

# 3D-QSAR predictions for $\alpha$ -cyclodextrin binding constants using quantum mechanically based descriptors

Lukas Linden, Kai-Uwe Goss, Satoshi Endo

<b>Citation</b>	Chemosphere, 169; 693-699
<b>Issue Date</b>	2017-02
<b>Type</b>	Journal Article
<b>Textversion</b>	author
<b>Rights</b>	© 2016 Elsevier Ltd. This manuscript version is made available under the CC-BY-NC-ND 4.0 License. <a href="https://creativecommons.org/licenses/by-nc-nd/4.0/">https://creativecommons.org/licenses/by-nc-nd/4.0/</a> . The article has been published in final form at <a href="https://doi.org/10.1016/j.chemosphere.2016.11.115">https://doi.org/10.1016/j.chemosphere.2016.11.115</a>
<b>DOI</b>	10.1016/j.chemosphere.2016.11.115

Self-Archiving by Author(s)  
Placed on: Osaka City University



## 17 **Abstract**

18 Binding of organic chemicals to  $\alpha$ -cyclodextrin ( $\alpha$ CD) is a typical example for host-guest  
19 complexation that is influenced by the 3D-structure of both the binding site (host) and the  
20 solute (guest). Prediction of the binding constant is challenging and requires a successful  
21 representation of the binding site-solute interactions in the 3D-space. In this study, we  
22 tested if a 3D quantitative structure activity relationship (3D-QSAR) model with quantum  
23 mechanically based local sigma profiles (LSPs) derived from the COSMOsar3D method is  
24 capable of predicting  $\alpha$ CD binding constants from the most recent literature and how the  
25 model performs in comparison to a standard comparative molecular field analysis and to a  
26 reference 2D-QSAR. The results showed that the new 3D-QSAR model was more predictive  
27 than both reference models (RMSE 0.45 vs 0.53/0.52,  $R^2$  0.70 vs 0.53/0.68). Furthermore,  
28 only the new model captured the differences in the binding constants between structural  
29 isomers of aliphatic alcohols and allowed an extrapolation of the prediction to another  
30 literature data set. The high performance of the 3D-QSAR model with LSPs tested in this  
31 study and its theoretical robustness suggest that this modeling approach should be  
32 applicable to other binding processes including protein binding.

## 33 **Keywords:**

34  $\alpha$ -Cyclodextrin (CD); Binding constant; Inclusion complex; Prediction

## 35 **1 Introduction**

36 Binding of organic chemicals to macromolecules is of high relevance in environmental  
37 science and related fields. For example, binding to macromolecular sorbents such as  
38 cyclodextrins (CDs) can be utilized for remediation of contaminated materials. Moreover,  
39 binding to proteins including binding proteins, enzymes, transporters, and receptors has  
40 strong impacts on toxicity of chemicals. Prediction of binding coefficients poses a major  
41 challenge, as the three-dimensional (3D) structure of both the solute and the binding site  
42 strongly influences the binding free energy, thus the binding constant (Herrmann, 2014).  
43 This is in contrast to the partition coefficients between liquids, for which the free energy is  
44 sufficiently well predicted by descriptors that characterize the interaction properties of the  
45 whole molecule without considering the molecular geometry (Karickhoff et al., 1991; Klamt,  
46 1995; Abraham et al., 2004; Endo and Goss, 2014).

47 3D quantitative structure activity relationships (3D-QSARs) attempt to establish a correlation  
48 between a macroscopic property (e.g., binding constant, receptor affinity) and 3D-structural  
49 features of the solute molecules. A widely used 3D-QSAR tool is comparative molecular field  
50 analysis (CoMFA) (Cramer et al., 1988). CoMFA uses 3D-discretized molecular field  
51 properties, called molecular interaction fields (MIFs), as descriptors for a statistical method  
52 (e.g., partial least square, PLS). Recently, Klamt et al. proposed the COSMOsar3D method  
53 (Klamt et al., 2012), which uses 3D-gridded COSMO surface polarization charge densities as a  
54 new set of MIFs. This extension of CoMFA emerges from the quantum mechanically-based  
55 COSMO-RS (conductor-like screening model for real solvent) method (Klamt, 1995; Klamt et  
56 al., 1998), which predicts the properties of a chemical by using the surface polarization  
57 charge densities (called sigma surface) of the molecule calculated quantum mechanically in a  
58 virtual conductor. For each molecule, the calculated sigma surface can be condensed into a

59 sigma profile, a histogram of all the 'partial' charges (or charge-patches) of the molecule.  
60 The sigma surface and the sigma profile of a chemical appear to accurately describe the  
61 abilities of the molecule to undergo intermolecular interactions including electrostatic,  
62 hydrogen-bond, and van der Waals interactions (Klamt, 2011). To extend this concept to 3D-  
63 QSARs, COSMOsar3D computes the sigma profiles at grid points within the 3D space to give  
64 the local sigma profiles (LSPs) (Thormann et al., 2012). The LSP is thus a histogram that  
65 contains information about the sigma surface of a specific part of the molecule. Considering  
66 the theoretical basis and the proven accuracy of COSMO-RS for partitioning between liquids,  
67 it is anticipated that the LSPs are ideal MIFs for 3D-QSAR modeling of the binding free energy  
68 that is strongly influenced by the molecular geometry of solutes. Nevertheless, the  
69 COSMOsar3D method has only been tested against standard sets of enzymatic inhibition  
70 activities by the developers and there has been no attempt to apply this method to  
71 equilibrium binding constants.

72 In this study, COSMOsar3D is used to model data sets of  $\alpha$ -cyclodextrin ( $\alpha$ CD) binding  
73 constants.  $\alpha$ CD is built of six 1-4-linked glucopyranose units that form a conic ring with a  
74 diameter of 5 Å. In water, all hydroxyl groups are positioned on the outside of the  $\alpha$ CD ring,  
75 resulting in a hydrophobic cavity inside (Cox et al., 1984), which enables  $\alpha$ CD to form host-  
76 guest complexes. Formation of such complexes (Connors, 1997) can improve the solubility of  
77 chemicals (Hedges, 1998), clean waste gas streams (Blach et al., 2008), remediate  
78 contaminated soils (Villaverde et al., 2005; Flaherty et al., 2013), and mask taste and odor  
79 compounds (Del Valle, 2004). Further, CDs can be used to enhance the bioavailability of  
80 organic pollutants (Liu et al., 2013), remove them from aqueous media (Sawicki and Mercier,  
81 2006), and extract dyes from sand (De Lisi et al., 2007). CDs are also considered a useful test  
82 material for investigating macromolecular binding because of their relatively simple and

83 well-studied structure as well as evidences of substantial molecular steric effects on the  
84 binding constants (Tabushi, 1982; Ishiwata and Kamiya, 1999; Schneider, 2009). In the  
85 common cyclodextrin family (i.e.,  $\alpha$ -,  $\beta$ -,  $\gamma$ -),  $\alpha$ CD may be the most suitable starting material  
86 for studying 3D-effects on binding, as it has the smallest cavity and thus the highest  
87 restriction for host-guest complexation.

88 The purpose of this study is to evaluate the LSP-based 3D-QSAR (i.e., COSMOsar3D) for  
89 predicting  $\alpha$ CD binding constants in comparison to a standard CoMFA model that uses steric  
90 and electrostatic fields as MIFs. In addition, these 3D-QSARs are compared to a well-  
91 established 2D-QSAR, namely the linear solvation energy relationship (LSER), which is a pp-  
92 LFER model using Abraham's descriptors (Abraham et al., 1994; Goss, 2005). Since the LSER  
93 does not explicitly include descriptors that describe molecular geometry, this comparison  
94 serves to evaluate whether taking into account the molecular 3D geometry improves the  
95 accuracy of predictions for  $\alpha$ CD binding constants.

## 96 **2 Methods**

### 97 **2.1 Data sets**

98 Two data sets of 1:1  $\alpha$ CD binding constants ( $K_{a1}$ ) [ $M^{-1}$ ] were considered in this study. The first  
99 has been measured in our laboratory under a consistent experimental condition, as reported  
100 previously (Linden et al., 2016). This data set, referred to as the "Linden data set", was used  
101 for the calibration and the first evaluation of the modeling approaches, because we consider  
102 these data of high quality and consistency. The second data set was from Suzuki (Suzuki,  
103 2001), who assembled literature data for  $\alpha$ CD binding constants. The Suzuki data set was  
104 used for an additional external validation of the modeling approaches.

105 The Linden data set (Linden et al., 2016) consists of 60 neutral aliphatic and aromatic  
106 chemicals (range of  $\log K_{a1}$ : 1.25–4.97, mean: 2.42, standard deviation (SD): 0.83). It contains  
107 several groups of isomers, e.g., 1-hexanol (i.e., end-substituted alcohol) and 3-hexanol (i.e.,  
108 middle-substituted alcohol) as well as homologous series of chemicals (e.g., alcohols,  
109 ketones, ether, chlorobenzenes). The Suzuki data set (Suzuki, 2001) includes 87 neutral  
110 aliphatic and aromatic chemicals (range of  $\log K_{a1}$ : -0.09–3.81, mean: 1.95, SD: 0.81). Ionic or  
111 partly ionic chemicals were not considered here to avoid uncertainty associated with the  
112 actual charge state of the bound molecule (i.e., ionic or neutral) and different descriptions of  
113 ionic molecules between MIFs. The chemicals and the respective  $\log K_{a1}$  values are listed in  
114 Table SI 1 and Table SI 2. Five alcohols, namely 1-butanol, 1-pentanol, 1-heptanol, 1-hexanol,  
115 and 1-octanol exist in both data sets. Their reported  $\log K_{a1}$  values are 0.28-0.51 log units  
116 higher in the Suzuki data set than in the Linden data set. The difference in  $\log K_{a1}$  might be,  
117 in part, caused by the different experimental temperatures (Suzuki data 25 °C, Linden data  
118 30 °C). Linden data were measured at 30 °C which was the lowest adjustable temperature in  
119 the experimental setting. This minor difference in temperature should be borne in mind  
120 when the results are evaluated (see below).

## 121 **2.2 Selection procedures for training and test sets**

122 For generation and evaluation of each model (i.e., 2D- and 3D-QSARs), the Linden data set  
123 was split into training and test sets. The training set was used for model calibration and  
124 selection, while the performance of the resulting model was validated with regard to the  
125 prediction of the test set. Prediction of data that were not part of the training set is essential  
126 as a control and should be considered the more important quality feature for 3D-QSARs  
127 (Gramatica, 2007).

128 For the general model evaluation, the training and test sets were generated with the  $\log K_{a1}$   
129 hierarchic bin system (Kauffman and Jurs, 2001) (procedure 1, see Fig. SI 3 for a scheme). In  
130 this system, the data set was sorted according to the  $\log K_{a1}$  values of the chemicals and  
131 then, from highest to lowest, four consecutive chemicals were placed in one bin. One  
132 chemical from each bin was selected randomly and placed in the test set. This classifies 25%  
133 chemicals of the data set to the test set. The rest of the chemicals formed the training set.  
134 The procedure was repeated five times, resulting in five random training sets and the  
135 corresponding test sets.

136 In order to evaluate varying steric effects within homologous series of chemicals and isomers,  
137 the following modified procedure was used to generate constructed test sets (procedure 2).  
138 As in the first procedure, the chemicals were sorted by  $\log K_{a1}$  and four chemicals in a row  
139 were grouped into one bin. Then, the numbers 1 to 4 were given randomly to the four  
140 chemicals of a bin. In the first run of chemical selection, the chemicals with the number 1  
141 embodied the test set, while the rest of the chemicals were used as the training set. In the  
142 second run, the chemicals with the number 2 were the test set, and so forth. This procedure  
143 resulted in four test and training set combinations. In comparison to procedure 1, the  
144 randomness of the selection is reduced, whereas each chemical is part of a test set once and  
145 the other three times it belonged to the training set.

### 146 **2.3 3D-QSARs**

147 The 3D-QSAR modeling followed the workflow shown in Fig. SI 1. Modeling generally takes  
148 the following steps: 3D-structure generation, alignment, MIFs generation, model calibration  
149 with PLS, and model evaluation using the test set. There are multiple options for each step,  
150 as explained below, and different combinations were tested in this work for comprehensive  
151 evaluation of the methods.



### 152 **2.3.1 3D structure generation**

153 The 3D structures of all chemicals were generated with Tinker or COSMOconfX13. Tinker  
154 (Marinescu and Bols, 2009) is a molecular modeling package implemented in Open3Dalign v.  
155 2.3 (O3A) (Tosco et al., 2011) and generates the structure-data files of the conformers for  
156 the O3A alignment. The quenched molecular dynamics conformational search of Tinker was  
157 performed with an implicit solvent calculation and a dielectric constant of 24, which is the  
158 dielectric constant of  $\beta$ CD (Yu et al., 2002), while for the rest of the parameters the default  
159 setting was chosen.

160 COSMOconfX13 is a tool box that uses Turbomole (Sijm et al., 2000) for the quantum  
161 mechanics calculations of COSMO files. The default COSMOconf procedure was modified so  
162 that it creates more conformers than usual (see SI). That is to say, the total number of  
163 possible conformers was increased, the energetic distance between conformers was  
164 reduced, and the clustering steps were loosened. These modifications were intended to  
165 account for the flexibility of the chemicals, which is more important for the  $\alpha$ CD binding than  
166 for bulk phase partitioning.

### 167 **2.3.2 Alignments**

168 The 3D structures of chemicals need to be aligned in the 3D space before performing  
169 statistical analysis. Ideally, the resulting position and orientation of a chemical in the 3D  
170 space corresponds to the optimal interaction possibility between the chemical and  $\alpha$ CD. In a  
171 target-based approach, the structure or a substructure of  $\alpha$ CD is used as the template to  
172 which all molecules are aligned. In a ligand-based approach, the template is generated with  
173 the help of chemicals that bind strongly to  $\alpha$ CD (i.e., with high  $\log K_{a1}$  values). For all  
174 approaches, up to ten conformers of each chemical were considered and the conformer with  
175 the highest alignment score and, if there are multiple conformers with the highest score,

176 then that with the lowest energy was chosen for the model. In this study, the following three  
177 alignment procedures were applied.

- 178 **1.** The O3A alignment maximizes the overlap of atoms of the template chemicals and of  
179 the remaining chemicals. This is a ligand-based method and a standard alignment for  
180 CoMFA approaches and was performed here by using O3A v. 2.3 (Tosco et al., 2011).  
181 The seven chemicals with the largest log  $K_{a1}$  values of the Linden data set, namely 1-  
182 dodecanol, 1-undecanol, 1-decanol, 1-nonanol, 2-undecanone, 2-decanone, and  
183 hexylbenzene were used as template chemicals. These chemicals were pre-aligned  
184 against each other and then each conformer of the remaining chemicals was aligned  
185 against the pre-aligned conformers of each template chemical. In the end, the  
186 position of the chemical/conformer with the highest score against any of the  
187 template chemicals was chosen.
- 188 **2.** The COSMOsim3D alignment (Thormann et al., 2012) maximizes the overlap between  
189 the sigma surfaces of the chemical and the template. Hereby, the template is an  
190 averaged sigma profile of the template chemicals. The template chemicals used were  
191 the same as in the previous alignment method.
- 192 **3.** The COSMOsim3D receptor alignment is a target-based approach that maximizes the  
193 overlap between the inverted sigma surface of  $\alpha$ CD (which is the sigma charge value  
194 of each surface patch multiplied with -1) and the sigma surface of the chemicals of  
195 the data set. The sigma surface of  $\alpha$ CD needs to be inverted because the alignment  
196 algorithm maximizes the overlap of like sigma charges in a ligand-based approach.  
197 The inversion therefore places the chemicals in a position where greatest interaction  
198 energies between both  $\alpha$ CD and the respective chemical occur, as the interaction  
199 energy is greatest when the difference between the sigma charges of two interacting

200 surface segments is maximal. This alignment already considers the steric restrictions  
201 of the  $\alpha$ CD cavity because the chemicals cannot be placed at the same position as the  
202  $\alpha$ CD. The input structure for the COSMOsim3D receptor alignment is the 3D structure  
203 of  $\alpha$ CD and the position of an exemplary ligand, the latter defines the starting  
204 position in the alignment procedure for all chemicals that need to be aligned. Two  
205 input structures were used in our approach to test the dependence of the  
206 COSMOsim3D receptor alignment on the input structure:

207 **(3a)** The 3D structure of  $\alpha$ CD and the position of a ligand (poly-*p*-phenylene rotaxane)  
208 were obtained from an X-ray measurement (Stanier et al., 2001) (three different  
209 views of the complex are shown in Fig. SI 4). The cosmo file of the  $\alpha$ CD structure was  
210 derived with a single point calculation using *COSMOconfX13*.

211 **(3b)** The 3D structure of  $\alpha$ CD and the position of a ligand (1-dodecanol) were  
212 estimated by a molecular dynamics simulation (MDsim), which was kindly provided  
213 by Sven Jakobtorweihen at Hamburg University of Technology. The complex with the  
214 smallest distance between the center of mass of  $\alpha$ CD and that of 1-dodecanol was  
215 chosen as the template for the alignment (Fig. SI 5). The cosmo file for the resulting  
216  $\alpha$ CD structure was derived with a single point calculation using *COSMOconfX13*.

### 217 2.3.3 MIFs

218 Two sets of MIFs were used as independent variables for the PLS regression analysis.

219 **1.** The van der Waals (vdW) and the electrostatic (ele) fields are the two standard  
220 CoMFA variables. Molecular mechanics calculations using the Merck force field  
221 (MMFF94) were performed with Open3DQSAR v. 2.3 (Tosco and Balle, 2011) to  
222 derive the vdW and ele fields. A  $sp^3$  carbon atom was used as the probe. A grid

223 spacing of 1 Å was used with a 5 Å gap, i.e., the minimal distance to the box, around  
224 the chemicals.

225 **2.** LSPs were derived from the cosmo files by COSMOsar3D (Klamt et al., 2012). For the  
226 3D-QSAR model used here the LSPs were split into several consecutive profiles, each  
227 covering a range of 0.006 e/Å<sup>2</sup>. Thus, MIFs 1, 2, ..., and 7 cover sigma values from -  
228 0.024 to -0.018 e/ Å<sup>2</sup>, -0.018 to -0.012 e/ Å<sup>2</sup>, ..., and, 0.012 to 0.018 e/ Å<sup>2</sup>, respectively  
229 (Fig. SI 2). In the end, the integral of each LSP serves as the value for the independent  
230 variable. A grid spacing of 2 Å was used in a box that leaves at least a 5 Å gap around  
231 the chemicals.

#### 232 **2.3.4 Statistical tool**

233 The independent variables, i.e., the MIFs, of the training set chemicals were correlated with  
234 the log  $K_{a1}$  values using PLS regression analysis. Prior to PLS regression analysis, the number  
235 of independent variables was reduced as following. An energy cutoff was set at  
236  $\pm 30$  kcal/mol (Kim, 1995), and variables that have a SD below a level of 0.1 among all  
237 training chemicals were excluded. The different MIFs were scaled before the PLS procedure  
238 using block unscaled weighting (Kastenholz et al., 2000). Moreover, fractional factorial  
239 design selection (Baroni et al., 1992; Baroni et al., 1993) was used to reduce the number of  
240 variables.

241 PLS analysis was performed to derive one to five PLS components. Thus, each run resulted in  
242 five different models that used one to five PLS components. Leave-two-out cross validation  
243 was performed with each model and then the model with the minimum of the root mean  
244 square error (RMSE) value was selected for further evaluation against the test set.

## 245 2.4 pp-LFER

246 The pp-LFER is among the most accurate and robust models to describe solute partitioning  
247 between liquids or liquid and gas phases, where molecular interactions are not sterically  
248 restricted. In a practical sense, a 3D-QSAR model may be considered meaningful only if it  
249 gives better predictions than the pp-LFER model, which is simple and quick as long as the  
250 solute descriptors are known. The pp-LFER used here appears,

$$251 \log K_{a1} = c + sS + aA + bB + vV + lL \quad (1)$$

252 where  $S$  is the polarizability/dipolarity parameter,  $A$  the solute H-bond acidity,  $B$  the solute  
253 H-bond basicity,  $V$  the McGowan characteristic volume ( $\text{cm}^3 \text{mol}^{-1}/100$ ) and  $L$  the logarithm  
254 of the hexadecane-air partitioning coefficient. In this work, the pp-LFER solute descriptors  
255 (capital letters in eq. 2) were obtained from the UFZ-LSER database (Endo et al., 2015) and  
256 the system parameters (lower case letters in eq. 2) were fitted with multiple linear  
257 regression analysis using the experimental data for  $\log K_{a1}$  of training chemicals.

## 258 3 Results & Discussion

259 Table 1 shows the statistical results for evaluation of the modeling approaches using the  
260 Linden data set. RMSE and  $R^2$  calculated with the test sets are considered more important  
261 evaluation criteria than  $q^2$ . Each value in the table represents the mean (+/- standard  
262 deviation) of five runs with five different training and test sets generated by test set  
263 selection procedure 1. In the following, the results of the pp-LFER approach are discussed  
264 first and then the results of the 3D-QSAR approach.

265 **Table 1. Comparison of the statistical results of the different modeling approaches for the**  
266 **prediction of  $\log K_{a1}$  of the Linden data set using test set selection procedure 1.**

Modeling	Method	Alignment	Field	$q^2 \pm \text{SD}$	RMSE $\pm$ SD	$R^2 \pm \text{SD}$
----------	--------	-----------	-------	---------------------	---------------	---------------------

approach							
M1	pp-LFER				0.52 ± 0.05	0.68 ± 0.07	
M2	3D-QSAR	O3A	LSP	0.63 ± 0.03	0.54 ± 0.08	0.56 ± 0.17	
M3	3D-QSAR	O3A	vdW ele	0.58 ± 0.08	0.53 ± 0.11	0.53 ± 0.11	
<b>M4</b>	<b>3D-QSAR</b>	<b>COSMOsim3D</b>	<b>LSP</b>	<b>0.83 ± 0.02</b>	<b>0.45 ± 0.06</b>	<b>0.70 ± 0.08</b>	
M5	3D-QSAR	COSMOsim3D	vdW ele	0.70 ± 0.01	0.56 ± 0.06	0.53 ± 0.12	
M6a	3D-QSAR	COSMOsim3D	LSP	0.66 ± 0.06	0.51 ± 0.06	0.61 ± 0.09	
		receptor X-ray					
M6b	3D-QSAR	COSMOsim3D	LSP	0.71 ± 0.04	0.49 ± 0.04	0.64 ± 0.07	
		receptor MDsim					
M7	3D-QSAR	COSMOsim3D	vdW ele	0.51 ± 0.08	0.55 ± 0.08	0.56 ± 0.13	
		receptor X-ray					

267 **O3A means open3DALIGN,  $q^2$  is the coefficient of determination for the leave-two-out**  
268 **cross validation using the training set, RMSE is the root mean square error of the test set in**  
269 **log units, and  $R^2$  is the coefficient of determination of the test set. LSP, vdW, and ele**  
270 **indicate the usage of local sigma profiles, van der Waals interaction field, and electrostatic**  
271 **interaction field as molecular interaction field, respectively, SD is standard deviation, and**  
272 **MDsim is molecular dynamics simulation.**

### 273 3.1 pp-LFER

274 First, the pp-LFER equation (eq. 2) was fitted to all experimental  $\alpha$ CD binding constants of  
275 the Linden data set (i.e., no test and training set selection) to have an idea to what extent  
276 the 2D model can describe the whole data set (Fig. SI 4). This fit resulted in the equation

$$277 \log K_{\alpha 1} = -0.32 (\pm 0.44) + 2.04 (\pm 0.63) S + 3.15 (\pm 0.63) A - 3.01 (\pm 0.50) B +$$

$$278 6.01 (\pm 0.88) V - 1.10 (\pm 0.21) L \quad (2)$$

279 The fit of the pp-LFER equation usually results in a standard deviation of 0.1 to 0.2 log units  
280 for homogeneous solvent-water partition systems, which are not influenced by steric effects,  
281 and a larger standard deviation for partitioning or binding to heterogeneous materials such  
282 as serum albumin and natural organic matter (Bronner and Goss, 2011; Endo and Goss,  
283 2011). The RMSE for the binding to  $\alpha$ CD (Fig. SI 4) is 0.48, being comparable to fits for other  
284 heterogeneous materials (Bronner and Goss, 2011).

285 The pp-LFER fits for training sets extracted from the Linden data set resulted in system  
286 parameters similar to those for the complete Linden data set (Table SI 3). The predictions for  
287 the corresponding test sets (Table 1, M1) were surprisingly accurate (RMSE =  $0.52 \pm 0.05$  and  
288  $R^2 = 0.68 \pm 0.07$ ). This result was unexpected because the experimental results do suggest  
289 strong steric effects, whereas the pp-LFER model does not capture such effects (Linden et al.,  
290 2016). A closer examination of the results revealed that systematic prediction errors do exist  
291 for binding constants, e.g.,  $\log K_{a1}$  values for end-substituted chemicals were systematically  
292 underestimated and those for middle-substituted chemicals were overestimated, which is an  
293 indication that the pp-LFER model is not able to cover the underlying steric effects. In  
294 addition, chemicals that are not expected to fit into the  $\alpha$ CD cavity due to the steric  
295 hindrance were over-predicted by the pp-LFER, e.g., the  $\log K_{a1}$  value of 1-chloronaphthalene  
296 is predicted as 2.13, while the experiment showed that it is  $< 1.3$  (Linden et al., 2016).

## 297 **3.2 3D-QSARs**

298 Seven 3D-QSAR model variants were constructed using different combinations of structure  
299 generation, alignment, and MIF methods and evaluated with the Linden data set, as  
300 explained in the method section (Fig. SI 1, Table 1). The results show the following trends: (i)  
301 RMSE and  $R^2$  of the 3D-QSAR model variants for test set predictions were 0.45–0.56 and  
302 0.53–0.70, respectively. While the best 3D-QSAR model (M4) performed slightly better than  
303 the pp-LFER, the statistics were similar on average. (ii) The models that used the LSPs (Klamt  
304 et al., 2012) as independent variables tended to result in better predictions than those using  
305 the vdW and ele MIFs for a given alignment (i.e., O3A, COSMOsim3D, or COSMOsim3d  
306 receptor). These outcomes suggest that LSPs are more suitable descriptors to describe the  
307 binding to  $\alpha$ CD than the tested CoMFA variables. This interpretation is in line with the claim

308 that LSPs are theoretically more relevant for linear regression models, like PLS, to describe  
309 the interaction energy (Klamt et al., 2012).

310 Of the 3D-QSARs tested, the model that uses the COSMOsim3D alignment with the LSP  
311 variables (M4, Table 1) was the best model variant (i.e., with the lowest RMSE). No  
312 improvement was observed for the use of the 3D-structure of  $\alpha$ CD as the template for the  
313 alignment (compare M6a and M6b to M4). Moreover, no difference was observed between  
314 the use of the two  $\alpha$ CD structures (M6a (X-Ray) vs. M6b (MDsim)) for the target-dependent  
315 alignment. The fact that no improvement was observed by the use of the target-dependent  
316 alignment suggests that the selected 7 template chemicals were sufficient for aligning the 60  
317 chemicals in the Linden set. This result, however, may not be general; alignments with a  
318 binding site structure are expected to be advantageous particularly if the data availability is  
319 limited. Note that, in principle, MDsim could directly calculate binding coefficients (Gebhardt  
320 and Hansen, 2016; Sancho et al., 2016) but such calculations would be time consuming for a  
321 larger number of chemicals, although these calculations are more and more automated and  
322 routinely performed.

323 The possibility of a chance correlation for the best modeling approach (M4) was evaluated  
324 by scrambling of the dependent  $\log K_{a1}$  values in two sorted bins (this means each chemical  
325 got a permuted  $\log K_{a1}$  value) (Tropsha et al., 2003; Rücker et al., 2007), which resulted in  
326 non-predictive models ( $R_{\text{training}}^2 = 0.40$ ,  $q_{\text{LTO}}^2 = -0.0030$ , the mean of 10 times evaluation) .

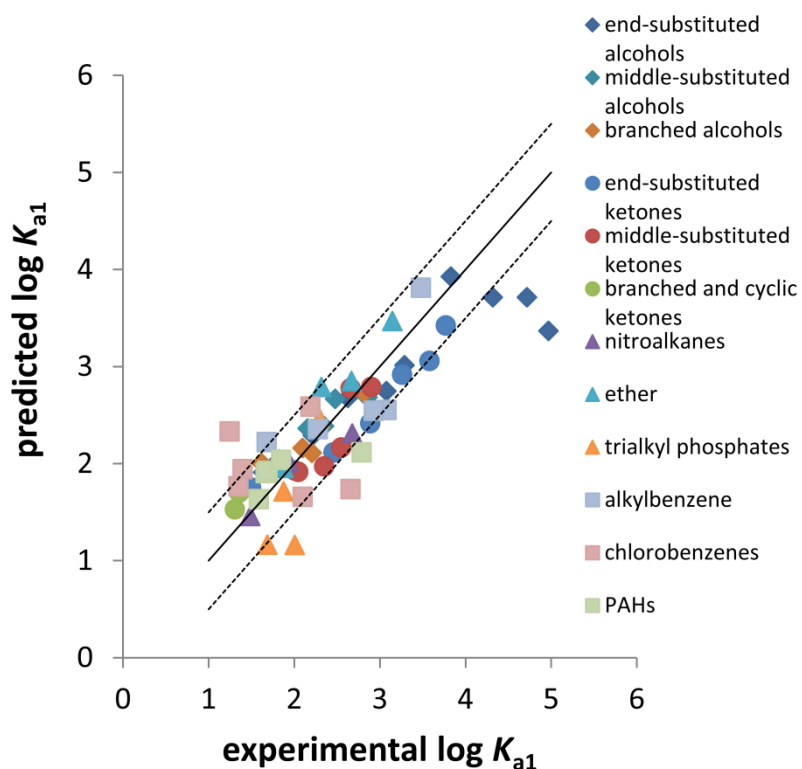
327 To infer binding mechanisms, the contributions of the MIFs (vdW and ele, or LSPs) to the PLS  
328 components are examined. The percentage contributions of the seven LSPs to the M4 PLS  
329 model are shown in Fig. SI 6. MIF 4 (-0.012 to 0 e/Å<sup>2</sup>, Fig. SI 2) had the highest contribution  
330 to the PLS components. This is an indication for the importance of vdW interactions and the



331 hydrophobic effect for the binding to  $\alpha$ CD (Marques, 2010). The contribution of MIF 4  
332 decreases slightly with increasing PLS component number, whereas the contributions of the  
333 other MIFs rather increased with increasing PLS component number. The PLS component 1  
334 in this example already explained 70% of the variance in the  $\log K_{a1}$  data, while the other  
335 four PLS components added up to an explained variance of 27%, i.e., the PLS components 2-  
336 5 serve for fine tuning of the model. The field contributions of model variants that used vdW  
337 and ele variables support the mechanistic interpretation obtained from the LSPs; the  
338 contribution of the vdW field is around 90% for the models.

### 339 **3.2.1 Predictions of specific molecular steric effects**

340 To evaluate the performance of the 3D-QSAR modeling approaches for predicting particular  
341 types of chemicals, four training and test sets were generated from the Linden data set  
342 according to test set selection procedure 2 (see the method section) and all prediction  
343 procedures were redone. Model approaches M3, M4, M5, and M6b were evaluated here  
344 because they performed best in the random evaluation above and allow comparison of the  
345 classical CoMFA approach and the new COSMO-based approach. The resulting statistics (i.e.,  
346  $q^2$ , RMSE,  $R^2$ ) were similar to those obtained above with test set selection procedure 1  
347 (Table 1), except for M3, for which the test set selection procedure 2 resulted in worse  
348 predictions (see Table SI 5). Fig. 1 compares the experimental data and the predictions by  
349 the best model variant (M4, with COSMOsim3D + LSPs) for individual chemicals.

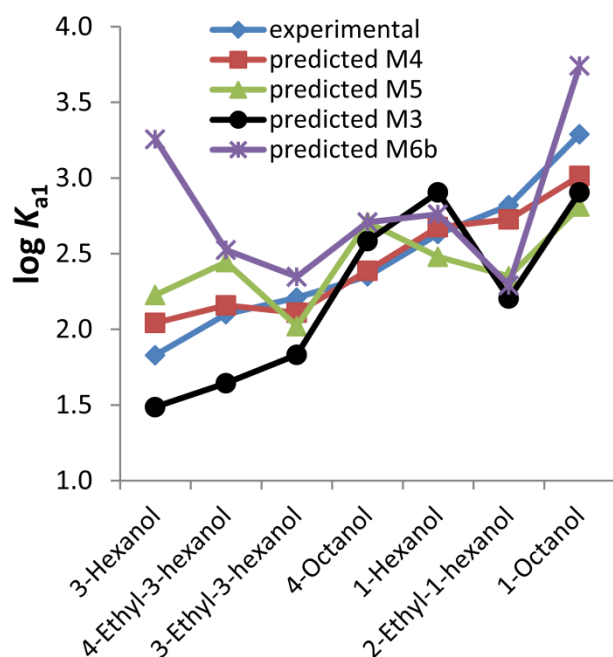


350

351 **Figure 1. Prediction of  $\log K_{a1}$  of 60 Linden's chemicals with COSMOsim3D alignment and**  
 352 **local sigma profiles as variables (M4). Test sets were selected with test set selection**  
 353 **procedure 2. The solid line indicates the 1:1 line and the dashed lines indicate a deviation**  
 354 **of 0.5 log units from the 1:1 line.**

355 Many trends of the data that are related to steric effects were quantitatively described in  
 356 the best 3D-QSAR model variant we found (M4). For example: experimental data show  
 357 relatively large differences in  $\log K_{a1}$  between isomeric chemicals with the functional group  
 358 at the terminal and the middle positions such as 1-heptanol and 4-heptanol. These chemicals  
 359 are predicted successfully by M4, e.g., 1-heptanol ( $\log K_{a1}$  exper. 3.08, pred. 2.75) and 4-  
 360 heptanol ( $\log K_{a1}$  exper. 2.16, pred. 2.36). Also, as is the case in the experimental data,  
 361 elongation of the alkyl chain in only one direction resulted in a higher increase of  $\log K_{a1}$  than  
 362 elongation in two or more directions (Fig. 2). The 3D-QSAR model variants M3, M5, and M6b  
 363 were not able to describe the differences between these alcohols so well as M4 (Fig. 2). The  
 364 comparison between M4 and M5 shows that the use of LSPs instead of vdW and ele not only  
 365 minimizes the overall prediction errors but helps distinguish structural isomers of alcohols.

366 The standard CoMFA model (M3) underestimates most of these alcohols and is not able to  
367 capture the steric effects. M6b uses LSPs as variables, but it appears that the target-based  
368 alignment cannot as accurately reproduce the trend of alcohol data as the ligand-based  
369 alignment in this case.



370

371 **Figure 2. Experimental and predicted  $\log K_{a1}$  for  $\alpha$ CD binding of two C6-alcohols and five**  
372 **C8-alcohols.**

373 Experimental data for chlorobenzenes showed a distinct substitution effect on the  $\alpha$ CD  
374 binding constant.  $K_{a1}$  increases with chlorine substitution up to two chlorine atoms, whereas  
375 a further substitution decreases  $K_{a1}$ , which can be explained by the size limitation of the  
376 cavity. This effect is not well described by any 3D-QSAR model tested here. For example,  
377 1,2,4,5-tetrachlorobenzene and 1,3-dichlorobenzene showed a prediction error larger than  
378 0.6 log units with the best model variant, M4. The use of the  $\alpha$ CD target structure  
379 (COSMOsim3D receptor alignment, M6b), the CoMFA variables vdW and ele (M5), and the  
380 standard CoMFA model (M3) did not improve the prediction of chlorobenzenes. A reason for  
381 the inaccurate predictions for chlorobenzenes could be the small number of data that

382 showed strong effects of steric restrictions. As shown in the previous work (Linden et al.,  
383 2016),  $K_{a1}$  for chemicals that undergo strong steric restrictions tend to have  $K_{a1}$  values that  
384 are too low to measure and thus such chemicals cannot be included in the data set for  
385 model calibration.

386 The end-substituted chemical 1-dodecanol was the biggest outlier in all predictions. A reason  
387 could be that 1-dodecanol has the longest alkyl chain and the largest  $K_{a1}$  in the data set.  
388 Therefore, the positive interaction between the long alkyl chain and  $\alpha$ CD may not be  
389 covered by the models. Additionally, the 3D-QSAR models in this work only consider one  
390 selected conformer of each chemical, which neglects the influence of different binding  
391 modes for predictions of flexible molecules like 1-dodecanol. Furthermore, a recent MDsim  
392 study showed that 1-dodecanol interacts substantially with the water surrounding  $\alpha$ CD and  
393 that the explicit consideration of the water molecules is necessary for a successful prediction  
394 of long chain alcohols (Gebhardt and Hansen, 2016). Note that, while the data we  
395 considered are for 1:1 binding constants, 2:1 binding can become more important for  
396 chemicals with long alkyl chain(s).

### 397 **3.3 Predictions of the Suzuki data set**

398 For a further evaluation of each modeling approach, models were generated using all Linden  
399 data as the training set and evaluated with the Suzuki data as an external test set. The  
400 prediction of the Suzuki data by the pp-LFER calibrated with the Linden data (Table SI 6, M1)  
401 was substantially worse (RMSE = 1.09,  $R^2$  = 0.13), as compared to the test set predictions of  
402 the Linden data set (Table 1, M1). This RMSE is even greater than the SD of the Suzuki data.  
403 It is notable that the pp-LFER, which does not include steric terms, does show promising  
404 statistics when evaluated with the Linden set alone (Table 1, M1), whereas the model  
405 calibrated with the Linden set does not extrapolate well to the external Suzuki set. We have

406 tried the reversed evaluation (i.e., using the Suzuki set as the training set and the Linden set  
407 as the test set, Table SI 6) and obtained similar statistics but substantially different  
408 regression coefficients.

409 The 3D-QSAR models handled the external prediction better than the pp-LFER model, but  
410 RMSE values for the predictions of the Suzuki data set (Table SI 6, M2-M7) were 0.13-0.19  
411 log units higher than the test set predictions for the Linden data set. The model variant that  
412 uses the COSMOsim3D alignment and LSPs (Table SI 6, M4) achieved an RMSE of 0.59 and an  
413  $R^2$  of 0.61, while all other models had  $RMSE > 0.68$  and  $R^2 < 0.5$ . For a given alignment, LSPs  
414 resulted in better or equivalent statistics as compared to vdW and ele. These results are in  
415 line with the findings we obtained from the model evaluation with the Linden data set only.  
416 Note that systematic under-predictions for the Suzuki data were not found; thus, the  
417 temperature difference is not a significant reason for the increased RMSE. We obtained  
418 similar statistics for the reversed evaluation (i.e., using the Suzuki set as training set and the  
419 Linden set as test set, Table SI 6). We also found that, if both Linden and Suzuki sets are  
420 combined and split to training and test sets, statistics for the test set prediction improves  
421 (RMSE,  $R^2$ ), which suggests that there are significant differences in the chemical domains  
422 that are covered by the two data sets. As an example, the Suzuki data set includes phenols  
423 and phenyl acetates, which are chemical classes not included in the Linden data set. On the  
424 other hand, only the Linden data set includes ethers and ketones. Moreover, the Suzuki data  
425 set is predominated by aromatic chemicals while the proportion of aromatic and aliphatic  
426 chemicals is comparable in the Linden data set.

427 We further tested if the steric restriction through the cavity can correctly be described by  
428 the model variant M4. The binding coefficients were predicted for the ten chemicals for  
429 which we were able to determine only the upper limit of  $\log K_{a1}$  ( $< 1.3$ ) in the previous work

430 (Linden et al., 2016). These chemicals are most likely too large to fit into the  $\alpha$ CD cavity.  
431 Eight of the ten chemicals had predicted  $\log K_{a1}$  values of  $1.3 \pm 0.4$ , which is in a semi-  
432 quantitative agreement with our experiments.  $\log K_{a1}$  values for 1-chloronaphthalene  
433 (predicted  $\log K_{a1}$  2.67) and acenaphthene (predicted  $\log K_{a1}$  2.42) were overestimated by >  
434 1 log unit. In contrast, the prediction of a similar chemical, acenaphthylene resulted in a  
435 predicted  $\log K_{a1}$  of 1.7. The COSMOsim3D alignment placed acenaphthene and  
436 acenaphthylene in different positions, which likely explains the deviation in the predictions.

#### 437 **4 Conclusions**

438 A 3D-QSAR model with COSMOsim3D (Thormann et al., 2012) for alignment and LSPs for  
439 independent variables in PLS regression analysis was capable of predicting  $\alpha$ CD binding  
440 constants for organic chemicals with an RMSE of 0.45 log units. This model can be used for  
441 the prediction of unknown  $\alpha$ CD binding constants for neutral organic chemicals and covers  
442 the most important steric effects that influence the binding to  $\alpha$ CD (Linden et al., 2016). As  
443 assumed, the description of the binding to  $\alpha$ CD needs to include the 3D-structure of the  
444 solutes because the 3D-QSAR model worked much better than the simple correlation with  
445  $\log K_{ow}$  (Linden et al., 2016) and better than the 2D-QSAR model (pp-LFER) considered here.  
446 Hence, it can be concluded that the LSPs are more suitable variables for 3D-QSAR modeling  
447 of the binding process to  $\alpha$ CD and probably for other binding processes as well, e.g., binding  
448 to other types of cyclodextrin with a different application range. Use of 7 out of 60 chemicals  
449 as templates for the alignment appeared to be sufficient, also with regard to the prediction  
450 for 84 external data (Suzuki, 2001). Consequently, the combination of COSMOsim3D and  
451 COSMOsar3D may be applicable to similar binding systems with an unknown or flexible  
452 target-structure, as far as data for some strongly binding chemicals are available. In an

453 upcoming study, we will apply the 3D-QSAR modeling approaches tested in this study to  
454 model the binding to serum albumin, which also showed specific 3D effects.

## 455 **Acknowledgements**

456 The authors thank the Helmholtz Interdisciplinary Graduate School for Environmental  
457 Research (HIGRADE) for financial support and Sven Jakobtorweihen at Hamburg University  
458 of Technology for providing the molecular dynamics simulations. SE acknowledges the  
459 financial support from the MEXT/JST Tenure Track Promotion Program. The authors thank  
460 Nadin Ulrich for helpful comments on an early version of the manuscript.

## 461 **Appendix A. Supplementary material**

462 Supplementary data associated with this article can be found, in the online version, at ...

- 464 Abraham, M.H., Andonian-Haftvan, J., Whiting, G.S., Leo, A., Taft, R.S., 1994. Hydrogen bonding. Part  
465 34. The factors that influence the solubility of gases and vapours in water at 298 k, and a new  
466 method for its determination. *Journal of the Chemical Society, Perkin Transactions 2*, 1777-1791.
- 467 Abraham, M.H., Ibrahim, A., Zissimos, A.M., 2004. Determination of sets of solute descriptors from  
468 chromatographic measurements. *J. Chromatogr. A* 1037, 29-47.
- 469 Baroni, M., Clementi, S., Cruciani, G., Costantino, G., Riganelli, D., Oberrauch, E., 1992. Predictive  
470 ability of regression models. Part ii: Selection of the best predictive pls model. *J. Chemom.* 6, 347-356.
- 471 Baroni, M., Costantino, G., Cruciani, G., Riganelli, D., Valigi, R., Clementi, S., 1993. Generating optimal  
472 linear pls estimations (golpe): An advanced chemometric tool for handling 3d-qsar problems. *Quant.*  
473 *Struct.-Act. Relat.* 12, 9-20.
- 474 Blach, P., Fourmentin, S., Landy, D., Cazier, F., Surpateanu, G., 2008. Cyclodextrins: A new efficient  
475 absorbent to treat waste gas streams. *Chemosphere* 70, 374-380.
- 476 Bronner, G., Goss, K.-U., 2011. Predicting sorption of pesticides and other multifunctional organic  
477 chemicals to soil organic carbon. *Environ. Sci. Technol.* 45, 1313-1319.
- 478 Connors, K.A., 1997. The stability of cyclodextrin complexes in solution. *Chem. Rev.* 97, 1325-1357.
- 479 Cox, G.S., Turro, N.J., Yang, N.C.C., Chen, M.J., 1984. Intramolecular exciplex emission from aqueous  
480  $\beta$ -cyclodextrin solutions. *J. Am. Chem. Soc.* 106, 422-424.
- 481 Cramer, R.D., Patterson, D.E., Bunce, J.D., 1988. Comparative molecular field analysis (comfa). 1.  
482 Effect of shape on binding of steroids to carrier proteins. *J. Am. Chem. Soc.* 110, 5959-5967.
- 483 De Lisi, R., Lazzara, G., Milioto, S., Muratore, N., 2007. Adsorption of a dye on clay and sand. Use of  
484 cyclodextrins as solubility-enhancement agents. *Chemosphere* 69, 1703-1712.
- 485 Del Valle, E.M., 2004. Cyclodextrins and their uses: A review. *Process Biochem.* 39, 1033-1046.
- 486 Endo, S., Goss, K.U., 2011. Serum albumin binding of structurally diverse neutral organic compounds:  
487 Data and models. *Chem. Res. Toxicol.* 24, 2293-2301.
- 488 Endo, S., Goss, K.U., 2014. Applications of polyparameter linear free energy relationships in  
489 environmental chemistry. *Environmental Science & Technology* 48, 12477-12491.
- 490 Endo, S., Watanabe, N., Ulrich, N., Bronner, G., Goss, K.-U., 2015. Ufz-Iser database v 2.1. Helmholtz  
491 Centre for Environmental, Leipzig, Germany.
- 492 Flaherty, R.J., Nshime, B., DeLaMarre, M., DeJong, S., Scott, P., Lantz, A.W., 2013. Cyclodextrins as  
493 complexation and extraction agents for pesticides from contaminated soil. *Chemosphere* 91, 912-920.
- 494 Gebhardt, J., Hansen, N., 2016. Calculation of binding affinities for linear alcohols to  $\alpha$ -cyclodextrin by  
495 twin-system enveloping distribution sampling simulations. *Fluid Phase Equilib.* 422, 1-17.
- 496 Goss, K.-U., 2005. Predicting the equilibrium partitioning of organic compounds using just one linear  
497 solvation energy relationship (Iser). *Fluid Phase Equilib.* 233, 19-22.
- 498 Gramatica, P., 2007. Principles of qsar models validation: Internal and external. *QSAR &*  
499 *Combinatorial Science* 26, 694-701.
- 500 Hedges, A.R., 1998. Industrial applications of cyclodextrins. *Chem. Rev.* 98, 2035-2044.
- 501 Herrmann, A., 2014. Dynamic combinatorial/covalent chemistry: A tool to read, generate and  
502 modulate the bioactivity of compounds and compound mixtures. *Chem. Soc. Rev.* 43, 1899-1933.
- 503 Ishiwata, S., Kamiya, M., 1999. Cyclodextrin inclusion: Catalytic effects on the degradation of  
504 organophosphorus pesticides in neutral aqueous solution. *Chemosphere* 39, 1595-1600.
- 505 Karickhoff, S.W., McDaniel, V.K., Melton, C., Vellino, A.N., Nute, D.E., Carreira, L.A., 1991. Predicting  
506 chemical reactivity by computer. *Environ. Toxicol. Chem.* 10, 1405-1416.
- 507 Kastenholz, M.A., Pastor, M., Cruciani, G., Haaksma, E.E.J., Fox, T., 2000. Grid/cpca: A new  
508 computational tool to design selective ligands. *J. Med. Chem.* 43, 3033-3044.
- 509 Kauffman, G.W., Jurs, P.C., 2001. Qsar and k-nearest neighbor classification analysis of selective  
510 cyclooxygenase-2 inhibitors using topologically-based numerical descriptors. *J. Chem. Inf. Comput.*  
511 *Sci.* 41, 1553-1560.



512 Kim, K.H., 1995. Comparative molecular field analysis (comfa). in: Dean, P.M. (Ed.). Molecular  
513 similarity in drug design. Springer Netherlands, Dordrecht, pp. 291-331.

514 Klamt, A., 1995. Conductor-like screening model for real solvents: A new approach to the  
515 quantitative calculation of solvation phenomena. *J. Phys. Chem.* 99, 2224-2235.

516 Klamt, A., 2011. The cosmo and cosmo-rs solvation models. *Wiley Interdisciplinary Reviews:  
517 Computational Molecular Science* 1, 699-709.

518 Klamt, A., Jonas, V., Bürger, T., Lohrenz, J.C.W., 1998. Refinement and parametrization of cosmo-rs. *J.  
519 Phys. Chem. A* 102, 5074-5085.

520 Klamt, A., Thormann, M., Wichmann, K., Tosco, P., 2012. Cosmosar3d: Molecular field analysis based  
521 on local cosmo  $\sigma$ -profiles. *Journal of Chemical Information and Modeling* 52, 2157-2164.

522 Linden, L., Goss, K.-U., Endo, S., 2016. Exploring 3d structural influences of aliphatic and aromatic  
523 chemicals on  $\alpha$ -cyclodextrin binding. *J. Colloid Interface Sci.* 468, 42-50.

524 Liu, H., Cai, X., Chen, J., 2013. Mathematical model for cyclodextrin alteration of bioavailability of  
525 organic pollutants. *Environmental Science & Technology* 47, 5835-5842.

526 Marinescu, L., Bols, M., 2009. Cyclodextrin derivatives that display enzyme catalysis. *Trends Glycosci.  
527 Glycotechnol.* 21, 309-323.

528 Marques, H.M.C., 2010. A review on cyclodextrin encapsulation of essential oils and volatiles. *Flavour  
529 Fragrance J.* 25, 313-326.

530 Rücker, C., Rücker, G., Meringer, M., 2007. Y-randomization and its variants in qspr/qsar. *Journal of  
531 Chemical Information and Modeling* 47, 2345-2357.

532 Sancho, M.I., Andujar, S.A., Porasso, R.D., Enriz, R.D., 2016. Theoretical and experimental study of  
533 inclusion complexes of  $\beta$ -cyclodextrins with chalcone and 2',4'-dihydroxychalcone. *J. Phys. Chem. B.*

534 Sawicki, R., Mercier, L., 2006. Evaluation of mesoporous cyclodextrin-silica nanocomposites for the  
535 removal of pesticides from aqueous media. *Environmental Science & Technology* 40, 1978-1983.

536 Schneider, H.J., 2009. Binding mechanisms in supramolecular complexes. *Angew. Chem. Int. Ed.* 48,  
537 3924-3977.

538 Sijm, D., Kraaij, R., Belfroid, A., 2000. Bioavailability in soil or sediment: Exposure of different  
539 organisms and approaches to study it. *Environ. Pollut.* 108, 113.

540 Stanier, C.A., O'Connell, M.J., Clegg, W., Anderson, H.L., 2001. Synthesis of fluorescent stilbene and  
541 tolan rotaxanes by suzuki coupling. *Chem. Commun.*, 493-494.

542 Suzuki, T., 2001. A nonlinear group contribution method for predicting the free energies of inclusion  
543 complexation of organic molecules with  $\alpha$ - and  $\beta$ -cyclodextrins. *J. Chem. Inf. Comput. Sci.* 41, 1266-  
544 1273.

545 Tabushi, I., 1982. Cyclodextrin catalysis as a model for enzyme action. *Acc. Chem. Res.* 15, 66-72.

546 Thormann, M., Klamt, A., Wichmann, K., 2012. Cosmosim3d: 3d-similarity and alignment based on  
547 cosmo polarization charge densities. *Journal of Chemical Information and Modeling* 52, 2149-2156.

548 Tosco, P., Balle, T., 2011. Open3dqsar: A new open-source software aimed at high-throughput  
549 chemometric analysis of molecular interaction fields. *J. Mol. Model.* 17, 201-208.

550 Tosco, P., Balle, T., Shiri, F., 2011. Open3dalign: An open-source software aimed at unsupervised  
551 ligand alignment. *J. Comput. Aided Mol. Des.* 25, 777-783.

552 Tropsha, A., Gramatica, P., Gombar, V.K., 2003. The importance of being earnest: Validation is the  
553 absolute essential for successful application and interpretation of qspr models. *QSAR &  
554 Combinatorial Science* 22, 69-77.

555 Villaverde, J., Pérez-Martínez, J.I., Maqueda, C., Ginés, J.M., Morillo, E., 2005. Inclusion complexes of  
556  $\alpha$ - and  $\gamma$ -cyclodextrins and the herbicide norflurazon: I. Preparation and characterisation. II.  
557 Enhanced solubilisation and removal from soils. *Chemosphere* 60, 656-664.

558 Yu, J.S., Wei, F.D., Gao, W., Zhao, C.C., 2002. Thermodynamic study on the effects of beta-  
559 cyclodextrin inclusion with berberine. *Spectrochim. Acta, Part A* 58, 249-256.

560

561 All tables:

562 **Table 2. Comparison of the statistical results of the different modeling approaches for the**  
563 **prediction of  $\log K_{a1}$  of the Linden data set using test set selection procedure 1.**

Modeling approach	Method	Alignment	Field	$q^2 \pm SD$	RMSE $\pm$ SD	$R^2 \pm SD$
M1	pp-LFER				$0.52 \pm 0.05$	$0.68 \pm 0.07$
M2	3D-QSAR	O3A	LSP	$0.63 \pm 0.03$	$0.54 \pm 0.08$	$0.56 \pm 0.17$
M3	3D-QSAR	O3A	vdW ele	$0.58 \pm 0.08$	$0.53 \pm 0.11$	$0.53 \pm 0.11$
<b>M4</b>	<b>3D-QSAR</b>	<b>COSMOsim3D</b>	<b>LSP</b>	<b><math>0.83 \pm 0.02</math></b>	<b><math>0.45 \pm 0.06</math></b>	<b><math>0.70 \pm 0.08</math></b>
M5	3D-QSAR	COSMOsim3D	vdW ele	$0.70 \pm 0.01$	$0.56 \pm 0.06$	$0.53 \pm 0.12$
M6a	3D-QSAR	COSMOsim3D	LSP	$0.66 \pm 0.06$	$0.51 \pm 0.06$	$0.61 \pm 0.09$
M6b	3D-QSAR	receptor X-ray COSMOsim3D	LSP	$0.71 \pm 0.04$	$0.49 \pm 0.04$	$0.64 \pm 0.07$
M7	3D-QSAR	receptor MDsim COSMOsim3D	vdW ele	$0.51 \pm 0.08$	$0.55 \pm 0.08$	$0.56 \pm 0.13$

564 **O3A means open3DALIGN,  $q^2$  is the coefficient of determination for the leave-two-out**  
565 **cross validation using the training set, RMSE is the root mean square error of the test set in**  
566 **log units, and  $R^2$  is the coefficient of determination of the test set. LSP, vdW, and ele**  
567 **indicate the usage of local sigma profiles, van der Waals interaction field, and electrostatic**  
568 **interaction field as molecular interaction field, respectively, SD is standard deviation, and**  
569 **MDsim is molecular dynamics simulation.**