# Quick and Efficient Quantitative Predictions of Androgen Receptor Binding Affinity for Screening Endocrine Disruptor Chemicals Using 2D-QSAR and Chemical Read-Across

**Arkaprava Banerjee[a], Priyanka De[a], Vinay Kumar[a], Supratik Kar[b], Kunal Roy[a,*]**

[a]Drug Theoretics and Cheminformatics Laboratory, Department of Pharmaceutical Technology, Jadavpur University, Kolkata700032, India

[b]Interdisciplinary Center for Nanotoxicity, Department of Chemistry, Physics and Atmospheric Sciences, Jackson State University, Jackson, Mississippi 39217, United States

*For Correspondence

E-mail: kunalroy_in@yahoo.com; kunal.roy@jadavpuruniversity.in;

URL: https://sites.google.com/site/kunalroyindia;

Tel: +91 98315 94140

24 **Abstract**

25 Endocrine Disruptor Chemicals are synthetic or natural molecules in the environment that

26 promote adverse modifications of endogenous hormone regulation in humans and/or in

27 animals. In the present research, we have applied two-dimensional quantitative structure-

28 activity relationship (2D-QSAR) modeling to analyze the structural features of these

29 chemicals responsible for binding to the androgen receptors (logRBA) in rats. We have

30 collected the receptor binding data from the EDKB database (https://www.fda.gov/science-

31 research/endocrine-disruptor-knowledge-base/accessing-edkb-database) and then employed

32 the **DTC-QSAR** tool, available from https://dtclab.webs.com/software-tools, for dataset

33 division, feature selection, and model development. The final partial least squares was

34 evaluated using various stringent validation criteria. From the model, we interpreted that

35 hydrophobicity, steroidal nucleus, bulkiness and a hyrdrogen bond donor at an appropriate

36 position contribute to the receptor binding affinity, while presence of electron rich features

37 like aromaticity and polar groups decrease the receptor binding affinity.  Additionally we

38 have also performed chemical Read-Across predictions using **Read-Across-v3.1** available

39 from https://sites.google.com/jadavpuruniversity.in/dtc-lab-software/home, and the results for

40 the external validation metrics were found to be better than the QSAR-derived predictions.

41 To explore the essential features responsible for the receptor binding, pharmacophore

42 mapping, molecular docking along with molecular dynamics simulation were also performed,

43 and the results are in accordance with the QSAR findings.

44

45 **Keywords:** Endocrine disruptors; Androgen receptor binding affinity; QSAR; Read-across;

46 docking; Pharmacophore

47

48

**1. Introduction**

It is fascinating that our brain is responsible for almost every physiological function that our body performs. The hypothalamus, also known as our "built-in thermostat" is the control centre for the endocrine system, which comprises various ductless chemical messengers commonly termed as hormones. In nature, there is existence of molecules which can potentially mimic these chemical messengers and bring about "disruption" in the normal physiological functioning of the body. Such compounds are classified as Endocrine Disrupting Chemicals (EDCs) as they mimic the natural hormones, bind to the specific receptors and bring about endocrine disruption in humans and wildlife [1-4]. In 2011, *Schug et al.* [5] reported that EDCs show various neurological, reproductive and cardiovascular adverse effects by interfering with the synthesis, transport, metabolism and release of hormones. However, it has also been observed that EDCs can act on transcriptional coactivators, synthesis and metabolism of steroids, non-steroidal receptors and various other mechanisms that ultimately converge to endocrine and reproductive systems [5]. The complexity in the mechanism of disruption in endocrine functions and activation of signaling pathways probably explains the reason for the lack of experimental toxicity data of EDCs [6]. As compared to estrogenic mode of disruption, little is known about how EDCs adversely affect the androgen receptors and hinders the male reproductive tract health [7]. Among various other targets, chemicals like DDTs, industrial chemical phthalates, organophosphate insecticides like parathion and herbicides of phenylurea derivatives like linuron can potentially bind to the androgen receptor and bring about disruption thus resulting in the toxicity [1].

Development and maintenance of male sexual characteristics is controlled by Androgen Receptors (AR), a class of ligand-activated transcriptional regulatory protein [8]. Most androgenic EDCs perform activation of transcription through receptor mediated mechanism

74  [9]. Using this information, it is possible to identify the potential EDCs through the

75  competitive binding assay at the AR. **Figure S1 (Supplementary Material SI-1)** represents

76  the potential of EDCs in inhibition of the androgen receptor inside the mammalian cell.

77

78  The Organization for Economic Co-operation and Development (OECD) promotes the use of

79  *in silico* approaches wherever applicable. As the resources are limited, it is highly impractical

80  to perform toxicity assessment of all EDCs against all possible end points in the exploration

81  of different disruption mechanisms experimentally [6]. Thus, with the aim of data gap filling,

82  efficient *in silico* approaches with scientifically well defined algorithms are adopted. In

83  recent times, there has been an increase in non-testing methods which comply with the 3Rs

84  (Reduction, Replacement and Refinement in animal experiments) in scientific

85  experimentations [10]. Among various other non-testing methods, Quantitative Structure-

86  Activity Relationship (QSAR) and Chemical Read-Across are two of the most widely used

87  methods for prediction of toxicity associated with chemicals [10-11]. The advantages

88  associated with *in silico* approaches in general are: a) they reduce experimental time, cost and

89  b) they speed up obtaining the desired results. The basic concept behind regression-based

90  QSAR lies in the development of a model consisting of the dependent variable (response) and

91  one or more features in the molecules (independent variables) which contribute to the

92  response values either positively or negatively and is expressed in numerical terms. Read-

93  Across, on the other hand, is performed by extrapolating the outcome of hazard identification

94  from certain source chemicals to one or more target chemicals based on "similarity" between

95  the source compound(s) and the target compound [11] and it does not involve the

96  development of supervised learning models.  Both of these approaches are mainly used for

97  two purposes: 1) to predict end point values of a completely new set of chemicals for the

98  purpose of filling data gaps (predictive models) and 2) mechanistic and physicochemical

99      interpretation of the structural features in a molecule which are responsible to elicit the

100     response [12].

101     In the recent past, efforts have been made to predict the binding affinity of various EDCs to

102     the androgen receptors using computational approach. *Hong et al.* [13] in 2003 studied the

103     binding affinity of natural, synthetic and environmental chemicals to the androgen receptor

104     by Comparative Molecular Field Analysis (CoMFA) (a 3D-QSAR approach), and they

105     inferred that the steric and electronic properties of the training compounds are essential in

106     describing the binding affinity of EDCs to the androgen receptor. In 2002, *Serafimova et al.*

107     [14] studied the active formulation ingredients of pesticides and their ability to bind to the

108     androgen receptor and performed their evaluation using COREPA method. They have

109     utilized stereochemical properties like the inter-atomic distances between the nucleophilic

110     sites and their charges and used them to predict the binding affinity in terms of $pK_i$. *Piir et*

111     *al.* [15] in 2020 performed binary and multi-class classifications for antagonists, agonists and

112     binders to the AR by implementing random forest classification models. They stated that the

113     accuracy obtained in their multi-class classification was good considering the large size of the

114     training set that they have utilized.

115     3D-QSAR methods involve computational complexity of conformational analysis and

116     alignment and inherit the property of being non-reproducible in nature. The novelty of the

117     current work is predicting the binding affinity of endocrine disruptors to the androgen

118     receptors in a quantitative and reproducible manner. The data was obtained from Endocrine

119     Disruptor Knowledge Base (EDKB) database (https://www.fda.gov/science-

120     research/bioinformatics-tools/endocrine-disruptor-knowledge-base) thus avoiding personal,

121     systemic or instrumental error in data collection. It was then divided into a modeling set and a

122     validation set based on the availability of experimental response values in terms of log RBA,

123     where RBA stands for Receptor Binding Affinity. A regression-based 2D-QSAR model was

124    generated using the modeling set, and subsequently similarity-based chemical Read-Across

125    was also performed.  The reliability of both of the approaches was evaluated using various

126    strict validation metrics. The physicochemical interpretation of different possible mechanisms

127    influencing the binding of EDCs to the androgen receptor were also discussed and reported

128    which can ultimately help a chemist recognize the features in a molecule that has potential to

129    cause androgen receptor toxicity. In support of this theory, pharmacophore mapping was also

130    performed to serve the purpose of screening of the features in a molecule which contribute to

131    AR binding affinity. Analysis of the binding of the ligand to the various amino acid residues

132    in the receptor was also done with the help of molecular docking and the stability of such

133    binding was evaluated using molecular dynamics (MD) simulation at 100 ns.

134

135    **2. Materials and methods**

136    *2.1 Collection of Androgen Receptor Binding Affinity data of EDCs and curation of their*

137    *structures*

138    The androgen Receptor Binding Affinity (RBA) data of various EDCs were collected from

139    the Endocrine Disruptor Knowledge Base (EDKB) database (https://www.fda.gov/science-

140    research/bioinformatics-tools/endocrine-disruptor-knowledge-base) obeying the strict OECD

141    guidelines.    The    chemical    structures    downloaded    from    PubChem    database

142    (https://pubchem.ncbi.nlm.nih.gov/) in .sdf format were represented in Marvin Sketch

143    (https://chemaxon.com/products/marvin) software. Chemical curation of our compounds was

144    performed by the application of a KNIME workflow (https://sites.google.com/site/dtclabdc/)

145    taking the single .sdf file as input. Further details are available in **Supplementary Material**

146    **SI-1.**

147

148    *2.2 Calculation of molecular descriptors and data pre-treatment*

149  Descriptors are certain properties in a molecule encoded in numerical terms which can be

150  handled statistically. Two molecules are said to be "identical" or 100% similar if they have

151  identical set of descriptor values. The descriptors for our curated compounds were calculated

152  using alvaDesc v2.0.6 [16]. To enhance simplicity in the interpretation of the developed

153  model, we have used only selected classes of descriptors **(Supplementary Material SI-1).**

154  The inter-correlated descriptors having correlation values >0.95 and variance cut-off 0.00001

155  were removed using the Java-based tool Data Pretreatment GUI 1.2 available from

156  https://dtclab.webs.com/software-tools.

157

158  *2.3 Dataset division and model development*

159  Dataset division into training and test sets during a QSAR model development ensures the

160  models' predictive ability. In the present study, the available data set was segregated into two

161  classes: 1) the modeling set which comprises the compounds having reported response values

162  in terms of log RBA and 2) the validation set consisting of compounds for which the response

163  values were not reported. We have eliminated six compounds from our modeling set due to

164  their aberrant nature of activity. The reduced modeling set was taken as an input for the java-

165  based software tool DTC-QSAR v1.0.5 (https://dtclab.webs.com/software-tools), where we

166  performed division into training and test sets in 70:30 ratios based on Euclidean Distance

167  method [17], and feature selection was done by employing Genetic Algorithm technique [18].

168  The descriptors obtained from the set of GA-MLR models were then pooled and the best

169  descriptor combinations from all possible models were obtained by using Best Subset

170  Selection (BSS) v2.1 available from https://dtclab.webs.com/software-tools. The Best Subset

171  Selection tool generates models based on all possible combination of descriptors, and one can

172  select the best models based on validation metrics like $r^2$, $Q^2_{LOO}$, $MAE_{95\%}$, $Q^2_{F1}$ and $Q^2_{F2}$. To

173  nullify the inter-correlation among descriptors, the final Partial Least Squares (PLS)

7

174  regression model was obtained with the best descriptor combination taking three latent

175  variables, and various internationally accepted validation metrics were calculated [19-20]

176  (**Supplementary Material SI-1**).

177

178  *2.4 DModX Applicability Domain Plots*

179  The Applicability Domain can be termed as a theoretical region in chemical space which

180  surrounds both the descriptors and response [21]. The distance to model in X-space (DModX)

181  approach was implemented to check the applicability domain of the model.

182

183  *2.5 Similarity based Read-Across prediction*

184  What differentiates Read-Across approach from classical QSAR is that Read-Across is

185  entirely a similarity-based approach which does not involve the development of a statitstical

186  model. QSAR models become statistically unreliable when there are limited number of data

187  points [11] and contrastingly, read-across approach not being a hardcore statistical approach

188  tends to yield better results even for small datasets and thus can be aimed for data gap filling.

189  In the present work, after performing feature selection, we have divided the training set

190  compounds into sub-training and sub-test sets based on Euclidean distance-based division.

191  These sets were further used for hyperparameter optimisation in the Read-Across-v3.1

192  (https://sites.google.com/jadavpuruniversity.in/dtc-lab-software/home) tool. The optimised

193  hyperparameters were then used for the original training and test set files as input.

194

195  *2.6 3D-Pharmacophore mapping*

196  In this investigation, 3D-Pharmacophore mapping was implemented to explore the potential

197  features that are crucial for the interaction at the active site of the androgen receptor. The

198  receptor binding affinity (RBA) expressed as logRBA was used as the dependent variable to

199  develop the pharmacophore models. Molecules prepared for the 2D-QSAR model

200  development were used for this study. The dataset was rationally divided into training (30

201  compounds for hypothesis development) and test (115 compounds for validation) sets based

202  on the logRBA values spanning four orders of magnitude. 3D-Pharmacophore modeling was

203  performed using HypoGen algorithm as embedded in Biovia Discovery Studio Client 4.1

204  client [22] following the protocol as discussed by *Kumar et al.* [23]. Details of the protocol

205  performed for 3D-Pharmacophore modeling is provided in **Supplementary Material SI-1**.

206  Validation of the obtained models was executed using different parameters such as cost

207  analysis, the Fischer randomization test (F-test), and test set prediction to evaluate the

208  robustness and predictive ability of models as discussed by *Kumar et al.* 2020 [23].

209

210  *2.7 Molecular docking study*

211  Molecular docking study was performed to predict the potential of complex formation and

212  explore the binding mode of the compounds showing the highest and lowest binding affinity

213  to the androgen receptor. The crystal structure of the protein was extracted from the protein

214  databank by the PDB ID: 3G0W [24] (available from https://www.rcsb.org/structure/3G0W).

215  A rigid docking approach was applied using the CDOCKER with a grid-based protocol [25]

216  for the aim of the receptor-ligand interaction, as prompted in Biovia Discovery Studio Client

217  4.1 client [22] following the protocol as discussed by *Kumar et al.* [23]. Details of the

218  protocol performed for molecular docking is provided in **Supplementary Material SI-1**.

219  After molecular docking, the docked inclusion complexes with the best ranked CDOCKER

220  interaction energy and bond formation between compounds and active amino acid residues

221  were chosen for the detailed interpretation and correlation. We have also validated the

222  docking protocol by redocking the bound ligand at the protein's active site (**Figure S2**)

223  (**Supplementary material SI-1**) and calculating the RMSD (**Figure S3**) (**Supplementary**

224    **material SI-1)** with the bound ligand and the redocked ligand. The ligplot shows the number

225    of interactions and active amino acids responsible for the important interaction in the crystal

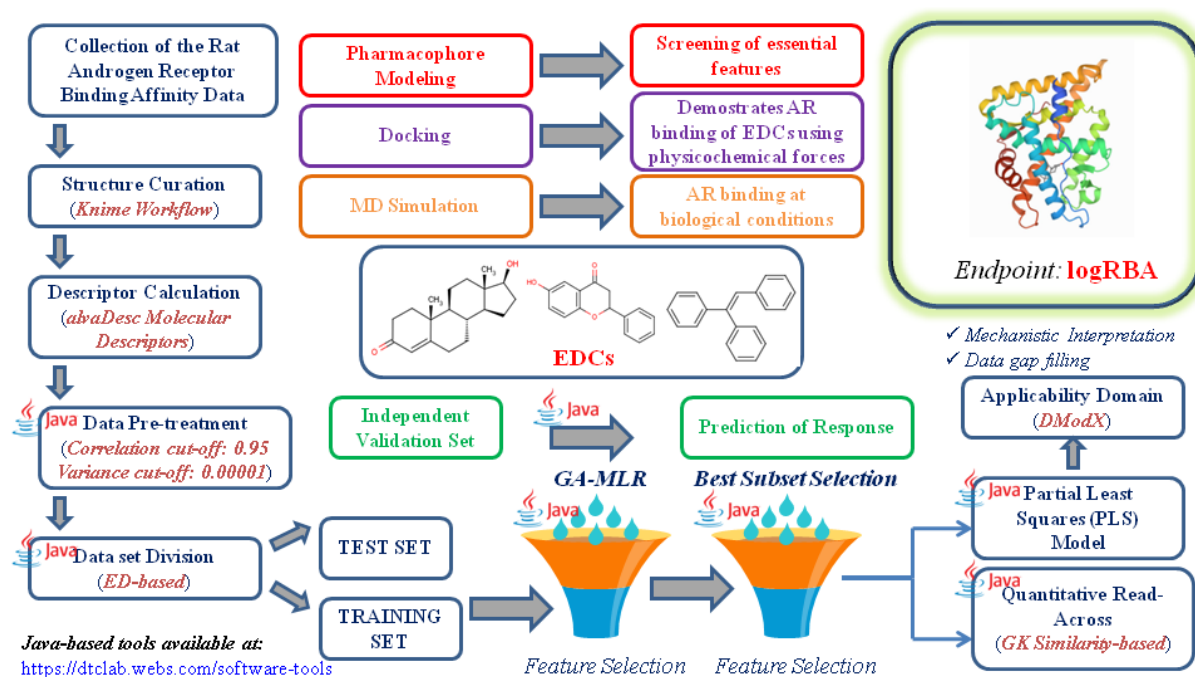226    structure of androgen receptor and with their bound ligand.

227

228    *2.8 Molecular dynamics (MD) simulation and MM/GBSA-Binding energy calculation*

229    Further to study the stability of ligand-receptor complex at biological conditions, molecular

230    dynamics simulation at 100ns was performed [26-29], and receptor binding affinity using

231    MM/GBSA [30] method was calculated.

232

233    The whole workflow of multiple cheminformatic applications applied to the ARB data set has

234    been depicted pictorially in **Figure 1**.



236    **Figure 1.** Schematic representation of the workflow of cheminformatic applications used in

237    this study.

238

239    **3. Results & Discussion**

240    **3.1 2D-QSAR analysis**

241 The modeling data set has been provided in an Excel sheet in the **Supplementary Material**

242 **SI-2**. The training set comprises 103 EDCs that were used for model development while the

243 test set comprises 44 EDCs that were used for prediction and external validation. The final

244 PLS equation with three Latent Variables is shown in Eq. (1). The descriptors have been

245 mentioned in the descending order of importance as per the Variable Importance Plot (**Figure**

246 **2**).

247

248 $LogRBA = -3.23 + 0.49 \times SsssCH - 0.41 \times MaxaaCH + 0.23 \times nCconj + 0.35 \times$

249 $LogP99 - 0.17 \times F10[C - O] + 0.06 \times minsOH + 0.06 \times N\% + 0.67 \times F08[O - F]$

250 (1)

251 $$R^2_{(TRAIN)} = 0.74, Q^2_{(LOO)} = 0.68, Q^2_{F1} = 0.58, Q^2_{F2} = 0.58$$

252 $$Scaled\ average\ r^2_m(Train) = 0.57, Scaled\ average\ r^2_m(Test) = 0.50$$

253 $$Scaled\ delta\ r^2_m(Train) = 0.18, Scaled\ delta\ r^2_m(Test) = 0.07$$

254 $$MAE_{(TRAIN)} = 0.46, MAE_{(TEST)} = 0.54, n_{(Training)} = 103, n_{(Test)} = 44$$

255

256 The statistical quality and internal and external validation metric values of the QSAR model

257 are satisfactory considering the diversity and heterogeneity of the data set. The descriptors

258 selected in the QSAR model are detailed below **(Figure S4 in Supplementary Materials SI-**

259 **1).** The different plots [10] related to the PLS model are provided in **Figures S5-S9 in**

260 **Supplementary Materials SI-1**.

261

262 **3.1.1 Descriptors contributing to Hydrophobicity**

263 In the final PLS model, we have obtained a total set of 6 descriptors contributing positively to

264 the response out of which some are responsible for directly influencing the hydrophobic

265 properties of the molecules (e.g., LOGP99, SsssCH, nCconj, F08[O-F]) while some induce

266   hydrophobicity indirectly (e.g., minsOH). The **SsssCH** descriptor stands for sum of E-states

267   of sssCH (tertiary carbon atoms) [31]. In this data set, the compounds containing a sterioidal

268   (cyclopentanoperhydrophenathrene) nucleus shows higher values for this descriptor. This

269   suggests that for a higher receptor binding affinity, presence of the steroidal nucleus is

270   preferred. The present dataset includes 5α-Androstan-17β-ol (**23**) which has a higher SsssCH

271   descriptor value and shows a higher receptor binding affinity, as compared to 4-

272   Hydroxybiphenyl (**148**) which is devoid of tertiary carbon atoms (**Figure 2**). The **nCconj**

273   descriptor signifies the number of non-aromatic conjugated carbons ($sp^2$), and it positively

274   correlates with the response values as in the case of Trenbolone (**157**), which has a higher

275   number of non-aromatic conjugated carbon atoms ($sp^2$) thus resulting in enhanced receptor

276   binding affinity while in case of Aldrin (No. **176**), where the sp2 carbons are not in

277   conjugation, exhibit a much lower receptor binding affinity. In the steroidal structures of the

278   data set, the descriptor nCconj actually signifies the importance of the conjugated enone

279   moiety in ring A (like **67**), as the keto group at 3 position serves as an hydrogen bond

280   acceptor (see molecular docking in a later section). The descriptor **LOGP99** stands for

281   Wildmann-Crippen octanol-water partition coefficient, and it positively contributes to the

282   response values, as an increase in the o/w partition coefficient value increases the lipid

283   solubility. For instance, Dihydrotestosterone benzoate (**134**) has a high LOGP99 value, and

284   thus has a higher receptor binding affinity compared to Diethyl phthalate (**34**) which has a

285   lower partition coefficient value. The descriptor **minsOH** stands for minimum E-state of the

286   sOH hydroxyl group [31]. This can be attributed to the inherent property of the hydroxyl

287   groups to be able to form hydrogen bond interactions with the receptor residues in an

288   appropriate location [32] and thus contributes to the enhancement in the receptor binding

289   affinity of the molecule. A higher minsOH value also signifies that there is a large

290   hydrophobic moiety attached to the hydroxyl group, thus bulkiness of the structure also

291     contributes to the overall hydrophobic property. The presence of OH group at a desired

292     location as well as its attachment to a bulky moiety in Norgestrel (**67**) is the reason for its

293     high receptor binding affinity whereas the molecule Aldrin (**176**) lacks the hydroxyl group

294     and results in lower receptor binding affinity. **N%** denotes the percentage of nitrogen present

295     in the molecular structure, and it shows a positive contribution to the response. In a previous

296     work, *Zhou et al.* stated that nitrogen in the form of primary amino group can be

297     accommodated in the same location as the hydroxyl group (probably due to the bio-isosteric

298     nature of O and NH) and thus can actively participate in hydrogen bonding with the receptor

299     residues like Asn705 [33] resulting in enhanced receptor binding affinity, as also

300     demonstrated in our model. Due to the presence of Nitrogen in Carbaryl (**72**), it exhibits

301     slightly higher receptor binding affinity than Bis(n-octyl) phthalate (**114**) which is devoid of

302     nitrogen atoms. The descriptor **F08[O-F]** stands for frequency of O and F atoms at the

303     topological distance of 8. The presence of F atoms can induce polarity, but the presence of O

304     at the topological distance of 8 suggests that the compounds are bulky in nature, thus

305     overshadowing the polar effects with the hydrophobic properties contributed due to bulkiness

306     of the structure. Presence of a lipophilic -CF3 group in Hydroxyflutamide (**187**) ensures

307     higher receptor binding affinity while 17α-Estradiol (**7**) is devoid of $CF_3$ atoms and does not

308     tend to bind well to the receptor.

309

310     **3.1.2 Descriptors contributing to Polarity and Electron Richness**

311     Out of the total 8 descriptors obtained in our model, two of them correlate negatively to the

312     response values and induce polarity and electron richness to the molecules. One of the

313     descriptors is **MaxaaCH,** which stands for maximum E-state of aaCH (aromatic CH groups)

314     [31]. This is probably due to the fact that aromatic compounds are comparatively more polar

315     than their alicyclic counterparts. This can be observed in 3-methyl-estriol (**102**) with an

316    aromatic ring showing reduced receptor binding affinity as compared to 3β-Androstanediol

317    (**183**) which is devoid of any aromatic ring and thus exhibiting higher receptor binding

318    affinity. The other descriptor is **F10[C-O]** which stands for frequency of C and O at the

319    topological distance 10. This descriptor depicts the presence of polar functionalities like

320    hydroxyl, ether or ester groups. It is to be noted that the hydroxyl group as minsOH

321    contributes positively to the receptor binding affinity due to its ability to form a hydrogen

322    bond at a desired location while attached to a bulky scaffold. Therefore, it can be concluded

323    that F10[C-O] descriptor actually acts to compensate that effect with the polar effects of OH

324    and this can be confirmed with the near-equal and opposite values of the standardized

325    coefficients of both these descriptors in our PLS model. Our dataset contains Dexamethasone

326    (**75**) which shows lower receptor binding affinity than Triphenylethylene (**4)** as the latter

327    lacks polar functionalities like hydroxy, ether or ester groups.  The hydrogen bond donor

328    group should be present at a specific position like 17 position of the steroidal nucleus as in

329    5α-Androstan-17β-ol (**23**) to participate in the hydrogen bonding interaction with the receptor

330    functionalities (see Molecular Docking in a later section). Presence of polar functionality at

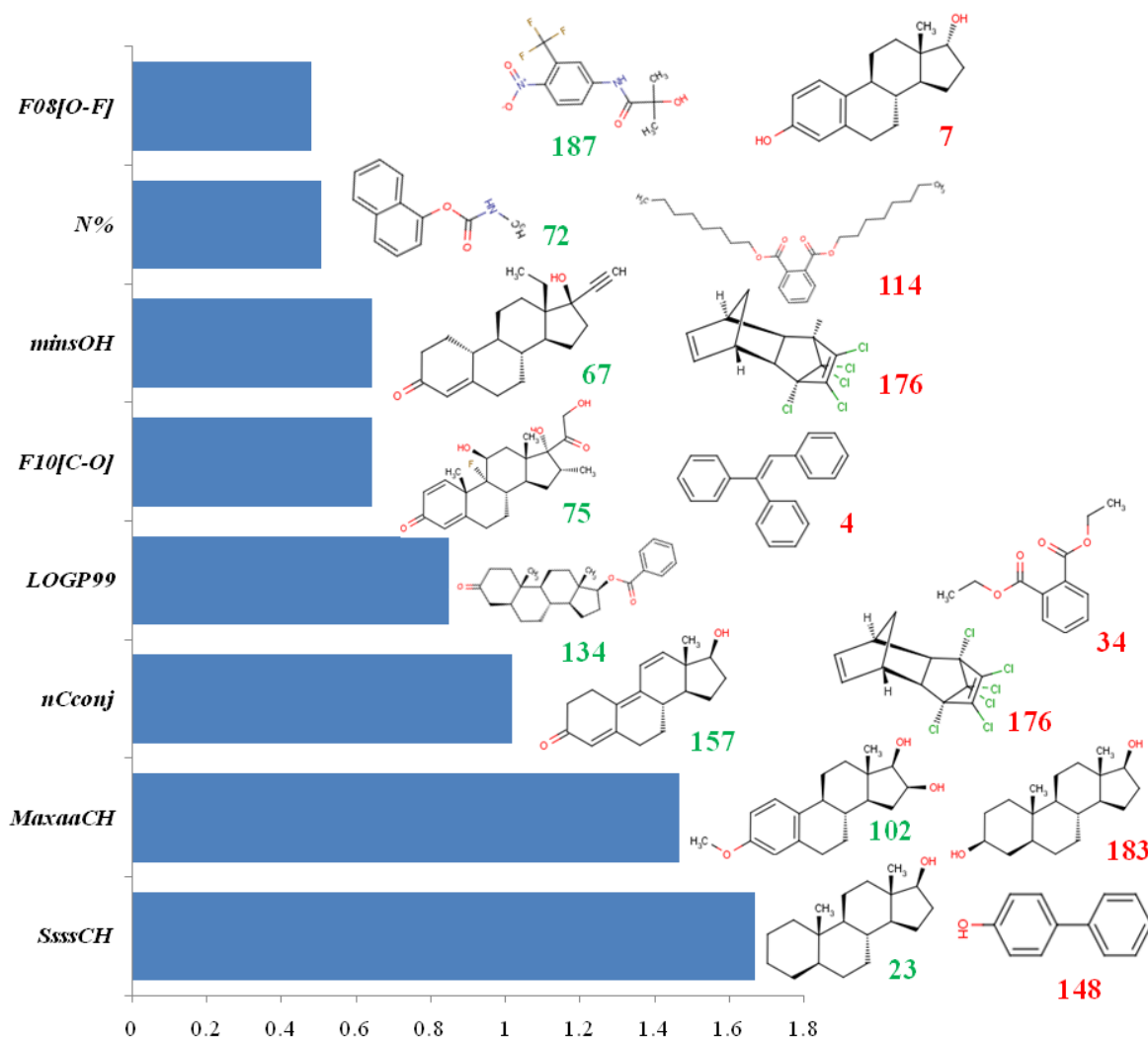331    any other locations decrease the RBA.

332

333

334

335

336

337

**Figure 2.** Variable Importance Plot. Structures of representative compounds having higher
and lower values of individual descriptors are also shown

### 3.1.3 Predictions for the validation set

Prediction of the receptor binding affinities for the compounds which constitute the validation
set was performed using a java-based software tool Prediction Reliability Indicator (PLS
Version) [34] available from https://dtclab.webs.com/software-tools. The results obtained

347 depicts that out of the 55 compounds, 12 were outside the applicability domain with

348 Bad/Unreliable prediction quality and among the remaining 43 compounds, two of them have

349 moderate prediction quality and the others have good prediction quality. The results of this

350 prediction is provided in an Excel sheet in the **Supplementary Material SI-2**.

351

352 **3.2 Chemical Read-Across results**

353 After QSAR model development, the same training and test set compounds were taken as

354 inputs for quantitative Read-Across-based predictions using the same input features as

355 descriptors, while implementing three different similarity functions: the Euclidean Distance-

356 based, the Gaussian Kernel Similarity-based and the Laplacean Kernel Similarity-based

357 predictions, and after optimization of the hyper-parameters, it was found that the external

358 validation results obtained from quantitative Read-Across algorithm using Gaussian Kernel

359 Similarity-based functions were better compared to the results obtained using QSAR and also

360 the other two read-across approaches **(Figure S10) (Supplementary Material SI-1)**. The

361 higher values of $Q_{F1}^2$ (0.64), $Q_{F2}^2$ (0.64) and lower $MAE_{(TEST)}$ (0.47) in Read-Across suggest

362 that predictive ability of the Read-Across algorithm was even better for predictions as

363 compared to the classical QSAR approach. It appears that the local similarity-based approach

364 gives better predictions over model-derived predictions obtained from the whole training data

365 set. The results of this prediction is provided in an Excel sheet in the **Supplementary**

366 **Material SI-2**.

367

368 **3.3 Comparison of present 2D-QSAR and Read-Across with previous models**

369 We have developed here an easily reproducible and transferable 2D-QSAR model using

370 simple interpretable descriptors. *Hong et al.* [13] employed Comparative Molecular Field

371 Analysis (CoMFA) (a 3D-QSAR approach) by taking similar number of data points and the

372 corresponding quality and validation metrics were $r^2 = 0.902$ and $q^2 = 0.571$ which

373 suggests that their model is less robust due to a high difference between $r^2$ and $q^2$ values.

374 Also it is important to note that CoMFA methodology requires conformation analysis and

375 alignment of the molecules making the results less reproducible. *Piir et al.* [15] applied

376 binary and multi-class classification techniques generating only qualitative results whereas

377 our model generates quantitative predictions. Thus, it can be concluded that our model is

378 robust, predictive (due to acceptable values of the external validation metrics) and

379 reproducible. **Table 1** depicts how our QSAR model and Read-Across based predictions

380 supersedes the previous results in the quantitative prediction quality.

381

382 **Table 1**: Comparison with the previous studies

| Authors | Method | $n_{(Train)}$ | $n_{(Test)}$ | End Point | $R^2$ | $Q^2$ | $Q_{F1}^2$ | $Q_{F2}^2$ | Inference |
|---|---|---|---|---|---|---|---|---|---|
| *Hong et al.* [13] | 3D-QSAR (CoMFA) (Regression) | 146 | 8 | logRBA | 0.90 | 0.57 | - | - | Less robust, non-reproducible |
| *Piir et al.* [15] | Classification-based QSAR | 1688 | 5273 | AR Activity | - | - | - | - | Graded predictions only |
| **Our work** | 2D-QSAR (Regression) | 103 | 44 | logRBA | 0.74 | 0.68 | 0.58 | 0.58 | Robust, Predictive, Reproducible |
| **Our work** | Quantitative Read-Across | 103 | 44 | logRBA | - | - | 0.64 | 0.64 | Predictive, Reproducible |

383

**3.4 3D Pharmacophore modeling analysis**

384

385 In this analysis, we have developed ten different 3D- pharmacophore hypotheses from a

386 training set of 30 compounds. The robustness of the generated models in terms of fitness,

387 stability, classical fitness metrics, and predictability was examined using stringent validation

388 metrics. In terms of internal validation, all the developed models were showing excellent

389 results, thus for the selection of the best hypothesis, we have checked the performance on the

390 test set. External validation of the developed models was implemented by mapping the test

391 set compounds with the same settings applied for the pharmacophore generation by the FAST

392 method. After analysis (**Table S1**) (**Supplementary Material SI-1),** Hypo-8 was found to be

393 the best one among the ten hypotheses with one Hydrogen bond acceptor (HBA), two

394 Hydrophobic (HYD), and one Hydrogen bond donor (HBD) features (**Figure S11**)

395 **(Supplementary Material SI-1).** In terms of internal validation, the best pharmacophore

396 model (Hypo 8) was obtained (**Table S1**) (**Supplementary Material SI-1)** in the cost

397 analysis with a higher correlation coefficient (R: 0. 757), total cost (329.866), maximum fit

398 (10.809), configuration cost (12.097) and higher cost difference (287.88). These values stated

399 that the selected model was appropriate in terms of internal quality metrics. After mapping,

400 we found that 27 compounds from the data set were correctly mapped and predicted, whereas

401 88 compounds were not mapped due to the absence of features found in the select

402 pharmacophore model. Out of these 88 compounds, 79 compounds have the ARB affinity

403 lower than the training set mean suggesting that these are low affinity compounds due

404 absence of the required pharmacophoric features (and hence not mapped).  The observed and

405 predicted values of the training and test set molecules obtained from the analysis using Hypo-

406 8 are given in **Sheets 2 and 3 (Supplementary Material SI-3).** We have developed a Java-

407 based software tool **Klassification1.0** for calculating the classification metrics and the tool is

408 now made available online at https://sites.google.com/jadavpuruniversity.in/dtc-lab-

409 software/home. The test set statistics are based only on the mapped compounds. The Fisher

410 validation test confirms the non-randomness of the selected pharmacophore (Hypo-8) model.

411 The total correlation and cost values obtained from the original and randomized models of the

412 hypothesis for the Fisher validation test are stated in **Sheets 4 and 5** in the **Supplementary**

413 **Material SI-3**. Additionally, the validated pharmacophore model was used to estimate the

414 affinity of the external dataset of 55 compounds, with no quantitative observed response

415 values in the source file. After prediction, we have found that only 13 compounds were

416 correctly mapped and predicted, whereas 42 compounds were not mapped due to the absence

417 of features found in the select pharmacophore model, out of the listed 55 compounds (see

418 **Sheet 6** in **Supplementary Material SI-3**). We have also predicted the 6 compounds omitted

419 from the original dataset because of their outlier behavior in the initial modeling (2D-QSAR)

420 exercises. After prediction, we have found that only 4 compounds were properly mapped and

421 estimated and whereas 2 compounds were not mapped due to the absence of features found in

422 the select pharmacophore model (see **Sheet 7** in **Supplementary Material SI-3**).
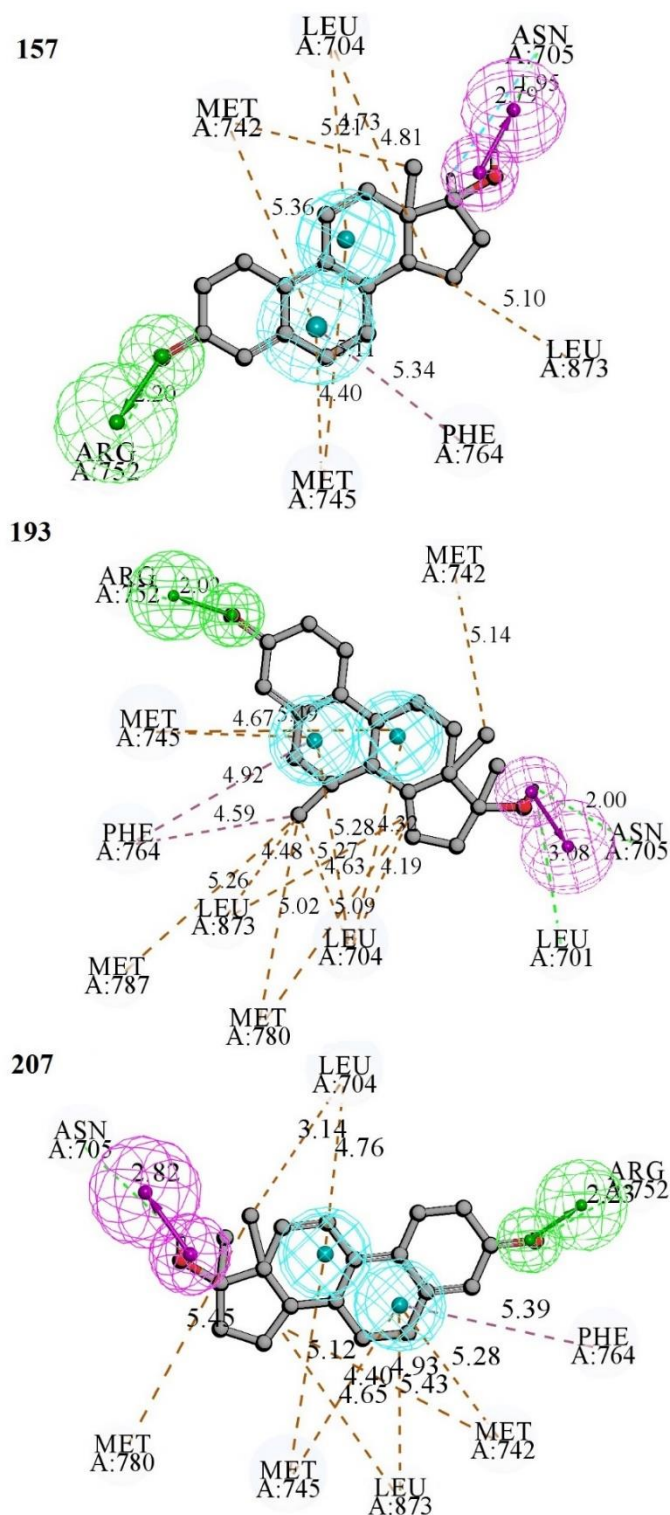
423

424 **3.5 Molecular docking analysis**

425 *3.5.1 Molecular docking analysis of the compounds with the highest and lowest binding*

426 **affinities from the dataset**

427 We have implemented the molecular docking using the three compounds with the highest

428 ARB (compound **157, 193,** and **207**) and three compounds with the lowest ARB (compound

429 **34, 87,** and **114**) from the whole dataset, to explore the potential interactions at the active

430 pocket of androgen receptor. The detailed information of docking interactions, CDOCKER

431 interaction energy, and their correlation with the features derived from the developed best

432 2D-QSAR model are illustrated in **Table S2** in **Supplementary Material SI-1.**

433

434 **3.5.1.1 Molecular docking analysis of the compounds with the highest binding affinity**

435 **from the dataset**

436 One of the highest ARB compounds from the dataset is compound **157,** which interacted with

437 the active site pocket of the receptor (**Figure 3**) *via* hydrogen bonding with the amino acid

438 residues ASN A: 705, ARG A: 752 in the distance of 1.95, 2.79 and 2.20 Å respectively, π-

439 alkyl hydrophobic bond with amino acid residue PHE A: 764 in the distance of 5.34 Å, and

440 alkyl hydrophobic bonding with the amino acid residues LEU A: 704, MET A: 742, MET A:

441 745, LEU A: 873 in the distance of 4.81, 4.73, 5.21, 5.36, 4.40, 5.11, 5.10 Å respectively.

**Figure 3.** Molecular docking interactions and correlation with pharmacophore model of the compound with the highest binding affinity (Compound **157, 193, 207**) from the dataset.

445

446 The next highest ARB compound in this series from the dataset is compound **193**, which

447 interacted with the active site pocket of the receptor (**Figure 3**) *via* hydrogen bonding with

448 the amino acid residues ARG A: 752, ASN A: 705, LEU A: 701 in the distance of 2.02, 2,

449 3.08 Å respectively, π-alkyl hydrophobic bond with amino acid residue PHE A: 764 in the

450 distance of 4.92, 4.59 Å, and alkyl hydrophobic bonding with the amino acid residues MET

451 A: 742, MET A: 745, MET A: 787, MET A: 780, LEU A: 873, LEU A: 704 in the distance of

452 5.14, 4.67, 5.49, 5.26, 5.02, 5.09, 4.48, 5.27, 4.63, 5.28, 4.32, 4.19 Å respectively.

453

454 The third highest ARB compound from the dataset is **207,** which interacted with the active

455 site pocket of the receptor (**Figure 3**) *via* hydrogen bonding with the amino acid residues

456 ASN A: 705, ARG A: 752 in the distance of 2.82, 2.23 Å respectively, π-alkyl hydrophobic

457 bond with amino acid residue PHE A: 764 in the distance of 5.39 Å, and alkyl hydrophobic

458 bonding with the amino acid residues LEU A: 704, MET A: 780, MET A: 745, LEU A: 873,

459 MET A: 742 in the distance of 3.14, 4.76, 5.45, 5.12, 4.40, 4.65, 4.93, 5.43, 5.28 Å

460 respectively.

461

462 The molecular docking analysis of the compounds with the lowest binding affinity from the

463 data set is given in **Figures S12-S14 in Supplementary Material SI-1**. The results of

464 molecular dynamic simulation are also given in **Supplementary Material SI-1**.

465

466 **3.6 Correlation of the 3D-pharmacophore model with the molecular docking analysis,**

467 **2D QSAR, and Read-across models**

468 We have mapped the highest and least ARB compounds from the data set using the selected

469 pharmacophore model (Hypo 8) and superimposed the mapped highest ARB compounds in

470     the pharmacophore with its docking interaction showing important amino acids (**Figure 3**).

471     From **Figures S15 and S16 (Supplementary Material SI-1)** we can see that the highest

472     ARB compounds of the dataset set **157** (logRBA: 2.05) and **193** (logRBA: 2.27) mapped

473     entirely on Hypo-8 with all of the three features appearing in the model. From **Figures S15,**

474     **S16, and 3** we can see that B and C rings of the steroid nucleus lie in the hydrophobic region

475     and interact with hydrophobic amino acids (MET A: 745, PHE A: 764, MET A: 742, LEU A:

476     704) via alkyl and π-alkyl bonding (hydrophobic bond), ketone group is in the hydrogen bond

477     acceptor region, interacting with the ARG A: 752 amino acid by hydrogen bond and hydroxy

478     group lies in the hydrogen bond donor region, interacting with ASN A: 705, LEU A: 701

479     amino acids via hydrogen bond. These features are well corroborated with the SsssCH,

480     nCconj, LOGP99, and minsOH descriptors of the 2D-QSAR models and Read-across

481     hypotheses. On the other hand, the least ARB compounds of the dataset set do not map

482     entirely due to the lack of hydrogen bond donor feature in the case of compound **34** (logRBA:

483     -3.44) (**Figure S17** in **Supplementary Materials SI-1**) and hydrogen bond acceptor in case

484     of compound **92** (logRBA: -3.15) (**Figure S18** in **Supplementary Materials SI-1**). Thus, we

485     can conclude from the above discussion that the absence of any of these three features in

486     compounds reduces the receptor binding affinity against androgen receptor.

487

488

489

490     .

491

492

493

494

495    **4. Overview and Conclusion**

496    This study reports a highly robust, reproducible, easily interpretable and sufficiently

497    predictive regression-based 2D-QSAR model which is developed in accordance to the OECD

498    guidelines. This model predicts that various structural features like o/w partition coefficient,

499    bulkiness of the structure, presence of a steroid (cyclopentanoperhydrophenanthrene)

500    nucleus, number of non-aromatic conjugated carbon ($sp^2$) and hydrogen bonding to the

501    specific receptor residues contribute positively to the receptor binding affinity leading to the

502    toxicity while features like aromaticity in a molecule and presence of polar functionalities

503    like hydroxyl, ether or ester groups at additional locations in the structures lower receptor

504    binding affinity. The similarity-based Quantitative Read-Across approach was also

505    implemented according to the Gaussian-kernel similarity function using an java-based

506    software tool, and it was found that the predictive ability of the Read-Across approach

507    supersedes that of the QSAR approach as the external validation metrics were slightly better

508    in the Read-Across based predictions. The response values of our validation set were

509    calculated using the Prediction Reliability Indicator tool (https://dtclab.webs.com/software-

510    tools) thus making a successful attempt to data gap filling. Pharmacophore mapping was done

511    to screen the essential features, and it was found that a hydrogen bond acceptor, two

512    hydrophobic and one hydrogen bond donor features are essential for receptor binding affinity.

513    This information was supported by performing molecular docking analysis and it was found

514    that the molecules having highest receptor binding affinity possess all the three different

515    features that our pharmacophore hypothesis suggested. Furthermore, the docking results

516    explained the possible amino acid residues present at the surface of the androgen receptor

517    interacting with the compounds resulting in greater receptor binding affinity of the ligand.

518    Additionally, to demonstrate the receptor binding at the biological conditions, Molecular

519    Dynamics Simulation was performed. We believe that our developed QSAR model and read-

520 across approach will be useful in the screening of compounds with lower androgen receptor

521 binding affinity and will possibly tend to reduce environmental hazards.

522

523

524 **Author contributions**

525 AB: computation, validation, software tool development, initial draft, PD: computation,

526 validation and editing, VK: computation, validation and initial draft, SK: computation, initial

527 draft, editing, KR: conceptualization, supervision and editing

528

529 **Declaration of Competing Interest**

530 The authors declare that there are no competing interests.

531

537

538

539 **References**

540   1) Fang, H., Tong, W., Branham, W.S, Moland, C.L., Dial, S.L., Hong, H., Xie, Q.,

541      Perkins, R., Owens, W., Sheehan, D.M., 2003. Study of 202 natural, synthetic, and

542      environmental chemicals for binding to the androgen receptor. Chem. Res. Toxicol.

543      16, 1338-1358

544   2)  Falco, M.D., Forte, M., Laforgia, V., 2015. Estrogenic and anti-androgenic endocrine

545        disrupting chemicals and their impact on the male reproductive system; Front.

546        Environ. Sci. 3,  1-12

547   3)  Tan, H., Wang, X., Hong, H., Benfenati, E., Giesy, J.P., Gini, G.C., Kusko, R., Zhang,

548        X., Yu, H., Shi, W., 2020. Structures of endocrine-disrupting chemicals determine

549        binding to and activation of the estrogen receptor α and androgen receptor. Environ.

550        Sci. Tech. 54, 11424-11433

551   4)  Kucheryavenko, O., Vogl, S., Marx-Stoelting, P., Endocrine disruptor effects on

552        estrogen, androgen and thyroid pathways: recent advances on screening and

553        assessment. In: Mantovani, A., Fucic, A. (Eds.), Challenges in Endocrine Disruptor

554        Toxicology and Risk Assessment, Royal Society of Chemistry, London, 2021, 1-24

555   5)  Schug, T.T., Janesick, A., Blumberg, B., Heindel, J.J., 2011. Endocrine disrupting

556        chemicals and disease susceptibility. J. Ster. Biochem. Mol. Bio.. 127, 204-215

557   6)  Khan, K., Roy, K., 2019. Ecotoxicological QSAR modeling of endocrine disruptor

558        chemicals. J. Hazard Mater. 369, 707-718

559   7)  Luccio-Camelo, D.C., Prins, G.S., 2011. Disruption of androgen receptor signaling in

560        males by environmental chemicals. J.  Ster. Biochem. Mol. Bio. 127, 74-82

561   8)  Zhao, C.Y., Zhang, R.S., Zhang, H.X., Xue, C.X., Liu, H.X., Liu, M.C., Hu, Z.D.,

562        Fan, B.T., 2005. QSAR study of natural, synthetic and environmental endocrine

563        disrupting compounds for binding to the androgen receptor; SAR QSAR Environ.

564        Res. 16 (4), 349-367

565   9)  Davey, R.A., Grossmann, M., 2016. Androgen receptor structure, function and

566        biology: from bench to bedside. Clin. Biochem. Rev. 37(1), 3-15

567    10) Seth, A., Roy, K., 2020. QSAR modelling of algal low level toxicity values of

568       different phenol and aniline derivatives using 2D descriptors. Aquat. Toxicol. 228, 1-

569       11

570    11) Chatterjee, M., Banerjee, A., De, P., Gajewicz, A., Roy, K., 2022. A novel

571       quantitative read-across tool designed purposefully to fill the existing gaps in

572       nanosafety data. Environ. Sci.: Nano 9, 189-203

573    12) Ambure, P., Gajewicz, A., Cordeiro, M.N.D.S., 2019. Roy, K., New workflow for

574       QSAR model development from small data sets: small dataset curator and small

575       dataset modeler. integration of data curation, exhaustive double cross-validation and a

576       set of optimal model selection techniques. J. Chem. Inf. Model. 59, 4070-4076

577    13) Hong, H., Fang, H., Xie, Q., Perkins, R., Sheehan, D.M., Tong, W., 2003.

578       Comparative molecular field analysis (CoMFA) model using a large diverse set of

579       natural, synthetic and environmental chemicals for binding to the androgen receptor.

580       SAR QSAR  Environ. Res. 14(5-6), 373-388

581    14) Serafimova, R., Walker, J., Mekenyan, O., 2002. Androgen receptor binding affinity

582       of pesticide "active" formulation ingredients. QSAR evaluation by COREPA method.

583       SAR QSAR in Environ. Res 13(1), 127-134

584    15) Piir, G., Sild, S., Maran, U., 2021. Binary and multi-class classification for androgen

585       receptor agonists, antagonists and binders. Chemosphere 262,  128313

586    16) Mauri, A., alvaDesc: A tool to calculate and analyze molecular descriptors and

587       fingerprints. In: Roy K. (Ed.), Ecotoxicological QSARs. Methods in Pharmacology

588       and Toxicology, Humana, New York, NY, 2020, pp. 801-820.

589    17) Martin, T.M., Harten, P., Young, D.M., Muratov, E.N., Golbraikh, A., Zhu, H.,

590       Tropsha, A., 2012. Does rational selection of training and test sets improve the

591       outcome of qsar modeling?. J.  Chem. Inf.  Model. 52, 2570-2578

592    18) Leardi, R., 2000. Application of genetic algorithm-PLS for feature selection in
593           spectral data sets. J. Chemom. 14, 643-655

594    19) Roy, K., Mitra, I., 2011. On various metrics used for validation of predictive qsar
595           models with applications in virtual screening and focused library design. Comb.
596           Chem. High Throughput Screen. 14, 450-474

597    20) Roy, K., Das, R.N., Ambure, P., Aher, R.B., 2016. Be aware of error measures.
598           Further studies on validation of predictive QSAR models. Chemom. Int. Lab. Sys.
599           152, 18-33

600    21) Roy, K., Kar, S., Das, R..N., Understanding The Basics Of QSAR For Applications In
601           Pharmaceutical Sciences And Risk Assessment, Elsevier Inc, NY, 2015

602    22) Discovery Studio Predictive Science Application | Dassault Systèmes BIOVIA.
603           https://www.3dsbiovia.com/products/collaborative-science/biovia-discovery-studio/

604    23) Kumar, V., Ojha, P.K., Saha, A.,Roy, K.,. 2020. Exploring 2D-QSAR for prediction
605           of beta-secretase 1 (BACE1) inhibitory activity against Alzheimer's disease, SAR
606           QSAR Environ. Res. 31(2), 87-133

607    24) Nirschl, A. A., Zou, Y., Krystek Jr, S. R., Sutton, J. C., Simpkins, L. M., Lupisella, J.
608           A., & Hamann, L. G., 2009. N-Aryl-oxazolidin-2-imine muscle selective androgen
609           receptor modulators enhance potency through pharmacophore reorientation. J. Med.
610           Chem. 52(9), 2794-2798.

611    25) Momany F.A., Rone R., 1992. Validation of the general purpose QUANTA®
612           3.2/CHARMm® force field. J. Comput. Chem. 13(7), 888-900.

613    26) Abraham, M. J., Murtola, T., Schulz, R., Páll, S., Smith, J. C., Hess, B., Lindahl, E.,
614           2015. GROMACS: High performance molecular simulations through multi-level
615           parallelism from laptops to supercomputers. SoftwareX. 1, 19-25.

616   27) Zoete, V., Cuendet, V., Grosdidier, A., Michielin, O., 2011. SwissParam: a fast force

617        field generation tool for small organic molecules. J. Comput. Chem. 32, 11, 2359-

618        2368.

619   28) Jorgensen, W. L., Chandrasekhar, J., Madura, J. D., Impey, R. M., Klein, M. L., 1983.

620        Comparison of simple potential functions for simulating liquid water. J. Chem. Phys.

621        79, 2, 926-935.

622   29) Chatterjee, S., Maity, A., Chowdhury, S., Islam, M. A., Muttinini, R. K., Sen, D.,

623        2021. In silico analysis and identification of promising hits against 2019 novel

624        coronavirus 3C-like main protease enzyme. J. Biomol. Struct. Dyn. 39(14), 5290-

625        5303.

626   30) Valdés-Tresanco, M. S., Valdés-Tresanco, M. E., Valiente, P. A., Moreno, E., 2021.

627        gmx_MMPBSA: a new tool to perform end-state free energy calculations with

628        GROMACS. J. Chem. Theory Comput. 17, 10, 6281-6291.

629   31) Butina, D., 2004. Performance of Kier-Hall E-state descriptors in quantitative

630        structure activity relationship (QSAR) studies of multifunctional molecules.

631        Molecules 9, 1004-1009

632   32) Wahl, J., Smieško, M.; 2018. Endocrine disruption at the androgen receptor:

633        employing molecular dynamics and docking for improved virtual screening and

634        toxicity prediction. Int. J. Mol. Sci. 19, 1784

635   33) Zhou, W., Duan, M., Fu, W., Pang, J., Tang, Q., Sun, H., Xu, L., Chang, S., Li, D.,

636        Hou, T., 2018. Discovery of novel androgen receptor ligands by structure-based

637        virtual screening and bioassays. Genom. Proteom. Bioinform. 16, 416-427

638   34) Roy, K., Ambure, P., Kar, S., 2018. How precise are our quantitative structure-

639        activity relationship derived predictions for new query chemicals? ACS Omega 3;

640        11392-11406