *Conferences Series*, 1-30 November, 2010, pp.1-25.
- 881
- 882 [78] A. R. Katritzky, R. Petrukhin, D. Tatham, S. Basak, E. Benfenati, M. Karelson, U. Maran,  
883 *J. Chem. Inf. Comput. Sci.*, 2001, **41**, 679-685.
- 884
- 885 [79] R. Todeschini, P. Gramatica, 3D QSAR in Drug Design – Vol. 2, H. Kubinyi, G. Folkers,  
886 Y. C. Martin (Eds.), Kluwer/ESCOM, Dordrecht (The Netherlands), 1998, 355–380.
- 887
- 888 [80] S. Parvez, C. Venkataraman, S. Mukherji, *Toxicol. in Vitro* 2008, **22**, 1806–1813.
- 889

- 890 [81] A. Z. Dudek, T. Arodzb, J. Gálvezc, *Comb.Chem. High Throughput Screening* 2006, **9**,  
891 213-228.
- 892
- 893 [82] P. Labute, *J. Mol. Graph. Model.* 2000, **18**, 464–477.
- 894
- 895 [83] Y. Z. Dai, D. S. Yang, F. Zhu, L. Y. Wu, X. Z. Yang, J. H. Li, *Chemosphere* 2006, **65**,  
896 2427–2433.

897 **Legend to the Figures**

898

899 **Figure 1:** Histogram of the (a) SMA data, (b) cellular uptake in PaCa2 cells, (c) pEC<sub>50</sub> of  
900 metallic oxide NPs in *E. coli*, (d) cell viability in CNT, and (e) pEC<sub>50</sub> of fullerenes  
901 in HIV-1 PR.

902

903 **Figure 2:** Radar plot of the distribution of selected descriptors used in (a) case study  
904 1, (b) case study 2, and (c) case study 3, (d) case study 4, and (e) case study 5 for  
905 nano-QSAR modeling.

906

907 **Figure 3:** Conceptual diagram of the (a) DTB-nano-QSAR and (b) DTF-nano-QSAR  
908 models.

909

910 **Figure 4:** Plot of the experimental and model predicted values of the biological activity in  
911 training and test data under (a) case study 1, (b) case study 2, (c) case study 3, (d)  
912 case study 4, and (e) case study 5, using DTB and DTF nano-QSARs.

913

914 **Figure 5:** Plot of the contribution of the selected descriptors in NPs biological activity  
915 prediction models for (a) case study 1, classification, (b) case study 1, regression,  
916 (c) case study 2, classification, (d) case study 2, regression, (e) case study 3,  
917 regression, (f) case study 4, regression, and (g) case study 5, regression.

918

919



920 **Table 1:** Optimal parameters of constructed EL-based nano-QSARs.

Case study/ Model	Classification			Regression		
	Number of trees	Max depth of any tree	No. of Average group splits	Number of trees	Max depth of any tree	No. of Average group splits
<b>Case Study-1</b>						
DTB	332	5	66.1	337	6	50.2
DTF	95	12	9.3	140	7	5.8
<b>Case Study-2</b>						
DTB	298	5	43.8	362	6	177.9
DTF	143	11	13.5	282	18	51.9
<b>Case Study-3</b>						
DTB	-	-	-	365	6	19.4
DTF	-	-	-	171	7	3.4
<b>Case Study-4</b>						
DTB	-	-	-	397	7	21.0
DTF	-	-	-	170	11	14.1
<b>Case Study-5</b>						
DTB	-	-	-	310	5	53.4
DTF	-	-	-	693	8	19.8

922

923

924 **Table 2:** Performance parameters for ensemble models for different case studies

925

Model	Sub-Sets	$Q_{F1}^2$	$Q_{F2}^2$	$Q_{F3}^2$	CCC	CV-RMSE
<b>Case Study-1 (n=31)</b>						
DTB	Test	0.898	0.898	0.890	0.946	3.83
	Complete	0.930	0.930	0.929	0.961	
DTF	Test	0.787	0.787	0.770	0.895	4.03
	Complete	0.849	0.849	0.847	0.918	
<b>Case Study-2 (n=109)</b>						
DTB	Test	0.785	0.785	0.724	0.843	0.31
	Complete	0.892	0.892	0.888	0.932	
DTF	Test	0.721	0.721	0.642	0.783	0.31
	Complete	0.859	0.859	0.853	0.906	
<b>Case Study-3 (n=17)</b>						
DTB	Test	0.915	0.913	0.930	0.957	0.16
	Complete	0.953	0.953	0.956	0.974	
DTF	Test	0.866	0.864	0.889	0.915	0.29
	Complete	0.865	0.865	0.874	0.917	
<b>Case Study-4 (n=29)</b>						

DTB	Test	0.791	0.777	0.933	0.926	5.25
	Complete	0.796	0.796	0.782	0.854	
DTF	Test	0.927	0.922	0.976	0.957	4.85
	Complete	0.850	0.850	0.867	0.898	
<b>Case Study-5 (n=48)</b>						
DTB	Test	0.853	0.849	0.946	0.912	1.02
	Complete	0.764	0.763	0.719	0.839	
DTF	Test	0.758	0.752	0.912	0.841	1.07
	Complete	0.877	0.877	0.894	0.920	

926

927

928

929 **Table 3:** Performance parameters for constructed nano-QSARs

930

Case study	Model	Sub-Sets	Mean	*SD	MAE	RMSE	R <sup>2</sup>
Case Study-1	Experimental	Training	-9.94	7.65	-	-	-
		Test	-9.65	8.52	-	-	-
		Complete	-9.89	7.68	-	-	-
	DTB	Training	-10.30	6.68	0.63	1.87	0.950
		Test	-10.31	7.99	1.86	2.48	0.906
		Complete	-10.30	6.81	0.87	2.00	0.939
	DTF	Training	-10.00	6.78	1.96	2.75	0.868
		Test	-10.75	8.57	3.06	3.59	0.817
		Complete	-10.14	7.01	2.17	2.93	0.851
Case Study-2	Experimental	Training	3.65	0.41	-	-	-
		Test	3.65	0.48	-	-	-
		Complete	3.65	0.42	-	-	-
	DTB	Training	3.65	0.33	0.10	0.12	0.947
		Test	3.64	0.29	0.17	0.22	0.905
		Complete	3.65	0.32	0.11	0.14	0.932
	DTF	Training	3.65	0.31	0.10	0.14	0.942
		Test	3.63	0.26	0.19	0.25	0.889
		Complete	3.65	0.30	0.11	0.16	0.923
Case Study-3	Experimental	Training	2.69	0.56	-	-	-
		Test	2.76	0.53	-	-	-
		Complete	2.71	0.53	-	-	-
	DTB	Training	2.69	0.52	0.08	0.09	0.974
		Test	2.83	0.52	0.11	0.14	0.936
		Complete	2.74	0.50	0.09	0.11	0.955
	DTF	Training	2.64	0.43	0.17	0.20	0.911
		Test	2.78	0.41	0.16	0.18	0.894
		Complete	2.69	0.41	0.16	0.19	0.896
Case Study-4	Experimental	Training	34.96	27.71	-	-	-
		Test	31.20	16.68	-	-	-
		Complete	34.31	25.93	-	-	-
	DTB	Training	35.12	21.77	6.88	8.61	0.931

		Test	32.50	23.64	6.03	7.04	0.971
		Complete	34.84	21.28	6.96	8.62	0.903
	DTF	Training	35.26	19.21	8.95	10.31	0.929
		Test	31.41	14.97	3.28	4.16	0.927
		Complete	34.73	17.96	8.34	9.88	0.922
<b>Case Study-5</b>	Experimental	Training	5.33	1.48	-	-	-
		Test	5.47	0.92	-	-	-
		Complete	5.36	1.38	-	-	-
	DTB	Training	5.36	1.16	0.32	0.39	0.970
		Test	5.55	0.78	0.28	0.34	0.863
		Complete	5.40	1.09	0.31	0.38	0.958
	DTF	Training	5.38	1.05	0.40	0.49	0.963
		Test	5.49	0.71	0.33	0.43	0.762
		Complete	5.40	0.99	0.39	0.48	0.943

931 \*standard deviation

932

933

934 **Table 4:** Classification performance parameters for EL nano-QSARs for different case studies

935

Model	Sub-Sets	Sensitivity (%)	Specificity (%)	Accuracy (%)	MCC
<b>Case study-1</b>					
DTB	Training	100.00	100.00	100.00	1.00
	Test	100.00	66.67	71.43	0.47
	Complete	100.00	91.67	95.45	0.91
DTF	Training	100.00	94.74	97.30	0.95
	Test	100.00	66.67	71.43	0.47
	Complete	100.00	88.00	93.18	0.87
<b>Case study-2</b>					
DTB	Training	100.00	97.44	98.78	0.98
	Test	88.24	100.00	92.59	0.86
	Complete	96.67	97.96	97.25	0.94
DTF	Training	100.00	100.00	100.00	1.00
	Test	87.50	90.91	88.89	0.78
	Complete	96.67	97.96	97.25	0.94

936

937

938

939

940

941

942

943

944 **Table 5a:** Applicability domain of the selected descriptors in classification nano-QSARs under  
 945 different case studies

Case Studies	Descriptors	Training Set		Test Set	
		Min	Max	Min	Max
Case study-1	Size	20.00	74.00	20.00	36.00
	R1	0.50	36.00	0.50	36.00
	R2	0.50	153.00	0.50	122.00
	ZP	-37.00	5.90	-21.90	3.24
Case study-2	SP-5	0.00	6.81	0.00	6.51
	VP-4	0.00	5.73	0.00	4.54
	WPSA-2	49.40	5128.96	43.77	665.06
	WNSA-3	-79.80	-1.50	-46.45	-2.80
	MOMI-XZ	1.40	215.31	1.04	77.56
	XLogP	-3.96	15.99	-4.08	7.04
	nRotB	0.00	32.00	0.00	15.00
	nHBDOn	0.00	6.00	0.00	3.00

946

947

948

949 **Table 5b:** Applicability domain of the selected descriptors in regression nano-QSARs under  
 950 different case studies

Case Studies	Descriptors	Training Set		Test Set	
		Min	Max	Min	Max
Case study-1	Size	20.00	74.00	20.00	31.00
	R1	0.00	36.00	0.00	32.00
	R2	0.00	153.00	0.00	62.00
	ZP	-37.00	5.90	-13.60	1.95
Case study-2	VP-4	0.00	5.73	0.00	3.39
	VPC-6	0.00	23.03	0.00	2.94
	IP	-1.00	9.05	-1.00	9.05
	nRotB	0.00	32.00	0.00	24.00
	nHBAcc	1.00	11.00	1.00	8.00
Case study-3	PSA	17.07	43.37	17.07	43.37
	MolRef	1.44	4.33	1.44	4.24
	OP	10.3	53.26	14.73	47.07
Case study-4	SP-5	1.74	8.17	3.29	3.99
	Kier3	3.62	7.80	4.59	7.16
	MDEC-22	5.06	14.14	5.08	10.30
	WTU	9.86	43.82	14.64	24.23
	MOMI-YZ	1.77	14.39	4.22	12.58
	XLogP	0.70	4.59	2.14	4.59

<b>Case study-5</b>	Mpol	95.80	160.83	97.61	153.70
	Mpa	76.52	137.82	75.93	111.66
	AC	76.00	132.00	80.00	115.00
	BC	107.00	166.00	111.00	148.00
	CBC	0.00	24.00	0.00	16.00
	AAC	18.00	58.00	22.00	48.00
	HRC	0.00	2.00	0.00	0.00
	PI	388.00	496.00	384.00	456.00
	BI	0.54	0.89	0.60	0.88
	WI	11438.00	35826.00	11488.00	27807.00

---

951

952

953

954