

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29

DR. LEI HUANG (Orcid ID : 0000-0002-7846-9760)

DR. PETER FANTKE (Orcid ID : 0000-0001-7148-6982)

Article type : Original Article

## **A Quantitative Property-Property Relationship for the Internal Diffusion Coefficients of Organic Compounds in Solid Materials**

Lei Huang<sup>1\*</sup>, Peter Fantke<sup>2</sup>, Alexi Ernstoff<sup>2</sup> and Olivier Jolliet<sup>1</sup>

<sup>1</sup>Department of Environmental Health Sciences, School of Public Health, University of Michigan, Ann Arbor, MI, USA

<sup>2</sup>Division for Quantitative Sustainability Assessment, Department of Management Engineering, Technical University of Denmark, 2200 Kgs. Lyngby, Denmark

\*Corresponding author, [huanglei@umich.edu](mailto:huanglei@umich.edu)

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/ina.12395](https://doi.org/10.1111/ina.12395)

This article is protected by copyright. All rights reserved

30

## 31 Abstract

32 Indoor releases of organic chemicals encapsulated in solid materials are major contributors to  
33 human exposures and are directly related to the internal diffusion coefficient in solid materials.  
34 Existing correlations to estimate the diffusion coefficient are only valid for a limited number of  
35 chemical-material combinations. This paper develops and evaluates a quantitative property-  
36 property relationship (QPPR) to predict diffusion coefficients for a wide range of organic  
37 chemicals and materials. We first compiled a training dataset of 1103 measured diffusion  
38 coefficients for 158 chemicals in 32 consolidated material types. Following a detailed analysis of  
39 the temperature influence, we developed a multiple linear regression model to predict diffusion  
40 coefficients as a function of chemical molecular weight (MW), temperature, and material type  
41 (adjusted  $R^2$  of 0.93). The internal validations showed the model to be robust, stable and not a  
42 result of chance correlation. The external validation against two separate prediction datasets  
43 demonstrated the model has good predicting ability within its applicability domain ( $R^2_{\text{ext}} > 0.8$ ),  
44 namely MW between 30 and 1178 g/mol and temperature between 4 and 180 °C. By covering a  
45 much wider range of organic chemicals and materials, this QPPR facilitates high-throughput  
46 estimates of human exposures for chemicals encapsulated in solid materials.

## 47 Keywords

48 Diffusion, Solid materials, Consumer products, Indoor release, Organic chemicals, Correlation

## 49 Practical implications

50 The quantitative property-property relationship developed by the present study provides a more  
51 comprehensive correlation method to estimate the diffusion coefficients, as it covers a wide  
52 range of organic chemicals and solid materials, and also considers the effect of temperature. This  
53 model provides the basis for facilitating high-throughput estimates of indoor human exposures  
54 for chemicals encapsulated in solid materials relevant for several science-policy fields, such as  
55 chemical alternatives assessment (CAA), risk assessment (RA) and life cycle assessment (LCA).

56

## 57 1. Introduction

58 Chemicals encapsulated in solid materials have been identified as a major source of passive  
59 emissions to indoor air<sup>1-3</sup> and of transfers into food<sup>4</sup> and onto skin<sup>5</sup>. Typical examples include

60 chemicals used as flame retardants in furniture and plasticizers in food contact materials. To  
61 estimate the release of these chemicals from solid materials, and eventually consumer exposures,  
62 the diffusion coefficient,  $D$  ( $\text{m}^2/\text{s}$ ), for chemicals encapsulated in solid materials, is essential  
63 information.  $D$  describes the transport of a molecule through a material, which is specific for a  
64 chemical-material combination and is also influenced by ambient temperature. Experimental  
65 techniques such as chamber tests for building materials<sup>6,7</sup>, and sorption/desorption experiments  
66 for polymer materials<sup>8-10</sup> have enabled measurement of a limited number of chemical diffusion  
67 coefficients for building materials such as vinyl flooring, gypsum board, particle board, plywood,  
68 carpet and cement<sup>11-14</sup>, as well as polymer materials including polyethylene (PE), polystyrene  
69 (PS), polypropylene (PP), and polyvinyl chloride (PVC)<sup>4, 15, 16</sup>. However, given the limited  
70 number of chemical-material combinations with measured  $D$ s, and the costly and time-  
71 consuming nature of experiments, quantitative relationships are needed to complement existing  
72 measurements by predicting the diffusion coefficients from known physicochemical properties for  
73 chemicals without experimental data. This is especially important for high-throughput  
74 approaches where a large number of chemical-material combinations need to be evaluated and  
75 for which it is unrealistic to perform experiments on all relevant combinations.

76 Several correlation methods have been developed to estimate the diffusion coefficients from  
77 physicochemical properties of chemicals<sup>8, 12, 17-19</sup>. For example, Berens and Hopfenberg  
78 correlated the  $D$  to the mean molecular diameter of the diffusing molecule, using data on more  
79 than 20 chemicals in 3 glassy materials including PVC, PS and polymethyl methacrylate  
80 (PMMA)<sup>8</sup>. Zhao et al. found a correlation between  $D$  and vapor pressure for water and 8  
81 aromatic hydrocarbons in polyurethane foam (PUF)<sup>19</sup>. Furthermore, both Bodalal et al. and Cox  
82 et al. estimated the  $D$  as a function of molecular weight<sup>12, 18</sup>. The former study considered  
83 measured  $D$  data on 5 aromatics and 5 aldehydes in several building materials<sup>12</sup>, while the latter  
84 study considered data on 4 alkanes in vinyl flooring<sup>18</sup>. For each of these aforementioned  
85 approaches, the main limitation is that the correlations are specific to certain chemical classes  
86 and materials; for example aldehydes in plywood, which limits their application for other  
87 materials and chemical classes. Addressing this research gap to facilitate wider applicability,  
88 Guo developed a method which estimates the diffusion coefficient as a function of the  
89 chemical's molar volume for mixed chemical classes<sup>17</sup>. However, this approach is limited to 6

90 building materials and are developed based on a small dataset of limited chemical classes ( $\leq 3$   
91 chemical classes for 5 of the 6 building materials).

92 The aforementioned correlation methods consider experiments for building materials at room  
93 temperature, and therefore temperature is not relevant and thus not considered in the correlation  
94 model. For other exposure scenarios, such as transfer of chemicals from food contact materials  
95 (FCMs) into food, ambient temperature is highly relevant because FCMs can be heated,  
96 refrigerated, or frozen. Accordingly, Begley et al. presented a correlation method to estimate the  
97 diffusion coefficient in 9 polymer materials as a function of molecular weight and temperature<sup>4</sup>,  
98 which is not applicable beyond the considered polymers.

99 In all, the currently available correlation methods to estimate  $D$  do not provide sufficient  
100 coverage of chemicals encapsulated in consumer products in different use scenarios (i.e. ambient  
101 temperatures). Developing low-tier, high-throughput methods to estimate exposure to chemical  
102 in consumer products across a variety of chemical-material combinations is a recent focus in  
103 various science-policy fields such as computational exposure science and life cycle assessment  
104 (LCA)<sup>20-25</sup>. Addressing the lack of methods to estimate  $D$  for a variety of chemical-product  
105 scenarios, the present study aims to develop a more comprehensive correlation method to  
106 estimate  $D$  for wide range of organic compounds in multiple solid materials. More specifically,  
107 we aim to:

- 108 1) Carry out a comprehensive and extensive literature review to collect experimental diffusion  
109 coefficient data on a wide range of materials and chemicals.
- 110 2) Use multiple linear regression techniques to establish the relationship between the diffusion  
111 coefficient and various predictor variables including physiochemical properties, material  
112 properties and environmental characteristics.
- 113 3) Perform internal and external validations to characterize the validity and predictive power of  
114 the developed correlation.

115 Since the material type is a categorical property variable and is not related to the chemical's  
116 molecular structure, we call this correlation a quantitative property-property relationship (QPPR)  
117 instead of a quantitative structure-activity relationship (QSAR). This QPPR provides a more  
118 advanced correlation method to estimate the diffusion coefficients of organic compounds  
119 compared to previous studies, as it covers a wide range of solid materials and physiochemical  
120 properties, and also considers the effect of temperature. By providing reliable estimates of this

121 key diffusion parameter for a large number of chemicals, this method will facilitate high-  
122 throughput assessments of chemical emissions and human exposures for chemicals encapsulated  
123 in solid materials relevant for chemical alternatives assessment (CAA), risk assessment and LCA.

## 124 2. Materials and methods

### 125 2.1 Dataset

126 Experimental diffusion coefficient data were compiled from 68 references from the peer-  
127 reviewed scientific literature. The initial dataset contained a total of 1124 records covering 161  
128 unique chemicals and 88 distinct solid materials (provided in Supporting Info). Experimental  
129 data expressed in  $\text{cm}^2/\text{s}$  were converted to  $\text{m}^2/\text{s}$ . There are different types of diffusion  
130 coefficients reported in the literature, so harmonization of these data was performed to develop a  
131 consistent dataset. For diffusion coefficients measured in liquid sorption experiments, the  
132 ‘intrinsic’ diffusion coefficients, corrected for the swelling of materials were collected <sup>10</sup>.  
133 Sorption of the liquid molecules inside the solid material may cause swelling of the material,  
134 which would lead to decreased observed diffusion coefficients and thus need to be corrected <sup>10</sup>.  
135 For porous materials consisting of pore space and solid material, two types of models can be  
136 used to describe the chemical transport through these materials. The one-phase model considers  
137 the porous material as an assumed homogeneously mixed material, so an ‘apparent’ diffusion  
138 coefficient is used to describe the chemical diffusion through such imaginary material <sup>7</sup>. In  
139 contrast, the multi-phase model considers the material as a mixture of pores and solid parts, and  
140 the chemical diffuses mainly through the pores if the pores are interconnected, or through the  
141 pores and solid parts alternately if the pores are isolated from each other. The gas-phase diffusion  
142 through the pores, which can be described by an ‘effective’ diffusion coefficient, is assumed to  
143 be much faster than the diffusion through the solid parts <sup>7</sup>. Haghghat et al., <sup>7</sup> has demonstrated  
144 that the ‘apparent’ diffusion coefficient is equivalent to the ‘effective’ diffusion coefficient ( $D_e$ )  
145 divided by the material phase-gas phase partition coefficient ( $K_{ma}$ ). Thus, for porous materials  
146 the ‘apparent’ diffusion coefficients reported in studies were collected <sup>26</sup>. For studies where only  
147 the  $D_e$  and  $K_{ma}$  were reported <sup>27-29</sup>, they were converted to ‘apparent’ diffusion coefficients using  
148 the aforementioned method. Data were excluded for studies where only the ‘effective’ diffusion  
149 coefficients were reported.

150 From the initial dataset, 21 records were excluded from further analyses because they involve  
151 chemicals that are inorganic, chemicals for which no CAS number could be identified, or

152 chemicals that are polymer chains with varying molecular weights. The final considered dataset  
153 thus includes 1103 records for 158 unique chemicals and 87 materials.

## 154 2.2 Modeling methods

### 155 2.2.1 Multiple linear regression

156 A multiple linear regression (MLR) analysis was performed to identify and quantify the effect of  
157 different parameters on the diffusion coefficient. The MLR model takes the following general  
158 form:

$$159 \log_{10}D = \alpha + \beta_1 \cdot X_1 + \dots + \beta_n \cdot X_n + b_1 \cdot M_1 + \dots + b_m \cdot M_m \quad (1)$$

160 where  $\log_{10}D$  is the logarithm of the diffusion coefficient ( $\text{m}^2/\text{s}$ ),  $\alpha$  is the intercept;  $X_1$  to  $X_n$  are  
161 independent variables related to physiochemical properties, such as molecular weight, molar  
162 volume, and vapor pressure, and/or environmental characteristics like temperature;  $\beta_1$  to  $\beta_n$  are  
163 regression coefficients for the respective independent variables  $X_1$  to  $X_n$ ; and  $M_1$  to  $M_m$  are  
164 dummy variables for the solid materials, with one dummy variable per type of material. A  
165 dummy variable equals 1 for the material type it represents, and equals 0 for all other materials;  
166 for example,  $M_1 = 1$  for material type = 1,  $M_1 = 0$  for material types 2 to  $m$ .  $b_1$  to  $b_m$  are  
167 regression coefficients for the respective dummy variables  $M_1$  to  $M_m$ . The number of  $m$  is equal  
168 to the number of material types considered minus 1, since the material type with the highest  
169 number of measured  $D$  data is used as the reference material type and does not require a dummy  
170 available in the MLR. Note that the MLR model gives one coefficient for each material type,  
171 while a material type can represent a single pure substance such as calcium silicate, a composite  
172 material such as vinyl flooring and gypsum board, or a group of similar materials such as  
173 wooden boards. Details of the material types will be discussed later. This regression equation  
174 also implies that the material coefficients ( $b_1$  to  $b_m$ ) and the physiochemical property coefficients  
175 ( $\beta_1$  to  $\beta_n$ ) are independent of each other, which if corroborated by internal and external  
176 validations (Section 2.3), allow for the maximum prediction coverage in terms of chemical-  
177 material combinations. All regression coefficients were estimated by the least squares (LS)  
178 method. All regression analyses were performed using IBM SPSS Statistics version 23 (IBM  
179 corporation, Armonk, New York).

### 180 2.2.2 Grouping of materials and initial regressions

181 To reduce the number of dummy variables, to avoid over-fitting of the MLR model, and to have  
182 a minimum of 10 records and 3 different chemicals per material type to ensure enough variability,

183 the 87 original materials were grouped into 32 consolidated material types, based on the  
184 similarity of the regression coefficients and the material types (see Supporting Information (SI),  
185 Section S1). Thus  $m = 31$  in Eq. 1, with PET as the 32<sup>nd</sup> and reference material, since it is the  
186 material with most reported diffusion coefficients.

187 In previous studies, either the chemical's molecular weight ( $MW$ ), molar volume ( $MV$ ) or vapor  
188 pressure ( $VP$ ) has been used as predictor of the diffusion coefficient in a given material<sup>12, 17-19</sup>.  
189 Begley et al.<sup>4</sup> also suggested that the logarithm of the diffusion coefficient varies linearly with  
190 the inverse of the absolute temperature ( $1/T$ ). Thus, the initial regression was performed to  
191 identify which of the above variables ( $MW$ ,  $MV$ ,  $VP$  and  $1/T$ ) are best predictors of the diffusion  
192 coefficients of compounds encapsulated in the 32 material types, i.e., to identify  $X_1$  to  $X_n$  in Eq.  
193 (1). Details of the initial regression process are presented in SI, Section S2. Results of the initial  
194 regression model suggest that the log-molecular weight and the inverse of the absolute  
195 temperature are the most important predictors, and therefore the employed MLR model takes the  
196 following form:

$$197 \quad \log_{10} D = \alpha + \beta_{\log_{10} MW} \cdot \log_{10} MW + \beta_{1/T} \cdot \frac{1}{T} + b_1 \cdot M_1 + \dots + b_m \cdot M_m \quad (2)$$

198 where  $MW$  is the chemical's molecular weight (g/mol) and  $T$  is the absolute temperature (K).  
199 The model performance of using log-molecular weight and molecular weight as predictors were  
200 very close when using the training dataset (1103 records,  $m=31$ ), but the model using log-  
201 molecular weight as predictor was finally selected since it performs better for high-molecular-  
202 weight chemicals (Section 3.3.3).

### 203 2.2.3 Temperature dependence

204 Studies have shown that the activation energy of diffusion is a contributor to the temperature  
205 dependence of the diffusion coefficient and varies as function of both the material and the  
206 chemical properties<sup>4, 30, 31</sup>. Thus, ideally a specific temperature correction coefficient should be  
207 used for each chemical-material combination. Since data availability is not sufficient to  
208 determine chemical-specific temperature coefficients for each of the 32 materials, and since  
209 chemical properties seem to have limited influence on the activation energy<sup>4, 30</sup>, we followed the  
210 strategy of Begley et al.<sup>4</sup>, differentiating temperature coefficients for a limited number of  
211 material groups, applying one generic temperature coefficient for all chemicals within each  
212 material group. Begley et al.<sup>4</sup> have introduced a variable  $\tau$  to adjust the temperature coefficient  
213 for two groups of materials, where  $\tau$  equals 0 or 1577 for 9 different polymers, which

214 corresponds to activation energy of 86.9 kJ/mol for e.g. LDPE or 100 kJ/mol for e.g. HDPE. To  
215 analyze the temperature dependency of the diffusion coefficients in our dataset, we first plotted  
216  $\log_{10}D$  against  $1/T$  for each of the 32 material types (SI Section S3). The plots generally show as  
217 expected<sup>4</sup> an inverse relationship in which  $\log_{10}D$  is decreasing with increasing  $1/T$ , different  
218 materials exhibiting different slopes. Since variability in diffusion coefficient is higher between  
219 than within given studies, we first determined a temperature coefficient for each chemical-  
220 material-study combination, and then calculated an average temperature coefficient for each  
221 material type by averaging all temperature coefficients belonging to the same material type. The  
222 analysis of the material-specific temperature coefficients showed that the materials can be  
223 grouped into three categories: (1) high-, (2) medium- and (3) low-coefficient categories, with  
224 three corresponding values for the temperature coefficient adjustment factor  $\tau$ , which are given in  
225 Section 3.1. Details are presented in SI Section S3.3. The adjusted MLR model takes the  
226 following form accordingly:

$$227 \quad \log_{10}D = \alpha + \beta_{\log MW} \cdot \log_{10}MW + \frac{\beta_{1/T+\tau}}{T} + b_1 \cdot M_1 + \dots + b_m \cdot M_m, \quad (3)$$

228

#### 229 2.2.4 Final regression

230 To avoid multicollinearity problems in the MLR model and to avoid the influence of the material  
231 type “Limited-data material group” on the temperature coefficients, we fixed the temperature  
232 coefficients determined using Eq. 3 and thus the final regression takes the following form:

$$233 \quad \log_{10}D - \frac{\beta_{1/T+\tau}}{T} = \alpha + \beta_{\log MW} \cdot \log_{10}MW + b_1 \cdot M_1 + \dots + b_m \cdot M_m, \quad (4)$$

234 where the dependent variable is  $\log_{10}D - (\beta_{1/T+\tau})/T$  instead of  $\log_{10}D$ , with the values of  $\beta_{1/T}$   
235 and  $\tau$  obtained from Eq. 3 and presented later in Section 3.1. In this final regression, all 1103  
236 records of measured  $D$  data were utilized including the material type “Limited-data material  
237 group”, leading to  $m=31$  material types, plus one reference material type, PET, with  $b_{PET} = 0$ .

#### 238 2.3 Model validation

239 Validation of the final MLR model (Eq. 4) was performed using the QSARINS software, version  
240 2.2.1 ([www.qsar.it](http://www.qsar.it)) which is developed by Gramatica et al.<sup>32, 33</sup>.

##### 241 2.3.1 Internal validation

242 The MLR model’s capacity to predict portions of the training dataset was evaluated in an internal  
243 validation process, using two techniques for internal validation in QSARINS. The first one is the

244 leave more out (LMO) cross-validation technique, which iteratively and randomly exclude a  
245 certain percentage of the measured diffusion coefficient data, and then computes the regression  
246 coefficients with the remaining data and uses those coefficients to make predictions for the  
247 excluded ones<sup>33</sup>. We used 1000 iterations and the percentage of the excluded elements was set  
248 as 20%.

249 The second technique for internal validation is the Y-scrambling procedure, which demonstrates  
250 that the model is not the result of chance correlation. In this procedure, the experimental  
251 responses (in our study, the temperature-adjusted diffusion coefficients) are shuffled at random  
252 and used with the original predictors to establish an MLR model. If the original MLR model is  
253 internally valid, the performances of the scrambled models should be much worse than the  
254 original model<sup>33</sup>. We used 1000 iterations for the Y-scrambling.

### 255 2.3.2 External validation

256 We also evaluated the model ability to provide reliable predictions on new datasets in a so-called  
257 external validation process, using the following two approaches.

258 The first approach was to split the existing dataset (1103 records) into one training dataset and  
259 one prediction datasets. The training dataset was used to generate regression coefficients of the  
260 MLR model, and then the MLR model was applied to the prediction set to examine the  
261 prediction performances of the model. Three kinds of splitting were performed using existing  
262 options in the QSARINS software (see SI, Section S5.1 for details) by random percentage (20%  
263 of the entire dataset randomly selected as the prediction set, 80% rest to the training set), by  
264 response and by structure (data first ordered by responses of the temperature-adjusted diffusion  
265 coefficient, or by the first axis of principal component analysis (PCA) of the descriptors,  
266 respectively). We introduced a fourth kind of splitting by studies, since variability across studies  
267 for a given material is in general larger than variability within a given study, yielding similar  
268 sample sizes of approximately 880 data for the training set and 220 data for the prediction set (SI,  
269 Table S3).

270 The second approach of external validation was to use the entire collected dataset (1103 records)  
271 as the training set and to use an entirely separate dataset as the prediction set. For the prediction  
272 set, two datasets were used. The first one is a database of diffusion coefficients from the United  
273 States Food and Drug Administration (FDA), which is a “database available upon request” for  
274 guidance for industry (<http://www.fda.gov/Food/ucm081818.htm>), and includes non-peer

275 reviewed diffusion coefficient data reported by industry. This dataset includes 191 records of  
276 experimental diffusion coefficients of 46 chemicals in 22 materials which are mainly polymers  
277 used for food contact materials (see SI, Section S5.1 for details). The quality and reliability of  
278 these data are not characterized by FDA. The second prediction dataset is constructed from  
279 several studies conducted before 1982<sup>34-36</sup>, referenced in<sup>37</sup>. This dataset, designated as “Data by  
280 1982”, includes 281 records of measured diffusion coefficients of 92 chemicals in 8 polymer  
281 materials, also including self-diffusion (see SI, Section S5.1 for details). Data for both prediction  
282 sets are provided in Supporting Info.

### 283 2.3.3 Applicability domain (AD)

284 The analysis and definition of the applicability domain (AD) of models is a fundamental issue  
285 that must be addressed in QSAR and QPPR studies. The study of AD can provide information on  
286 the reliability of the model predictions, i.e., if the chemicals are inside the AD, the predictions  
287 are interpolated and are more reliable; if the chemicals are outside the AD, the predictions are  
288 extrapolated and less reliable, because effects can occur outside the AD that do not exist within  
289 the AD<sup>38</sup>. Three complementary methods were applied to define the AD of the diffusion  
290 coefficient QPPR: the range of model predictors, the leverage approach, and the PCA of the  
291 model predictors<sup>39</sup>. More explanation of these methods is provided in SI, Section S4. In our  
292 analysis, chemicals are considered inside the AD if they are viewed inside AD by all three  
293 methods, whereas chemicals are considered outside AD if they are viewed outside AD by all  
294 three methods, and finally chemicals that fall inside the AD for only one or two methods are  
295 considered as ‘borderline.’

296

## 297 3. Results and discussion

### 298 3.1 Temperature dependence of the diffusion coefficient

299 The compiled dataset of 1103 records including 158 chemicals and 32 material types shows that  
300 the diffusion coefficient in solid materials decreases with decreasing temperature, as  
301 demonstrated by the highly significant negative regression coefficient for the variable  $1/T$ , with  
302  $\beta_{1/T} = -4440$  (K) with a standard error (SE) of 164 (K) and  $p < 0.001$  in Eq. 2 (SI, Section  
303 S3.1). This is in agreement with previous studies<sup>4, 30, 31</sup>. This general tendency of decreasing  
304 diffusion with increasing  $1/T$  is well illustrated by the example of PET, the material with the  
305 most data available (Figure 1A – see SI, Figure S1 for other materials). To further refine the

306 coefficient for the temperature variable into specific materials groups, Figure 1B illustrates well  
307 for methyl methacrylate (MMA) homopolymer the importance of first determining a temperature  
308 coefficient for each separate study and material-chemical combination (Section 2.2.3) and then  
309 averaging the temperature coefficients across studies. The molecular weight-normalized  
310 diffusion coefficients show a negative linear relationship with  $1/T$  within each of the three  
311 experimental studies of Figure 1B<sup>40-42</sup>, with similar regression coefficients of -4530 (K), -5704  
312 (K), -3415 (K), averaging -4550 (K) with an SE of 305 (K). However, since the absolute  
313  $\log_{10}MW$ -normalized diffusion coefficients reported by Hennebert et al.<sup>42</sup> are much higher than  
314 those reported by the other two studies, doing one regression with all data from the three studies  
315 would result in a non-significant temperature coefficient (p-value of 0.19), thus demonstrating  
316 the importance to first perform temperature regressions using data from the same study and for  
317 the same chemical.

318 Table 1 presents the average temperature coefficients and their standard errors for each of the 32  
319 consolidated material types. Based on the values of the temperature coefficients (unit in K), the  
320 32 material types can be grouped into three categories: (1) high-coefficient category with  
321 relatively high (absolute value) temperature coefficients ( $< -5000$ ), i.e., materials in which  
322 diffusion coefficients are highly sensitive to the change in temperature, (2) medium-coefficient  
323 category with temperature coefficients in between ( $-5000 < (\beta_{1/T} + \tau) < -3000$ ), and (3) low-  
324 coefficient category with relatively low (absolute value) temperature coefficients ( $> -3000$ ), i.e.,  
325 materials in which diffusion coefficients are least sensitive to the change in temperature. Details  
326 for the grouping of temperature coefficients can be found in SI, Section S3.3.

327 The temperature coefficients  $\beta_{1/T}$  and  $\tau$  used in Eq. 4 for each of the three temperature-  
328 dependency material categories are obtained from the regression using the MLR model of Eq.  
329 S3-2 (SI, Section S3.3), yielding values of  $\beta_{1/T} = -3486 \pm 299$  (K) and  $\tau_{\text{high}} = -2391 \pm$   
330  $356$  (K),  $\tau_{\text{medium}} = 0$  (K) and  $\tau_{\text{low}} = +1676 \pm 510$  (K). Thus, for the High-, Medium- and  
331 Low-coefficient categories, the final temperature coefficients ( $\beta_{1/T} + \tau$ ) are -5877 (K), -3486 (K),  
332 and -1810 (K), corresponding to activation energy of 113, 66.7 and 34.7 (kJ/mol), respectively.  
333 Begley et al.<sup>4</sup> also aggregated 9 types of polymer materials into two temperature categories, with  
334 activation energy of 100 and 86.9 (kJ/mol), which have similar values with the high- and  
335 medium-coefficient categories in the present paper, to which these 9 polymer materials are  
336 assigned. These results indicate that the categorization of the temperature coefficient in the

337 present paper is consistent with previous studies, while extending the QPPR to a wider range of  
338 materials.

### 339 3.2 Final QPPR and model fitting

340 Using the full dataset (1103 records) and Eq. 4, the final MLR model for predicting the diffusion  
341 coefficient in solid materials is as follows:

$$342 \log_{10}D - \frac{\tau-3486}{T} = 6.39 - 2.49 \cdot \log_{10}MW + b \quad (5)$$

$$343 N = 1103, R^2 = 0.932, R^2_{\text{adj}} = 0.930, SE = 1.17, RMSE = 1.15$$

$$344 \text{ANOVA: } F = 457, df = 32, p < 0.0001$$

345 where  $D$  is the diffusion coefficient ( $\text{m}^2/\text{s}$ ),  $MW$  is molecular weight ( $\text{g/mol}$ ),  $T$  is absolute  
346 temperature (K),  $b$  and  $\tau$  (K) are the material-specific coefficients presented in Table 2. This  
347 model is provided as an excel model in Supporting Info to facilitate application. The standard  
348 errors for the intercept (6.39) and the coefficient of  $\log_{10}MW$  (-2.49) are 0.29 and 0.13,  
349 respectively. An SE of 1.17 of the final model (Eq. 5) indicates that the 95% confidence interval  
350 (CI) of the predicted response,  $\log_{10}D - (\tau-3486)/T$ , is the predicted value  $\pm 2.30$ . The 95% CI of  
351 the  $\log_{10}D$  cannot be directly calculated, but the average absolute difference between predicted  
352 and measured  $\log_{10}D$  is 0.83 across the whole dataset (1103 records), and 95% of this absolute  
353 difference is below 2.54.

354 This MLR model shows excellent fitting of the experimental data, with an adjusted R-square of  
355 0.932 and a root mean square error (RMSE) of 1.15. The model fit is highly significant with an  
356 ANOVA p-value smaller than 0.0001. Figure 2 shows the scatter plot of experimental versus  
357 predicted responses, which aligns well with the 1:1 line. In this MLR model, the response  
358 (dependent variable) is the temperature-adjusted log diffusion coefficient, i.e.,  $\log_{10}D - (\tau-3486)/T$ ,  
359 instead of  $\log_{10}D$ , in order to fix the temperature coefficients and to avoid multicollinearity  
360 problems, as mentioned in Section 2.2.4. The residual plot (Figure 3) shows that the residuals are  
361 distributed evenly throughout the dataset, again indicating the good fit of the linear model for the  
362 data.

363 The key predictors other than temperature in the MLR model are the material type and the  
364 molecular weight of the diffusing chemical. The regression coefficient when considering log-  
365 molecular weight is equal to -2.49, indicating that the diffusion coefficient decreases with  
366 increasing molecular weight. This implies that larger molecules diffuse more slowly compared to  
367 smaller molecules in solid materials, which is intuitive and consistent with findings from

368 previous studies<sup>4, 12, 17, 18</sup>. However, although the molecular weight is a highly significant  
369 predictor ( $p < 0.0001$ ), it explains less than 10% of the total variance of the diffusion coefficient  
370 (SI, Section S4).

371 The 31 dummy variables for the material types reflect the material dependency and account for  
372 most of the total variance of the diffusion coefficient, indicating that the diffusion coefficient in  
373 solid materials is strongly dependent on the material type. Since “Polyethylene terephthalate  
374 (PET)” was used as the reference material in the regression, the value of its coefficient  $b$  is zero  
375 (Table 2). For each of the other material types, the coefficient  $b$ , combined with the temperature  
376 coefficient  $\tau$ , i.e.  $b+(\tau+2391)/T$ , determines the difference in log-diffusion coefficient between  
377 that material type and PET, since PET has a temperature coefficient  $\tau$  of -2391 (K) (Table 2, last  
378 column). Chemicals in material types with high values of  $b+(\tau+2391)/T$  diffuse quicker than in  
379 material types with low values. Therefore, under room temperature ( $T = 298.15$  K), the values of  
380  $b+(\tau+2391)/T$  and the corresponding diffusion coefficients tend to be lower in dense, rigid  
381 materials such as glass, stainless steel, methyl methacrylate (MMA) polymers, polyethylene  
382 naphthalate (PEN), and rigid polymers including polyether ether ketone (PEEK), rigid PVC,  
383 polytetrafluoroethylene (PTFE), and polycarbonate (Table 2). In contrast, the values of  
384  $b+(\tau+2391)/T$  and the corresponding diffusion coefficients can be up to 13 orders of magnitude  
385 higher in flexible or porous materials, such as gypsum, wood, rubber, and polyurethane foam-  
386 based materials (Table 2). It should be noted that the composition and properties of a given  
387 material type may vary considerably depending on the intended use, as well as over time as  
388 material substitutions are made and production procedures differ. Thus, the material type  
389 coefficients in Table 2 actually represent an average composition and diffusion behavior for the  
390 specific material types.

391 The significance of the material type coefficient only indicates that the coefficients  $bs$  of these  
392 material types are significantly different from the reference material type, PET, but if another  
393 material type was selected as the reference material, the regression coefficients and statistical  
394 significance of all materials would change. Thus, the insignificance of the regression coefficients  
395 for material type variables does not indicate that those material types do not have a relevant  
396 influence on the diffusion coefficient. As a result, we keep all 31 material type dummy variables  
397 in the final regression to retain as much information as possible.

398 The MLR model given in Eq. 5 contains material-specific variables, so it is only valid for the 32  
399 material types presented in Table 2. For materials that do not belong to those 32 types, we built  
400 another generic QPPR to predict the diffusion coefficients, which is presented in SI, Section S4,  
401 which should be used with caution because of higher uncertainties.

### 402 3.3 Model validation results

#### 403 3.3.1 Internal validation

404 For the 20% leave-more-out (LMO) cross validation, the correlation coefficient,  $Q^2_{LMO}$  for the  
405 1000 iterations ranges from 0.89 to 0.95, with an average of 0.93, and a root mean square error  
406 for cross validation ( $RMSE_{cv}$ ) average of 1.19. Both the  $Q^2_{LMO}$  and  $RMSE_{cv}$  are similar to the  $R^2$   
407 and  $RMSE$  computed using the full dataset, which is 0.93 and 1.15, respectively. These results  
408 indicate that when fitted to a random 80% of the dataset the model is still able to predict the  
409 remaining 20% of the dataset, meaning that the model is internally stable.

410 For the Y-scrambling, the average  $R^2_{Yscr}$  and  $Q^2_{Yscr}$  for the 1000 iterations are 0.029 and -0.033,  
411 respectively, which are much smaller than the  $R^2$  and  $Q^2_{LMO}$  of the original model. The  $RMSE$   
412 for Y-scrambling,  $RMSE_{Yscr}$ , is 4.36 which is much higher than the  $RMSE$  and  $RMSE_{cv}$  of the  
413 original model. These results demonstrate that no correlation exists between the scrambled  
414 responses and the predictors. Thus, chance correlation for the original model can be ruled out.  
415 Overall, the internal validation demonstrates that the MLR model represented by Eq. 5 is robust  
416 and stable, and is not a result of chance correlation.

#### 417 3.3.2 External validation

418 As described in Section 2.3.2, the first method of external validation was to split the full dataset  
419 (1103 records) into training set and prediction set, and four types of splitting were performed,  
420 including splitting by a random 20%, by ordered response, by ordered structure, and by studies.  
421 Six criteria for external validation were computed and are presented in Table 3. The  $R^2_{ext}$  is the  
422 determination coefficient of the prediction set data using the model calculated using the training  
423 set data. The other five criteria,  $Q^2_{F1}$ <sup>43</sup>,  $Q^2_{F2}$ <sup>44</sup>,  $Q^2_{F3}$ <sup>45</sup>,  $r_m^2$ <sup>46</sup>, and  $CCC$ <sup>47</sup>, are external  
424 validation criteria proposed by different studies, which evaluate various aspects of the model's  
425 external prediction ability. These criteria are usually in accordance with each other but can  
426 sometimes give contradictory results<sup>47</sup>, so they need to be evaluated together. Chirico and  
427 Gramatica have proposed threshold values for these different criteria<sup>48</sup>, which are presented in  
428 Table 3. For the first three types of splitting (by random 20%, by ordered response, and by

429 ordered structure), the  $R^2_{\text{ext}}$  are higher than 0.9, and all of the other five criteria pass the  
430 threshold values and are also higher than 0.9, indicating good prediction ability of the model  
431 calculated using only the training set data. In these three types of splitting, the data were assigned  
432 to the training and prediction data sets either randomly or alternately (by ordered response or  
433 structure), so it is likely that a portion of the data from each study was assigned to the training set  
434 while the remaining portion of the data was assigned to the prediction set. As the result, the  
435 prediction set is well within the applicability domain (AD) defined by the training set (SI,  
436 Figures S2-S7), so it is expected that the model calculated using the training set can well predict  
437 the prediction set.

438 For the fourth type of splitting, splitting by studies, data from 30 studies were selected as the  
439 prediction set, while data from the remaining 48 studies constituted the training set. Thus, all  
440 data from one study and for one particular material will be either in the training or in the  
441 prediction set, so the validation using this splitting is close to a truly “external” validation. Most  
442 of the prediction set is inside the AD defined by the training set except for two data points (SI,  
443 Figures S8-S9). As a result, the  $R^2_{\text{ext}}$  dropped to 0.85, and the values of the other five validation  
444 criteria are apparently lower than those for the above three types of splitting, reflecting that  
445 variability is higher between than within studies. The five validation criteria nevertheless all pass  
446 the threshold values (Table 3), indicating that the model calculated using the training set has  
447 good prediction ability.

448 As a second method of external validation, the 1103 data points from the 68 studies were used as  
449 the training set, and additional data from an FDA database and from studies before 1982 were  
450 used as two separate prediction sets. As presented in Table 3, when using FDA dataset as the  
451 prediction set, the  $R^2_{\text{ext}}$  is reduced to 0.80 which is lower than the  $R^2_{\text{ext}}$  for the above four types  
452 of splitting. Four of the five validation criteria pass the threshold values, while  $Q^2_{\text{F3}}$  does not  
453 pass the threshold. In contrast, when using data by 1982 as the prediction set, the  $R^2_{\text{ext}}$  is 0.93,  
454 which is very close to the  $R^2$  of the training dataset (Section 3.2). The absolute difference  
455 between predicted and measured  $\log_{10}D$  averages 2.20 (95<sup>th</sup> percentile of 5.53) for the FDA  
456 dataset, and averages 1.08 (95<sup>th</sup> percentile of 2.68) for the data by 1982. Figure 3 presents the  
457 comparison between model predicted and experimental responses for these two prediction sets.  
458 Data from both prediction sets are generally distributed close to the 1:1 line, but the FDA data  
459 are more dispersed compared to the training set data while the data by 1982 are almost as

460 compact as the training set data. The FDA data lack documentation of experimental details, so  
461 their quality may not be as good as the data reported in peer-reviewed literature. Also, when the  
462 FDA polymer types were linked to our consolidated material types, mismatches may have  
463 occurred due to lack of description of the polymers in the FDA dataset, which may lead to  
464 inaccuracies in model predictions. Overall, however, our QPPR performs reasonably well on  
465 these two fully external datasets, demonstrating its good predictive ability.

### 466 3.3.3 Applicability domain (AD)

467 We performed the analysis of the model's applicability domain (AD) using the three approaches  
468 explained in Section 2.3.3. The model being evaluated is the final MLR model presented in Eq. 5,  
469 which was calculated using the training set of 1103 data points collected from 68 studies  
470 obtained from the peer-reviewed literature. For the analysis of AD, we focus on the two external  
471 prediction datasets: the FDA dataset (189 data points) and the data by 1982 (239 data points).  
472 Detailed results of the AD analysis are presented in SI, Section S6.1.

473 Combining the three methods, none of the data points in both prediction sets fell out of the AD.  
474 For the FDA dataset, the majority of the data points were inside the AD, while 15 data points  
475 were on borderline of AD. Similarly, only 35 data points from the data by 1982 were on  
476 borderline of AD. Thus, it is valid to use the present QPPR to make reliable estimates of  
477 diffusion coefficients for all data points in the two prediction sets. The physiochemical property  
478 space covered by the QPPR is mainly determined by the chemical's molecular weight, which  
479 ranges from 30 to 1178 g/mol. The vapor pressure at 25 °C may also be a relevant property,  
480 which ranges from  $9.8 \cdot 10^{-29}$  to  $5.2 \cdot 10^5$  Pa. The range of  $\log_{10}D$  covered by the QPPR ranges  
481 from -22.1 to -5.2 where  $D$  is measured in  $\text{m}^2/\text{s}$ .

482 As mentioned in Section 2.2.2, the model performances of using log-molecular weight and  
483 molecular weight as predictors were very close to each other when using the training dataset.  
484 However, residual analysis and external validation showed that  $\log_{10}MW$  is a more stable  
485 predictor than  $MW$  when handling high-molecular-weight chemicals, which becomes prominent  
486 for the FDA dataset which includes certain chemicals with molecular weight higher than 1500  
487 g/mol. While none of the data points in the FDA dataset fell out of the AD using the  $\log_{10}MW$   
488 model, 11 data points would be outside AD using the  $MW$  model. Details are presented in SI,  
489 Section S6.2. Thus,  $\log_{10}MW$  instead of  $MW$  was selected as a predictor in the final QPPR (Eq.  
490 5).

491 Schwope et al.<sup>37</sup> suggested that the linear relationship between  $\log_{10}D$  and  $\log_{10}MW$  may only  
492 be valid for a certain range of molecular weight, and there may be a saturation of diffusion  
493 coefficients for small molecular weights, i.e., for a given material and a given temperature, the  
494 diffusion coefficient does not continue to increase for chemicals with molecular weight lower  
495 than a certain value, which is likely determined by the material type. To further examine the  
496 effect of molecular weight on model applicability, we analyzed the model residuals versus the  
497 log of molecular weight for the training dataset and the two prediction sets (Figure 4). For the  
498 three datasets, the residuals are distributed evenly on both sides of zero in the MW range of the  
499 training dataset of 30 and 1178 g/mol ( $\log_{10}MW$  of 1.48 to 3.07). For methane (MW=16 g/mol),  
500 most of the predictions overestimate diffusivity, suggesting that diffusivity may indeed not  
501 further decrease below MW 30 g/mol. Since methane was the only chemical with data available  
502 for MW lower than 30 g/mol, data for additional chemicals and materials are therefore needed to  
503 further test this hypothesis of saturation at low MW. Similarly, additional data are needed to  
504 provide more accurate estimates for chemicals with very high molecular weights.

505 Overall, the performance of the final model (Eq. 5) in this external validation indicates that it has  
506 the ability to provide reliable predictions, as long as the considered chemicals are within the  
507 model's applicability domain. With the log-molecular weight as a predictor, our model is able to  
508 make reliable extrapolations on chemicals with molecular weights up to about 2500 g/mol, but  
509 caution still needs to be taken when applying the model on extremely-high-molecular-weight  
510 chemicals. Ideally, the model should be applied to predict diffusion coefficients for chemicals  
511 with molecular weights lower than 1178 g/mol which is the maximum within the training dataset.  
512 Caution also needs to be taken when applying the model on very-low-molecular-weight  
513 chemicals due to the possible saturation effect. Both the FDA dataset and the data by 1982 were  
514 used for the external validation but not combined with the original training dataset to calculate a  
515 more comprehensive MLR model, because these data are somewhat outdated; the FDA data are  
516 not published in literature, so there is a lack of experimental details, making these undocumented  
517 data less reliable than the data collected from peer-reviewed literature.

518

### 519 3.4 Limitations and future work

520 While the extension to 32 different consolidated material types is a major progress, the present  
521 model is still not fully comprehensive. First, the model may not be valid for very high or very

522 low molecular weight (MW) chemicals. It may not be valid for ionizing organic chemicals either,  
523 since ionizing chemicals such as acids, alcohols/phenols and amines are not well represented in  
524 the training dataset, as they only account for less than 10% of the data points, and the model does  
525 not consider chemical ionization or interaction within a material, which may make the  
526 chemical's diffusivity lower than that predicted by the model. Second, the present model is not  
527 applicable for materials types other than the 32 types in the training set, e.g. for material such as  
528 resin and textiles, due to the lack of experimental data. Although a more general MLR model (SI,  
529 Section S4) was developed which does not require material type as the predictor, it gives much  
530 less accurate predictions of the diffusion coefficient. Third, the present model does not consider  
531 any interaction between MW and material type, i.e., it assumes the effect of MW is the same  
532 across different materials. Although model validations show that this assumption may be  
533 reasonable for the existing data, ideally it needs to be further verified using data spanning the  
534 whole MW range (30 to 1178 g/mol) for each material. Therefore, more experimental diffusion  
535 coefficient data need to be obtained, or more advanced experimental methods to measure  
536 diffusion coefficients need to be developed, for other material types and chemical sizes and  
537 classes to make the model more comprehensive.

538 There are also large variations in the experimental diffusion coefficients between some of  
539 different studies for three material types, namely "MMA homopolymer", "Natural rubber" and  
540 "Rigid polymers", even after correcting for molecular weight and temperature, as shown in  
541 Figure 1 and SI, Figure S1. This means that the regression coefficients  $b$  and  $\tau$  for these material  
542 types should be taken with care. The variations could be due to three causes. First, experimental  
543 variation; for example, Franz et al.<sup>40</sup> used desorption experiments to measure the diffusion  
544 coefficients in MMA homopolymer, while Hennebert et al.<sup>42</sup> used sorption experiments. Second,  
545 the swelling of polymers during liquid sorption experiments, which generally occurs for  
546 crosslinked polymers in low-molecular weight solvents<sup>49</sup>, may not always be accounted for, and  
547 can lower the diffusion coefficients by orders of magnitude<sup>10</sup>. Third, the properties of the same  
548 material can vary between studies depending on how it was made and which additives were used.  
549 This may also be the case for some other materials such as vinyl flooring, carpet, synthetic  
550 rubber, etc., for which the material type coefficients in Eq. 5 can only represent some sort of  
551 average composition and diffusion behavior for the specific materials. Ideally, quantitative,  
552 continuous properties of the solid materials, such as density, porosity and crystalline state of the

553 material as well as other descriptors of the material's composition and molecular structure,  
554 instead of qualitative material types could be measured and entered into the model as predictors,  
555 so that the model can be more accurate and can be extrapolated to various material types outside  
556 the training dataset.

557

#### 558 4. Conclusions

559 A multiple linear regression model has been developed to predict the internal diffusion  
560 coefficients of organic compounds in various solid materials (excel model provided in SI).  
561 Experimental diffusion coefficient data collected from 68 studies of the peer-reviewed literature  
562 were used as the training set for the regression. The model uses two continuous variables,  
563 molecular weight and inversed absolute temperature, and one categorical variable, material type,  
564 as predictors. The model has been internally validated to be robust, stable and not a result of  
565 chance correlation. External validation using two prediction sets demonstrates that the model  
566 predictions are most reliable within the model's applicability domain, namely molecular weight  
567 between 30 and 1178 g/mol temperature between 4 and 180 °C, and material type belonging to  
568 the 32 consolidated types.

569 The main advantage of the present model is that it is applicable for chemicals with a wide range  
570 of molecular weights (but only up to about 16 to 2500 g/mol, with special treatment for  
571 molecular weight lower than 30 g/mol) in various materials. This is advantageous compared to  
572 the correlation methods developed in previous studies often specific for certain chemical classes  
573 or materials. The present model is able to provide reliable estimates of diffusion coefficients for  
574 a large number of chemical-material combinations, making it suitable for high-throughput  
575 assessments of the releases and human exposures to chemicals encapsulated in solid materials,  
576 particularly building materials and food contact materials. To make the model comprehensive,  
577 more experimental diffusion coefficient data need to be obtained for other material types, or  
578 quantitative and continuous parametrization of various solid materials needs to be further  
579 developed.

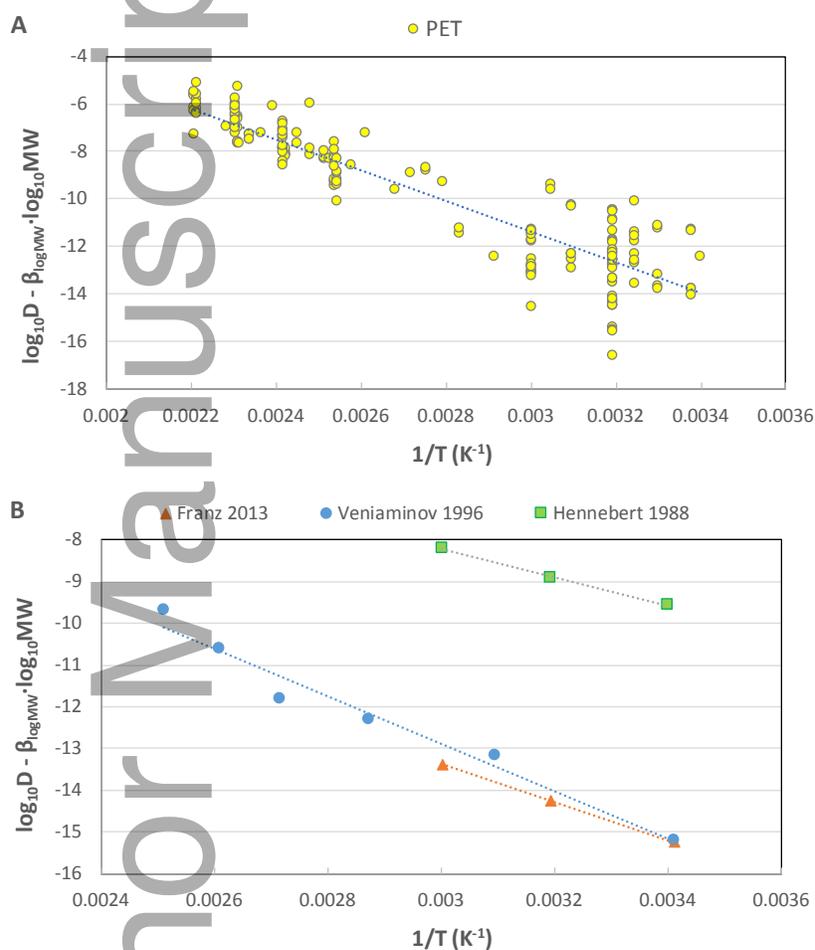
580

#### 581 Acknowledgements

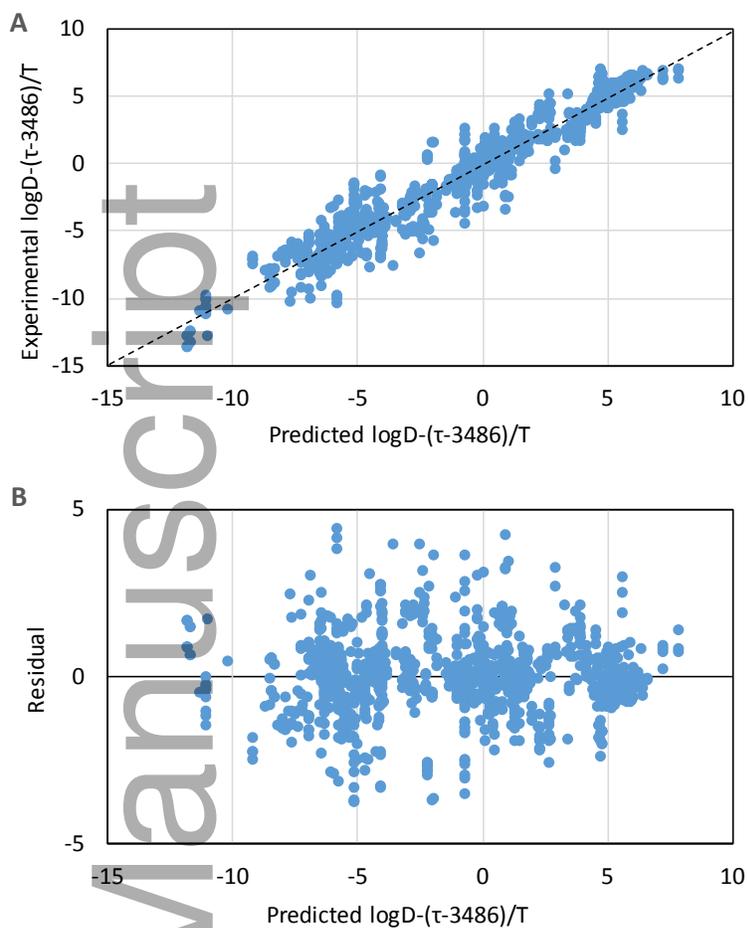
582 The authors thank Prof. Ester Papa, Dr Alessandro Sangion, and Prof. Paola Gramatica from the  
583 University of Insubria, Italy for advice on MLR modeling and validation, as well as support for

584 the QSARINS software. Funding for this research was provided by US EPA contract EP-16-C-  
585 000070 and by the Long Range Research Initiative of the American Chemistry Council. P.  
586 Fantke was supported by the Marie Curie project Quan-Tox (GA No. 631910) funded by the  
587 European Commission under the Seventh Framework Programme.

588 Tables and Figures

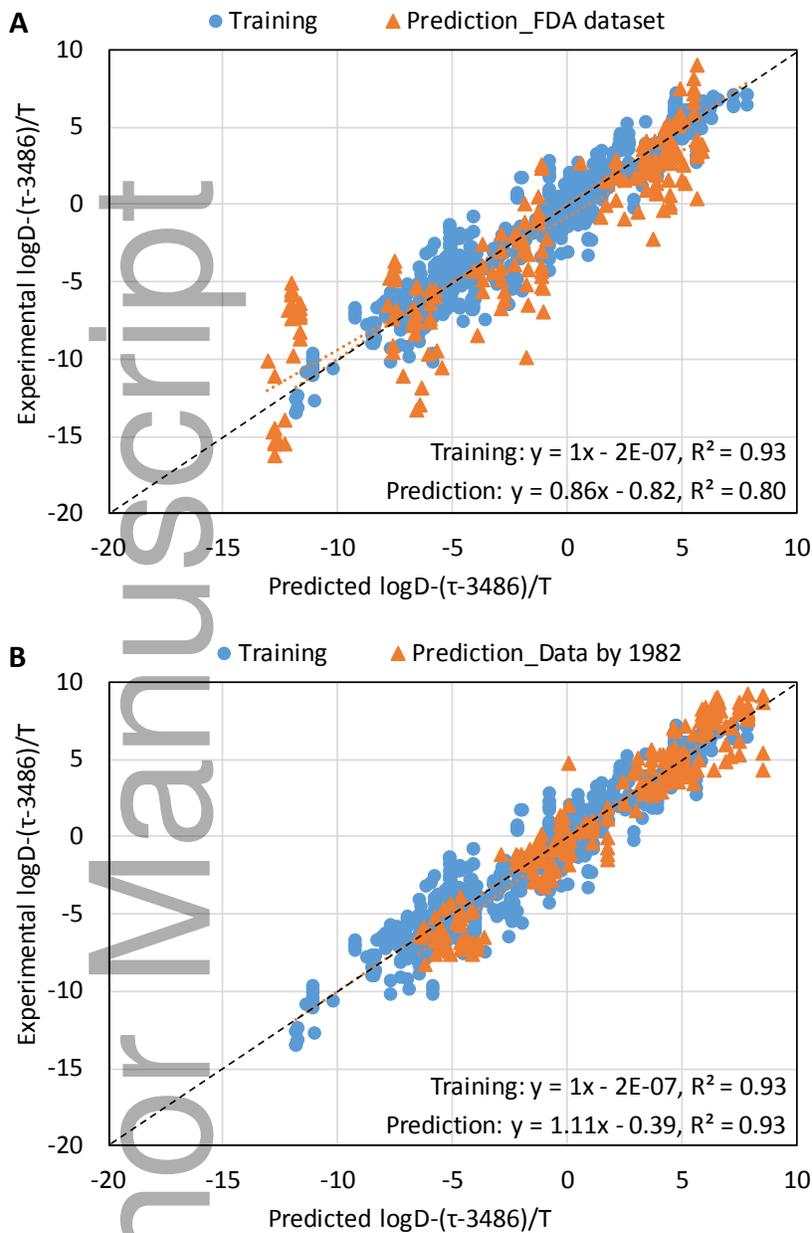


589  
590 Figure 1. Relationship between the diffusion coefficient  $D$  (corrected for  $\log_{10} MW$ ) and the  
591 inverse of temperature for (A) PET, and (B) methyl methacrylate (MMA) homopolymer. The  
592 units of  $D$  and  $MW$  are  $m^2/s$  and  $g/mol$ , respectively.



593

594 Figure 2. Values of  $\log_{10} D - (\tau - 3486)/T$  predicted by the final QPPR (Eq. 5) vs. (A) experimental  
 595 values, and (B) residuals. The dotted line in (A) indicates the 1:1 line. The units of  $D$  and  $T$  are  
 596  $\text{m}^2/\text{s}$  and  $\text{K}$ , respectively.



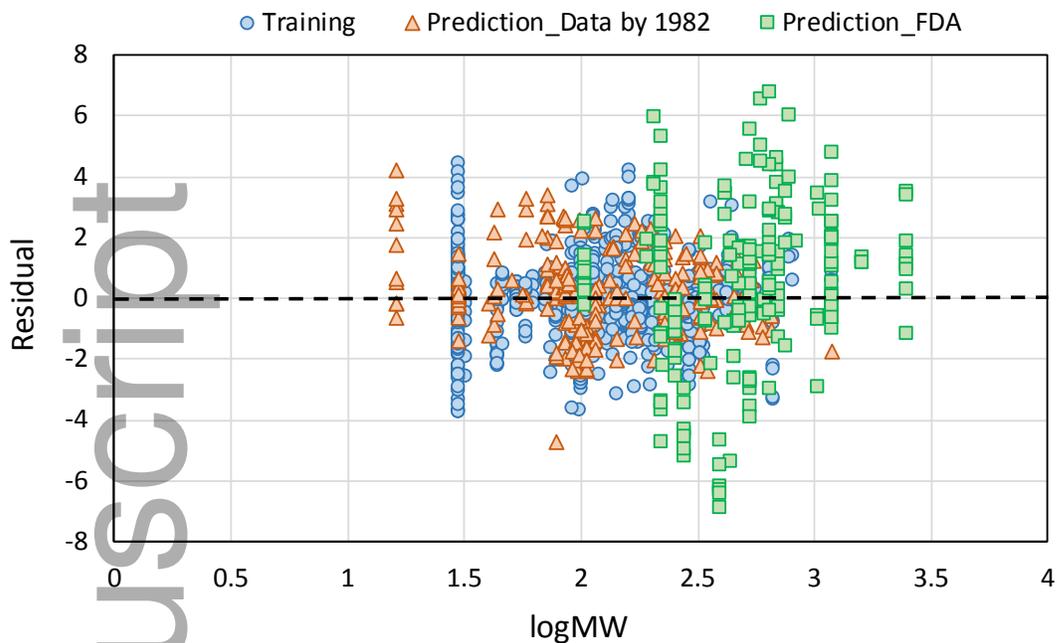
597

598 Figure 3. Values of  $\log_{10}D-(\tau-3486)/T$  predicted by the final QPPR (Eq. 5) vs. experimental  
 599 values when using (A) FDA dataset and (B) Data by 1982 as the prediction sets. The black  
 600 dotted line indicates the 1:1 line. The units of  $D$  and  $T$  are  $m^2/s$  and  $K$ , respectively.

601

602

603



604  
 605 Figure 4. Residual between the present QPPR and observed data as a function of  $\log_{10}MW$  for  
 606 the training dataset, the FDA dataset, and the data by 1982 set. The unit of MW is g/mol.

607  
 608  
 609  
 610  
 611  
 612  
 613  
 614  
 615  
 616  
 617  
 618 Table 1. Temperature dependence of diffusion coefficient in the 32 consolidated material types  
 619 (all numbers are in the unit of K)

Category	Material	Mean coefficient of 1/T	SD between studies	Coefficient value for Eq. 5		
				$\beta_{1/\pi}$	$\tau$	$\beta_{1/\pi} + \tau$
High-coefficient category	PP homopolymer	-6665	2354	-3486	-2391	-5877
	Polyethylene terephthalate (PET)	-6567	2399			
	General polystyrene (PS)	-5713	3560			
	Polyethylene naphthalate (PEN)	-5449	1940			
	PP copolymer	-5384	1194			
	High-density polyethylene (HDPE)	-5294	1124			
Medium-coefficient category	MMA homopolymer	-4549	1145	-3486	0	-3486
	ABS, EVOH	-4222	n/a			
	High-impact polystyrene (HIPS)	-4215	n/a			
	Polyamide (PA)	-4179	1854			
	MMA copolymer-medium or low density	-4056	1272			
	Polyethylene (PE, LDPE, LLDPE)	-3713	536			
	Limited-data material group	n/a	n/a			
	Calcium silicate	n/a	n/a			
	Carpet	n/a	n/a			
	Glass, Stainless steel	n/a	n/a			
	Vinyl acetate-based polymers	n/a	n/a			
	Cement	n/a	n/a			
	Low-coefficient category	Gypsum board	n/a			
Plywood		n/a	n/a			
Flexible PVC		-2917	2618			
Other wooden boards		-2411	888			
Polychloroprene (CR)		-2127	286			
Vinyl flooring		-1951	n/a			
Polystyrene foam (XPS, EPS)		-1806	n/a			
Polyurethane foam-based materials*		-1705	699			
Synthetic rubber		-1326	205			
Ethylene-propylene rubbers		-1145	300			
Natural rubber (NR)		-939	337			
Rigid polymers		-510	1552			
Paper		-312	n/a			
Gypsum and cellulose ceiling tile		331	294			

\*This material type refers to low-density polyurethane foams with a density of 0.005 to 0.03 g/cm<sup>3</sup>.

620  
621  
622  
623  
624  
625

Table 2. Material-specific coefficients for Eq. 5

Material	Coefficient <i>b</i>			$\tau$ (K)	$b + (\tau + 2391.15)/T$ at 25 °C
	Coefficient	SE <sup>f</sup>	p-value		
Calcium silicate	1.17	0.29	< 0.0001	0	9.19
Carpet	-1.23	0.28	< 0.0001	0	6.79
Cement	0.330	0.226	0.15	0	8.35
Ethylene-propylene rubbers	-6.32	0.29	< 0.0001	1676	7.32
Flexible PVC	-8.51	0.31	< 0.0001	1676	5.13
General polystyrene (PS)	2.04	0.30	< 0.0001	-2391	2.04
Glass, Stainless steel	-8.57	0.38	< 0.0001	0	-0.550
Gypsum and cellulose ceiling tile	-1.24	0.31	< 0.0001	1676	12.4
Gypsum board	-5.77	0.30	< 0.0001	1676	7.87
High density polyethylene (HDPE)	5.11	0.20	< 0.0001	-2391	5.11
High-impact polystyrene (HIPS)	-7.11	0.27	< 0.0001	0	0.907
Methyl methacrylate (MMA) copolymer-medium or low density	-7.73	0.21	< 0.0001	0	0.294
Methyl methacrylate (MMA) homopolymer <sup>h</sup>	-7.84	0.31	< 0.0001	0	0.175
Natural rubber (NR) <sup>h</sup>	-3.60	0.27	< 0.0001	1676	10.0
Other wooden boards <sup>a</sup>	-6.72	0.21	< 0.0001	1676	6.92
Paper	-8.53	0.34	< 0.0001	1676	5.11
Plywood	-5.61	0.34	< 0.0001	1676	8.03
Polyamide (PA)	-5.40	0.16	< 0.0001	0	2.62
Polyacrylonitrile butadiene styrene (ABS), Ethylene vinyl alcohol (EVOH)	-4.97	0.23	< 0.0001	0	3.05
Polychloroprene (CR)	-6.31	0.35	< 0.0001	1676	7.33
Polyethylene (PE, LDPE, LLDPE)	-1.65	0.16	< 0.0001	0	6.37
Polyethylene naphthalate (PEN)	-1.16	0.28	< 0.0001	-2391	-1.16
<b>Polyethylene terephthalate (PET)<sup>g</sup></b>	<b>0.00</b>	<b>0.15</b>	<b>n/a</b>	-2391	0.00
Polystyrene foam (XPS, EPS)	-8.32	0.29	< 0.0001	1676	5.32
Polyurethane foam-based materials <sup>b</sup>	-7.35	0.25	< 0.0001	1676	6.30
PP copolymer	4.79	0.28	< 0.0001	-2391	4.79
PP homopolymer	4.53	0.15	< 0.0001	-2391	4.53
Rigid polymers <sup>c,h</sup>	-11.9	0.25	< 0.0001	1676	1.70
Synthetic rubber	-5.93	0.32	< 0.0001	1676	7.71
Vinyl acetate-based polymers <sup>d</sup>	-0.459	0.326	0.16	0	7.56
Vinyl flooring	-6.77	0.21	< 0.0001	1676	6.87
Limited-data material group <sup>e</sup>			see footnotes		

<sup>a</sup> Includes Particleboard, Oriented strand board (OSB), Medium-density fiberboard (MDF), High-density board, and Wood chamber wall.

<sup>b</sup> This material type refers to low-density polyurethane foams with a density of 0.005 to 0.03 g/cm<sup>3</sup>.

<sup>c</sup> Includes Polyether ether ketone (PEEK), Rigid PVC, Polytetrafluoroethylene (PTFE), and Polycarbonate.

<sup>d</sup> Includes Ethyl vinyl acetate (EVA), Polyvinyl acetate (PVA), and Polyvinyl acetate polyacrylic acid copolymer.

<sup>e</sup> The coefficient *b* for this group is -2.26 with an SE of 0.18, and the coefficient  $\tau$  is 0. "Limited-data material group" includes data from 20 different materials, so the accuracy of the coefficients is low and they are not recommended for use in predicting diffusion coefficients. This group includes Alginate film, Balance, Decorative and Overlay layers of wooden flooring, Cellulose, Epichlorhydrin-dimethylamine polymer (EDP), Epoxy/acrylic copolymer, latex, MMA/Butyl methacrylic (BMA) copolymer -very low density, Nanocomposite polyamide, Paint, Pectin film, Pectin/Alginate composite film, Polydimethylsiloxane (PDMS) membrane, Polyisoprene (PI) membrane, Polyoctenamer (PO) membrane, Polyoxymethylene, Polytrimethylene terephthalate (PTT), Polyvinylidene chloride (PVDC), and Silicone.

<sup>f</sup> Standard error.

<sup>g</sup> Reference material.

<sup>h</sup> Coefficients should be taken with care due to large variations between studies.

626

627 Table 3. External validation results

External validation criteria	$R^2_{ext}$	$Q^2_{F1}$	$Q^2_{F2}$	$Q^2_{F3}$	$\overline{r^2}_m$	CCC
Threshold		> 0.70	> 0.70	> 0.70	> 0.65	> 0.85
Splitting by random percentage	0.92	0.92	0.92	0.92	0.90	0.96
Splitting by ordered response	0.94	0.94	0.94	0.95	0.93	0.97
Splitting by ordered structure	0.94	0.94	0.94	0.94	0.91	0.97
Splitting by studies	0.85	0.85	0.84	0.85	0.78	0.92
FDA dataset as prediction set	0.80	0.77	0.77	0.60	0.71	0.89
Data by 1982 as prediction set	0.93	0.93	0.92	0.90	0.85	0.95

$R^2_{ext}$ : determination coefficient of the prediction set external data.

$Q^2_{F1}$ : correlation coefficient proposed by Shi et al.

$Q^2_{F2}$ : correlation coefficient proposed by Schuurmann et al.

$Q^2_{F3}$ : correlation coefficient proposed by Consonni et al.

$\overline{r^2}_m$ : determination coefficient proposed by Ojha et al.

CCC: concordance correlation coefficient proposed by Chirico and Gramatica.

628

629

## 630 References

- 631 1. Little JC, Weschler CJ, Nazaroff WW, et al. Rapid methods to estimate potential exposure to  
632 semivolatile organic compounds in the indoor environment. *Environ Sci Technol.* 2012; 46(20): p.  
633 11171-11178.
- 634 2. Xu Y, Cohen Hubal EA, Clausen PA, et al. Predicting residential exposure to phthalate plasticizer  
635 emitted from vinyl flooring: a mechanistic analysis. *Environ Sci Technol.* 2009; 43(7): p. 2374-  
636 2380.
- 637 3. Guo Z. Review of indoor emission source models. Part 1. Overview. *Environ Pollut.* 2002; 120(3):  
638 p. 533-549.
- 639 4. Begley T, Castle L, Feigenbaum A, et al. Evaluation of migration models that might be used in  
640 support of regulations for food-contact plastics. *Food Addit Contam.* 2005; 22(1): p. 73-90.
- 641 5. Xie M, Wu Y, Little JC, et al. Phthalates and alternative plasticizers and potential for contact  
642 exposure from children's backpacks and toys. *J Expo Sci Env Epid.* 2016; (26): p. 119-124.
- 643 6. Liu Z, Ye W, Little JC. Predicting emissions of volatile and semivolatile organic compounds from  
644 building materials: a review. *Build Environ.* 2013; 64: p. 7-25.
- 645 7. Haghghat F, Huang H, Lee C-S. Modeling approaches for indoor air VOC emissions from dry  
646 building materials—a review. *ASHRAE Trans.* 2005; 111(1): p. 635-645.

- 647 8. Berens A and Hopfenberg H. Diffusion of organic vapors at low concentrations in glassy PVC,  
648 polystyrene, and PMMA. *J Membrane Sci.* 1982; 10(2-3): p. 283-303.
- 649 9. Hickey AS and Peppas NA. Solute diffusion in poly (vinyl alcohol)/poly (acrylic acid) composite  
650 membranes prepared by freezing/thawing techniques. *Polymer.* 1997; 38(24): p. 5931-5936.
- 651 10. John J, Kunchandy S, Kumar A, et al. Transport of methyl methacrylate monomer through  
652 natural rubber. *J Mater Sci.* 2010; 45(2): p. 409-417.
- 653 11. Luo R and Niu J. Determining diffusion and partition coefficients of VOCs in cement using one  
654 FLEC. *Build Environ.* 2006; 41(9): p. 1148-1160.
- 655 12. Bodalal A, Zhang J, Plett E, et al. Correlations between the internal diffusion and equilibrium  
656 partition coefficients of volatile organic compounds (VOCs) in building materials and the VOC  
657 properties. *ASHRAE Trans.* 2001; 107: p. 789.
- 658 13. Bodalal A, Zhang J, Plett E. A method for measuring internal diffusion and equilibrium partition  
659 coefficients of volatile organic compounds for building materials. *Build Environ.* 2000; 35(2): p.  
660 101-110.
- 661 14. Little JC, Hodgson AT, Gadgil AJ. Modeling emissions of volatile organic compounds from new  
662 carpets. *Atmos Environ.* 1994; 28(2): p. 227-234.
- 663 15. Dole P, Feigenbaum AE, Cruz CDL, et al. Typical diffusion behaviour in packaging polymers—  
664 application to functional barriers. *Food Addit Contam.* 2006; 23(2): p. 202-211.
- 665 16. Reynier A, Dole P, Humbel S, et al. Diffusion coefficients of additives in polymers. I. Correlation  
666 with geometric parameters. *J Appl Polym Sci.* 2001; 82(10): p. 2422-2433.
- 667 17. Guo Z. Review of indoor emission source models. Part 2. Parameter estimation. *Environ Pollut.*  
668 2002; 120(3): p. 551-564.
- 669 18. Cox SS, Zhao D, Little JC. Measuring partition and diffusion coefficients for volatile organic  
670 compounds in vinyl flooring. *Atmos Environ.* 2001; 35(22): p. 3823-3830.
- 671 19. Zhao D, Cox S, Little J. *Source/sink characterization of diffusion controlled building materials.* in  
672 *Proceedings of the 8th International Conference on Indoor Air Quality and Climate-Indoor Air.*  
673 1999.
- 674 20. Jolliet O, Ernstoff AS, Csiszar SA, et al. Defining Product Intake Fraction to Quantify and Compare  
675 Exposure to Consumer Products. *Environ Sci Technol.* 2015; 49: p. 8924-8931.
- 676 21. Shin H-M, Ernstoff A, Arnot JA, et al. Risk-based high-throughput chemical screening and  
677 prioritization using exposure models and in vitro bioactivity assays. *Environ Sci Technol.* 2015;  
678 49(11): p. 6760-6771.

- 679 22. Shin H-M, McKone TE, Bennett DH. Intake fraction for the indoor environment: a tool for  
680 prioritizing indoor chemical sources. *Environ Sci Technol.* 2012; 46(18): p. 10063-10072.
- 681 23. Ernstoff AS, Fantke P, Csiszar SA, et al. Multi-pathway exposure modelling of chemicals in  
682 cosmetics with application to shampoo. *Environ Int.* 2016; 92-93: p. 87-96.
- 683 24. Csiszar SA, Ernstoff AS, Fantke P, et al. Stochastic modeling of near-field exposure to parabens in  
684 personal care products. *J Expo Sci Env Epid.* 2017; (27): p. 152-159.
- 685 25. Egeghy PP, Sheldon LS, Isaacs KK, et al. Computational exposure science: An emerging discipline  
686 to support 21st-century risk assessment. *Environ Health Persp.* 2016; 124(6): p. 697.
- 687 26. Deng Q, Yang X, Zhang J. Study on a new correlation between diffusion coefficient and  
688 temperature in porous building materials. *Atmos Environ.* 2009; 43(12): p. 2080-2083.
- 689 27. Xu J and Zhang JS. An experimental study of relative humidity effect on VOCs' effective diffusion  
690 coefficient and partition coefficient in a porous medium. *Build Environ.* 2011; 46(9): p. 1785-  
691 1796.
- 692 28. Xu J, Zhang JS, Liu X, et al. Determination of partition and diffusion coefficients of formaldehyde  
693 in selected building materials and impact of relative humidity. *J Air Waste Ma.* 2012; 62(6): p.  
694 671-679.
- 695 29. Park J-S, Little JC, Kim S-D, et al. The Determination of Diffusion and Partition Coefficients of PUF.  
696 *J Korean Soc Atmos Environ.* 2010; 26(1): p. 77-84.
- 697 30. Welle F and Franz R. Diffusion coefficients and activation energies of diffusion of low molecular  
698 weight migrants in Poly(ethylene terephthalate) bottles. *Polym Test.* 2012; 31(1): p. 93-101.
- 699 31. Ewender J and Welle F. Determination of the activation energies of diffusion of organic  
700 molecules in poly (ethylene terephthalate). *J Appl Polym Sci.* 2013; 128(6): p. 3885-3892.
- 701 32. Gramatica P, Cassani S, Chirico N. QSARINS - chem: Insubria datasets and new QSAR/QSPR  
702 models for environmental pollutants in QSARINS. *J Comp Chem.* 2014; 35(13): p. 1036-1044.
- 703 33. Gramatica P, Chirico N, Papa E, et al. QSARINS: A new software for the development, analysis,  
704 and validation of QSAR MLR models. *J Comp Chem.* 2013; 34(24): p. 2121-2132.
- 705 34. Flynn JH. A collection of kinetic data for the diffusion of organic compounds in polyolefins.  
706 *Polymer.* 1982; 23(9): p. 1325-1344.
- 707 35. Park G. The diffusion of some organic substances in polystyrene. *Transactions of the Faraday*  
708 *Society.* 1951; 47: p. 1007-1013.
- 709 36. Park G. The diffusion of some halo-methanes in polystyrene. *Transactions of the Faraday Society.*  
710 1950; 46: p. 684-697.

- 711 37. Schwope A, Goydan R, Reid R, *Methods for assessing exposure to chemical substances Volume*  
712 *11: Methodology for Estimating the Migration of Additives and Impurities from Polymeric*  
713 *Materials* 1990, U.S.EPA: Washington, D.C.
- 714 38. Gramatica P. Principles of QSAR models validation: internal and external. *QSAR & combinatorial*  
715 *science*. 2007; 26(5): p. 694-701.
- 716 39. Cassani S and Gramatica P. Identification of potential PBT behavior of personal care products by  
717 structural approaches. *Sustain Chem Pharm*. 2015; 1: p. 19-27.
- 718 40. Franz R and Brandsch R. Migration of acrylic monomers from methacrylate polymers–  
719 establishing parameters for migration modelling. *Packag Technol Sci*. 2013; 26(8): p. 435-451.
- 720 41. Veniaminov A and Sedunov YN. Diffusion of phenanthrenequinone in poly (methyl  
721 methacrylate): holographic measurements. *Polym Sci Ser A*. 1996; 38: p. 59-63.
- 722 42. Hennebert P. Solubility and diffusion coefficients of gaseous formaldehyde in polymers.  
723 *Biomaterials*. 1988; 9(2): p. 162-167.
- 724 43. Shi LM, Fang H, Tong W, et al. QSAR models using a large diverse set of estrogens. *J Chem Inf*  
725 *Comp Sci*. 2001; 41(1): p. 186-195.
- 726 44. Schüürmann G, Ebert R-U, Chen J, et al. External validation and prediction employing the  
727 predictive squared correlation coefficient - Test set activity mean vs training set activity mean. *J*  
728 *Chem Inf Model*. 2008; 48(11): p. 2140-2145.
- 729 45. Consonni V, Ballabio D, Todeschini R. Comments on the definition of the  $Q^2$  parameter for QSAR  
730 validation. *J Chem Inf Model*. 2009; 49(7): p. 1669-1678.
- 731 46. Ojha PK, Mitra I, Das RN, et al. Further exploring  $r_m^2$  metrics for validation of QSPR models.  
732 *Chemometr Intell Lab*. 2011; 107(1): p. 194-205.
- 733 47. Chirico N and Gramatica P. Real external predictivity of QSAR models: how to evaluate it?  
734 Comparison of different validation criteria and proposal of using the concordance correlation  
735 coefficient. *J Chem Inf Model*. 2011; 51(9): p. 2320-2335.
- 736 48. Chirico N and Gramatica P. Real external predictivity of QSAR models. Part 2. New  
737 intercomparable thresholds for different validation criteria and the need for scatter plot  
738 inspection. *J Chem Inf Model*. 2012; 52(8): p. 2044-2058.
- 739 49. Nandi S and Winter HH. Swelling behavior of partially cross-linked polymers: a ternary system.  
740 *Macromolecules*. 2005; 38(10): p. 4447-4455.
- 741
- 742