1
2  •  **Article title:** Lipophilicity in drug design: An overview of lipophilicity descriptors in 3D-QSAR studies
3
4  •  **Short running title:** Novel lipophilicity descriptors in 3D-QSAR.
5
6  •  **Author names:** Tiziana Ginex,[1,*] Javier Vazquez,[1,2] Enric Gilbert,[2] Enric Herrero,[2] & F. Javier Luque[1,*]
7
8  •  **Author affiliations**
9  1 Department of Nutrition, Food Sciences and Gastronomy, Faculty of Pharmacy and Food Sciences, Campus
10 Torribera, Institute of Biomedicine (IBUB), and Institute of Theoretical and Computational Chemistry (IQTC-UB),
11 University of Barcelona, Av. Prat de la Riba 171, Santa Coloma de Gramenet E-08921, Spain.
12 2 Pharmacelera, Plaça Pau Vila, 1, Sector 1, Edificio Palau de Mar, Barcelona 08039, Spain.
13
14 •  **Corresponding author details:** Tiziana Ginex (tiziana.ginex@ub.edu), F. Javier Luque (fjluque@ub.edu)
15
24 •  **Information pertaining to writing assistance:** N/A.
25
26 •  **Ethical disclosure:** N/A
27
28 •  **Author Contributions:** TG and JV performed computations. All authors contributed to the analysis and discussion
29    of results. TG and FJL wrote the manuscript.
30
31 **Word count:  6262**
32 **Figure number:  3**
33 **Table number:  2**
34

**Abstract**

The pharmacophore concept is a fundamental cornerstone in drug discovery, playing a critical role in determining the success of *in silico* techniques, such as virtual screening and 3D-QSAR studies. The reliability of these approaches is influenced by the quality of the physicochemical descriptors used to characterize the chemical entities. In this context, a pivotal role is exerted by lipophilicity, which is a major contribution to host-guest interaction and ligand binding affinity. Several approaches have been undertaken to account for the descriptive and predictive capabilities of lipophilicity in 3D-QSAR modelling. Recent efforts encode the use of quantum mechanical-based descriptors derived from continuum solvation models, which open novel avenues for gaining insight into structure-activity relationships studies.

70    *1. The pharmacophore concept and its application in drug design.*

71    Almost all processes of life are determined by the recognition between biomolecules, a process dictated by

72    the chemical complementarity between the interacting partners [1]. An effective characterization of the

73    chemical features associated to the structure of both "host" and "guest" is necessary for disclosing the key

74    molecular determinants implicated in the formation of the host-guest complex. In drug discovery studies

75    addressing the interaction of small molecules (ligands) with macromolecular receptors, these determinants

76    are generally encoded under the concept of pharmacophore. A simple and intuitive definition can be

77    attributed to Paul Ehrlich, since this concept can be related to "a molecular framework that carries (*phoros*)

78    the essential features responsible for a drug's (*pharmacon*) biological activity" [2]. Nevertheless, Ehrlich did

79    not use the term *pharmacophore* in his papers, where the terms *haptophore* and *toxophore* were adopted [3].

80    Instead, the modern concept of pharmacophore evolved from the identification of "chemical groups" to the

81    definition as "patterns of abstract features in space" by Schueler [4], reflected in early models depicting key

82    features for biological activity that must satisfy certain geometrical relationships [5, 6], and the development

83    of the first pharmacophore pattern recognition programs [7]. Thus, according to the International Union of

84    Pure and Applied Chemistry (IUPAC), a pharmacophore "does not represent a real molecule or a real

85    association of functional groups, but a purely abstract concept that accounts for the common molecular

86    interaction capacities of a group of compounds towards their target structure", being the largest common

87    denominator shared by a set of active molecules [8].

88    This evolution has been accompanied by the progressive refinements triggered by advances in molecular

89    descriptors and computational methods seen in the last 30 years, since a variety of *in silico* techniques have

90    exploited the pharmacophore concept. This is exemplified by virtual screening (VS) studies of large

91    molecular databases performed to identify new promising compounds according to their similarity to a given

92    privileged template, which should contain reference physicochemical features relevant for biological activity

93    [9-11]. Molecular/chemical (global/local) similarity is a subjective concept since it depends on the specific

94    details of the methodological approach, the nature of the molecular features relevant for similarity

95    assessment, and the definition of the similarity function [12]. A sensitive and effective estimation of

96    molecular similarity is a fundamental pre-requisite for the identification of potential leads starting from a

97    chemical reference, which represents the paradigm of virtual screening.

98    Another successful application of the pharmacophore concept is linked to 3D quantitative structure-activity

99    relationships (3DQSAR) [13], such as CoMFA [14], CoMSIA [15] and GRID/GOLPE [16]. These methods

100   permit to identify a pharmacophore from the relationships between the biological activities of a set of aligned

101   molecules and the projection of selected physicochemical descriptors into the surrounding space, leading to

102     the disclosure of regions favourable or not to the bioactivity of compounds. 3D-QSAR approaches are also

103     used to model ADME(T) properties in the attempt to predict whether a molecular candidate would be able to

104     achieve its biological target [17]. Optimization of both ligand potency and ADME(T) profile is absolutely

105     required to translate promising molecular candidates to successful low-dose therapeutics. However, the

106     success of this operation is not trivial, since the final result depends on factors such as the quality of the input

107     data, as well as the adequacy and level of description of the physicochemical parameters used in the analysis.

108     In fact, Gleeson and collaborators [18] have observed the existence of a diametrically opposed relationship

109     between descriptors that efficaciously model drug potency and ADME(T) properties, making more

110     challenging the drug discovery process.

111

112     *2. Lipophilicity in drug design*

113     The relevance of lipophilicity in understanding the pharmacological profile of drug-like compounds is

114     widely recognized [19], as a broad variety of biodistribution and toxicological processes are ultimately

115     related to the differential solubility of solutes in aqueous and non-aqueous environments. This is illustrated

116     by Lipinski's rule-of-five [20], which relates the drug-likeness of oral compounds with molecular weight,

117     hydrogen bonding, and lipophilicity. Being a key property for the prediction of ADME(T) properties, this

118     has stimulated the development of experimental and computational approaches to quantify the lipophilicity

119     of a (bio)organic molecule.

120     Experimentally, the lipophilicity of a molecule can be quantified by its partition coefficient (*P*), as this

121     equilibrium thermodynamic property measures the ratio of concentrations of the compound between two

122     immiscible solvents, generally water and *n*-octanol. In turn, the partition coefficient can be expressed in

123     terms of the transfer free energy ($\Delta G_{tr}^{o/w}$) between the two solvents (Eq 1).

124

125     $$\Delta G_{tr}^{o/w} = -2.303\ RT\ logP \hspace{8cm} \text{Eq 1}$$

126     Lipophilicity reflects the complex interplay between the intermolecular forces that dictate the differential

127     solvation in the aqueous and organic phases. Accordingly, it can be factorized in terms of selected physico-

128     chemical properties of the compound that may be relevant for the preferential solvation in aqueous and non-

129     aqueous solvents, as shown in Eq 2 [21, and references therein].

130

131     $$logP = \upsilon V - \Lambda + I + IE \hspace{8cm} \text{Eq 2}$$

132     where *v* is a constant, *V* is the molar volume, which encompass the ability of the solute to elicit nonpolar

133     interactions, $\Lambda$ is related to the polarity of the compound, and finally *I* and *IE* accounts for the solute capacity

134 to form ionic interactions, which favor partitioning into the aqueous phase, and for the contribution due to
135 intramolecular effects, respectively.

136 Let us note that lipophilicity and hydrophobicity, which are often used as equivalent concepts, are not strictly

137 synonymous, the latter being in fact one of the contributions to molecular lipophilicity [22]. Thus, while

138 hydrophobicity can be defined as the tendency of non-polar groups of a molecule to aggregate in order to

139 minimize the unfavourable exposition to the surrounding polar (water) solvent, lipophilicity is a measure of

140 the affinity of the molecule for the non-polar solvent in a biphasic system constituted by a polar and a non-

141 polar solvent.

142 Lipophilicity affects a number of pharmacokinetic parameters (Figure 1). Low lipophilicity is responsible of

143 high aqueous solubility, which is a key factor for drug-likeness, but an excessively low lipophilicity could

144 compromise the ability of the drug to achieve the biological target. On the opposite site, highly soluble

145 compounds possess poor permeability through biological membranes, limiting absorption along the

146 gastrointestinal tract, or the transport across the blood-brain barrier. Therefore, optimal requirements for

147 efficient solubility and permeability properties are inevitably enclosed in a very narrow range of

148 lipophilicity. Another key aspect for drug-likeness is bioavailability, which is inversely correlated to low

149 first-pass clearance. Once again, lipophilicity is crucial since high lipophilicity is associated to high

150 clearance and low metabolic stability. Overall, a careful handling of lipophilicity is required to optimize

151 compound availability at the biological target.

152

153 Figure 1 here

154

155

156 On the other hand, lipophilicity has rarely been used as the primary descriptor in ligand-receptor recognition.

157 Indeed, following the IUPAC recommendation for the definition of a pharmacophore, it is defined as "the

158 ensemble of steric and electronic features that is necessary to ensure the optimal supramolecular interactions

159 with a specific biological target structure" [8]. This definition hides the key role played by (de)solvation in

160 the recognition and binding of a drug-like compound to its macromolecular target [23], especially keeping in

161 mind that the maximal achievable affinity that can be attained for target binding sites is largely influenced by

162 nonpolar desolvation [24]. This is consistent with the concept that favourable drug binding is largely driven

163 not only by the global lipophilicity of a compound, but more importantly by the spatial distribution of polar

164 and apolar regions along the chemical skeleton. Thus, while apolar regions determine the binding affinity

165 with complementary lipophilic regions of the binding site, polar interactions would provide 'anchor points'

166    contributing to ligand specificity and/or directionality in the binding pocket, as well as to modulate binding

167    kinetics of the ligand [25-30].

168    Taken together, these data suggest that a concomitant optimization of both pharmacokinetic profile and drug

169    potency have to be done to obtain successful drug products. This is encoded in the concept of lipophilicity

170    efficiency (LipE), which provides a metric that normalizes the potency (generally measured as $K_i$ or $IC_{50}$) of

171    the ligand against a protein target for the lipophilicity of the compound [31-33]. This is achieved by

172    substracting the $logP$ (or the distribution coefficient for ionizable molecules, $logD$) from the negative

173    logarithm of the potency (Eq 3).

174

175    $LipE = -log(\text{potency}) - logP$                                      Eq 3

176

177    LipE can be useful to provide guidelines to study the simultaneous effects exerted by structural changes on

178    potency and lipophilicity, which is central for drug design and lead optimization programmes, thus giving

179    support to the formulation of the "lipophilic pharmacophore" concept.

180

181    *3. From empirical fragment/atom-based approaches to 3D structure-based methods to estimate lipophilicity*

182    Numerous efforts have been done to assess lipophilicity by means of experimental methods [34-36].

183    Similarly, a plethora of computational approaches for estimating $logP$ have also been developed [37-42]. We

184    limit ourselves to remark selected fundamental concepts, while the reader is addressed to the previously

185    quoted reviews for detailed comparative analysis.

186    Within the framework of *substructure-based* methods for $logP$ estimation, fragmental and atom-based

187    techniques follow a general additive scheme as shown in Eq 4,

188

189    $logP = \sum_{i=1}^{n} a_i f_i + \sum_{j=1}^{m} b_j F_j$                                      Eq 4

190

191    where $logP$ is the sum of the weighted ($a_i$) contribution of each fragment/atom ($f_i$) and a correction factor

192    ($b_j F_j$).

193    Fragmental methods are illustrated by the work of Leo, Hansch and Elkins [43] as well as Nys and Rekker

194    [44]. The former relies on the concept of substituent constant, which encodes the lipophilicity contribution of

195    a chemical group or atom when it replaces an hydrogen atom in a reference compound, and the theoretical

196    estimation of $logPo/w$ follows an additivity scheme, named cLOGP. This method permits to extrapolate the

197    partition coefficients starting from a list of experimentally fitted fragmental contributions to lipophilicity. An

198  arbitrary set of interfragmental rules was then used to compile a database library of *fragment-weighted*

199  lipophilicity contributions,. On the other hand, Nys and Rekker [44] introduced the concept of hydrophobic

200  fragmental constant (*f*), which represents the lipophilicity contribution of a constituent part of a structure to

201  the total lipophilicity of a given compound. Fragments range from atoms to heterocyclic rings, so that

202  functional groups with direct contribution to resonance interactions were left intact, and are differentitated

203  upon linkage to aliphatic and aromatic structures. The differences between experimental *logP* and the

204  additive value estimated from the $\sum f$ approach was accounted for by correction rules, reflecting factors such

205  as the presence of vicinal electronegative centres in the chemical structure, aromatic condensation, cross-

206  conjugation or hydrogen-bonding [45].

207  An example of atom-based partitioning strategy was undertaken by Ghose and Crippen, who developed a

208  procedure that combines lipophilicity contributions at an atomic level leading to the ALOGP method. This

209  method encompassed a list of 120 atom types for carbon, hydrogen, oxygen, nitrogen, sulfur, and halogens

210  [46-48]. An alternative strategy is the XLOGP method [49], which is based on the summation of atomic

211  contributions derived from experimental lipophilicity data of 1831 organic molecules, and includes

212  correction factors for some intramolecular interactions.

213  In the last decades, the evolution of computer performances enabled the development of whole molecule-

214  based strategies to predict the lipophilicity by taking into account the three-dimensional structure of

215  compounds, and thus the effect of molecular conformation. Among all the available techniques, the

216  molecular lipophilicity potential (MLP) [51] offers an empirical quantitative 3D description of the

217  lipophilicity potential from all the molecular fragments on the surrounding space of a compound. The MLP

218  approach is then intended to model the lipophilic interactions between ligand and receptor as noted in Eq 5,

219

220  $MLP_k = \sum_{i=1}^{N} F_i \, f(d_{ik})$ 　　　　　　　　　　　　　　　　　　　　Eq 5

221

222  where $F_i$ is the lipophilic fragmental contribution and $f(d_{ik})$ is a distance function which depends on the

223  separation between a given fragment (*i*) and any point on the molecular surface or volume (*k*).

224  Molecular fields derived from the MLP potential have found a wide range of pharmaceutical applications,

225  including the prediction of skin permeation and distribution of new chemical entities [50], modeling of

226  peptides and proteins [52, 53], and structure-activity relationships studies [54].

227  The Hydrophobic INTeraction (HINT) method represents an alternative, promising strategy for the study of

228  lipophilicity in biomolecular interactions [55, 56]. This method exploits a scale of hydrophobic fragments

229  constants at the atomic level by means of an adaptation of the CLOGP method, which are then used to

230  evaluate a pairwise interaction energy term ($b_{ij}$) between atoms $i$ and $j$ in the interacting partners according

231  to Eq 6.

232

233  $$b_{ij} = a_i S_i a_j S_j T_{ij} R_{ij} + r_{ij} \qquad\qquad\qquad\qquad\qquad \text{Eq 6}$$

234

235  where $a_i$ and $S_i$ are respectively the hydrophobic constant and the accessible surface area of the atom $i$, $T_{ij}$ is a

236  logic function describing the character of interacting pairs (attraction or repulsion), and $R_{ij}$ and $r_{ij}$ denote

237  functions of the distance between atoms $i$ and $j$, the former following an exponential form and the latter a

238  Lennard-Jones implementation.

239  Eq. 5 encodes the formalism of the "natural" HINT force-field, which has been used to explore a variety of

240  applications in ligand-protein and protein-protein interactions [57-61].

241  Other approaches have relied on molecular properties derived from quantum mechanical treatments of

242  molecules. An early attempt is the work by Roger and Cammarata [62, 63], who related the logP of aromatic

243  compounds with the charge density of both π and σ electron frameworks and the induced polarization. In a

244  distinct approach, the BLOGP method relied on semiempirical AM1 calculations to derive geometrical and

245  quantum chemical descriptors for the prediction of logP [64, 65]. In a similar approach, Clark and coworkers

246  performed AM1 and PM3 calculations to derive a series of descriptors, including electrostatic potentials,

247  total dipole moments, mean polarizabilities, surfaces, volumes and charges, which were used in the

248  prediction of partition coefficients [66, 67].

249  These efforts can also be exemplified with the concept of heuristic molecular lipophilic potential (HMLP)

250  [68, 69]. In this approach, the lipophilic/hydrophilic features of a compound are determined from the

251  analysis of the electrostatic potential computed at the molecular surface. To this end, a dimensionless

252  distance-dependent screening function is used to compare the local electron density at the surface of a given

253  atom with the electrostatic potential generated on the rest of atoms. The screening function, which was

254  derived from statistical mechanical treatment of polar solvent molecules as dipoles, accounts for the

255  influence exerted by the atomic descriptors of the electrostatic potential from surrounding atoms. Ultimately,

256  such a comparison leads to the definition of an atomic lipophilicity index, which can adopt positive or

257  negative values, reflecting the lipophilic and hydrophilic nature, respectively, of such an atom.

258  Finally, a distinct approximation comes from the usage of solute-solvent correlation functions derived by

259  using the Reference Interaction Site Model (RISM) as descriptors for QSAR studies. By using a classical

260  statistical mechanics-based solvent model combined with machine learning, 1D solute-solvent correlation

261  functions were used to predict Caco-2 cell permeabilities [70]. As an extension of this approach, Güssregen

et al. proposed the Comparative Analysis of 3D-RISM Maps (CARMa) methodology [71]. In this computational strategy, the classical electrostatic and steric fields generally used in CoMFA are replaced by solute–solvent distribution functions determined from 3D-RISM computations, which are subsequently treated as descriptors to perform QSAR analysis. The method was validated using a set of serine protease inhibitors as a test system.

Even though CARMa uses a statistical mechanics solvent model, the electrostatic and steric effects implemented in CoMFA cannot be directly captured. This issue has been recently addressed by solving 3D-RISM equations for a solvent comprising CoMFA probes in aqueous solution, this extension being referred to as CARMa(electrolyte) [72]. The analysis performed for six protein–ligand systems reveals a small but consistent increase in prediction accuracy compared to CoMFA.

*4. Lipophilicity from QM continuum solvation methods.*

More elaborate methods for estimating the partition coefficients have been proposed in the framework of QM-based continuum solvation models [73, 74], which were developed with the aim of predicting the solvation free energy of solutes treating the solvent as a continuum polarizable medium. In spite of this rather crude approximation, these methods have proved to be a promising strategy that combines well established physical formalisms, a straightforward mathematical implementation, and a reduced computational cost, while predicting solvation free energies of (bio)organic compounds with chemical accuracy after a careful parameterization against experimental data [75-77]. Since a broad review of these formalisms and their applications exceeds the aims of this review, we limit ourselves to stress a selected set of recent studies addressing the potential impact of QM-based continuum methods in drug design.

*4.1 COSMO and COSMO-RS-based approaches*

In this context, the Continuum Solvation Model for Real Solvents (COSMO-RS) has been recently utilized to evaluate the similarity between molecules within the so-called COSMO*sim* method [78]. This method relies on the conductor-like screening model (COSMO) calculations to derive the so-called σ-profile of a given compound. The σ-profile collects the set of polarization charge densities generated on the surface patches of the molecule immersed in the solvent, which is treated as an ideal conductor. The one-dimensional histogram distribution of the σ values for the whole set of surface elements enclosed in the molecular surface gives rise to a characteristic signature of the solute, which can be used to measure a σ-profile-based similarity between compounds with application for the detection of bioisosteric fragments or

293    molecules. In order to enhance the computational efficiency, the σ-profile of a new compound can be

294    replaced with a composition of partial σ-profiles taken from similar fragments of precalculated molecules

295    stored in a database using COSMOfrag [79].

296    Since the σ-profile does not contain information about the spatial distribution of the polarization charge

297    density, COSMOsim3D has been recently proposed to alleviate this limitation [80]. To this end,

298    COSMOsim3D projects the surface charge density of each surface segment onto a regular 3D grid, so that

299    each point of the grid has an associated local σ-profile. In other words, instead of generating a single 1D σ-

300    profile for the entire molecule, COSMO*sim3D* creates a local 1D σ-profile at each position of a regular 3D

301    grid. This process leads to a four-dimensional histogram defined by the three Cartesian dimensions of the

302    grid point and the local σ-profile as the fourth dimension. If calculated for two molecules, this strategy can

303    be ultimately used to estimate their overall similarity. Furthermore, these local σ-profiles have been also

304    used to generate molecular interactions fields for 3D-QSAR studies [81].

305

306    *4.2 Fragmental lipophilicity model from the MST method: The Hyphar approach*

307    The Miertus-Scrocco-Tomasi (MST) solvation model has been used to develop 3D distribution patterns of

308    lipophilicity, which in turn have been exploited in predicting molecular overlays and 3D-QSAR studies [82-

309    83]. The MST model is a parametrized version of the polarizable continuum model developed by Tomasi

310    and coworkers [85, 86] at both semiempirical, Hartree-Fock and B3LYP levels [87-90] (for a review see

311    [91]). From the solvation free energies in water and *n*-octanol, one can derive the *n*-octanol/water partition

312    coefficient (Eq 1), which is a property of the whole molecule. Nevertheless, by decomposing the solvation

313    free energy into atomic contributions, one can obtain the 3D profile of lipophilicity from the corresponding

314    atomic contributions to the logP. For a molecule (M) containing $N$ atoms, this is achieved by decomposing

315    the logP (or the corresponding transfer free energy, $\Delta G_{tr,M}^{o/w}$) into electrostatic ($logP_{ele,i}$), cavitation ($logP_{cav,i}$)

316    and van der Waals ($logP_{vW,i}$) components, which can be derived from the polar ($\Delta G_{ele,i}^{o/w}$) and non-polar

317    ($\Delta G_{cav,i}^{o/w}$, $\Delta G_{vW,i}^{o/w}$) contributions to the solvation free energy (Eqs 7 and 8).

318

319    $\Delta G_{tr,M}^{o/w} = \sum_{i=1}^{N} \Delta G_{tr,i}^{o/w} = \sum_{i=1}^{N} \left( \Delta G_{ele,i}^{o/w} + \Delta G_{cav,i}^{o/w} + \Delta G_{vW,i}^{o/w} \right)$         Eq 7

320    $logP_M = \sum_{i=1}^{N} logP_i = \sum_{i=1}^{N} \left( logP_{ele,i} + logP_{cav,i} + logP_{vW,i} \right)$         Eq 8

321

322  Partitioning of the electrostatic term into atomic contributions can be made resorting to a perturbation

323  approximation of the coupling between the solute charge distribution and the solvent reaction field [92],

324  leading to Eq 9.

325

326  $$logP^{o/w}_{ele,i} = \frac{1}{2}\langle\Psi^0 \left| \sum_{\substack{k=1 \\ k\in i}}^{K} \frac{q_k^w}{|r_k^w - r|} - \sum_{\substack{l=1 \\ l\in i}}^{L} \frac{q_l^o}{|r_l^o - r|} \right| \Psi^0\rangle$$       Eq 9

327

328  where $\Psi^o$ is the solute wave function in the gas phase, and $K$ and $L$ stand for the total number of reaction

329  field charges in water ($q_k^w$) and $n$-octanol ($q_l^o$), located at positions $r_k^w$ and $r_l^o$.

330  The atomic decomposition of the cavitation and van der Waals terms takes advantage of the linear

331  dependence with the solvent-exposed surface of the atoms in the molecule (Eqs 10 and 11).

332

333  $$logP^{o/w}_{cav,i} = \sum_{i=1}^{N} \frac{S_i}{S_T} \Delta G^{o/w}_{P,i}$$       Eq 10

334  $$logP^{o/w}_{vW,i} = \sum_{i=1}^{N} S_i \Delta\xi^{o/w}$$       Eq 11

335

336  where $\Delta G^{o/w}_{P,i} = \Delta G^{w}_{P,i} - \Delta G^{o}_{P,i}$, $\Delta G_{P,i}$ being the cavitation free energy of atom $i$, $\Delta\xi^{o/w} = \xi^w - \xi^o$, with

337  $\xi_i$ being the atomic surface tension, and $S_i$ denotes the contribution of atom $i$ to the total molecular surface

338  ($S_T$).

339  In contrast to the COSMO-RS-based approaches, which rely on the concept of σ-profile (see above), the

340  MST-derived applications use the atomic contributions to the thermodynamic components of the differential

341  solvation free energy in water and $n$-octanol, which are encoded under the partition coefficient between these

342  two solvents. Accordingly, they take into account the effect of specific chemical features of the molecule,

343  such as the existence of specific tautomers or conformational species, or the formation of specific

344  intramolecular interactions (i.e., hydrogen bond), in the computation of the 3D distribution pattern of

345  molecular lipophilicity.

346  These patterns have been exploited to predict the chemical similarity between compounds [84]. By using the

347  MST-based hydrophobic descriptors $logP^{o/w}_{ele,i}$ and $logP^{o/w}_{cav,i}$, a computational procedure has been proposed to

348  identify the molecular overlay that maximizes the lipophilic similarity. To this end, molecular similarity was

349  achieved by comparing the hydrophobic fields generated by the molecules, which were pre-aligned

following multipole expansions of the atomic lipophilic contributions. On the other hand, simple descriptors of the hydrogen-bond (HB) donor/acceptor character of atoms were used to complement the information about the chemical nature of polar atoms in a molecule (briefly, the current implementation assigns an arbitrary value of +1 to hydrogen atoms in HB donors, and -1 to N and O atoms that may act as acceptors). This choice obeys to the fact that the polar nature of hydrophilic groups cannot distinguish the HB donor/acceptor character, as this information is not implicitly encoded by the $log P_{ele,i}^{o/w}$ term. Hydrophobic and HB properties are then projected into a 3D grid using the exponential function (Eq 12) implemented in CoMSiA [15], and then compared by means of the Tanimoto coefficient.

$$p_q = \sum_{i=1}^{N} w_i\, e^{-\alpha r_{iq}^2} \hspace{6cm} \text{Eq 12}$$

The method was implemented in PharmScreen software [83,93] and was successfully used to evaluate the molecular overlay for a collection of 121 molecular systems compiled by AstraZeneca, denoted as the AstraZeneca Overlays Validation Test Set [94]. This set contains molecular overlays experimentally characterized for 119 targets, which were grouped in four categories according to the expected difficulty in predicting the experimental overlay: easy, moderate, hard, and unfeasible. The results pointed out that correct overlays were predicted for 94% (easy), 79% (moderate), and 54% (hard) of the cases. Moreover, the overall performance obtained from classical electrostatic/steric descriptors and from Hyphar ones was fairly similar for easy and moderate subsets, but the accuracy obtained with Hyphar for the subset of hard cases exceeded the performance obtained with electrostatic/steric properties. Finally, it was found that the similar performance of Hyphar and electrostatic/steric descriptors does not imply that they lead to identical overlays. Rather, the analysis of the predicted poses revealed that the degree of identity in molecular overlays was reduced with the increase in the difficulty of the target. Overall, these findings point out that Hyphar descriptors may be a valuable alternative for molecule superposition and virtual screening of chemical libraries, especially for targets that may be challenging for predictive molecular similarity techniques.

On the other hand, the atom-centered MST-derived hydrophobic contributions have also been used as physicochemical descriptors to derive 3D-QSAR models using PharmQSAR [82]. MST/IEFPCM calculations were performed for 5 sets of compounds, including dopamine D2/D4 receptor antagonists, antifungal chromanones, glycogen synthase kinase-3 inhibitors, cruzain inhibitors, and thermolysin inhibitors. The compounds in these sets covered a wide range of variance in selected physicochemical properties (molecular weight, hydrogen-bond donor/acceptor, clogP, and number of rotatable bonds). The 3D-QSAR models obtained with the hydrophobic pharmacophore (HyPhar) were found to have a predictive

382    accuracy comparable to standard CoMFA and CoMSiA techniques. Moreover, Hyphar descriptors were also

383    valuable to discriminate the selectivity of compounds acting as inhibitors of thrombin, trypsin, and factor Xa

384    [83].

385    Overall, these findings support the usefulness of the MST-derived lipophilic descriptors as a valuable

386    alternative to electrostatic/steric properties to carry out virtual screening of chemical libraries for molecular

387    similarity, as well as to derive 3D lipophilic pharmacophores, thus providing valuable complementary

388    information to gain insight into the molecular determinants of bioactivity.

389

390    5. *A comparative analysis between Hyphar and electrostatic/steric properties*

391    The strength of Hyphar descriptors in 3D-QSAR studies may be attributed to two major features. First, the

392    concept of lipophilicity is very intuitive and widely accepted in medicinal chemistry. Second, the partitioning

393    of lipophilicity, which reflects a property of the whole molecule, into atomic or fragmental contributions

394    permits to obtain a graphical representation of the distribution pattern of polar and apolar regions adapted to

395    the 3D structure of a given compound. In turn, this paves the way to rationalize the recognition between a

396    small compound and its macromolecular target from the complementarity between hydrophilic and lipophilic

397    groups of the ligand and the polar and apolar nature of the side chains of residues that shape the binding

398    pocket. As an additional remark, let us note that resorting to Hyphar descriptors benefits from the accurate

399    description of the molecular charge distribution that can be attained by QM methods, which may take into

400    account the influence arising from the chemical features of the bioactive compound, such as the ionization

401    state, the preference for a tautomeric species, and the adoption of a given conformational state representative

402    of the binding mode of the ligand.

403    Given the novelty of MST-based atomic lipophilicity contributions, it is nevertheless necessary to explore

404    their suitability for 3D-QSAR studies. In this context, this section reports the results of a comparative

405    analysis performed to calibrate the performance of Hyphar descriptors through comparison with

406    electrostatic/steric ones. This analysis has been carried out using the comprehensive benchmark data set

407    compiled by Sutherland and coworkers [95], which comprises 113 angiotensin converting enzyme (ACE)

408    inhibitors, 111 acetylcholinesterase (AChE) inhibitors, 147 ligands for benzodiazepine receptors (BZR), 282

409    cyclooxygenase-2 (COX-2) inhibitors, 361 dihydrofolatereductase (DHFR) inhibitors, 66 glycogen

410    phosphorylase b (GPB) inhibitors, 74 thermolysin (THER) inhibitors, and 87 thrombine (THR) inhibitors.

411    Accordingly, the CoMFA/CoMSiA results reported in ref. 95 were compared with the 3D-QSAR models

412    obtained using Hyphar descriptors, which combine both "polar" ($logP_{ele,i}$) and "non-polar" ($logP_{cav,i}$)

413    hydrophobic contributions (see above). To this end, the atomic electrostatic and non-electrostatic

components of the lipophilicity were used to generate the molecular fields through projection into a grid that encloses the set of aligned compounds using a similarity index function (see [82] for further details). For the sake of comparison, the original molecular geometries and protonation states of compounds were kept in this study. All the details about models generation, grid dimensions and points, training/test sets, and related activity ranges for the eight sets compiled by Sutherland are reported in Supplementary Material (Tables S1-S3).

As a preliminary step, the effect of the QM method selected to derive the hydrophobic contributions on the performance of the 3D-QSAR Hyphar models was evaluated for a subset of four systems (D2 inhibitors, antifungal chromanones, GSK3-β and cruzain inhibitors) taken from our previous study [82]. To this end, Hyphar descriptors were derived from continuum computations performed with the MST version parametrized for the semiempirical RM1 method [96], and alternatively with the version parametrized at the B3LYP/6-31G(d) level [90]. Comparison of the statistical parameters obtained for the subset of training and test compounds defined for each molecular system is shown in Table 1.

**Table 1**. Statistical parameters of the 3D-QSAR HyPhar models obtained from MST/B3LYP and MST/RM1 calculations for the four sets of compounds.[a]

| System | Training set | | | | Test set | | Nc | Field (%) | |
|---|---|---|---|---|---|---|---|---|---|
| | $r^2$ | $q^2$ | $S$ | $S_{press}$ | $r^2$ | $S$ | | Elec | Non-elec |
| **D2** | | | | | | | | | |
| MST/B3LYP | 0.94 | 0.77 | 0.31 | 0.60 | 0.78 | 0.57 | 3 | 68.6 | 31.4 |
| MST/RM1 | 0.93 | 0.74 | 0.28 | 0.65 | 0.71 | 0.63 | 3 | 70.9 | 29.1 |
| **Chromanones** | | | | | | | | | |
| MST/B3LYP | 0.77 | 0.51 | 0.49 | 0.29 | 0.81 | 0.20 | 3 | 34.3 | 65.7 |
| MST/RM1 | 0.76 | 0.42 | 0.51 | 0.32 | 0.66 | 0.82 | 3 | 42.1 | 57.9 |
| **GSK3** | | | | | | | | | |
| MST/B3LYP | 0.91 | 0.80 | 0.12 | 0.19 | 0.79 | 0.21 | 3 | 54.5 | 45.5 |
| MST/RM1 | 0.91 | 0.82 | 0.30 | 0.18 | 0.79 | 0.21 | 5 | 64.7 | 35.3 |
| **Cruzain** | | | | | | | | | |
| MST/B3LYP | 0.81 | 0.50 | 0.31 | 0.51 | 0.69 | 0.47 | 2 | 53.0 | 47.0 |
| MST/RM1 | 0.91 | 0.65 | 0.31 | 0.44 | 0.70 | 0.46 | 3 | 58.4 | 41.6 |

[a] See [92] for a proper description of the molecular sets. Nc denotes the number of PLS components in the best 3D-QSAR model, and the terms Elec and Non-elec stand for the fraction (in percentage) of electrostatic ($logP_{ele,i}$) and non-electrostatic ($logP_{cav,i}$) hydrophobic contributions to the final model.

The results reveal that there is large resemblance in the overall performance of the 3D-QSAR models obtained from MST/RM1 and MST/B3LYP Hyphar descriptors for all data sets. This finding is remarkable, since 3D-QSAR models derived from the RM1 hydrophobic descriptors compare well with the performance

obtained at the B3LYP level, but at a much lower computational cost, making the usage of semiempirical

440 methods highly attractive for the study of large libraries of drug-like compounds. Accordingly, the

441 computationally less demanding RM1 method seems to be a promising choice for 3D-QSAR studies with

442 Hyphar parameters.

443 On the basis of these results, the benchmark data set reported by Sutherland and coworkers [95] was

444 examined using the MST/RM1 Hyphar descriptors. The 3D-QSAR Hyphar models were compared with the

445 CoMFA/CoMSIA results reported in [95], which were obtained by using electrostatic potential-fitted charges

446 at the MNDO level, but for the THER set, where Gasteiger-Marsili charges were used. For the sake of

447 comparison, an additional model, denoted CoMFA (RM1), which exploits RM1 electrostatic-potential fitted

448 partial charges in conjunction with an steric field obtained from the Lennard-Jones potential with a positively

449 charged C.3 atom probe, was also examined. This model, therefore, is intended to explore the efficiency of

450 RM1-based partial charges in defining electrostatic features of molecules at the atomic level.

451 Table 2 shows the statistical parameters of the 3D-QSAR models. In general, similar performances were

452 obtained for the different 3D-QSAR models determined for molecules in the training test included in a given

453 system, as noted in the large resemblance between the statistical values of the regression ($r^2$) and cross-

454 validation ($q^2$) models. The same trend can be observed for the test set compounds, although a small

455 improvement was found for CoMFA (RM1) and Hyphar models in GPB and THERM systems compared to

456 reference CoMFA/CoMSiA models. In addition, a higher level of accuracy was also achieved by the models

457 derived from RM1 calculations since the number of outliers in the test set was lower than in classical

458 CoMFA/CoMSIA (Supplementary Material, Table S4). On the other hand, both BZR and COX2 were

459 confirmed to be challenging systems for QSAR modelling, as already noted by Sutherland and coworkers

460 [95]. For instance, in case of COX2, part of the reason for the poor predictive behaviour may probably be

461 ascribed to the fact that training and test set cover different ranges of in the property space.

462

463 **Table 2.** Statistical parameters obtained for CoMFA and CoMSiA models reported in [95] with the results

464 determined by using COMFA (RM1) and Hyphar models in this study for the eight molecular systems

465 (ACE, AChE, BZR, COX2, DHFR, GPB, THERM and THR).[a]

| System | Training set | | | | Test set | | | Field (%) | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $r^2$ | $q^2$ | S | Spress | $r^2$ | S | Nc[c] | Ele | N-Ele | HB |
| **ACE** [b] | | | | | | | | | | |
| CoMFA | 0.80 | 0.68 | 1.04 | - | 0.49/0.55 | 1.54/1.47 | 3 | - | - | - |
| CoMSiA | 0.76 | 0.65 | 1.15 | - | 0.52/0.58 | 1.48/1.41 | 3 | - | - | - |
| CoMFA (RM1) | 0.82 | 0.67 | 0.42 | 1.37 | 0.54/0.61 | 1.45/1.32 | 3 | 29.4 | 70.6 | - |
| Hyphar | 0.75 | 0.64 | 0.51 | 1.43 | 0.42/0.62 | 1.62/1.35 | 2 | 28.8 | 53.5 | 17.7 |
| **AChE** | | | | | | | | | | |
| CoMFA | 0.88 | 0.52 | 0.41 | - | 0.47/0.56 | 0.95/0.87 | 5 | - | - | - |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| CoMSiA | 0.86 | 0.48 | 0.45 | - | 0.44/0.60 | 0.98/0.81 | 6 | - | - | - |
| CoMFA (RM1) | 0.90 | 0.54 | 0.32 | 0.85 | 0.35/0.52 | 1.07/0.86 | 6 | 20.0 | 80.0 | - |
| Hyphar | 0.76 | 0.45 | 0.50 | 0.92 | 0.65 | 0.78 | 4 | 64.1 | 18.7 | 17.2 |
| **BZR** | | | | | | | | | | |
| CoMFA | 0.61 | 0.32 | 0.41 | - | 0.00/0.18 | 0.97/0.81 | 3 | - | - | - |
| CoMSiA | 0.62 | 0.41 | 0.41 | - | 0.08/0.30 | 0.93/0.75 | 3 | - | - | - |
| CoMFA (RM1) | 0.60 | 0.36 | 0.64 | 0.53 | 0.21/0.21 | 0.81/0.80 | 3 | 30.5 | 69.5 | - |
| Hyphar | 0.67 | 0.37 | 0.58 | 0.54 | 0.00/0.02 | 0.91/0.86 | 6 | 48.8 | 16.7 | 34.5 |
| **COX2** | | | | | | | | | | |
| CoMFA | 0.70 | 0.49 | 0.56 | - | 0.29/0.37 | 1.24/1.09 | 5 | - | - | - |
| CoMSIA | 0.69 | 0.43 | 0.56 | - | 0.03/0.22 | 1.44/1.20 | 6 | - | - | - |
| CoMFA (RM1) | 0.74 | 0.51 | 0.52 | 0.72 | 0.19/0.34 | 1.20/1.07 | 5 | 28.6 | 71.4 | - |
| Hyphar | 0.60 | 0.52 | 0.63 | 0.71 | 0.26/0.40 | 1.15/0.99 | 3 | 85.4 | 4.3 | 10.3 |
| **DHFR** | | | | | | | | | | |
| CoMFA | 0.79 | 0.65 | 0.59 | - | 0.59/0.70 | 0.89/0.73 | 5 | - | - | - |
| CoMSiA | 0.76 | 0.63 | 0.62 | - | 0.52/0.63 | 0.96/0.81 | 5 | - | - | - |
| RM1 CoMFA | 0.81 | 0.67 | 0.44 | 0.73 | 0.42/0.55 | 1.04/0.91 | 4 | 17.7 | 82.3 | - |
| Hyphar | 0.72 | 0.63 | 0.53 | 0.78 | 0.53/0.56 | 0.94/0.89 | 5 | 36.2 | 38.8 | 25.0 |
| **GPB** | | | | | | | | | | |
| CoMFA | 0.84 | 0.42 | 0.43 | - | 0.42/0.37 | 0.94/0.70 | 4 | - | - | - |
| CoMSiA | 0.78 | 0.43 | 0.50 | - | 0.46/0.34 | 0.90/0.82 | 4 | - | - | - |
| CoMFA (RM1) | 0.88 | 0.43 | 0.36 | 0.85 | 0.51 | 0.89 | 4 | 24.4 | 75.6 | - |
| Hyphar | 0.83 | 0.54 | 0.42 | 0.75 | 0.71 | 0.68 | 3 | 52.0 | 2.7 | 45.3 |
| **THERM** [c] | | | | | | | | | | |
| CoMFA | 0.94 | 0.51 | 0.55 | 1.54 | 0.60 | 1.26 | 7 | - | - | - |
| CoMSiA | 0.85 | 0.54 | 0.73 | - | 0.36/0.46 | 1.87/1.60 | 6 | - | - | - |
| CoMFA (RM1) | 0.90 | 0.46 | 0.33 | 1.57 | 0.51/0.66 | 1.39/1.18 | 5 | 25.5 | 74.5 | - |
| Hyphar | 0.84 | 0.49 | 0.41 | 1.51 | 0.67 | 1.13 | 4 | 37.9 | 25.5 | 36.6 |
| **THR** [d] | | | | | | | | | | |
| CoMFA | 0.86 | 0.59 | 0.36 | - | 0.54/0.73 | 1.59/0.56 | 4 | - | - | - |
| CoMSiA | 0.88 | 0.62 | 0.34 | - | 0.55/0.62 | 0.76/0.66 | 5 | - | - | - |
| CoMFA (RM1) | 0.89 | 0.59 | 0.33 | 0.64 | 0.45/0.58 | 0.86/0.82 | 5 | 16.0 | 84.0 | - |
| Hyphar | 0.87 | 0.64 | 0.37 | 0.59 | 0.53/0.56 | 0.79/0.74 | 4 | 37.5 | 41.7 | 20.8 |

466
467 [a] For test sets compounds, statistical parameters ($r^2$ and $S$) with (left) and without (right) outliers (i.e.,
468 compounds with residuals higher than 2.5-fold the standard deviation) are indicated. The number of outliers
469 for each system is reported in Supplementary Material (Table S4).
470 [b] mol0088 (original file name mol_17) was excluded because it contains iodine atom.
471 [c] Partition between training and test sets made as indicated in [15].
472 [d] mol0088 (original file name 82) was excluded due to problems with the input geometry.
473

474 The predictive performance of the models was also examined by analyzing their capacity to discriminate

475 between active and inactive compounds. To this end, for each molecular system the compounds in the test

476 set were ranked according to their experimental potency: "active/positive" (P) and "inactive/negative" (N)

477 were categorized by applying a threshold value of 6.0 (in $pIC_{50}/pK_i$ units). Then, test set compounds with a

478 predicted $pIC_{50}/pK_i$ value larger than the threshold value were considered "actives/positives" (TP), whereas

479 compounds with a predicted $pIC_{50}/pK_i$ value lower than the threshold were considered "inactives/negatives"

480 (TN). For each molecular system, the number of P, N, TP and TN compounds, as well as false positives (FP)

481 and false negatives (FN) are compiled in Supplementary Material (Table S5). In turn, these values were used

482 to identify correctly negative (specificity or TNR; in green in Figure 2) and positive (sensitivity or TPR; in

483    blue in Figure 2) compounds, and to reduce the false negative rate ("fall-out" or FPR; in red in Figure 2) by

484    applying Eqs. 13-15.

485

486    $Specificity\ (TNR) = \dfrac{TN}{N} = \dfrac{TN}{(TN+FP)}$                    Eq. 13

487    $Sensitivity\ (TPR) = \dfrac{TP}{P} = \dfrac{TP}{(TP+FN)}$                    Eq. 14

488    $Fall-out\ (FPR) = \dfrac{FP}{N} = \dfrac{FP}{(FP+TN)} = 1 - TNR$            Eq. 15

489

490    Figure 2 here

491

492    These parameters, which can vary from 0 to 1, can be considered a measure of the predictive performance of

493    the model. According to this classification, a model can be considered good if it has high

494    specificity/sensitivity and low fall-out values. Nevertheless, this analysis requires a balanced partition of

495    active and inactive compounds in the set of compounds, a requirement that is not fulfilled in the case of BZR

496    and GPB systems, since only one inactive and one active compound are present in these two sets,

497    respectively. Accordingly, the results obtained for BZR and GPB should be excluded from the analysis. For

498    the rest of molecular systems, both CoMFA (RM1) and Hyphar models exhibit generally similar trends

499    (Figure 2). The Hyphar model has a slightly better performance in sensitivity/specificity and fall-out values

500    for AchE, THERM and THR systems, whereas the opposite trend in found for CoMFA (RM1) in ACE and

501    COX2.

502    Finally, the ability of CoMFA (RM1) and Hyphar models to rank the compounds according to their potency

503    was also examined (Figure 3). To this end, the Spearman ($Rs$) coefficient for the first (Q1; in green), second

504    (Q2; in blue) and third (Q3; in red) quartiles, which would encompass molecules with highest, medium and

505    low activity/affinity, were determined for the test set compounds in each system. Although there is a notable

506    resemblance in the general trends obtained for CoMFA (RM1) and Hyphar models, slightly better

507    performances (higher $Rs$ values) are observed for Hyphar models, especially for compounds of higher

508    activity/affinity (Q1/Q2), whereas the differences are less pronounced for compounds in Q3, probably due to

509    the larger noise associated to the biological activity low active compounds.

510

511    Figure 3 here

512

513  Overall, the results obtained for the benchmark systems reveal that the Hyphar descriptors yield 3D-QSAR

514  models with an overall performance that compares with the results obtained using standard

515  CoMFA/CoMSiA. Hyphar models also seem to be more effective in locating (high sensibility) and ranking

516  (high $Rs$) true positives, especially in regions of high and medium activity/affinity.

517

518  6. *Final consideration and perspectives*.

519  The concept of pharmacophore is essential to disclose the key features that dictate the interaction between

520  ligand and receptor. Hence, it represents an important tool to identify guidelines valuable in computer-aided

521  drug design, covering a variety of applications such as molecular similarity, virtual screening, ligand

522  optimization, scaffold hopping, as well as modeling of ADME(T) properties and target identification. The

523  descriptive and predictive power of pharmacophores depends on the quality and adequacy of molecular

524  properties used to disclose the hidden relationship between activity and chemical structure. In the last

525  decades several strategies were developed to derive descriptors capable of capturing the chemical features

526  relevant for drug design, including the application of descriptors derived from QM methods coupled to

527  continuum solvation models.

528  Although fundamental for the activity of drug-like compounds, inclusion of lipophilicity as a major

529  descriptor has revealed more elusive, possibly due to the complexity of the chemical processes encompassed

530  by this concept, or the difficulty to find a rigorous formalism to reduce it to atomic contributions since

531  lipophilicity reflects a property of the whole molecule. In this context, it is worth stressing the efforts in

532  deriving tools such as MLP [50] and HINT [55, 56], where the molecular lipophilicity was treated by means

533  of empirical atomic contributions, and hence enabling the analysis of the 3D distribution of polar/apolar

534  regions along the chemical scaffold to provide a novel interpretation to the molecular determinants

535  responsible of biological activity.

536  QM-based continuum solvation methods are a promising strategy for deriving 3D descriptors, such as

537  COSMO-RS-based σ-profiles [78-81] or MST-derived 3D lipophilicity patterns [82-84,97-99], which in turn

538  may be exploited in computer-aided drug design. The set of studies reported up to now for a variety of

539  benchmark datasets, covering both measurements of molecular similarity for aligned compound or the

540  derivation of 3D-QSAR models, are encouraging. In general, the statistical performance of these QM-based

541  descriptors compares well with the results obtained from classical approaches, generally combining

542  electrostatic and steric fields, as illustrated in the comparative analysis reported here for the sets of

543  compounds considered by Sutherland and coworkers [95]. At least in part, this may be due to the limitations

544  of electrostatic/steric descriptors for describing enthalpy and entropy contributions to the binding affinity. On

545 the other hand, QM-based approaches permit to account directly for the specific features of the bioactive

546 species of the ligand, including effects attributable to ionization, tautomerism, or the specific conformation,

547 which may be advantageous compared to generic descriptors derived from empirical contributions. These

548 computational approaches benefit from the usage of lipophilicity, a property widely used in drug design,

549 easy to interpret by medicinal chemists, and linked to a physicochemical property that can be measured

550 experimentally. Through partitioning of the molecular lipophilicity into atomic contributions, novel

551 fractional models that account for the 3D lipophilicity pattern of compounds can then be exploited in

552 computer-assisted drug design.

553 Overall, the analysis of structure-activity relationships in terms of the lipophilic/hydrophilic balance may

554 provide a useful signature to complement studies performed with electrostatic/steric properties. In this sense,

555 the QM MST-based hydrophobic descriptors are valuable in predicting molecular overlays and elucidating

556 molecular similarity patterns. The higher descriptive quality of these descriptors could thus offer interesting

557 clues in searching for novel bioactive compounds, especially for challenging targets.

558

559 *Executive summary*.

560 ▪ All biological and biochemical processes are driven by the general concept of host-guest

561 complementarity. Accordingly, an essential but effective description of the "guest" is required for a

562 successful prediction of "host" recognition.

563 ▪ The pharmacophore concept is a fundamental cornerstone in drug discovery, as it accounts for the

564 common interaction features of a group of compounds towards their target structure, playing a

565 critical role in determining the success of *in silico* techniques.

566 ▪ Optimized descriptors able to model both pharmacokinetics and pharmacodynamics properties in

567 drug design are not easily achievable, and the use of sub-optimal physicochemical parameters may

568 be a more effective strategy.

569 ▪ Besides the relevance in predicting ADME(T) properties, lipophilicity exerts a pivotal role in

570 accounting for the maximal achievable affinity that can be attained between ligand and receptor.

571 ▪ The usage of lipophilicity descriptors may offer novel opportunities to disclose the underlying

572 relationships between chemical features and biological activity. In this context, the availability of

573 refined version of QM-based continuum solvation models may be an effective strategy for deriving

574 novel descriptors well suited for drug design.

575    ▪    In 3D-QSAR studies, the MST-derived Hyphar descriptors have been shown to provide models for

576         structure-activity relationships with a predictive accuracy comparable to CoMFA/CoMSiA

577         techniques based on electrostatic/steric parameters.

578    ▪    The Hyphar descriptors are also a valuable alternative for molecule superposition and virtual

579         screening of chemical libraries, especially for targets that may be challenging for predictive

580         molecular similarity techniques.

581    ▪    The availability of "polar" and "non-polar" fractional descriptors obtained from MST-based

582         continuum solvation models may be valuable to explore the molecular determinants of bioactivity,

583         providing complementary interpretations to classical descriptors in the rational design of novel

584         compounds.

585

586    References

587

588    1    Gohlke H, Klebe, G. Approaches to the description and prediction of the binding affinity of small-

589         molecule ligands to macromolecular receptors. *Angew. Chem. Int. Ed.* 41, 2644-2676 (2002).

590    2    Khedkar SA, Malde AK, Coutinho EC, Srivastava S. Pharmacophore modeling in drug discovery

591         and development: An overview. *Med. Chem.* 3, 187-197 (2007).

592    3    Güner OF, Bowen JP. Setting the record straight: The origin of the pharmacophore concept. *J.*

593         *Chem. Inf. Model.* 54, 1269-1283 (2014).

594    4    Schueler FW. *Chemobiodynamics and Drug Design.* McGrawHill, New York, (1960).

595    5    Beckett AH, Harper NJ, Clitherow JW. The impact of stereoisomerism in muscarinic activity. *J.*

596         *Pharm. Pharmacol.* 15, 362-371 (1963).

597    6    Kier LB. Receptor mapping using molecular orbital theory. In: *Fundamental Concepts in Drug-*

598         *Receptor Interactions.* Academic Press, New York, 15-46 (1970).

599    7    Gund P, Wipke WT, Langridge R. Computer searching for molecular structure file for

600         pharmacophoric patterns. In: *Computers in Chemical Research and Education (Volume 3).* Hadzi D,

601         Zupan J (Eds.), Elsevier Scientific, Amsterdam, 5-33 (1973).

602    8    Wermuth CG, Ganellin CR, Lindberg P, Mitscher LA. Glossary of terms used in medicinal

603         chemistry (IUPAC recommendations 1998). *Pure Appl. Chem.* 70, 1129–1143 (1998).

604    9    Bender A, Glen RC. Molecular similarity: A key technique in molecular informatics. *Org. Biomol.*

605         *Chem.* 2, 3204-3218 (2004).

606    10   Wolber G, Seidel T, Bendix F, Langer T. Molecule-pharmacophore superpositioning and pattern

607         matching in computational drug design. *Drug Discov. Today* 13, 23-29 (2008).

608    11   Kaserer T, Beck KR, Akram M, Odermatt A, Schuster D. Pharmacophore models and

609         pharmacophore-based virtual screening: Concepts and applications exemplified on hydroxysteroid

610         dehydrogensases. *Molecules* 20, 22799-22832 (2015).

611    12   Maggiora G, Vogt M, Stumpfe D, Bajorath J. Molecular similarity in medicinal chemistry. *J. Med.*

612         *Chem.* 57, 3186-3204 (2013).

613    13   Verma J, Khedkar VM, Coutinho EC. 3D-QSAR in drug design - A review. *Curr. Top. Med.*

614         *Chem.* 10, 95-115 (2010).

615    14   Cramer RD, III, Patterson DE, Bunce JD. Comparative molecular field analysis (CoMFA). 1. Effect

616         of shape on binding of steroids to carrier proteins. *J. Am. Chem. Soc.* 110, 5959–5967 (1988).

617    15  Klebe G, Abraham U, Mietzner T. Molecular similarity indices in a comparative analysis (CoMSIA)
618          of drug molecules to correlate and predict their biological activity. *J. Med. Chem.* 37, 4130-4146
619          (1994).

620    16  Nilsson J, Wikström H, Smilde A, Glase S, Pugsley T, Cruciani G, Pastor M,Clementi S.
621          GRID/GOLPE 3D quantitative structure-activity relationship study on a set of benzamides and
622          naphthamides, with affinity for the dopamine D3 receptor subtype. *J Med Chem.* 40, 833-40 (1997).

623    17  Winiwarter S, Ridderström M, Ungell A-L, Andersson T, Zamora I. Use of molecular descriptors for
624          absorption, distribution, metabolism, and excretion predictions. In: *Comprehensive Medicinal*
625          *Chemistry II. (Volume 5)*. Testa B, van de Waterbeemd H (Eds.), Elsevier, Amsterdam, 531-554
626          (2006).

627    18  Gleeson MP, Hersey A, Montanari D, Overington J. Probing the links between in vitro potency,
628          ADMET and physicochemical parameters. *Nat. Rev. Drug Discov.* 10, 197 (2011).

629    19  Testa B, Carrupt PA, Gaillard P, Tsai RS. Intramolecular interactions encoded in lipophilicity: Their
630          nature and significance. In: *Lipophilicity in Drug Action and Toxicology*. Pliska V, Testa B, van de
631          Waterbeemd H (Eds.), VCH, Weinheim, 49–71 (1996).

632    20  *Drug bioavailability: Estimation of solubility, permeability, absorption and bioavilability*. van de
633          Waterbeemd H, Lennernäs H, Artursson P (Eds.), Wiley-VCH, Weinheim, (2003).

634    21  Caron G, Ermondi G, Scherrer RA. Lipophilicity, Polarity, and Hydrophobicity. In: *Comprehensive*
635          *Medicinal Chemistry II*. Taylor JB, Triggle DJ (Eds.), Elsevier Science, Oxford, 5, 425-452 (2007).

636    22  Van de Waterbeemd H, Carter RE, Grassy G, Kubinyi H, Martin YC, Tute MS, Willett P. Glossary
637          of terms used in computational drug design (IUPAC Recommendations 1997). *Pure Appl. Chem.* 69,
638          1137-1152 (1997).

639    23  Spyrakis F, Ahmed MH; Bayden AS, Cozzini P, Mozzarelli A, Kellog GE. The roles of water in the
640          protein matrix: A largely untapped resource for drug discovery. *J. Med. Chem. 60*, 6781–6827
641          (2017).

642    * This contribution provides an updated perspective on the roles of water molecules in protein structure,
643    function and dynamics, with a particular focus on the applications in drug discovery and design.

644    24  Cheng AC, Coleman RG, Smyth KT, Cao Q, Soulard P, Caffrey DR, Salzberg AC, Huang ES.
645          Structure-based maximal affinity model predicts small-molecule druggability. *Nat. Biotechnol. 25*,
646          71–75 (2007).

647    ** This study reports a model-based approach to predict druggable binding sites and estimate the
648    maximal affinity acievable by a small compound that relies on the hydrophobic desolvation, and the
649    nonpolar surface and curvatuve of the target binding site.

650    25  Davis AM, Teague SJ. Hydrogen bonding, hydrophobic interactions, and failure of the rigid receptor
651        hypothesis. *Angew. Chem. Int. Ed. 38*, 736–749 (1999).

652    26  Hajduk PJ, Huth JR, Fesik SW. Druggability indices for protein targets derived from NMR-based
653        screening data. *J. Med. Chem. 48*, 2518–2525 (2005).

654    27  Egner U, Hillig RC. A structural biology view of target drugability. *Expert Opin. Drug Discov. 3*,
655        391–401 (2008).

656    28  Schmidtke P; Barril X. Understanding and predicting druggability. A high-throughput method for
657        detection of drug binding sites. *J. Med. Chem. 53*, 5858–5867 (2010).

658    29  Schmidtke P, Luque FJ, Murray JB, Barril X. Shielded hydrogen bonds as structural determinants of
659        binding kinetics: Application in drug design. *J. Am. Chem. Soc. 133*, 18903–18910 (2011).

660    30  Tsopelas F, Giaginis C, Tsantili-Kakoulidou A. Lipophilicity and biomimetic properties to support
661        drug discovery. *Expert Opin. Drug Discov. 12*, 885–896 (2017).

662    31  Freeman-Cook KD, Hoffman, RL, Johnson TW. Lipophilic efficiency: The most important
663        efficiency metric in medicinal chemistry. *Future Med. Chem.* 5, 113-115 (2013).

664    32  Jopkins AL, Keserü GM, Leeson PD, Ress DC, Reynolds CH. The role of ligand efficiency metrics
665        in drug discovery. *Nat. Rev. Drug Discov.* 13, 105-121 (2014).

666    33  Johnson TW, Gallego RA, Edwards MP. Lipophilic efficiency as an important metric in drug design.
667        *J. Med. Chem.* 61, 6401-6420 (2018).

668    * An updated overview of the role of lipophilic efficiency as a metric with increasing impact in guiding
669    drug discovery.

670    34  Chen Z, Weber SG. A high-throughput method for lipophilicity measurement. *Anal. Chem.* 79,
671        1043-1049 (2007).

672    35  Giaginis C, Tsantili-Kakoulidou A. Alternative measures of lipophilicity: From octanol-water
673        partitioning to IAM retention. *J. Pharm. Sci.* 97, 2984-3004 (2008).

674    36  Andrés A, Rosés M, Ràfols C, Bosch E, Espinosa S, Segarra V, Huerta JM. Setup and validation of
675        shake-flask procedures for the determination of partition coefficients (logD) from low drug amounts.
676        *Eur. J. Phar. Sci.* 76, 181-191 (2015).

677    37  Mannhold R, Dross K. Calculation procedures for molecular lipophilicity: A comparative study.
678        *Quant-Struct.-Act. Relat.* 15, 403–409 (1996).

679     38    Ghose AK, Viswanadhan VN, Wendoloski JI. Prediction of hydrophobic (lipophilic) properties of

680           small organic molecules using fragmental methods: An analysis of ALOGP and CLOGP methods *J.*

681           *Phys*. *Chem*. *A*. 19, 172–178 (1998).

682     39    Mannhold R, van de Waterbeemd H. Substructure and whole molecule approaches for calculating

683           logP. *J*. *Comput*. *Aided Mol*. *Des*. 15, 337–354 (2001).

684     40    Caron G, Ermondi G, Scherrer RA. Lipophilicity, polarity and hydrophobicity. In: *Comprehensive*

685           *Medicinal Chemistry II*. *(Volume 5.18)*. Elsevier (Ed.), 425-452 (2006).

686     41    Chen HF. In silico logP prediction for a large data set with support vector machines, radial basis

687           neutral networks and multiple linear regression. Chem. Biol. Drug Des 74, 142-147 (2009).

688     42    Mannhold R, Poda GI, Ostermann C, Tetko IV. Calculation of molecular lipophilicity: state-of-the-

689           art and comparison of logP methods on more than 96,000 compounds. J. Pharm. Sci. 98, 861-893

690           (2009).

691     43    Leo A, Hansch C, Elkins D. Partition coefficients and their uses. *Chem*. *Rev*. 71, 525–616 (1971).

692     44    Nys GG, Rekker RF. The concept of hydrophobic fragmental constants (f values). II. Extension of

693           its applicability to the calculation of lipophilicities of aromatic and heteroaromatic structures. *Eur*. *J*.

694           *Med*. *Chem*. 9, 361–375 (1974).

695     45    Mannhold R, Rekker RF. The hydrophobic fragmental constant approach for calculating logP in

696           octanol/water and aliphatic hydrocarbon/water systems. *Perspect*. *Drug Discovery Des*. 18, 1–18

697           (2000).

698     46    Ghose AK, Crippen GM. Atomic physicochemical parameters for three-dimensional-structure-

699           directed quantitative structure-activity relationships. 2. Modeling dispersive and hydrophobic

700           interactions. *J*. *Chem*. *Inf*. *Comput*. *Sci*. 27, 21–35 (1987).

701     47    Viswanadhan VN, Ghose AK, Revankar GR, Robins RK. An estimation of the atomic contribution

702           to octanol-water partition coefficient and molar refractivity from fundamental atomic and structural

703           properties: Its uses in computer-aided drug design. *Math*. *Comput*. *Model*. 14, 505-510 (1990).

704     48    Wildman SA, Crippen GM. Prediction of physicochemical properties by atomic contributions. *J*.

705           *Chem*. *Inf*. *Comput*. *Sci*. 39, 868–873 (1999)

706     49    Wang R, Fu Y, Lai L. A new atom-additive method for calculating partition coefficients. *J*. *Chem*.

707           *Inf*. *Model*. 37, 615-621 (1997).

708     50    Gaillard P, Carrupt PA, Testa B, Boudon A. Molecular Lipophilicity Potential, a Tool in 3D QSAR:

709           Method and Applications. *J*. *Comput*. *Aided Mol*. *Des*. 8, 83–96 (1994).

710     51    Ottaviani G, Martel S, Carrut P-A, In silico and in vitro filters for the fast estimation of skin

711      permeation and distribution of new chemical entities. *J. Med. Chem.* 50, 742–748 (2007).

712    52   Laguerre M, Saux M, Dubost J, Carpy A, MLPP: A program for the calculation of molecular
713      lipophilicity potential in proteins. *Pharm. Pharmacol. Commun.* 3, 217–222 (1997).

714    53   Efremov RG, Chugunov AO, Pyrkov TV, Priestle JP, Arseniev AS, Jacoby E. Molecular
715      lipophilicity in protein modeling and drug design. *Curr. Med. Chem.* 14, 393–415 (2016).

716    54   Bitam S, Hamadache M, Hanini S. QSAR model for prediction of the therapeutic potency of N-
717      benzylpiperidine derivatives as AChE inhibitors. *SAR QSAR Environ. Res.*, 28, 471-489 (2017).

718    55   Kellogg GE, Semus SF, Abraham DJ. HINT: A new method of empirical hydrophobic field
719      calculation for CoMFA. *J. Comput-Aided Mol. Des.* 5, 454–552 (1991).

720    56   Kellogg GE, Abraham DJ. Hydrophobicity: Is Log P(o/w) more than the sum of its parts? *Eur. J.*
721      *Med. Chem.* 35, 651–661 (2000).

722    57   Fornabaio M, Spyrakis F, Mozzarelli A, Cozzini P, Abraham DJ, Kellogg GE. Simple, intuitive
723      calculations of free energy of binding for protein-ligand complexes. 3. The free energy contribution
724      of structural water molecules in HIV-1 protease complexes. *J. Med. Chem.* 47, 4507–4516 (2004).

725    58   Amadasi A, Spyrakis F, Cozzini P, Abraham DJ, Kellogg GE, Mozzarelli A. Mapping the energetics
726      of water-protein and water-ligand interactions with the ″natural″ HINT forcefield: Predictive tools
727      for characterizing the roles of water in biomolecules. *J. Mol. Biol.* 358, 289–309 (2006).

728    59   Marabotti A, Spyrakis F, Facchiano A, Cozzini P, Alberti S, Kellogg GE, Mozzarelli A. Energy-
729      based prediction of amino acid-nucleotide base recognition. J. *Comput. Chem.* 29, 1955–1969
730      (2008).

731    60   Amadasi A, Surface JA, Spyrakis F, Cozzini P, Mozzarelli A, Kellogg GE. Robust classification of
732      ″relevant″ water molecules in putative protein binding sites. *J. Med. Chem.* 51, 1063–1067 (2008).

733    61   Ahmed MH, Spyrakis F, Cozzini P, Tripathi PK, Mozzarelli A, Scarsdale JN, Safo MA, Kellogg
734      GE. Bound water at protein-protein interfaces: partners, roles and hydrophobic bubbles as a
735      conserved motif. *PLoS One* 6, e24712 (2011).

736    62   Rogers KS, Cammarata A. A molecular orbital description of the partitioning of aromatic
737      compounds between polar and non-polar phases. *Biochim. Biophys. Acta* 193, 22–29 (1969).

738    63   Rogers KS, Cammarata A. Superdelocalizability and charge density. A correlation with partition
739      coefficients. *J. Med. Chem.* 12(4), 692–693 (1969).

740    64   Bodor N, Gabanyi Z, Wong C. A new method for the estimation of partition coefficient. *J. Am.*
741      *Chem. Soc.* 111, 3783–3786 (1989).

742    65   Bodor N, Huang MJ. An extended version of a novel method for the estimation of partition

743        coefficients. *J. Pharm. Sci.* 81, 272–281 (1992).

744    66  Breindl A, Beck B, Clark T, Glen RC. Prediction of the n-octanol/water partition coefficient, logP,

745        using a combination of semiempirical MO-calculations and a neural network. J. Mol. Model. 3, 142–

746        155 (1997).

747    67  Beck B, Breindl A, Clark T. QM/NN QSPR models with error estimation: Vapor pressure and Log

748        P. *J. Chem. Inf. Comput. Sci.* 40, 1046–1051 (2000).

749    68  Du Q, Liu PJ, Mezey PG. Theoretical derivation of heuristic molecular lipophilicity potential: A

750        quantum chemical description for molecular solvation. *J. Chem. Inf. Model.* 45, 347–353 (2005).

751    69  Du Q-S, Li D-P, He W-Z, Chou K-C. Heuristic molecular lipophilicity (HMLP): Lipophilicity and

752        hydrophilicity of amino acid side chains. *J. Comput. Chem.* 27, 685–692 (2006).

753    70  Palmer DS, Mišin M, Fedorov MV, Llinas A. Fast and general method to predict the

754        physicochemical properties of druglike molecules using the integral equation theory of molecular

755        liquids. *Mol. Pharm.* 12(9), 3420–3432 (2015).

756    71  Güssregen S, Matter H, Hessler G, Lionta E, Heil J, Kast SM. Thermodynamic characterization of

757        hydration sites from integral equation-derived free energy densities: Application to protein binding

758        sites and ligand series. *J. Chem. Inf. Model.* 57, 1652–1666 (2017).

759    72  Ansari SM, Palmer DS. Comparative molecular field analysis using molecular integral equation

760        theory. *J. Chem. Inf. Model.* 58 (6), 1253–1265 (2018).

761    73  Orozco M, Luque FJ. Theoretical methods for the description of the solvent effect in biomolecular

762        systems. *Chem. Rev.* 100, 4187–4226 (2000).

763    74  Tomasi J, Mennucci B, Cammi R. Quantum mechanical continuum solvation models. *Chem. Rev.*

764        105, 2999–3094 (2005).

765    75  Cramer CJ, Truhlar DG. A universal approach to solvation modeling. *Acc. Chem. Res.* 41, 760–768

766        (2008).

767    76  Klamt A, Mennucci B, Tomasi J, Barone V, Curutchet C, Orozco M, Luque FJ. On the performance

768        of continuum solvation methods. A comment on "Universal approaches to solvation modeling". *Acc.*

769        *Chem. Res.* 42, 489–492 (2009).

770    77  Klamt A. The COSMO and COSMO-RS solvation models. *WIRES Comput. Mol. Sci.* 8, e1338

771        (2018).

772    78  Thormann M, Klamt A, Hornig M, Almstetter M. COSMO*sim*: Bioisosteric similarity based on

773        COSMO-RS $\sigma$-profiles. *J. Chem. Inf. Model.* 64, 1040–1053 (2006).

** This study reports the application of the σ-profiles derived from Conductor-like Screening Model for Realistic Solvation (COSMO-RS) calculations in drug similarity measurements.

79 Hornig M, Klamt A. COSMOfrag: A novel tool for high-throughput ADME property prediction and similarity screening based on quantum chemistry. *J. Chem. Inf. Model.* 45, 1169–1177 (2005).

80 Thormann M, Klamt A, Wichmann K. COSMO*sim3D*: 3D-Similarity and alignment based on COSMO polarization charge densities. *J. Chem. Inf. Model.* 52, 2149–2156 (2012).

81 Klamt A, Thormann M, Wichmann K, Tosco P. COSMO*sar3D*: Molecular field analysis based on local COSMO σ-profiles. *J. Chem. Inf. Model.* 52, 2157–2164 (2012).

** The contribution examines the usage of local σ-profiles in molecular field analysis, providing an interpretation about the features of the virtual free energy field generated from the target binding pocket.

82 Ginex T, Muñoz‑Muriedas J, Herrero E, Gibert E, Cozzini P, Luque FJ. Development and validation of hydrophobic molecular fields derived from the quantum mechanical IEF/PCM‑MST solvation models in 3D‑QSAR. *J. Comput. Chem.* 37, 1147–1162 (2016).

** This study examines the performance of QM-based MST lipophilic (Hyphar) descriptors for calculation of molecular fields in the derivation of structure-activity relationships models.

83 Ginex T, Muñoz-Muriedas J, Herrero E, Gibert E, Cozzini P, Luque FJ. Application of the quantum mechanical IEF/PCM-MST hydrophobic descriptors to selectivity in ligand binding. *J. Mol. Model.* 22, 136 (2016).

84 Vázquez J, Deplano A, Herrero A, Ginex T, Gilbert E, Rabal O, Oyarzabal J, Herrero E, Luque FJ. Development and validation of molecular overlays derived from 3D Hydrophobic Similarity with PharmScreen. *J. Chem. Inf. Model.* 58, 1596–1609 (2018).

** A comparative analysis of electrostatic/steric and QM-based lipophilici (Hyphar) descriptors for predicting molecular overlays from three-dimensional similarity measurements.

85 Miertus S, Scrocco E, Tomasi J. Electrostatic interaction of a solute with a continuum. A direct utilization of ab initio molecular potentials for the prevision of solvent effects. *Chem. Phys.* 55, 117–129 (1981).

86 Cancès E, Mennucci B, Tomasi J. A new integral equation formalism for the polarizable continuum

model: Theoretical background and applications to isotropic and anisotrpic dielectrics. *J. Chem. Phys.* 107, 3032 (1997).

87 Bachs M, Luque FJ, Orozco M. Optimization of solute cavities and van der Waals parameters in *ab initio* MST‑SCRF calculations of neutral molecules. *J. Comput. Chem.* 15, 446–454 (1994).

88 Luque FJ, Bachs M, Orozco M. An optimized AM1/MST method for the MST‑SCRF representation of solvated systems. *J. Comput. Chem.* 15, 847–857 (1994).

89 Curutchet C, Orozco M, Luque FJ. Solvation in octanol: parametrization of the continuum MST model. *J. Comput. Chem.* 22, 1180–1193 (2001).

90 Soteras I, Curutchet C, Bidon-Chanal A, Orozco M, Luque FJ. Extension of the MST model to the IEF formalism: HF and B3LYP parametrizations. *THEOCHEM-J Mol. Struct.* 727, 29–40 (2005).

91 Luque FJ, Curutchet C, Muñoz-Muriedas J, Bidon-Chanal A, Soteras I, Morreale A, Gelpí JL, Orozco M. Continuum solvation models: Dissecting the free energy of solvation. *Phys. Chem. Chem. Phys.* 5, 3827–3836 (2003).

92 Luque FJ, Bofill JM, Orozco M. New strategies to incorporate the solvent polarization in self-consistent reaction field and free-energy perturbation simulations. *J. Chem. Phys.* 103, 10183–10191 (1995).

93 PharmScreen - PharmQSAR, Pharmacelera, Barcelona, Spain.

94 Giangreco I, Cosgrove DA, Packer MJ. An extensive and diverse set of molecular overlays for the validation of pharmacophore programs. *J. Chem. Inf. Model.* 53, 852–866 (2013).

95 Sutherland JJ, O'Brien LA, Weaver DF. A comparison of methods for modeling quantitative structure–activity relationships. *J. Med. Chem.* 47, 5541–5554 (2004).

96 Forti F, Barril X, Luque FJ, Orozco M. Extension of the MST continuum solvation model to the RM1 semiempirical Hamiltonian. J. Comput. Chem. 29, 578–587 (2008).

97 Muñoz J, Barril X, Hernandez B, Orozco M, Luque FJ. Hydrophobic similarity between molecules: A MST-based hydrophobic similarity index. *J. Comput. Chem.* 23, 554–563 (2002).

98 Muñoz-Muriedas J, Perspicace S, Bech N, Guccione S, Orozco M, Luque FJ. Hydrophobic molecular similarity from MST fractional contributions to the octanol/water partition coefficient. *J. Comput. Aided Mol. Des.* 19, 401–419 (2005).

99 Muñoz-Muriedas J, Barril X, Lopez JM, Orozco M, Luque FJ. A hydrophobic similarity analysis of solvation effects on nucleic acid bases. *J. Mol. Model.* 13, 357–365 (2007).