# Chemical Similarity to Identify Potential Substances of Very High Concern – an Effective Screening Method

Pim N.H. Wassenaar[1,2] *, Emiel Rorije[1], Nicole M.H. Janssen[1], Willie J.G.M. Peijnenburg[1,2], Martina G. Vijver[2]

[1] *National Institute for Public Health and the Environment (RIVM), Centre for Safety of Substances and Products, P.O. Box 1, 3720 BA Bilthoven, The Netherlands*

[2] *Institute of Environmental Sciences (CML), Leiden University, P. O. Box 9518, 2300 RA Leiden, The Netherlands*

* <u>Corresponding author:</u> Pim N.H. Wassenaar. <u>Email:</u> pim.wassenaar@RIVM.nl. <u>Address:</u> National Institute for Public Health and the Environment (RIVM), Centre for Safety of Substances and Products, P.O. Box 1, 3720 BA Bilthoven, The Netherlands.
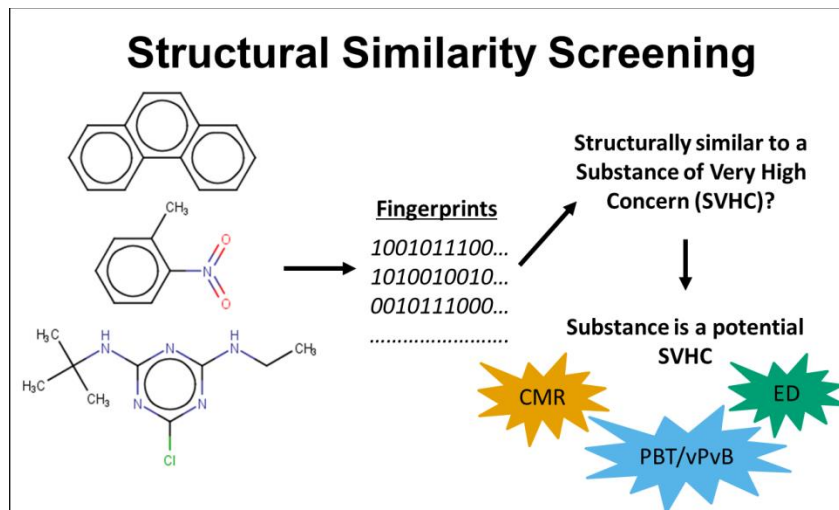
<u>Abbreviations</u>[1]

---

[1] *SVHC = Substances of Very High Concern ; CMR = Carcinogenic, Mutagenic or Reprotoxic ; PBT = Persistent, Bioaccumulative and Toxic ; vPvB = very Persistent and very Bioaccumulative ; ED = Endocrine Disruption; SPOKs = Single Point of Knowledge structures ; $K_{ow}$ = octanol/water partition coefficient ; UVCB = Substances of Unknown or Variable composition, Complex reaction products or Biological materials ; ECFP = Extended Connectivity Fingerprints ; FCFP = Functional-Class Fingerprints ; JT = Jaccard-Tanimoto coefficient ; HL = Harris-Lahey coefficient ; CT4 = Consonni-Todeschini 4 coefficient ; SS3 = Sokal-Sneath 3 coefficient ; Coh = Cohen coefficient ; SM = Simple Matching coefficient ; Yu2 = Yule 2 coefficient; TP = True Positives ; FP = False Positives ; FN = False Negatives ; TN = True Negatives.*

19     **Graphical Abstract**



20

21

22     **Highlights**

23     • Potential Substances of Very High Concern can be identified by chemical similarity.

24     • High balanced accuracies (≥0.8) were obtained for all SVHC-subgroup models.

25     • Improvement of the ED model by extending the database is considered necessary.

26     • The best performing similarity models can be used for screening and prioritization.

27

28  **Abstract**

29  There is a strong demand for early stage identification of potential substances of very high concern

30  (SVHC). SVHCs are substances that are classified as carcinogenic, mutagenic or reprotoxic (CMR);

31  persistent, bioaccumulative and toxic (PBT) or very persistent and very bioaccumulative (vPvB); or as

32  substances with an equivalent level of concern, like endocrine disruption (ED). The endeavor to improve

33  the identification of potential SVHCs is also acknowledged by the European Commission, in their long-

34  term vision towards a non-toxic environment. However, it has been shown difficult to identify substances

35  as potentially harmful.

36      With this goal in mind, we have developed a methodology that predicts whether a substance is a

37  potential SVHC based on chemical similarity to chemicals already identified as SVHC. The approach is

38  based on the structural property principle, which states that structurally similar chemicals are likely to

39  have similar properties.

40      We systematically analyzed the predictive performance of 112 similarity measures (i.e. all

41  different combinations of 16 binary fingerprints and 7 similarity coefficients) classifying the substances in

42  the dataset as (potential) SVHC or non-SVHC. The outcomes were analyzed for 546 substances that we

43  collected within the Dutch SVHC database – with identified CMR, PBT/vPvB and/or ED properties - and

44  411 substances that lack these hazardous properties. The best similarity measures showed a high

45  predictive performance with a balanced accuracy of 85% correct identifications for the whole dataset of

46  SVHC substances, and 80% for CMR, 95% for PBT/vPvB and 99% for ED subgroups.

47      This effective screening methodology showed great potential for early stage identification of potential

48  SVHCs. This model can be applied within regulatory frameworks and safe-by-design trajectories, and

49  hence can contribute to the EU goal of achieving a non-toxic environment.

50  **<u>Keywords:</u>** Substances of Very High Concern, Screening, Chemical similarity, Classification model.

## 1. Introduction

In recent decades, exposure to specific chemicals appeared of greater concern than previously anticipated, including concerns for polychlorinated biphenyls (PCBs), dichlorodiphenyltrichloroethane (DDT) and perfluorooctanesulfonic acid (PFOS) [1]. In many cases, when safety concerns are raised, widespread exposure has often already occurred, and typically the set of available toxicity data is inadequate to introduce risk management measures immediately. Consequently, chemicals of potential concern continue to be emitted, with the risk of significant effects on human and environmental health in the long-term. Therefore, it is important to signal emerging concerns and improve the early stage identification of hazardous chemicals before widespread exposure occurs. This endeavor is also acknowledged by the European Commission in their long-term vision towards a non-toxic environment [2,3]. In particular, high priority is given to so-called substances of very high concern (SVHC), which include substances with carcinogenic, mutagenic or reprotoxic (CMR) properties, substances with persistent, bioaccumulative and toxic (PBT) or very persistent and very bioaccumulative (vPvB) properties, or substances with endocrine disrupting (ED) properties [4]. Substances can be identified as SVHC following a regulatory decision process in which all available data is evaluated.

To improve the identification of potential SVHCs, it is essential to make efficient use of the limited amount of available (fate and toxicity) data. Several models have been described in the literature that predict hazard properties of chemicals from simple properties, like aquatic toxicity based on the octanol/water partition coefficient ($K_{ow}$) and/or structural alerts [5–7], or based on more complex algorithms [8–13]. Many of these models are (at least partially) based on the structural property principle, which assumes that (structurally) similar chemicals are likely to have similar properties [14]. Although these models are very useful to predict the effect of a chemical on a specific endpoint, their applicability to identify potential SVHC substances is limited. This is a consequence of the fact that the group of SVHC substances covers a broad range of different toxicological endpoints and mode of actions - and are only identified following a regulatory decision process. Within current models it is difficult to simulate

76  such a regulatory weight-of-evidence approach. Potentially, total chemical similarity to known SVHC

77  substances can be a useful way to estimate (potential) SVHC status, as such a method might be able to

78  cover more information on SVHC identification properties.

79      To our knowledge, only two models, both with the aim of prioritization, attempt to identify

80  potential SVHCs directly based on structural similarity to substances already identified as being SVHCs,

81  including the SINimilarity tool developed by ChemSec [15], and screening scenarios as applied by the

82  European Chemical Agency (ECHA) within the SVHC Roadmap program [16]. However, these methods

83  do not provide optimized and cross-validated methodologies, resulting in an unknown predictive

84  performance. If a high predictive accuracy could be achieved using only chemical similarity information,

85  the lack of toxicity information can be bypassed, and those substances of potential SVHC concern, that

86  are currently deemed "safe" in the absence of toxicity information, can be prioritized for further follow-up

87  action. In addition, the chemical similarity information also provides a clear follow-up direction, as the

88  potential concern is directly related to the concern of the most similar SVHC substance.

89      The aim of the present study was to evaluate the efficiency of a broad set of similarity measures

90  for the identification of potential SVHCs, with a specific focus on separately identifying CMR,

91  PBT/vPvB and ED concerns. We built upon the knowledge gained (see e.g. [17]) for calculating chemical

92  similarity, that generally consists of two main elements: a descriptor (or representation) of the chemical

93  structure and a similarity coefficient. First, descriptors are used to characterize the molecules that are

94  compared by assigning numerical values to structures [17–19]. These values are in most methods related

95  to the absence or presence of specific chemical substructures and are often encoded in fixed-length bit-

96  strings (consisting of zeros and ones) [20]. These bit-strings are also known as fingerprints. Secondly,

97  similarity coefficients are used to quantitatively express the similarity between two chemical descriptors

98  [17,19,21]. For our purpose, the similarity between two fingerprints can be used to quantify the structural

99  overlap between a chemical with unknown hazardous properties and known SVHCs. Many types of

100 descriptors and similarity coefficients are available and there is no similarity measure that consistently is

101 most effective (i.e. there is no single best "fingerprint - coefficient" combination for all applications)

102  [17,20,22]. Our study outcome provides the most optimal set of similarity measures as a first screening

103  model to identify substances of potential SVHC concern.

## 2. Methods

The study approach consists of four general steps (Figure 1). First, a dataset of substances with and without CMR, PBT/vPvB and/or ED properties was constructed (paragraph 2.1). Secondly, binary fingerprints were generated for all substances in the datasets (paragraph 2.2). Thirdly, similarity values (i.e. quantitative values of chemical similarity) were calculated between substances by comparing the fingerprints with similarity coefficients (paragraph 2.3). Only the extent of similarity to substances with identified CMR, PBT/vPvB and/or ED properties leading to the SVHC status was investigated. Finally, we determined an optimal similarity threshold and the predictive performance of each "fingerprint - coefficient" combination (paragraph 2.4). Steps two to four were reiterated for multiple "fingerprint - coefficient" combinations, as well as for different SVHC subgroups (i.e. for CMR, PBT/vPvB and ED separately and together), in order to identify the optimal model(s) based on balanced accuracy. A more elaborate description of these steps is provided in the following paragraphs.

*2.1 Dataset*

In order to identify chemicals of (potential) concern based on structural similarity to known toxicants, a set of known CMR, PBT/vPvB and ED substances is required. For this purpose, a Dutch list of substances of very high concern  was selected, as all substance on this list have CMR, PBT/vPvB and/or ED properties (see [23]; extracted on 01-03-2018). This list covers a broader range of chemicals than the EU-SVHC list under REACH, but are identified based on the same hazard criteria as the EU-SVHC substances (i.e. REACH article 57 [4]). The generation and composition of this list of substances is more elaborately described in Supplemental Material S.1.

In addition, for modelling purposes we also compiled a list of substances that are known not to have CMR, PBT/vPvB and/or ED properties. All substances on the REACH Annex IV - which lists chemicals that are considered to be inherently safe - were selected for this purpose, as well as all approved biocides and pesticides (see [24,25]; extracted on 23-05-2018). The list of biocides and

129  pesticides is suited for our purpose as all substances approved for introduction on the European market

130  have been tested experimentally and are negative for CMR, PBT/vPvB and ED endpoints, according to

131  the SVHC criteria.

132  Several adjustments were made to the compiled substance lists, as chemical similarity searches

133  require a specific and unambiguous chemical structure as input information. In cases that a group of

134  substances was included in one of the above-mentioned lists (e.g. polychlorinated naphthalenes),

135  representative chemical structures were generated and selected for inclusion in order to ensure that the

136  structures represent the varying types of branching and/or substituents (e.g. tri- up till octachloro

137  naphthalene, with two isomers per chlorine-atom count). When a substance is a mixture or a UVCB

138  (Substances of Unknown or Variable composition, Complex reaction products or Biological materials),

139  only the (representative) chemical structures of those components causing the concern were included (e.g.

140  benzene in some of the UVCBs). When a substance is considered a non-SVHC substance, the main

141  constituent(s) were included. Each unique chemical structure was included once in the final list. In

142  addition, specific metal-complexes (i.e. based on arsenic, beryllium, cadmium, chromium, lead, mercury,

143  nickel and cobalt) and fibers were excluded. For these metal-based complexes, it is generally the metal

144  atom causing the concern, irrespective of the organic counterparts. In case of fibers, the toxicity is (also)

145  determined by physical aspects other than their chemical structure (e.g. diameter, length and shape). In

146  addition, all inorganic substances were removed from the list of non-SVHC substances.

147  In total, a dataset of 546 SVHC and 411 non-SVHC single chemical structures was compiled (see

148  Supplemental Material Excel). Of the 546 SVHC substances, 306 are known to have CMR properties,

149  209 to have PBT/vPvB properties, and 52 are known to have ED properties. All chemical structures were

150  represented by a (single) SMILES code [26] and all charged structures were converted to their neutral

151  counterparts, where possible (Supplemental Material S.2). These SMILES codes were used for the

152  analyses.

153

154  *2.2 Fingerprints*

155    We restricted this study to binary fingerprints based on 2D-fragments, as they tend to be more selective

156    than whole molecule descriptors. Moreover, 2D-fragments descriptors are (computationally) easier to

157    handle than 3D-fragment descriptors [17]. The fingerprints were selected in such a way to ensure

158    maximum diversity and include dictionary-based, path-based, circular-based and pharmacophore-based

159    fingerprints (Table 1) [27]. The fingerprints were generated using freely available resources, including the

160    software packages RDkit and PaDEL-Descriptor (based on the Chemistry Development Kit (CDK)

161    libraries) [28,29]. For all non-dictionary based fingerprints, a string length of 1024 bits was used. More

162    details on the generation of the fingerprints are given in Supplemental Material S.3.

163    *2.3 Similarity coefficients*

164    The similarity between two 2D-binary fingerprints of known SVHCs and non-SVHC substances can be

165    computed by using various formulas, the so-called similarity coefficients. When comparing two binary

166    fingerprints, four different bit-combinations could be identified - denoted as *a*, *b*, *c* and *d*. *A*, *b*, *c* and *d*

167    represent the counts that a feature is present in one structure and absent in the other ("x=1 and y=0"),

168    absent in the first and present in the second structure ("x=0 and y=1"), present in both ("x=1 and y=1")

169    and absent in both ("x=0 and y=0"), respectively. These four numbers are combined in similarity

170    coefficients to quantify chemical similarity. In total, 44 different similarity coefficients are available to

171    calculate similarity values between binary fingerprints [21]. We selected seven coefficients for our

172    analysis based on diversity and based on their performance as observed by Todeschini et al. (2012) and

173    Floris et al. (2014) [21,30] (see Table 2). Similarity coefficients "SS1", "Ja" and "Gle" all showed a high

174    performance within Todeschini et al. 2012, but have an exactly similar performance as the JT-coefficient.

175    Therefore, it has been decided to only include the JT-coefficient within this study. All included similarity

176    coefficients were rescaled to provide similarity values between 0 and 1 using Equation 1, similar to

177    Todeschini et al. (2012) [21].

178

179
$$s' = \frac{s + \alpha}{\beta}$$
*Equation* 1

180    Where s is the original similarity value (Table 2), s' is the rescaled function in the range [0, 1], and α and

181    β are numerical parameters whose values are reported in Table 2. When α = 0 and β = 1, this means that

182    no transformation has been applied [21].

183

184    *Table 1: Binary fingerprints included in this study.*

| Name | Number of bits | Type of fingerprint | Source |
|---|---|---|---|
| Substructure Fingerprints | 307 | Dictionary based fingerprints | PaDEL-Descriptor [29] |
| MACCS Fingerprints | 166 | | |
| E-State Fingerprints | 79 | | |
| PubChem Fingerprints | 881 | | |
| Klekota-Roth Fingerprints | 4860 | | |
| CDK Extended Fingerprints | 1024 | Topological or Path-based fingerprints | |
| Atom Pairs Fingerprints | 1024 | | |
| Topological Torsion Fingerprints | 1024 | | |
| Extended Connectivity Fingerprints (diameter = 0) (ECFP0) | 1024 | Circular fingerprints * | RDkit [28] |
| Extended Connectivity Fingerprints (diameter = 2) (ECFP2) | 1024 | | |
| Extended Connectivity Fingerprints (diameter = 4) (ECFP4) | 1024 | | |
| Extended Connectivity Fingerprints (diameter = 6) (ECFP6) | 1024 | | |
| Functional-Class Fingerprints (diameter = 0) (FCFP0) | 1024 | Circular/pharmacophore fingerprints * | |
| Functional-Class Fingerprints (diameter = 2) (FCFP2) | 1024 | | |
| Functional-Class Fingerprints (diameter = 4) (FCFP4) | 1024 | | |
| Functional-Class Fingerprints (diameter = 6) (FCFP6) | 1024 | | |

185    *\*Morgan fingerprints were calculated using RDkit with radius of 0, 1, 2 and 3; which is roughly equivalent to*

186    *ECFP and FCFP0, 2, 4, and 6.*

187

188    *Table 2: Similarity coefficients included in this study (obtained from* [21]*).*

| Name | Formula | A | β | Class | Conditions |
|---|---|---|---|---|---|
| Jaccard-Tanimoto (JT) | $s = \dfrac{c}{c + a + b}$ | 0 | 1 | A | c=0 → s=0 |
| Harris-Lahey (HL) | $s = \dfrac{c(2d + a + b)}{2(c + a + b)} + \dfrac{d(2c + a + b)}{2(a + b + d)}$ | 0 | P | S | c=p or d=p → s=1; den=0 → s=0 |
| Consonni-Todeschini 4 (CT4) | $s = \dfrac{\ln(1 + c)}{\ln(1 + c + a + b)}$ | 0 | 1 | A | None |
| Sokal-Sneath 3 (SS3) | $s = \dfrac{1}{4}\left[\dfrac{c}{c + a} + \dfrac{c}{c + b} + \dfrac{d}{a + d} + \dfrac{d}{b + d}\right]$ | 0 | 1 | S | c=p or d=p → s=1; c=0 and d=0 → s=0 |
| Cohen (Coh) | $s = \dfrac{2(cd - ab)}{(c + a)(a + d) + (c + b)(b + d)}$ | +1 | 2 | Q | c=p or d=p → s=1; den=0 → s=0 |

| | | | | | |
|---|---|---|---|---|---|
| Simple Matching (SM) | $s = \dfrac{c + d}{c + a + b + d}$ | 0 | 1 | S | None |
| Yule 2 (Yu2) | $s = \dfrac{\sqrt{cd} - \sqrt{ab}}{\sqrt{cd} + \sqrt{ab}}$ | +1 | 2 | Q | c=p, d=p or ab=0 → s=1 |

189    *Names of the coefficients are provided as in accordance to Todeschini et al. 2012 [21], though the definition of a*

190    *and c are switched in Todeschini et al. 2012 [21]. The column "Class" represents the type of coefficient: S =*

191    *symmetric coefficient (counts a and d are considered equally); A = asymmetric coefficient (only count a is*

192    *considered); Q = correlation based coefficients that are transformed to obtain a value between zero and one. The*

193    *column "conditions" represents conditions that were assumed in order to avoid singularities. Den = denominator; p*

194    *= a + b + c + d.*


195    *2.4 Performance assessment*


196    *2.4.1 Performance statistics*

197    In total, 112 different similarity measures were selected (i.e. all different combinations of 16 fingerprints

198    and 7 similarity coefficients) and we analyzed their predictive performance on classifying the substances

199    in the dataset as (potential) SVHC or non-SVHC. For non-SVHC substances, similarities were calculated

200    to all substances in the SVHC set based on the fingerprint-coefficient combination. Similarities for SVHC

201    substances were calculated to all other substances on the SVHC set. Iteratively, one SVHC molecule at a

202    time was left out of the dataset and compared to the other SVHC substances. For each substance, only the

203    highest similarity value was retained.

204            For each fingerprint-coefficient combination, we determined the maximum balanced accuracy

205    (Equation 2), by selecting the optimal threshold (i.e. a value between 0 and 1) to predict (potential) SVHC

206    status versus non-SVHC status. Substances with a similarity value equal to or above this threshold are

207    predicted to be structurally similar to a substance with CMR, PBT/vPvB or ED properties to such an

208    extent that they are potential CMR, PBT/vPvB or ED themselves (and vice versa). When using a

209    threshold value, the number of 'True Positives (TP)', 'False Positives (FP)', 'False Negatives (FN)' and

210    'True Negatives (TN)' predictions can be determined for a fingerprint-coefficient combination, as well as

211    the balanced accuracy (Equation 2). By iteratively assessing the fingerprint-coefficient performance for

212    all distinguishing threshold values (ranging from 0-1), the optimal threshold, with maximum balanced

213    accuracy could be determined. The optimal threshold was selected for each specific fingerprint-coefficient

214    combination to ensure equal model comparisons.

215

216    $$Balanced\ Accuracy\ = \frac{Sensitivity + Specificity}{2} = \frac{\frac{TP}{TP + FN} + \frac{TN}{TN + FP}}{2} \qquad Equation\ 2$$

217

218    *2.4.2 Best model selection*

219    In addition to the overall performance (with all CMR, PBT/vPvB and ED substances together in the

220    reference set), also the predictive performance of all fingerprint-coefficient combinations for specific

221    subgroups were analyzed (i.e. for the subgroups of CMR, PBT/vPvB and ED substances separately). The

222    whole set of non-SVHC substances was used as truly negative data in each case. The best performing

223    model was selected based on the balanced accuracy.

224

225    *2.4.3 Best model evaluation*

226    Within the best performing models, we analyzed whether potential bias was introduced by the optimal

227    similarity coefficient. Specifically, symmetric similarity coefficients may tend to predict small substances

228    - with many '0-bits' - as similar to small SVHC substances, because of common absence of many features

229    (i.e. *d*-fragments*).* Although such a model could be considered most optimal based on statistical

230    performance of the dataset, the occurrence of this type of similarities is undesirable, as upon application

231    many small substances will incorrectly be classified as (potential) SVHC. Therefore, when potential

232    symmetric coefficient bias was identified in a best performing model, we decided to use an asymmetric

233    similarity coefficient for substances with a low number of '1-bits' (i.e. JT or CT4, which only considers *c*-

234    fragments as similar). The most optimal fragment count cut-off was analyzed based on balanced accuracy.

235         Furthermore, we analyzed the robustness of the best performing models by assessing the

236    performance after two different robustness checks. Within the first robustness check, we extended the

237 non-SVHC dataset by adding the substances of the "non-relevant" SVHC subgroup to the non-SVHC

238 dataset. To illustrate, for the CMR-model, all PBT/vPvB and ED SVHC substances that do not have CMR

239 properties were considered as not-CMR, and thus added to the non-SVHC set for this robustness check.

240 This robustness check could not have been conducted on the overall model, as in this case all SVHC

241 subgroups are relevant. Within a second robustness check, we reduced the number of representative

242 structures of group entries that were included within the SVHC as well as within the non-SVHC set to

243 generally two structures (see Supplemental Material Excel). In addition, some structurally similar

244 substances are represented various times in the SVHC or non-SVHC datasets, including a large number of

245 individual PCB isomers, chlorinated dibenzofurans, chlorinated dibenzodioxins and polybrominated

246 diphenyl ethers on the PBT/vPvB dataset. To determine the robustness of the best performing models,

247 such groups have also been reduced to a representation of generally two representative structures (see

248 Supplemental Material Excel). The performance of the adjusted datasets within the different robustness

249 checks was assessed similarly as described above, using the optimal threshold of the best-performing

250 model.

251 In addition, hierarchical cluster diagrams were generated for the different SVHC subgroups in

252 order to analyze the diversity within the subgroups. Hierarchical clusters were based on the similarity

253 matrix of the subgroup, using single-linkage method.

254 The performance of the best predictive models was also compared to existing methodologies –

255 using the SVHC dataset – including Toxtree (i.e. Benigni/Bossa rulebase for mutagenicity and

256 carcinogenicity), DART and the PB-score tool [6,7,31]. For this analysis, the presence of a structural alert

257 from Toxtree and/or DART was interpreted as a prediction of SVHC status based on CMR properties.

258 Besides performance evaluation, also applicability domain was analyzed by determining the 95th

259 percentile of molecular weight, log $K_{ow}$ [5], number of atoms, number rings and number of aromatic rings

260 within the applied datasets.

261 All data was analyzed in R (version 3.5.1) [32], using *caret*, *ChemmineR*, *caTools*, *ROCR* and *rcdk*

262 [33–37].

263  **3.  Results**

264  *3.1 Best model selection*

265  *3.1.1 Overall model performance*

266  Table 3 shows the ten best performing models when all CMR, PBT/vPvB and ED substances are taken

267  together in a single SVHC dataset. A wide variety of fingerprints was identified in the top ten models,

268  including dictionary-based, path-based, circular-based and pharmacophore-based fingerprints. In contrast,

269  one similarity coefficient, the Simple Matching (SM), is dominating the top ten models. Furthermore, it

270  can be observed that relatively high optimal similarity thresholds are determined. The height of the

271  threshold is highly related to the used similarity coefficient, and is specifically high for the SM coefficient

272  (Figure S.1). This is a consequence of the fact that $c$ and $d$ variables are treated as similar in this

273  coefficient (Table 2).

274      The overall best performing model, PubChem-SM combination, has an overall balanced accuracy

275  of 0.846. However, this specific combination is not the most optimal for the specific subgroups, having

276  different (toxicological) concerns. Therefore, we also analyzed model performances for the CMR,

277  PBT/vPvB and ED groups separately.

278

*Table 3: Ten best performing fingerprint-coefficient combinations for the dataset with all CMR, PBT/vPvB and ED substances included. Also specific subgroup performances – in balanced accuracy - are provided based on the optimal overall threshold values. The numbers represent the number of SVHC substances, 411 non-SVHC substances were included. Highest balanced accuracies are given in italic bold. AUC is the area under the curve of ROC-plot.*

| Model | | Threshold | Overall model performance (n=546 SVHC) | | | | | Balanced accuracy of subgroups using overall threshold value | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Fingerprint | Coefficient | | Sensitivity | Specificity | Precision | AUC (ROC) | Balanced accuracy | CMR (n=306 SVHC) | PBT/vPvB (n=209 SVHC) | ED (n=52 SVHC) |
| Pubchem | SM | 0.985 | 0.810 | 0.883 | 0.902 | 0.904 | *0.846* | 0.801 | 0.929 | 0.988 |
| Extended | SM | 0.957 | 0.806 | 0.878 | 0.898 | 0.897 | 0.842 | *0.811* | 0.889 | 0.981 |
| MACCS | SM | 0.970 | 0.734 | 0.946 | 0.948 | 0.897 | 0.840 | 0.760 | *0.951* | 0.960 |
| FCFP4 | SM | 0.991 | 0.835 | 0.842 | 0.875 | 0.893 | 0.839 | 0.802 | 0.911 | *0.990* |
| KlekotaRoth | SM | 0.998 | 0.773 | 0.898 | 0.909 | 0.889 | 0.835 | 0.777 | 0.921 | 0.942 |
| ECFP2 | SM | 0.992 | 0.852 | 0.813 | 0.858 | 0.900 | 0.832 | 0.798 | 0.925 | 0.987 |
| ECFP4 | SM | 0.984 | 0.832 | 0.832 | 0.868 | 0.882 | 0.832 | 0.791 | 0.900 | *0.990* |
| Extended | SS3 | 0.895 | 0.714 | 0.942 | 0.942 | 0.888 | 0.828 | 0.775 | 0.902 | 0.971 |
| Extended | Coh | 0.884 | 0.711 | 0.934 | 0.935 | 0.887 | 0.822 | 0.769 | 0.899 | 0.981 |
| MACCS | SS3 | 0.923 | 0.716 | 0.922 | 0.924 | 0.875 | 0.819 | 0.739 | 0.924 | 0.969 |

*3.1.2 Subgroup model performance*

284     The best performing similarity models optimized for the separate CMR, PBT/vPvB and ED subgroups are

285     shown in Table 4 (in row one till three, respectively). For the ED subgroup, 30 out of the 112 tested

286     different similarity measures showed similar predictive performance, but the rank of the fingerprints and

287     coefficients separately shows a highest rank for the FCFP4 fingerprint and the SS3 similarity coefficient.

288     The best performing combination of fingerprint and similarity coefficient is different for the different

289     subgroups, and a (slightly) higher balanced accuracy is obtained when compared to the best performing

290     overall model (Table 3).

291

292 *Table 4: Best performing fingerprint-coefficient combination for the CMR, PBT/vPvB and ED subgroups, including balanced accuracies after*

293 *robustness checks (see section 3.2). The CMR model was improved by combining a symmetric and asymmetric coefficient in order to prevent*

294 *symmetric coefficient bias (see section 3.2). In robustness check 1, the SVHC substances that did not belong to the subgroup of concern were*

295 *added to the dataset as non-SVHCs. In robustness check 2, the number of representative structures for group entries and structurally similar*

296 *substances were reduced to generally two structures in the SVHC and non-SVHC set. The numbers represent the number of SVHC substances. The*

297 *number of non-SVHC substances varies between the full model assessment (n=411) and the robustness checks (see 3.2.2). '-' means that it is not*

298 *possible to calculate a single AUC for a combination of two models. AUC is the area under the curve of ROC-plot.*

| Subset | Model | | Threshold | Sensitivity | Specificity | Precision | AUC (ROC) | Balanced accuracy | Robustness check | |
| | Fingerprint | Coefficient | | | | | | | 1 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|
| CMR (n=306) | Extended | SM | 0.944 | 0.784 | 0.854 | 0.800 | 0.859 | 0.819 | 0.735 | 0.799 |
| PBT/vPvB (n=209) | MACCS | SM | 0.970 | 0.919 | 0.983 | 0.965 | 0.971 | 0.951 | 0.942 | 0.911 |
| ED (n=52) | FCFP4 | SS3 | 0.866 | 0.981 | 1.000 | 1.000 | 0.984 | 0.990 | 0.969 | 0.917 |
| CMR improved (n=306) | Extended | CT4 (<85) SM (≥85) | 0.851 0.944 | 0.650 | 0.949 | 0.905 | - | 0.800 | 0.742 | 0.769 |

299 *3.2 Best model evaluation*

300 *3.2.1 Symmetric coefficient bias*

301 By applying the "Extended fingerprint – SM coefficient" combination for the CMR dataset, with a 0.944

302 similarity threshold, all substances with less than 63 fingerprint bits were considered to be similar to

303 CMR-SVHCs (Figure 2A). This coefficient bias is also observed upon visual inspection of the FP-

304 substances, perceiving a better similarity assessment with increased number of fingerprint bits (e.g.

305 'methyl octanoate' and '3-propanolide'; or 'Captan' and 'Captafol'; Figure 2B).

306 Based on our assessment, finding an optimal cut-off within the range of 63 to 100 fingerprint bits,

307 the combination of the CT4 coefficient for substances with less than 85 fingerprint bits and the SM

308 coefficient for substances with 85 or more fingerprint bits is most optimal, with a balanced accuracy of

309 0.800 and threshold values of 0.851 and 0.944, respectively (Table 4, row 4). The statistical performance

310 of the CT4-SM combination is lower than the SM coefficient only (when looking at the balanced

311 accuracy), due to an increase in FN-classified substances. On the contrary, also more substances are

312 correctly classified as negative, including structures with a relative low number of fingerprint bits, like

313 methyl octanoate and the terpenoid blend QRD-460 (Figure 2B; Figure S.2). This results in a much better

314 specificity and precision (Table 4; Table S.1). The PBT/vPvB and ED models do not require a

315 combination of asymmetric and symmetric coefficients as no symmetric coefficient bias was observed

316 (Supplemental Material S.4; Figure S.2).

317

318 *3.2.2 Robustness checks*

319 The robustness of the best-performing subgroup models was investigated via two robustness checks

320 (Table 4). Within the first robustness check, the SVHC substances that did not belong to the subgroup of

321 concern were added to the dataset as non-SVHCs (i.e. 'robustness check 1'). For the best performing

322 CMR model, 651 non-SVHC substances were included, for the best PBT/vPvB model 748 non-SVHC

323 substances and for the best ED model 905 non-SVHC substances. Within the second robustness check,

324     we reduced the number of representative structures for group entries and structurally similar substances of

325     the SVHC and non-SVHC set to generally two structures (i.e. 'robustness check 2'). In total, 30

326     substances were excluded from the non-SVHC set, 35 from the CMR subset, 96 from the PBT/vPvB

327     subset, and 34 from the ED subset.

328           Adding the non-target SVHC-substances to the non-SVHC set lowered the balanced accuracy and

329     hence the predictive performance, specifically for the CMR similarity model. Conversely, removal of

330     close structural analogues resulted in a larger decrease in predictive performance for the PBT/vPvB and

331     ED specific models.

332

333     *3.2.3 Single-point-of-knowledge*

334     The CMR and PBT/vPvB subgroup have a quite broad basis with 306 and 209 substances, respectively,

335     whereas the ED subgroup only consists of 52 substances. Within the PBT/vPvB and ED subgroups, some

336     groups of very similar structures can be identified, and only a few single-point-of-knowledge structures

337     (SPOKs) are included (Figure 3). SPOKs are substances that are not comparable to any other substance in

338     the subgroup and thus are single-point-of-knowledges within the dataset (i.e. the FN). Within the ED

339     substances, four groups and one distinct substance are present; in the PBT/vPvB subgroup, 15 groups and

340     17 distinct substances were identified (giving 1 and 17 false negatives, respectively). On the contrary, the

341     CMR-SVHC dataset is much more diverse in chemical structures and contains much more SPOKs,

342     reflected in the high number of FN-classified substances (n=107). For the CMR subgroup, no

343     unambiguous hierarchical clustering can be generated as the CT4-SM coefficient combination does not

344     fulfill the mathematical conditions for all substances (i.e. similarity between substance x and y is not

345     necessarily similar to the similarity between y and x). Nevertheless, some groups can be identified,

346     including polycyclic aromatic hydrocarbons, haloalkanes, cyclic and acyclic ethers, alkyl phenols,

347     phthalates, aromatic amines, nitroaromatics and chloroaromatics. As a consequence of the high structural

348     diversity, the calculated balanced accuracy is also lower for the CMR subgroup compared to the

349    PBT/vPvB and ED groups. It should be noted that the SPOK false negatives will be included in the full

350    dataset of SVHC substances when applying the model to a new substance.

351    *3.2.4 Performance of existing models*

352    The performance of a CMR model (i.e. the sum outcome from Toxtree and DART [7,31]) on the used

353    SVHC-set was analyzed. Substances were considered as CMR by the model when a Toxtree or DART

354    alert was identified. A balanced accuracy of 0.62 was determined, with a sensitivity of 0.78 and a

355    specificity of 0.47. Furthermore, the performance of a PBT model was evaluated (i.e. PB-score tool [6]).

356    For four substances no PB-score could be calculated as no log $K_{aw}$ could be estimated. For the used

357    dataset, a balanced accuracy of 0.73 was determined, with a sensitivity of 0.53 and a specificity of 0.93.

358    No ED model was analyzed because of the limitations identified in the ED-similarity model (see

359    discussion).

360

361

## 4. Discussion

As ever-increasing amounts of substances are produced, applied and emitted, it is important to focus attention on assessing the risks of those substances that are most likely to actually cause problems. Therefore, there is a need for efficient screening and prioritization methods to identify chemicals with a high potential of being hazardous. Within this study we evaluated the efficiency of a set of similarity measures for the identification of (potential) SVHCs. Based on our approach, we identified the three best performing models for CMR, PBT/vPvB and ED subgroups, that all show a promising balanced accuracy ($\geq$0.8) based on the used dataset.


### *4.1 Model performance*

The three subgroup-specific models showed a better performance than one single overall model. This is likely related to a difference in mode(s) of action between CMR, PBT/vPvB and ED substances, and is also reflected in the most optimal fingerprints. In addition, predictive performance appeared reasonably robust with less than 10% reduction of balanced accuracy following the two robustness checks for all best performing models.

For the PBT/vPvB substances, the MACCS fingerprint performed best. The MACCS fingerprint contains only 166 predefined bits and was particularly developed to categorize substances in functional groups [38]. The PBT/vPvB dataset has a low structural diversity, with many substances sharing common structural features (Figure 3), including aromatic-rings and high levels of halogenation. In addition, small substances are often not considered PBT/vPvB, as in general a lower octanol-water-partitioning is observed for smaller substances, and this in turn is related to the bioaccumulation potential [39]. Apparently, the MACCS fingerprint is very effective in making a distinction between PBT/vPvB and non-PBT/vPvB substances based on these common features. Consequently, a high predictive performance is observed for this dataset (0.951).

386        The CMR substances are structurally much more diverse, with 107 SPOKs in the SVHC dataset.

387        This diversity is also reflected in the most optimal fingerprint, the Extended Fingerprint. This path-based

388        fingerprint, which is based on the well-known Daylight fingerprint [40], recognizes all paths within a

389        structure consisting of 1-9 atoms (i.e. search depth of 8 bonds) and also includes some additional bits that

390        describe ring features [29]. Compared to dictionary-based fingerprints, it is assumed that this method is

391        more suitable to capture the broad diversity in CMR substances, as it characterizes all possible fragments

392        within a structure.

393        As the balanced accuracy for the CMR subgroup was relatively low (compared to the PBT/vPvB

394        and ED groups), we added an extra fingerprint that encodes for the presence of CMR-specific fragments

395        identified in expert-models like Toxtree and DART [7,31]. Nonetheless, the inclusion of the

396        mechanistically based substructures in the fingerprint did not lead to any improvement in the predictive

397        performance (Supplemental Material S.5). Apparently, the size of the dataset and the fragments present in

398        the optimal fingerprint already cover the specific structural features that have been linked to our collective

399        knowledge of mechanisms of action leading to CMR effects. The additional fingerprint is therefore

400        excluded again.

401        For ED substances, the FCFP-4 is identified as best performing fingerprint. FCFP-4 identifies

402        fragments based on functional group patterns. It recognizes atoms as hydrogen donors, hydrogen

403        acceptors, aromatics, halogens, basic-atoms and acidic-atoms, and it identifies fragments based on

404        patterns between these atoms (e.g. hydrogen donor – hydrogen acceptor – hydrogen donor) [28].

405        Endocrine disruptors generally interact with specific hormone receptors or interact with proteins in the

406        hormone pathway [41], and such (receptor) binding properties are potentially identified best by the

407        features covered in the FCFP-fingerprint. Furthermore, the diameter of 4 (FCFP-4) scored slightly better

408        for the similarity search than a diameter of 2 or 6, which is in line with earlier findings [42]. Rogers and

409        Hahn (2010) [42] concluded that a diameter of four is typically sufficient for similarity searches whereas

410        a diameter of six or eight is best for activity learning methods.

411        Despite the very high performance for the ED subgroup (0.990), prediction results from this

412    model should be interpreted with caution. The currently used ED-SVHC dataset is limited as it only

413    consists of a few number of substances that have a large structural overlap (Figure 3) and consequently

414    results in higher uncertainty around the optimal threshold value compared to the other models (Figure

415    S.3). In addition, there is only one substance on the ED-list with a hormone backbone (i.e. Diosgenin).

416    The reason for the low number of identified ED-SVHC substances is partially related to the fact that only

417    those substances are identified as ED for which SVHC-identification is of added regulatory value. In

418    addition, only recently guidance and criteria are developed for the identification of ED substances [43]. It

419    is recommended to further develop the ED model when more substances are classified as ED-SVHC, or

420    by including known endocrine disrupting substances such as the natural substrates (and synthetic variants

421    derived thereof) interacting with estrogen/androgen/thyroid and steroidogenic pathways. With a broader

422    dataset, a more sophisticated screening model will be possible. Based on the current dataset the ED-

423    SVHC similarity model is expected to miss many (potential) ED substances.

424        A higher performance is observed for the best-scoring CMR and PBT/vPvB similarity models

425    compared to existing models [6,7,31], when using the SVHC dataset. This indicates the value and

426    relevance of the structural property principle for identifying potential SVHC substances. For the ED

427    model, no comparison was made with existing models because of the limitations as mentioned above.

428

429    *4.2 Focus and restriction of the modelling*

430    We limited our assessment to the performance of 2D-binary fingerprints, and the presence or absence of

431    2D-fragments. More sophisticated fingerprints are also available, including count-based fingerprints,

432    taking into account how many times a fragment is present, or 3D-fingerprints that consider chemical

433    conformation. Particularly, 3D-fingerprints could be relevant to identify potential ED substances, as

434    receptor-binding properties are highly important for this group. In general, however, 2D-binary

fingerprints are most popular as they are an acceptable trade-off between the wealth of (possible) information and simplicity, enabling an easy and quick comparison [17,30]. Especially for the proposed screening activities, the currently evaluated methodology is considered adequate.

In principle, all non-SVHC substances that have been used for modelling purposes within this study are tested on CMR, PBT/vPvB and ED properties. Nevertheless, it is possible that some substances are currently not identified as such, but will become a SVHC substance in future, when new information becomes available or when new evaluations are conducted. For instance, glyphosate is included in the non-SVHC list used in this study, although its carcinogenicity is currently extensively discussed [44,45]. Furthermore, as shown in Figure 2, Captafol is considered as CMR substance whereas its close structural analogue Captan is not (see Supplemental Material S.1). Captafol is classified as a carcinogen category 1B (leading to SVHC status), and Captan as a carcinogen category 2 [46]. Although the model identifies Captan as a false positive, the results could be very useful and may provide further arguments for (de)-classification of these substances. For instance, within European regulatory frameworks, a category 2 classification (for carcinogenicity but also for mutagenicity and reproductive toxicity) is often the highest classification that can be agreed upon when there are insufficient (experimental) data to support a category 1B classification [47].

Despite the conductance of a performance analysis, including robustness checks, we were not able to conduct a proper external validation in order to analyze the performance on an external dataset. As SVHCs are identified after a regulatory decision process in which all available data is evaluated, we are not in the position to mark substances as SVHC for external validation purposes. Similarly, non-SVHC substances are challenging to assign, as many substances are not extensively evaluated on all SVHC endpoints (i.e. CMR, PBT/vPvB and ED). A proper external validation set can therefore only be developed in future, when new SVHC and non-SVHC substances are identified. Future work will focus on the application of the developed methodology to large sets of substances to obtain a better idea of the application performance.

460

461   *4.3 Use and applicability domain of the model*

462   The assumption, that structurally similar substances are likely to have similar properties, seems valid

463   based on our analysis and model performances. The proposed similarity models focus on multiple

464   endpoints (i.e. CMR, PBT/vPvB and ED) and could be applied as a first screening model, enabling to

465   prioritize further follow-up analyses. The model directly highlights the most similar SVHC substance(s),

466   which could provide additional information on the specific concerns. The absolute results should not be

467   interpreted as a conclusive outcome. The methodology is framed to give systematic and transparent ways

468   to identify relations that would not manually be identified. Based on the follow-up, it could be concluded

469   that 1) the substance is likely to have similar effects, 2) that further data is required to substantiate the

470   outcome, or 3) that the substance is not expected to have CMR, PBT/vPvB or ED properties.

471   Furthermore, it should also be highlighted that the developed model considers a screening model

472   to identify whether new chemicals are structurally similar to known SVHC substances. It should be kept

473   in mind that SVHCs are identified based on a regulatory decision process in which available data is

474   evaluated. Consequently, a negative model results (i.e. not structurally similar to a SVHC substance) does

475   not necessarily means that the substance for instance has no carcinogenic, or persistent properties. What it

476   does mean is that the chemical is not structurally similar to a SVHC and that related regulatory

477   consequence may - at the moment - not be applicable for the new chemical.

478   A short guide on the application of the methodology is provided in Supplemental Material S.3.

479   With respect to the applicability domain, an increase in reliability is observed with an increase in structure

480   complexity for all three models, especially for the CMR model (i.e. number of atoms and different atom

481   types). The structure similarity models are not applicable to arsenic, beryllium, cadmium, chromium,

482   lead, mercury, nickel and cobalt-metal derivatives. For these chemicals, the metal atoms (or ions) are

483   thought to be the cause of concern, irrespective of the (organic) groups present in the inorganic molecule.

484    These metal-based complexes are by definition predicted to be SVHC substances. However, the models

485    can be used to generate a first prediction for non-dissociating metals (e.g. organotin substances). In

486    principle, the chemical similarity itself is an applicability domain descriptor. If the new substance is

487    sufficiently similar to an existing SVHC, the substance is clearly within the applicability domain of the

488    model. Furthermore, physicochemical boundaries (i.e. $95^{th}$ percentiles) have been calculated for the

489    different models based on molecular weight, log $K_{ow}$, number of atoms, number of rings and the number

490    of aromatic rings (Table S.2). The similarity methodology does not discriminate between pristine

491    substances or environmental and/or metabolic breakdown products; this model is applicable to both. Risk

492    assessors, we therefore advise not only to apply the predictive model to the parent substance, but also to

493    the breakdown products as well as possible tautomers, as these may give different similarity outcomes.

494        This effective screening method can particularly be applied during product development and

495    chemical synthesis. By enhancing attention on chemicals of potential SVHC concern as early as possible

496    within regulatory frameworks and safe-by-design trajectories, this methodology contributes to the

497    transition towards a non-toxic environment.

498

## 5. Conclusions

Within this study, a systematic and transparent methodology was established that could identify potential SVHCs based on structural similarity to a known set of SVHCs. We have analyzed the influence of selected similarity characterizations (fingerprints and coefficients) on the identification of chemicals of potential SVHC concern. A good statistical performance was obtained for CMR, PBT/vPvB and ED substances, but nevertheless further work is considered necessary to improve the ED part due to the small reference dataset for this SVHC concern.

Application of the developed methodology is considered useful to identify chemicals of potential concern as early as possible, and as such may ensure that up-front more adequate risk management measures can be applied to contribute towards a non-toxic environment. It is foreseen that this scientifically-based model is beneficial to (environmental) risk assessors, industrial partners and academia.

## 6. Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## 7. Acknowledgements

## 8. References

519     [1]     S. Sauvé, M. Desrosiers, A review of what is an emerging contaminant, Chem. Cent. J. (2014).

520            doi:10.1186/1752-153X-8-15.

521     [2]     European Parliament, Decision No 1386/2013/EU on a General Union Environment Action

522            Programme to 2020 'Living well, within the limits of our planet,' 2013.

523     [3]     European Commission, Fact Sheet: 7th Environmental Action Plan, (2013).

524            http://ec.europa.eu/environment/pubs/pdf/factsheets/7eap/en.pdf (accessed February 2, 2019).

525     [4]     European Parliament, REACH Regulation EC/1907/2006, 2006.

526     [5]     US Environmental Protection Agency (US EPA), Estimation Programs Interface Suite for

527            Microsoft Window, v4.1, (2012).

528     [6]     E. Rorije, E.M.J. Verbruggen, A. Hollander, T.P. Traas, M.P.M. Janssen, Identifying potential

529            POP and PBT substances : Development of a new Persistence/Bioaccumulation-score (RIVM

530            Report 601356001), (2011) 1–88.

531     [7]     R. Benigni, C. Bossa, N. Jeliazkova, T. Netzeva, A. Worth, The Benigni / Bossa rulebase for

532            mutagenicity and carcinogenicity – a module of Toxtree, 2008.

533     [8]     P. Banerjee, A.O. Eckert, A.K. Schrey, R. Preissner, ProTox-II: A webserver for the prediction of

534            toxicity of chemicals, Nucleic Acids Res. (2018). doi:10.1093/nar/gky318.

535     [9]     H. Yang, C. Lou, L. Sun, J. Li, Y. Cai, Z. Wang, W. Li, G. Liu, Y. Tang, AdmetSAR 2.0: Web-

536            service for prediction and optimization of chemical ADMET properties, Bioinformatics. (2019).

537            doi:10.1093/bioinformatics/bty707.

538    [10]     BIOVIA, TOPKAT (Toxicity Prediction by Komputer Assisted Technology), (n.d.).

539            https://www.3dsbiovia.com/products/datasheets/qsar-admet-and-predictive-toxicology-with-

540       ds.pdf.

541   [11]  MultiCase Inc, MultiCase, (n.d.). http://www.multicase.com/.

542   [12]  Leadscope, Leadscope, (n.d.). http://www.leadscope.com/.

543   [13]  DTU, Danish (Q)SAR Database, (2015). http://qsar.food.dtu.dk/.

544   [14]  M.A. Johnson, G.M. Maggiora, Concepts and Applications of Molecular Similarity, Wiley, New
545       York, 1990.

546   [15]  ChemSec, Methodology for grouping the SIN List and development of the SINimilarity tool,
547       (2015). https://chemsec.org/publication/sin-list/methodology-for-grouping-the-sin-list-and-
548       development-of-the-sinimilarity-tool/.

549   [16]  European Chemicals Agency (ECHA), Screening Definition Document - Methodology for
550       identifying (groups of) potential substances of concern for (further) regulatory action, 2019.

551   [17]  P. Willett, The calculation of molecular structural similarity: Principles and practice, Mol. Inform.
552       (2014). doi:10.1002/minf.201400024.

553   [18]  R.D. Brown, Descriptors for Diversity Analysis, Perspect. Drug Discov. Des. 7 (1997) 31–49.

554   [19]  R. Todeschini, V. Consonni, Molecular descriptors for chemoinformatics: volume I: alphabetical
555       listing/volume II: appendices, references., John Wiley and Sons, 2009.

556   [20]  P. Willett, J.M. Barnard, G.M. Downs, Chemical similarity searching, J. Chem. Inf. Comput. Sci.
557       (1998). doi:10.1021/ci9800211.

558   [21]  R. Todeschini, V. Consonni, H. Xiang, J. Holliday, M. Buscema, P. Willett, Similarity coefficients
559       for binary chemoinformatics data: Overview and extended comparison using simulated and real
560       data sets, J. Chem. Inf. Model. (2012). doi:10.1021/ci300261r.

561 [22] C.L. Mellor, R.L. Marchese Robinson, R. Benigni, D. Ebbrell, S.J. Enoch, J.W. Firman, J.C.

562 Madden, G. Pawar, C. Yang, M.T.D. Cronin, Molecular fingerprint-derived similarity measures

563 for toxicological read-across: Recommendations for optimal use, Regul. Toxicol. Pharmacol.

564 (2019). doi:10.1016/j.yrtph.2018.11.002.

565 [23] National Institue for Public Health and the Environment (RIVM), List of Dutch Substances of

566 Very High Concern [in Dutch], (2018). https://rvszoeksysteem.rivm.nl/ZZSlijst/Index (accessed

567 March 1, 2018).

568 [24] European Chemicals Agency, Biocidal Active Substances., (2018).

569 http://echa.europa.eu/web/guest/information-on-chemicals/biocidal-active-substances (accessed

570 May 23, 2018).

571 [25] European Commission, EU Pesticides Database, (2018).

572 http://ec.europa.eu/food/plant/pesticides/eu-pesticides-

573 database/public/?event=activesubstance.selection&language=EN (accessed May 23, 2018).

574 [26] Daylight, SMILES - A Simplified Chemical Language, (2008).

575 http://www.daylight.com/dayhtml/doc/theory/theory.smiles.html (accessed February 15, 2019).

576 [27] S. Riniker, G.A. Landrum, Open-source platform to benchmark fingerprints for ligand-based

577 virtual screening, J. Cheminform. (2013). doi:10.1186/1758-2946-5-26.

578 [28] G. Landrum, RDKit: Open-source Cheminformatics and machine-learning,

579 Http://Www.Rdkit.Org/. (2019). doi:10.2307/3592822.

580 [29] C.W. Yap, PaDEL-descriptor: An open source software to calculate molecular descriptors and

581 fingerprints, J. Comput. Chem. (2011). doi:10.1002/jcc.21707.

582 [30] M. Floris, A. Manganaro, O. Nicolotti, R. Medda, G.F. Mangiatordi, E. Benfenati, A generalizable

583      definition of chemical similarity for read-across, J. Cheminform. (2014). doi:10.1186/s13321-014-

584      0039-1.

585   [31]   S. Wu, J. Fisher, J. Naciff, M. Laufersweiler, C. Lester, G. Daston, K. Blackburn, Framework for

586      identifying chemicals with structural features associated with the potential to act as developmental

587      or reproductive toxicants, Chem. Res. Toxicol. (2013). doi:10.1021/tx400226u.

588   [32]   R Core Team, R: A Language and Environment for Statistical Computing, Vienna, Austria.

589      (2018).

590   [33]   M. Kuhn, R Package: caret, Ver. 6.0-81, CRAN. (2018).

591   [34]   Y. Cao, A. Charisi, L.C. Cheng, T. Jiang, T. Girke, ChemmineR: A compound mining framework

592      for R, Bioinformatics. (2008). doi:10.1093/bioinformatics/btn307.

593   [35]   J. Tuszynski, caTools: Tools: moving window statistics, GIF, Base64, ROC AUC, etc. R package

594      version 1.17.1.1, URL Http//CRAN. R-Project. Org/Package= CaTools. (2018).

595   [36]   T. Sing, O. Sander, N. Beerenwinkel, T. Lengauer, ROCR: Visualizing classifier performance in

596      R, Bioinformatics. (2005). doi:10.1093/bioinformatics/bti623.

597   [37]   R. Guha, Chemical Informatics Functionality in R, J. Stat. Softw. (2007).

598      doi:10.18637/jss.v018.i05.

599   [38]   MACCS, Molecular ACCess System (MACCS) keys. MDL Information Systems - As interpreted

600      by CDK, (n.d.).

601   [39]   European Chemicals Agency (ECHA), Guidance on information requirements and chemical safety

602      assessment (Chapter R.11: PBT/vPvB assessment), 2017. doi:10.2823/128621.
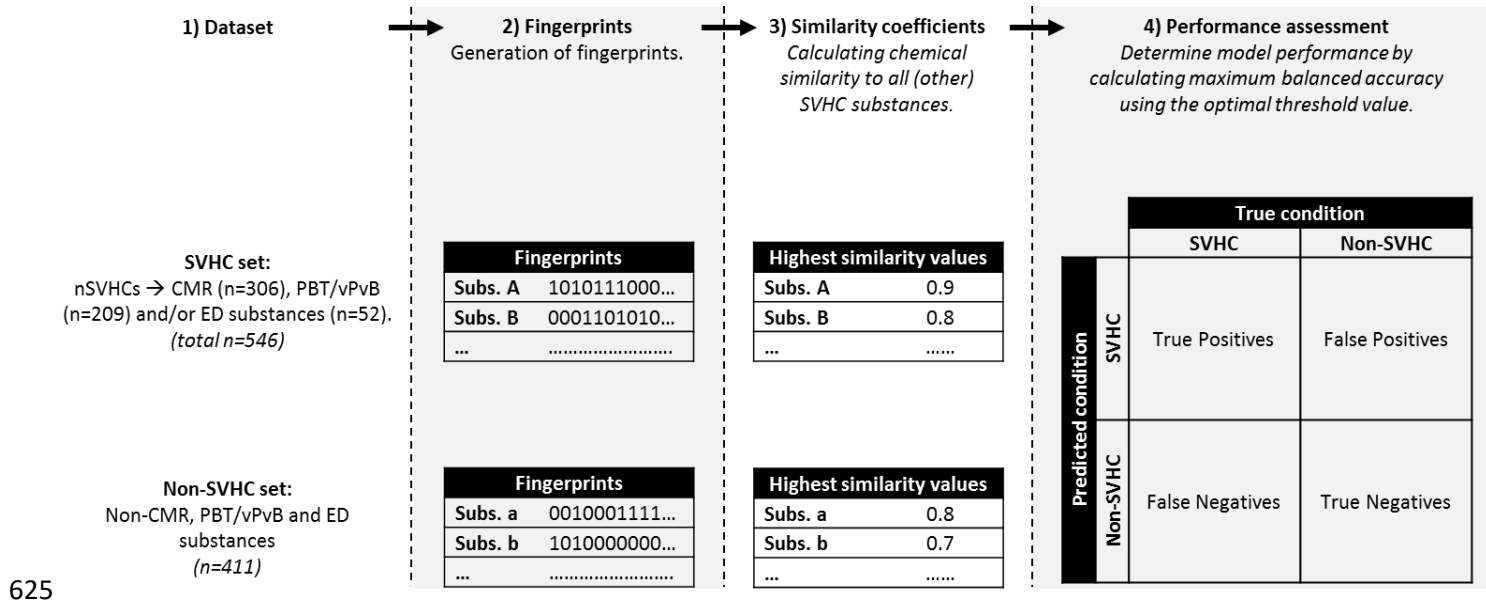
603   [40]   Daylight, Fingerprints - Screening and Similarity, (2008).

604      http://www.daylight.com/dayhtml/doc/theory/theory.finger.html (accessed February 15, 2019).

605   [41]   Å. Bergman, J. Heindel, S. Jobling, K. Kidd, R. Zoeller, Endocrine Disrupting Chemicals - 2012,

606          2013.

607   [42]   D. Rogers, M. Hahn, Extended-connectivity fingerprints, J. Chem. Inf. Model. (2010).

608          doi:10.1021/ci100050t.

609   [43]   European Chemicals Agency (ECHA), European Food Safety Authority (EFSA), Joint Research

610          Centre (JRC), Guidance for the identification of endocrine disruptors in the context of Regulations

611          (EU) No 528/2012 and (EC) No 1107/2009, EFSA J. (2018). doi:10.2903/j.efsa.2018.5311.

612   [44]   International Agency for Research on Cancer, Glyphosate Monograph, 2015.

613   [45]   J. V. Tarazona, D. Court-Marques, M. Tiramani, H. Reich, R. Pfeil, F. Istace, F. Crivellente,

614          Glyphosate toxicity and carcinogenicity: a review of the scientific basis of the European Union

615          assessment and its differences with IARC, Arch. Toxicol. (2017). doi:10.1007/s00204-017-1962-

616          5.

617   [46]   European Chemicals Agency (ECHA), Search for Chemicals, (2019).

618          https://echa.europa.eu/en/information-on-chemicals (accessed September 20, 2002).

619   [47]   M. Woutersen, M. Beekman, M.E.J. Pronk, A. Muller, J.A. de Knecht, B.C. Hakkert, Does

620          REACH provide sufficient information to regulate mutagenic and carcinogenic substances?, Hum.

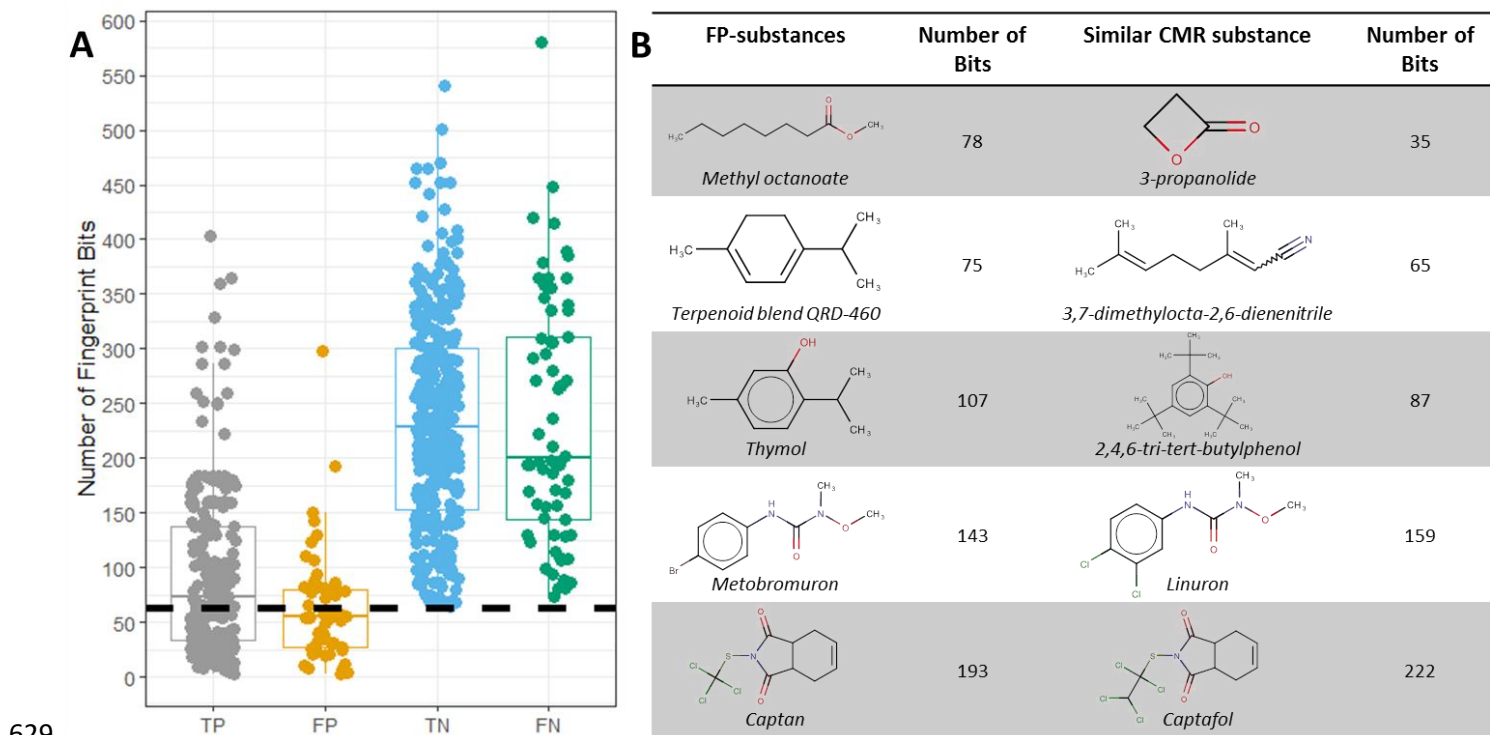621          Ecol. Risk Assess. (2018). doi:10.1080/10807039.2018.1480351.

622

623

624   **9.  Figures**



| | 1) Dataset | 2) Fingerprints<br>Generation of fingerprints. | 3) Similarity coefficients<br>*Calculating chemical similarity to all (other) SVHC substances.* | 4) Performance assessment<br>*Determine model performance by calculating maximum balanced accuracy using the optimal threshold value.* |

625

626   *Figure 1. Overview of the methodology divided into four steps. Steps two to four were reiterated for multiple*
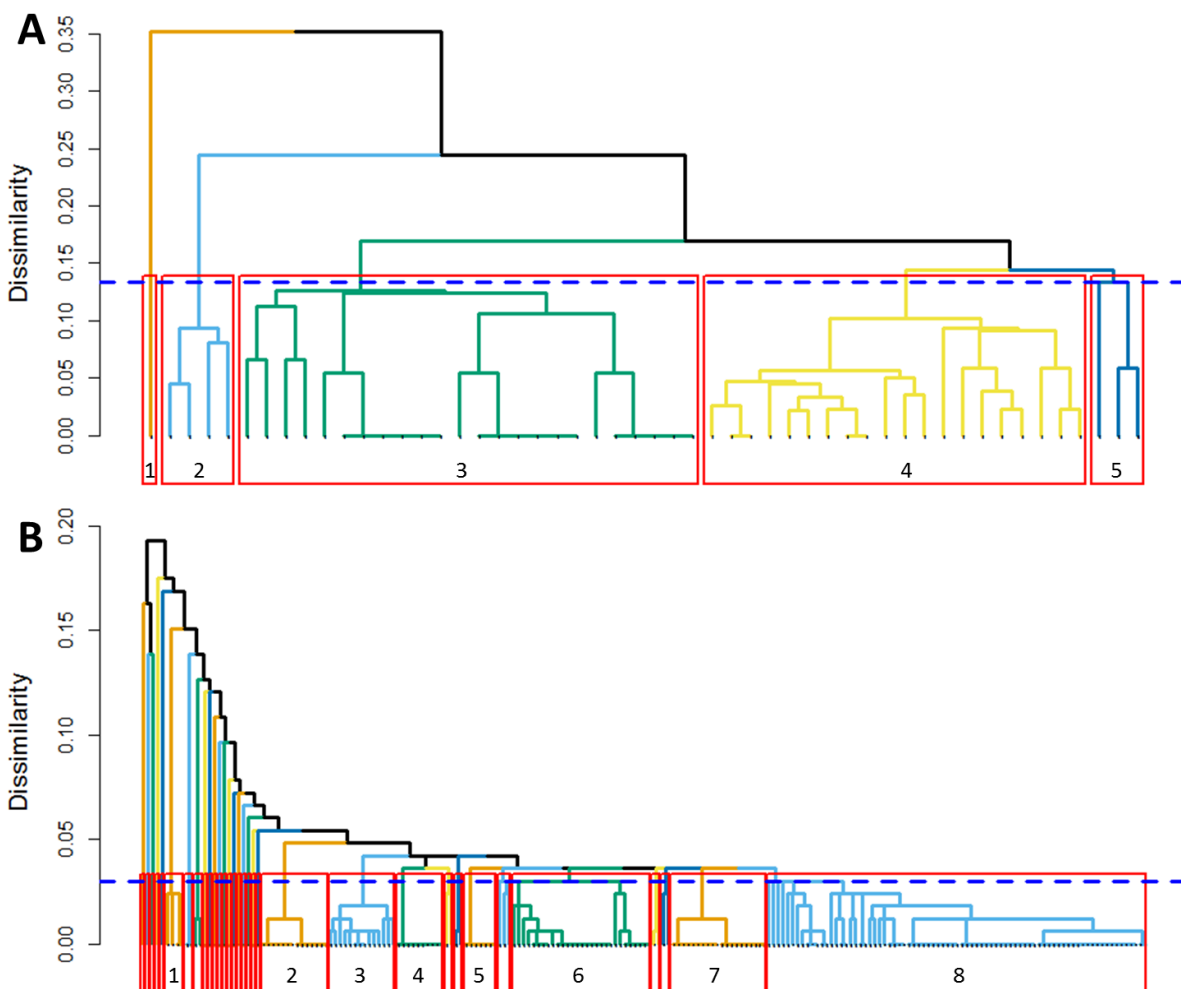
627   *fingerprint-coefficient combinations.*

628

629

*Figure 2: Classification of the CMR-SVHC and non-SVHC substances using the "Extended Fingerprint – SM coefficient" combination. A) Fingerprint bit count distributions across the different classifications: True Positive, False Positives, True Negatives and False Negatives. All substances with less than 63 fingerprint bits are classified as positive (dashed-line). B) Illustration of some False Positive classified substances and the most similar CMR substance. With an increase in the number of fingerprint bits, less ambiguous similarities are established.*

636
637 *Figure 3: Hierarchical clustering for the ED and PBT/vPvB subgroups based on single linkage method. For ED, the*

638 *FCFP4 fingerprint and SS3 coefficient are plotted, and for PBT/vPvB the MACCS fingerprint and SM coefficient.*

639 *The y-axis describes the dissimilarity between the SVHC structures and is equal to 1 minus the similarity. The blue*

640 *dotted line represents the used threshold (i.e. 1 minus threshold values). The red-colored boxes represent clusters of*

641 *similar substances. A) ED clusters. Five different clusters can be identified: 1 = Diosgenin, 2 = Phthalates, 3 =*

642 *Ethoxylated phenols, 4 = Nonyl and heptyl phenols, 5 = Octyl, pentyl and bi-phenols (Bisphenol A). B) PBT/vPvB*

643 *clusters. Thirty-two different clusters can be identified, including some large clusters: 1 = Phenolic benzotriazoles,*

644 *2 = Halogenated Dioxins, 3 = Chlorinated paraffins, 4 = Brominated diphenyl ethers, 5 = Perfluorinated*

645 *carboxylic acids, 6 = Polycyclic aromatic hydrocarbons, 7 = Halogenated dibenzofurans, 8 = Halogenated*

646 *aromatics and cycloalkanes.*

647