

# Machine learning models based on molecular descriptors to predict human and environmental toxicological factors in continental freshwater

Rémi Servien, Eric Latrille, Dominique Patureau, Arnaud Hélias

► **To cite this version:**

Rémi Servien, Eric Latrille, Dominique Patureau, Arnaud Hélias. Machine learning models based on molecular descriptors to predict human and environmental toxicological factors in continental freshwater. 2021. hal-03109818v2

**HAL Id: hal-03109818**

**<https://hal.inrae.fr/hal-03109818v2>**

Preprint submitted on 21 Jul 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Machine learning models based on molecular descriptors to predict**  
2 **human and environmental toxicological factors in continental**  
3 **freshwater**

4  
5 Rémi Servien<sup>a,b,\*</sup>, Eric Latrille<sup>a,b</sup>, Dominique Patureau<sup>a</sup>, Arnaud Hélias<sup>c,d</sup>

6  
7 <sup>a</sup>INRAE, Univ. Montpellier, LBE, 102 Avenue des étangs, F-11000 Narbonne, France

8 <sup>b</sup>ChemHouse Research Group, Montpellier, France

9 <sup>c</sup>ITAP, Univ Montpellier, INRAE, Institut Agro, Montpellier, France

10 <sup>d</sup>ELSA, Research group for environmental life cycle sustainability assessment and ELSA-  
11 Pact industrial chair, Montpellier, France

12 \*corresponding author : [remi.servien@inrae.fr](mailto:remi.servien@inrae.fr)

13  
14 **Highlights:**

- 15 • Characterization factors (for human health and ecotoxicological impacts) were  
16 predicted using molecular descriptors.
- 17 • Several linear or non-linear machine learning methods were compared.
- 18 • The non-linear methods tend to outperform the linear ones using a train and test  
19 procedure. Cluster-then-predict approaches often show the best performances,  
20 highlighting their usefulness.
- 21 • This methodology was then used to derive characterization factors that were missing  
22 for more than a hundred chemicals in USEtox®.

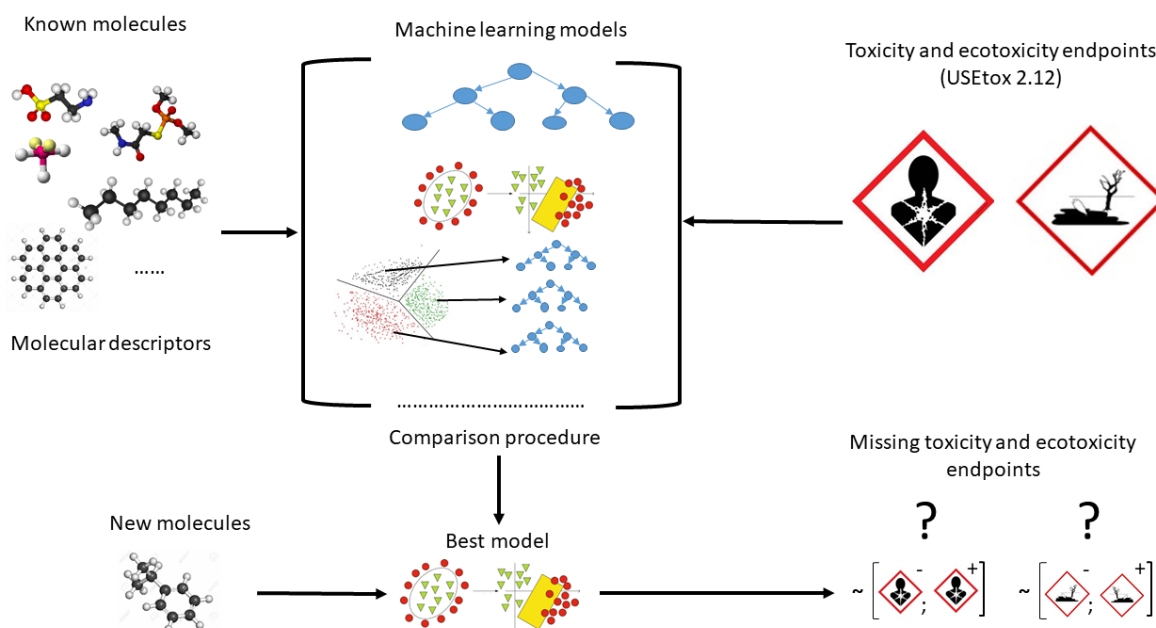
23  
24  
25 **Abstract:** It is a real challenge for life cycle assessment practitioners to identify all relevant  
26 substances contributing to the ecotoxicity. Once this identification has been made, the lack of

1 corresponding ecotoxicity factors can make the results partial and difficult to interpret. So, it is  
2 a real and important challenge to provide ecotoxicity factors for a wide range of compounds.  
3 Nevertheless, obtaining such factors using experiments is tedious, time-consuming, and made  
4 at a high cost. A modelling method that could predict these factors from easy-to-obtain  
5 information on each chemical would be of great value. Here, we present such a method, based  
6 on machine learning algorithms, that used molecular descriptors to predict two specific  
7 endpoints in continental freshwater for ecotoxicological and human impacts. The different  
8 tested machine learning algorithms show good performances on a learning database and the  
9 non-linear methods tend to outperform the linear ones. The cluster-then-predict approaches  
10 usually show the best performances which suggests that these predicted models must be  
11 derived for somewhat similar compounds. Then, predictions were derived from the validated  
12 model for compounds with missing toxicity/ecotoxicity factors.

13

#### 14 Graphical abstract:

15



16

17

1 **Keywords:** machine learning, Life Cycle Assessment, characterisation factors, toxicity,  
2 ecotoxicity, continental freshwater.

3

#### 4 **1. Introduction**

5

6 Recent legislations such as the Registration, Evaluation, Authorization and restriction of  
7 Chemicals (REACH) regulation in the EU requires that manufacturers of substances and  
8 formulators register to provide eco/toxicological data for substances with volume higher than  
9 one metric ton per year. As an example, the U.S. Environmental Protection Agency (EPA) has  
10 more than 85,000 chemicals listed under the Toxic Substances Control Act (Hinds and Weller,  
11 2016). So, robust (eco)-toxicological data are quickly needed to make informed decisions on  
12 how to regulate new chemicals. These data must also be coupled with environmental  
13 exposures and sources data, to better understand the impact on the environment.

14

15 To address the cause-effect relationships between the flow of molecules emitted by human  
16 activities and the consequences for ecosystems and humans, LCA offers a structured,  
17 operational, and standardized (Finkbeiner et al., 2006) methodological framework. Two main  
18 steps are at the core of this approach:

- 19 • Quantification of the masses of substances emitted into the environment through the  
20 Life Cycle Inventory (LCI). While it is possible to rely on databases that facilitate this  
21 inventory work for the background of the system under study, this task must  
22 nevertheless be carried out on a case-by-case basis to represent all the specificities  
23 of the foreground elements. To best describe human activities, their specificities must  
24 be represented on a case-by-case basis. This is the task of the LCA practitioner.
- 25 • Calculation of the impacts on ecosystems and human health of these emitted masses.  
26 Due to the complexity of environmental mechanisms, it is not possible to (re)model  
27 impact pathways on a case-by-case basis. Therefore, LCA uses characterization  
28 factors (CF) to assess the potential impacts of a compound. Concretely, if two

1 compounds are emitted with the same mass, the one with the higher CFs will have the  
2 higher impact. Then, CFs are multiplied by the emitted masses of each compound to  
3 determine the impacts. CFs are not recalculated for each study but provided within a  
4 Life Cycle Impact Assessment (LCIA) method.

5  
6 For a given impact, the LCIA method designer refers to the knowledge of the scientific  
7 community to model the mechanisms involved. For human toxicity and freshwater ecotoxicity,  
8 USEtox® (Rosenbaum et al., 2008), was developed by life cycle initiative under the United  
9 Nations Environmental Programme (UNEP) and the Society for Environmental Toxicology and  
10 Chemistry (SETAC) (Henderson et al. 2011) to produce a transparent and consensus  
11 characterization model. USEtox® is also used for the European Product Environmental  
12 Footprint (PEF) (Saouter et al., 2020). This model gathers in one single characterization factor  
13 the chemical fate, the exposure, and the effect for each of the several thousands of organic  
14 and inorganic compounds. Then, the USEtox® model results can be extended to determine  
15 endpoint effects expressed as total (i.e. cancer and non-cancer) disability-adjusted life years  
16 (DALY) for human health impacts and potentially disappeared fraction of species (PDF) for  
17 ecotoxicological impacts. The PDF represents an increase in the fraction of species potentially  
18 disappearing as a consequence of emission in a compartment while the DALY represents an  
19 increase in adversely affected life years. These endpoints are now consensual at an  
20 international level (Verones et al., 2017).

21 If the structure of the USEtox® multimedia model is always the same, to determine the CF of  
22 a molecule, numerous physico-chemical parameters (such as solubility, hydrophobicity,  
23 degradability) and detailed toxicological and ecotoxicological data must be provided. For  
24 example, EC50 values (i.e. the effective concentration at which 50% of a population died) for  
25 at least three species from three different trophic levels are required for the ecotoxicological  
26 effect factor.

27

1 Over the past few decades, thousands of tests (in laboratory and field) have been carried out  
2 to evaluate the potential hazard effects of chemicals (He et al., 2017). Usually, toxicity testing  
3 has relied on *in vivo* animal models, which is extremely costly and time-consuming (Xia et al.,  
4 2008). In recent years, under societal pressures, there has been a significant paradigm shift  
5 in toxicity testing of chemicals from traditional *in vivo* tests to less expensive and higher  
6 throughput *in vitro* methods (National Research Council, 2007). However, it is still extremely  
7 hard to test the number of existing and ever-increasing numbers of new chemicals, which  
8 leaves their impacts largely unknown. That's why more computational models are needed to  
9 complement experimental approaches to decrease the experimental cost and determine the  
10 prioritization for those chemicals which may need further *in vivo* studies. Such models already  
11 exist, like QSAR models that are mostly linear models based on the chemical structure of  
12 compounds (Danish QSAR database (DTU, 2015), ECOSAR (Mayo-Bean et al., 2011), VEGA  
13 (Benfenati et al., 2013)) and are used to predict ecotoxicological data (LC50) needed for  
14 REACH for example. Recently, machine learning algorithms have been used to predict  
15 hazardous concentration 50% (HC50) based on 14 physico-chemical characteristics (Hou et  
16 al., 2020a) or on 691 more various variables (Hou et al., 2020b). In the case of USEtox®,  
17 despite its wide use in LCA, it only offers characterization factors for approximately 3000  
18 chemicals and even for this limited number of compounds, 19% of ecotoxicity CFs and 67%  
19 of human toxicity CFs are missing.

20 The objective of this article is thus to propose a new way of calculating CFs using machine  
21 learning approaches to solve the problem of nonlinearity that could affect a linear QSAR  
22 method. This makes it possible, when the CFs are not determined due to lack of time or lack  
23 of data, to propose values based solely on easily identifiable molecular descriptors. Here, the  
24 main differences with the above-cited methods are twofold: first, our input variables are only  
25 molecular descriptors that could be easily collected for any newly available compounds;  
26 second, our output variables are directly the CFs that are closer to the endpoints (DALY and  
27 PDF) than the HC50 or the LC50 (i.e. the acute aquatic toxicity experimental threshold). These  
28 two specific endpoints will be studied in the present paper through the emission of compounds

1 in continental freshwater and will be named  $CF_{ET}$  for ecotoxicological impacts and  $CF_{HT}$  for  
2 human ones. To address this aim, we will test different methods (linear and non-linear) and  
3 assess their performances, to build a robust model that could predict CFs that are currently  
4 lacking.

5

6

## 7 **2. Materials & Methods**

### 8 **2.1. USEtox® database**

9

10 The last version of the USEtox® database was downloaded, namely the corrective release  
11 2.12 (USEtox®, 2020). The whole USEtox® 2.12 database contains 3076 compounds.

12

### 13 **2.2. TyPol database**

14

15 We recently developed TyPol (Typology of Pollutants), a classification method based on  
16 statistical analyses combining several environmental parameters (i.e., sorption coefficient,  
17 degradation half-life, Henry constant) and an ecotoxicological parameter (bioconcentration  
18 factor BCF), and structural molecular descriptors (i.e., number of atoms in the molecule,  
19 molecular surface, dipole moment, energy of orbitals). Molecular descriptors are calculated  
20 using an *in silico* approach (combining Austin Model1 and Dragon software). In the present  
21 paper, we only extract and use the molecular descriptors from the TyPol database, as this  
22 information could be easily collected for any new compound. The 40 descriptors included in  
23 the TyPol database have been selected based on a literature review on QSAR equations used  
24 to predict the main environmental processes as degradation, sorption, volatilization. These 40  
25 descriptors were the ones most frequently used in the equations, meaning describing the best  
26 the behaviour of organic compounds in the environment. By consequence, even if no  
27 environmental parameters are directly incorporated as input in our model, some information  
28 that are directly linked to them are included in the 40 molecular descriptors. These descriptors

1 are constitutional, geometric, topological, and quantum-chemical descriptors (see Table 1); 36  
 2 described the 2D-structure of the compound while the other four are linked to its 3D-structure.  
 3 An important advantage of the unique use of molecular descriptors is that they are easily and  
 4 quickly computable for not yet synthesized compounds. For more details, we refer the  
 5 interested reader to Servien et al. (2014) where the choice of the 40 molecular descriptors is  
 6 described in details. Now, TyPol gathers 526 compounds, including pesticides,  
 7 pharmaceuticals and their transformation products (Benoit et al. 2017, Traoré et al. 2018).

8

9 **Table 1** – List of the 40 molecular descriptors in TyPol

Category	Molecular descriptors		
Constitutional	Number of atoms	Number of non-H atoms	Number of hydrogen atoms
	Number of hydrogen atoms	Number of carbon atoms	Number of nitrogen atoms
	Number of oxygen atoms	Number of phosphorus atoms	Number of sulfur atoms
	Number of fluorine atoms	Number of chlorine atoms	Number of halogen atoms
	Number of bonds	Number of non-H bonds	Number of double bonds
	Number of triple bonds	Number of multiple bonds	Number of rotatable bonds
	Number of aromatic bonds	Sum of conventional bond order	Number of rings
	Number of circuits	Molecular weight	
Geometric	Connolly molecular surface area		



Topological	Connectivity index of order 0	Connectivity index of order 1	Connectivity index of order 2
	Connectivity index of order 3	Connectivity index of order 4	Connectivity index of order 5
	Valence connectivity index of order 0	Valence connectivity index of order 1	Valence connectivity index of order 2
	Valence connectivity index of order 3	Valence connectivity index of order 4	Valence connectivity index of order 5
Quantum-chemical	Polarizability	Electric dipole moment	HOMO energy
	LUMO energy	Total energy	

1

2

3

4

### 2.3. Machine learning methods

5

6 To predict the CFs using the molecular descriptors we use three modelling methods combined.

7 The first method is a linear well-known prediction method namely the Partial Least Squares

8 (PLS) (Wold, 1985). It finds the multidimensional directions in the observable variable

9 (molecular descriptor) space that explains the maximum multidimensional variance direction

10 in the predicted variable (CF) space. That provides a linear regression model based on the

11 observable variables to predict the predicted variable. We also choose to compare two

12 machine learning methods adapted to non-linear problems: the random forest (Breiman 2001)

13 and the support vector machines (SVM) (Drucker et al. 1996). Random forests are a machine

14 learning method, for classification or, in our case, regression, that operate by constructing a

15 multitude of decision trees that uses a random subset of the training data and limits the number

16 of variables used at each split and outputting the mean prediction (regression) of the individual

1 trees. SVM constructs a hyperplane or set of hyperplanes in a high- or infinite-dimensional  
2 space in which the problem is linearly separable.

3  
4 These choices allow us to compare several ideas. The PLS is a simple linear method that will  
5 not exhibit good performances if the underlying relationship is not linear. The SVM and RF  
6 methods are well-known non-linear machine learning algorithms that used to show good  
7 results in this kind of problem (Hou et al., 2020a).

8  
9 All the models were computed in the freeware R (R core team, 2019). The PLS has been  
10 computed using the package mixOmics (Rohart et al., 2017), the random forests using the  
11 package randomForest (Liaw et al., 2002), and the SVM using the package e1071 (Meyer et  
12 al., 2019). These 3 modelling methods have some parameters that needed to be fixed: the  
13 number of latent components for the PLS (fixed using the tune.pls function), the number of  
14 variables randomly sampled as candidates at each split for the random forests (selected using  
15 the tune.randomForest function) and, for the SVM, the gamma parameter of the radial kernel  
16 and the cost of constraints violation (using the tune.svm function). All these different tune  
17 functions are based on cross-validation (i.e. a training/test procedure to find the best value for  
18 the parameters) using default function values.

19

## 20 **2.4. Clustering-based model**

21

22 A recent popular way to make predictions is to use a cluster-then-predict approach. That is,  
23 clustering is used for pre-classification which is to arrange a given collection of input patterns  
24 into natural meaningful clusters. Then, the clustering results are used to construct a predictor  
25 in each cluster. The main idea of the cluster-then-predict approach is that if the clustering  
26 performs well the prediction will be easier by modelling only similar compounds. If a new  
27 compound with no  $CF_{ET}$  and/or  $CF_{HT}$  is investigated, the clustering can easily be applied to it  
28 before the prediction model itself. The cluster-then-predict approach has already been applied

1 with success in various domains such as sentiment prediction (Sony et al., 2015), finance  
2 (Tsai et al., 2014), chemometrics (Minh Mai Le et al., 2018). So we decided to use the  
3 clustering given by the TyPol application (more details in Servien et al., 2014) based on the  
4 whole database and the molecular descriptors. Note that the TyPol clustering has already  
5 been shown relevant on various occasion: in combination with mass spectrometry to  
6 categorize tebuconazole products in soil (Storck et al., 2016), to explore the potential  
7 environmental behaviour of putative chlordecone transformation products (Benoit et al., 2017)  
8 or to classify pesticides with similar environmental behaviours (Traore et al., 2018). So, the  
9 clustering procedure of TyPol was applied on the whole database of 526 compounds using  
10 the 40 molecular descriptors. This approach provides us a global clustering based on all the  
11 available information contained in the TyPol database. It is based on PLS, hierarchical  
12 clustering and an optimal choice of the number of clusters and is detailed in Servien et al.  
13 (2014). The obtained clustering is given in Supplementary Figure S1 and relies on 5 different  
14 clusters.

15 Based on this clustering, we then define three other competing methods. For these methods,  
16 a different model (with different parameters) will be derived for the compounds in each cluster.  
17 Consequently, six different models will be calibrated and tested for each CF prediction: global  
18 PLS, global SVM, global random forest, cluster-then-PLS, cluster-the-SVM and cluster-then-  
19 Random Forest.

20

## 21 **2.5. Comparison procedure**

22

23 To assess the performances of the different models we will use the following procedure:

24 1. Split each cluster between a training set (85% of the dataset) and a test set (15%).

25 The test set is not used for any step of the procedure (such as the imputation of the  
26 missing data, the calibration of the parameters ...).

27 2. Imputation of the NA values (less than 1%) in the descriptor matrix using the NIPALS  
28 algorithm (Wold, 1985).

- 1 3. Tune the parameters and train the specific models by performing cross-validation on
- 2 the training set. We have 3 global models to train (PLS, random forest, and SVM) and
- 3 the cluster-then-test models (PLS, random forest and SVM for each cluster).
- 4 4. Test the different models on the test set. Compute the absolute error.
- 5 5. Back to step 1.

6

7 The Typol clustering focused on the common compounds is plotted in Supplementary Figure

8 S2 and we see that, for cluster 5, the 3 global models are the only ones available as we can't

9 define a cluster-then-test model due to a lack of data (only one compound remaining). The

10 whole algorithm is repeated 200 times. All the performances are compared in terms of

11 absolute error. The absolute error is the absolute difference between the prediction and the

12 true value. It has been shown to be the most natural and unambiguous measure of error

13 (Willmott et Matsuura, 2005) and is chosen to be easily comparable to the assumed error on

14 the experimental CFs (2-3 logs, see Rosenbaum, 2008). For each cluster, we chose the model

15 with the lowest median absolute error.

16

17 Then, the best model is calibrated and computed on the whole cluster. Finally, it is applied to

18 the compounds, according to their clusters, with a  $CF_{ET}$  (or a  $CF_{HT}$ ) equals to NA to provide a

19 prediction. For the compounds in cluster 5, this best model cannot be a cluster-then-predict

20 one and, by consequence, is a global one. To assess the robustness of our prediction we

21 derive a 95% prediction interval for each prediction. The type of model and its corresponding

22 parameters are fixed during this process, according to the best model of the cluster. For

23 example, if the best model of cluster 1 was the random forest approach, random forest models

24 are used with the parameters optimized during the previous step. Then, we perform a leave-

25 one-out bootstrap on the dataset that was used to compute the model (the whole dataset if

26 the model is global, only the data lying in the dedicated cluster if that is a cluster-then-predict

27 model) and a new model is computed on this leave-one-out sample. A prediction is carried for

1 each leave-one-out model and the 2.5% and 97.5% quantile of these predictions are computed  
2 and considered as the prediction interval (Hou et al., 2020a).

3  
4 The five descriptors contributing the most to the prediction are then derived for each chosen  
5 model to assess the differences between models and to interpret their relevance. For a  
6 random forest model, these descriptors are calculated using variable permutations (Breiman,  
7 2001), for the SVM they are the descriptors with the higher coefficients in absolute value.

8

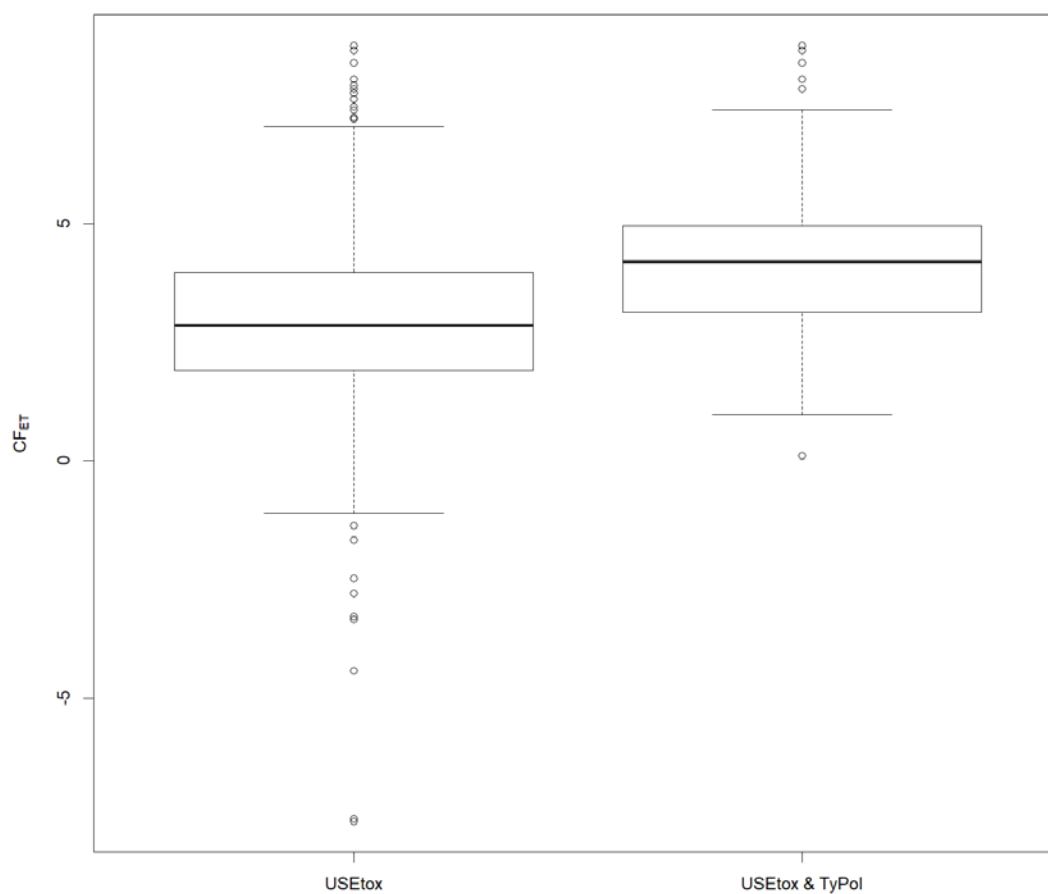
### 9 **3. Results**

10

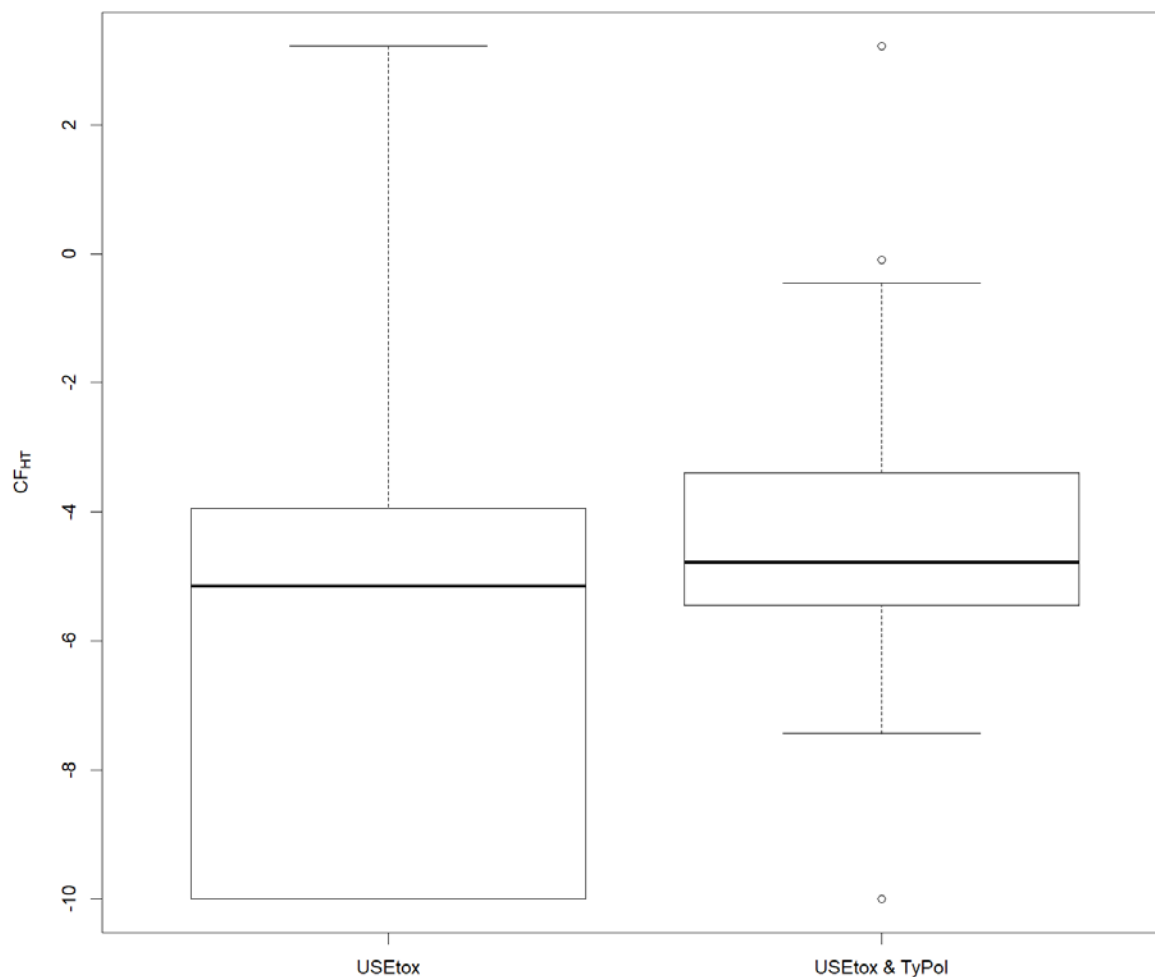
#### 11 **3.1. Descriptive analysis of the intersection of the TyPol and the USEtox®** 12 **databases**

13

14 As the objective of this proof-of-concept study is to predict USEtox®  $CF_{ET}$  and  $CF_{HT}$  using the  
15 molecular descriptors contained in TyPol, we could only use the compounds that are present  
16 in both databases. This results in 274 compounds that are detailed in Table S1 in  
17 supplementary material and the range of their  $CF_{ET}$  and  $CF_{HT}$  values are summarized in the  
18 boxplots in Figures 1 and 2. Note that for the 274 common compounds there are 15 NA values  
19 for the  $CF_{ET}$  and 102 for the  $CF_{HT}$ .



1  
 2 **Figure 1-** Boxplots of the  $CF_{ET}$  for the USEtox® database and the common molecules  
 3 between the USEtox® and the TyPol databases. This  $CF_{ET}$  is equal to the  $\log_{10}(\text{PDF} \cdot \text{m}^3 \cdot \text{d} \cdot \text{kg}^{-1})$ .  
 4 1).



1

2 **Figure 2** Boxplots of the  $CF_{HT}$  for the USEtox® database and the common molecules between  
 3 the USEtox® and the TyPol databases. This  $CF_{HT}$  is equal to  $\log_{10}((DALY+\epsilon).kg^{-1})$ . The  $\epsilon$  is  
 4 needed as some values of the DALY are exactly equal to zero.  $\epsilon$  has been chosen equal to  
 5  $1e-10$  to be below the minimum of the USEtox® database ( $5e-9$ ).

6

7 We could see on these two figures that the common compounds present higher  $CF_{ET}$  and  
 8  $CF_{HT}$  values than the one of the complete USEtox® database: it focuses on the more  
 9 dangerous compounds as their boxplots are above the USEtox® counterparts.

10

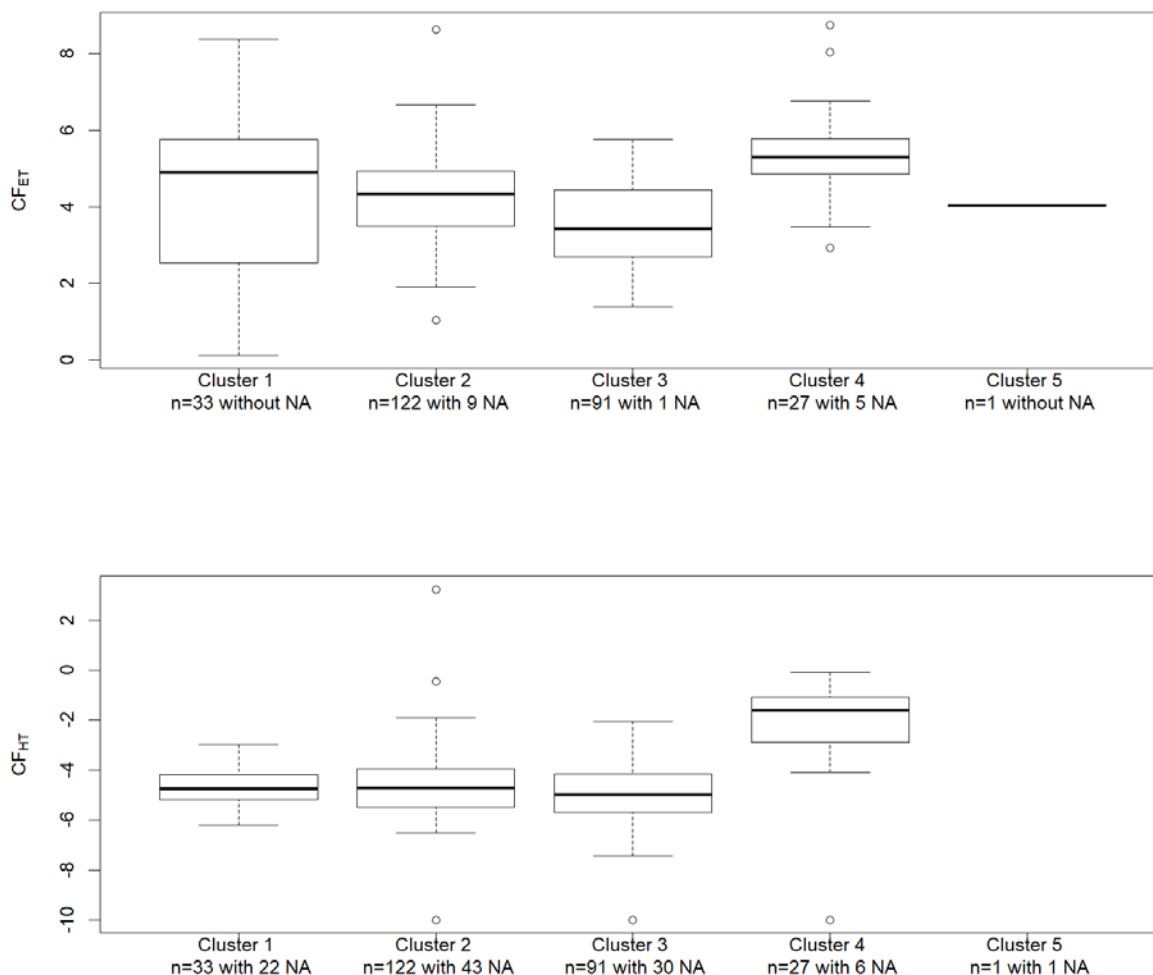
1 The Typol clustering focused on the common compounds is plotted in Supplementary Figure  
2 S2 and the boxplots of each molecular descriptor per cluster are given in Supplementary  
3 Figure S3 with different indicators in Table S2. We could see that they are clustered in 5 groups  
4 with different sizes (respectively 33 compounds in the first black cluster, 122 compounds in  
5 the second red cluster, 91 compounds in the third green cluster, 27 compounds in the fourth  
6 blue cluster, and one compound in the fifth brown cluster). Cluster 1 grouped compounds with  
7 a high number of aromatic bonds, double bonds, rotatable bonds, and multiple bonds. Cluster  
8 2 is an intermediate one between clusters 1 and 3, with less extreme values. Cluster 3 is made  
9 of compounds with the lowest molecular mass. Cluster 4 gathered compounds presenting a  
10 high number of halogens, rings, and circuits. The unique compound in the fifth cluster is  
11 erythromycin (highest molecular mass and number of H and C, lowest number of rings) and,  
12 obviously, no cluster-then-predict model could be built for this cluster

13

14 As a first analysis of the clustering given by TyPol, we could see in Figure 3 below the boxplots  
15 of the  $CF_{ET}$  and  $CF_{HT}$  within the 5 clusters.

16





1

2 **Figure 3-** Boxplot by cluster for the CF<sub>ET</sub> and CF<sub>HT</sub> values. Note that the unique compound of  
 3 Cluster 5 has no CF<sub>HT</sub> value. The size of the clusters and the numbers of NA are gathered in  
 4 the legend.

5

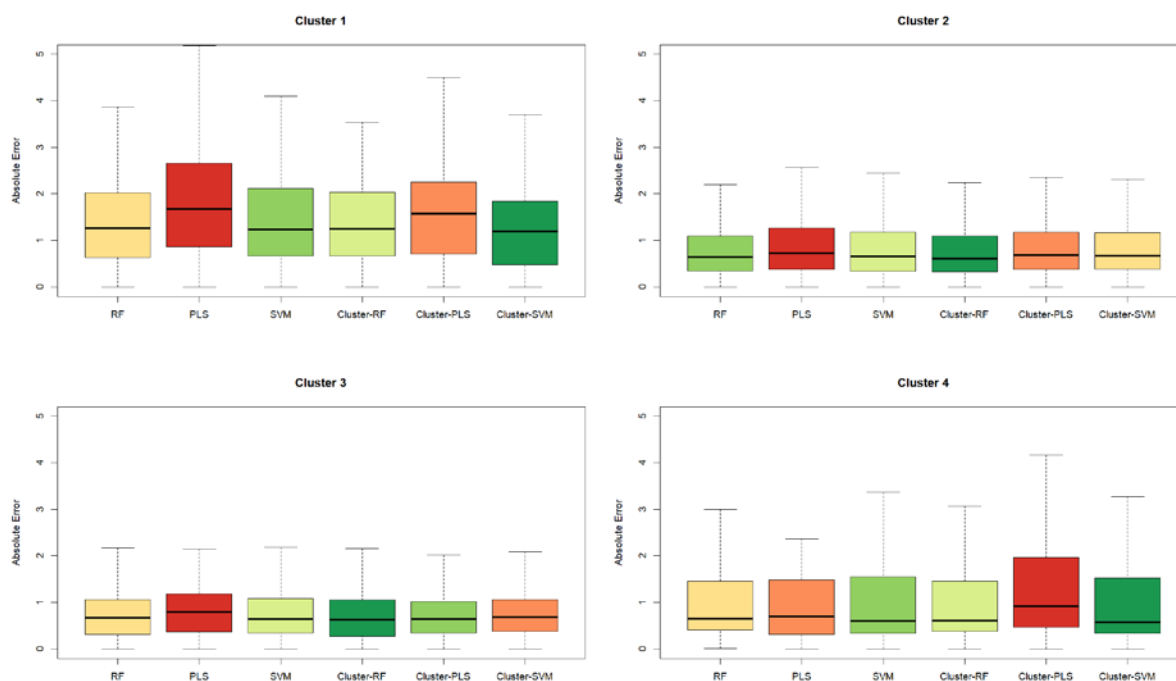
6 The predictions will be made difficult for the CF<sub>ET</sub> of cluster 1 as it covers a wide range whereas  
 7 it includes a relatively small number of compounds. On the contrary, cluster 3 covers a small  
 8 range with no extreme values and includes a high number of compounds, for this cluster the  
 9 cluster-then-predict approach could produce interesting results.

10

11 **3.2. Models and prediction of the CF<sub>ET</sub>**

### 1                    3.2.1.    Performances of the machine learning methods

2    The methodology described in the previous section was applied to our dataset and gave the  
3    results gathered in Figure S4 for the global results and in Figure 4 for the results detailed on  
4    each cluster.



5  
6  
7    **Figure 4** - Performances of the different methods in terms of the log of the absolute error of  
8    the  $CF_{ET}$  with respect to the different clusters. In each cluster, the models are coloured from  
9    green (best) to red (worst) according to their median of the absolute error.

10  
11    The performances are not similar in each cluster. For example, performances of all methods  
12    for cluster 1 are very poor (median absolute error above 1) whereas performances for cluster  
13    4 seem good despite its smallest size (median absolute error around 0.6). So, a future  
14    prediction of an unknown compound which lies in cluster 1 will be less reliable than in other  
15    clusters. Note that we could not test this in the next section as no NA value is present in this  
16    cluster 1.

17

1 The cluster-then-predict methods seem more appropriate in each cluster. The cluster-then-RF  
 2 approach has the best performances (with a global median absolute error equals to 0.64 and  
 3 the best performances on clusters 2 and 3), even if there is not a big difference between the  
 4 different methods. The cluster-then-SVM is also the best method for the two clusters 1 and 4.  
 5 The linear methods (PLS and cluster-then-PLS) have higher absolute errors but are  
 6 competitive. The individual predictions of the best method in each cluster are reported in  
 7 Figure S5.

### 9 3.2.2. Prediction with the best model

10  
 11 Then we apply the best model in each cluster: a cluster-then-predict approach using SVM for  
 12 clusters 1 and 4 and using random forest for clusters 2 and 3. To compare the different models  
 13 in each cluster and give an idea of what are the important molecular descriptors we provide  
 14 the five most important molecular descriptors for each cluster in the following table.

15  
 16 **Table 1**- The five most important molecular descriptors for each best model for each cluster.  
 17 The most important descriptors are in the first line of the table.

Cluster 1: cluster-then-SVM model	Cluster 2: cluster-then-RF model	Cluster 3: cluster-then-RF model	Cluster 4: cluster-then-SVM model
HOMO energy	Number of Chlorine atoms	Number of triple bonds	Number of double bonds
Molecular surface area	Number of halogen atoms	Molecular mass	Number of Nitrogen atoms
Number of Sulfur atoms	Number of Oxygen atoms	Number of Phosphorus atoms	HOMO energy
Connectivity index chi-5	Molecular mass	Number of Oxygen atoms	Number of triple bonds

Connectivity index chi- 3	Number of bonds	Number of halogen atoms	Electric dipole moment
------------------------------	-----------------	----------------------------	------------------------

1

2 We could see in this Table that the important molecular descriptors strongly differ from one  
3 cluster to another, highlighting the usefulness of the cluster-then-predict approaches.

4

5 Then the models were used to predict the missing CF<sub>ET</sub> of the common compounds between  
6 USEtox® and TyPol databases. These values are by consequence new estimations of the  
7 CF<sub>ET</sub> for compounds on which we have no information. The prediction intervals are relatively  
8 small: less than 0.5 log<sub>10</sub> in a log scale which highlights the robustness of the estimation. They  
9 are given in Table S3. No NA value was present in cluster 1 with no prediction for this cluster.  
10 For cluster 2 gathering molecules with intermediate molecular mass, 9 CF<sub>ET</sub> values were  
11 predicted for various kinds of compounds. One value concerns the antibiotic sulfamethazine  
12 and its value is quite near to the one of sulfamethoxazole and sulfadiazine of the same  
13 sulphonamide antibiotic family constituted of the sulphonamide group (-S(=O)<sub>2</sub>-NR<sub>2</sub>R<sub>3</sub>). Cluster  
14 3 grouped compounds with the lowest molecular mass and the lowest median CF<sub>ET</sub> like  
15 ibuprofen, phthalates, cresol constituted of monoaromatic ring substituted with methyl,  
16 carboxylic groups. The CF<sub>ET</sub> prediction for acetylsalicylic acid seemed coherent with the value  
17 of the nearest compounds (herbicides mecoprop) of this group. Cluster 4 gathered compounds  
18 with the highest median CF<sub>ET</sub> and that presented a high number of rings halogenated or not,  
19 like PAH and hormones. The 5 CF<sub>ET</sub> predicted concerned 4 PAHs and 1 hormone. By  
20 comparison to the 2 other PAHs present in this cluster, the 4 predicted CF<sub>ET</sub> are quite similar  
21 and higher. Concerning the prediction for the hormone, the CF<sub>ET</sub> is intermediate between the  
22 CF<sub>ET</sub> of the 3 other hormones in the cluster. It seems that all these 5 predicted values are very  
23 closed, falling near the median value of this cluster.

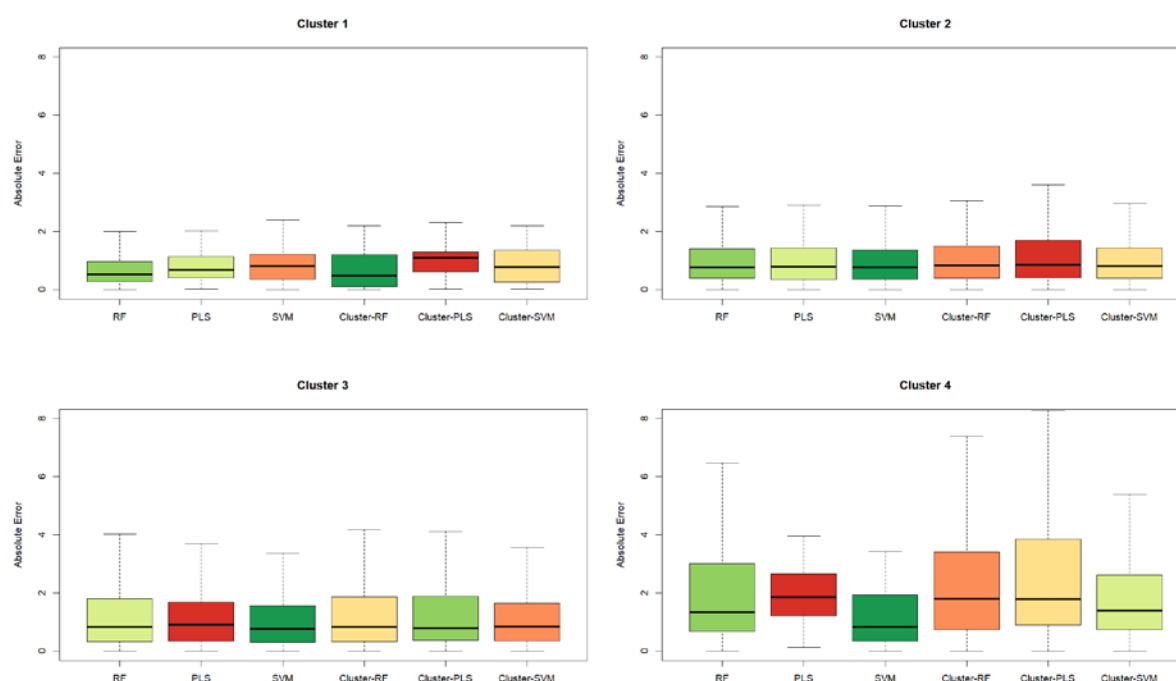
24

25

### 3.3. Models and prediction of the $CF_{HT}$

#### 3.3.1. Performances of the methods

Let us recall that we have more NA values for the  $CF_{HT}$  (102) than for the  $CF_{ET}$  (15). The performances of the methods are illustrated in the following figure.



**Figure 5-** Performances of the different methods in terms of the log of the absolute error of the  $CF_{HT}$  with respect to the different clusters. In each cluster, the models are coloured from green (best) to red (worst) according to their median of the absolute error.

We observe that, despite its small size (11 compounds), the  $CF_{HT}$  of the first cluster are well predicted (with the best performance for the cluster-then-RF approach). It could be explained by the small range of the  $CF_{HT}$  values of this cluster, as illustrated on the boxplot in Figure 3. The performances of all the methods are comparable on clusters 2 and 3 where the best method is the SVM. Cluster 4 seems to be the more difficult to predict: all the methods have their worst results on this cluster and, if the SVM has an acceptable median absolute error of

1 0.82, all the medians of the other methods are above 1.3. Global performances of the different  
2 methods are given in Supplementary Figure S6. Note that, as for  $CF_{ET}$ , the linear methods  
3 based on PLS are outperformed by the other ones.

4

### 5 **3.3.2. Prediction with the best model**

6

7 The global SVM model was then calibrated and computed on the whole dataset. It was then  
8 used to predict the compound of clusters 2, 3, 4, and 5. Let us recall that there is a lonely  
9 molecule in cluster 5 and, as it has a NA value for its  $CF_{HT}$ , the best global model (SVM) is  
10 used. For cluster 1, a cluster-then-RF model is computed. The more important descriptors of  
11 these two models are gathered in the following table.

12

13 **Table 2-** Five most important molecular descriptors for each best model for each cluster.

14 The most important descriptors are in the first line of the table.

Cluster 1 : cluster-then-RF model	Cluster 2, 3, 4 and 5: SVM model
Number of Fluorine atoms	Number of halogen atoms
Connectivity index chi-5	Electric dipole moment
Connectivity index chi-1	Number of double bonds
Number of circuits	Number of Chloride atoms
Number of rings	Number of Oxygen atoms

15

16 Then, this model was used to predict the  $CF_{HT}$  value for the 102 common compounds without  
17 a  $CF_{HT}$  value. These predictions are reported in Supplementary Table S4. As for the  $CF_{ET}$ , the  
18 small width of the prediction interval (less than a  $\log_{10}$  in a log scale) highlights the robustness  
19 of the approach even with a relatively small number like estimations made for compounds that  
20 lie in cluster 1. In this cluster 1,  $CF_{HT}$  for a phthalate (DEHP) is already known, but the one for

1 diisodecyl and diisononyl phthalate was predicted with value in the same range. The 3 cyclines  
2 (tetracycline, aureomycin, and oxytetracycline) present in cluster 1, presented also similar  
3 predicted  $CF_{HT}$ . This was also the case for triclosan and triclocarban in cluster 2. Similar  
4 predicted and known  $CF_{HT}$  were found for four herbicides from the substituted urea family  
5 (linuron, diuron, monolinuron, isoproturon) in cluster 3. Cluster 4 gathered a small number of  
6 molecules but with the highest median  $CF_{HT}$ , the predicted  $CF_{HT}$  of the organochlorine  
7 insecticide isodrin was similar to another congener of the same family, aldrin.

8

#### 9 **4. Discussion**

10

11 It is a real and important challenge to provide characterization factors for a wide range of  
12 compounds. Obviously, it is expected that these new calculated factors have an acceptable  
13 margin of error. As reported in UNEP/SETAC (2019), it is commonly assumed that the  
14 uncertainty of the characterization factors can vary by approximately 2-3 orders of log-  
15 magnitude (Rosenbaum et al. 2008) or significantly higher (up to 7 orders) if all sources of  
16 uncertainty are considered (Douziech et al. 2019). Using our methodology, we can exhibit a  
17 median absolute error of 0.62 log for the prediction of the  $CF_{ET}$  and 0.75 log for the prediction  
18 of the  $CF_{HT}$ . These results are very promising as they are below the level of uncertainty  
19 commonly assumed and as they are based on molecular descriptors that could be easily  
20 obtained for each compound without ecotoxicity factor. Based on this fact we could already  
21 provide 15 new  $CF_{ET}$  and 102 new  $CF_{HT}$  for the common molecules between USEtox® and  
22 TyPol without a previous value.

23

24 The idea of predicting ecotoxicity characterization factors for chemicals using machine  
25 learning algorithms has already been used (Hou et al., 2020a and 2020b). But, here, our  
26 findings go further. Indeed, we show that we could directly obtain accurate estimations of  
27 endpoint values from easy-to-obtain molecular descriptors. This will open the door to the fast  
28 characterization of each new unknown compound that appears, including transformation

1 products. We also show that the cluster-then-predict approach can give better performances  
2 than the usual ones. This local approach confirms that local models could be an efficient  
3 prediction method when heterogeneity of data generates nonlinear relations between the  
4 response and the explicative variables (Lesnoff et al., 2020).

5 Across the clusters and models, there is a general trend that the non-linear models tend to  
6 outperform the linear ones. This suggests that a linear model is not fully adequate to capture  
7 the complexity of the relationship between the molecular descriptors and the CFs. However,  
8 the use of linear model for e.g. a QSAR is likely due to the ease of interpreting its coefficients,  
9 while interpretation is much more challenging for machine learning approaches such as  
10 random forest or SVM. Thus, the advantages or drawbacks of linear/non-linear approaches  
11 must be balanced according to the final goal of each study. Here, as the main goal is to  
12 calculate the most accurate CFs, non-linear models seem more suited. We must also mention  
13 that a new emerging field is developing tools needed to help make black-box models (e.g.  
14 random forest) more interpretable (Bénard et al., 2021).

15 The difficult interpretability of the machine learning models used in this study can thus be  
16 viewed as a limitation. On another side, even if we already had an acceptable number of  
17 compounds on our training datasets, the model accuracies would benefit of the inclusion of  
18 new compounds. These compounds could be carefully chosen to improve the models where  
19 there is a clear need (i.e. where the performances of the models are not good enough), for  
20 example in the cluster 1 for  $CF_{ET}$  or in the cluster 4 for  $CF_{HT}$ .

21 One of the interests of USEtox® and its three-step structure (fate - exposure - effect) is that it  
22 can be adapted to some specific contexts (a more accurate and spatialised fate model, a  
23 different exposure...) while keeping the steps that are not modified. However, these  
24 adaptations of USEtox® are not widely used and are reserved for advanced users. Our  
25 approach does not allow this, with a direct one-step estimation of CFs. It was designed to  
26 provide default CF values for molecules where information is missing. We have chosen to  
27 directly predict the CF by simplicity, as the first tests revealed that doing three models (for the  
28 three steps) and then calculating the CFs produced less accurate results. It would however be



1 an interesting perspective to estimate only some of the stages by these learning approaches  
2 and to combine them with stages modelled in a classical way in USEtox®

3

## 4 **5. Conclusion**

5

6 In a recent study, Aemig et al. (2021) studied the potential impacts on Human health and  
7 aquatic environment of the release of 286 micropollutants (organic and inorganic) at the scale  
8 of France. One of their conclusion was that, due to a lack of characterization factors, these  
9 impacts could be assessed only for 1/3 of these molecules. This paper fills this gap by  
10 providing a new modelling method to derive characterization factors from easily obtainable  
11 molecular descriptors. The results presented here show that models that can handle non-  
12 linearity and that could be adapted to a small number of compounds (using the cluster-then-  
13 predict approaches) are the best suited. By consequence, the missing characterization factors,  
14 as well as those of new molecules, could now be quickly estimated with an overall good  
15 precision. More generally, one of the key factors in the evaluation of toxicity and ecotoxicity in  
16 LCA lies in the construction of the characterization factors: a task requiring a large amount of  
17 data and a consequent investment of time. The use of machine learning allows us to go  
18 beyond these constraints. This makes it possible to obtain characterization factor values in a  
19 fast and simple way, which can be used as long as conventionally established CFs are not  
20 available.

21

## 22 **Declaration of Competing Interest**

23

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

26

## 27 **Acknowledgments**

28

1 The authors are grateful to Pierre Benoit, Laure Mamy, and Virginie Rossard for their work  
2 on TyPol.

3

#### 4 **Funding**

5

6 This research did not receive any specific grant from funding agencies in the public,  
7 commercial, or not-for-profit sectors.

8

#### 9 **Supplementary materials**

10

11 Supplementary material associated with this article can be found, in the online version.

12

#### 13 **References**

14

15 Aemig, Q., Hélias, A., Patureau, D., 2021. Impact assessment of a large panel of organic and  
16 inorganic micropollutants released by wastewater treatment plants at the scale of France,  
17 *Water Research*, 188, 116524, <https://doi.org/10.1016/j.watres.2020.116524>.

18

19 Bénard, C., Biau, G., da Veiga, S., Scornet, E, 2021. Interpretable random forests via rule  
20 extraction. In *International Conference on Artificial Intelligence and Statistics*, vol. 130 of  
21 *Proceedings of Machine Learning Research*, 937–945 (PMLR, 2021).

22

23 Benfenati, E., Manganaro, A., Gini, G.C., 2013. VEGA-QSAR: AI Inside a Platform for  
24 Predictive Toxicology. *CEUR Workshop Proceedings*, 21-28.

25

26 Benoit, P., Mamy, L., Servien, R., Li, Z., Latrille, E., Rossard, V., Bessac, F., Patureau, D.,  
27 Martin-Laurent, F., 2017. Categorizing chlordecone potential degradation products to explore

1 their environmental fate, *Science of the Total Environment*, 574, 781–795.  
2 <https://doi.org/10.1016/j.scitotenv.2016.09.094>.

3

4 Breiman, L., 2001. Random Forests, *Machine Learning*, 45 (1), 5–32.  
5 <https://doi.org/10.1023/A:1010933404324>.

6

7 Cortes, C., Vapnik, V., 1995. Support-vector networks, *Machine Learning*, 20 (3), 273–297.  
8 <https://doi.org/10.1007/BF00994018>.

9

10 Douziech, M., Oldenkamp, R., van Zelm, R., King, H., Hendriks, A.J., Ficheux, A.-S.,  
11 Huijbregts, M.A.J., 2019. Confronting variability with uncertainty in the ecotoxicological impact  
12 assessment of down-the-drain products, *Environment International*, 126, 37-45,  
13 <https://doi.org/10.1016/j.envint.2019.01.080>.

14

15 Drucker, H., Burges, C.C., Kaufman, L., Smola, A.J., Vapnik, V., 1997, Support Vector  
16 Regression Machines, *Advances in Neural Information Processing Systems 9*, NIPS, 155–  
17 161, MIT Press. <https://dl.acm.org/doi/10.5555/2998981.2999003>.

18

19 DTU, 2015. Danish QSAR database. Danish QSAR group, National Food Institute, Technical  
20 University of Denmark.

21 Finkbeiner, M., Inaba, A., Tan, R., Christiansen, K., Klüppel, H.-J., 2006. The New  
22 International Standards for Life Cycle Assessment: ISO 14040 and ISO 14044. *The*  
23 *International Journal of Life Cycle Assessment*, 11 (2), 80–85.  
24 <https://doi.org/10.1065/lca2006.02.002>.

25 He, J., Tang, Z., Zhao, Y., Fan, M., Dyer, S. D., Belanger, S. E., Wu, F., 2017. The combined  
26 QSAR-ICE models: practical application in ecological risk assessment and water quality

1 criteria, Environmental Science & Technology, 51, 8877.  
2 <https://doi.org/10.1021/acs.est.7b02736>.

3

4 Henderson, A.D., Hauschild, M.Z., Van De Meent, D., Huijbregts, M.A.J., Larsen, H.F., Margni,  
5 M., McKone, T.E., Payet, J., Rosenbaum, R.K., Jolliet O., 2011. USEtox® fate and ecotoxicity  
6 factors for comparative assessment of toxic emissions in life cycle analysis: sensitivity to key  
7 chemical properties, The International Journal of Life Cycle Assessment, 16, pp. 701-709  
8 <https://doi.org/10.1007/s11367-011-0294-6>.

9

10 Hinds, R.d.C., Weller, J.L., 2016. Toxic Substances Control Act. Environmental Law Practice  
11 Guide, vol. 4.

12

13 Hou, P., Jolliet, O., Zhu, J., Xu, M., 2020a. Estimate ecotoxicity characterization factors for  
14 chemicals in life cycle assessment using machine learning models. Environment International,  
15 135, 105393. <https://doi.org/10.1016/j.envint.2019.105393>.

16

17 Hou, P., Zhao, B., Jolliet, O., Zhu, J., Wang, P., Xu, M., 2020b. Rapid Prediction of Chemical  
18 Ecotoxicity Through Genetic Algorithm Optimized Neural Network Models, ACS Sustainable  
19 Chemistry & Engineering, 8 (32), 12168-12176.  
20 <https://dx.doi.org/10.1021/acssuschemeng.0c03660>.

21

22 Lesnoff, M., Metz, M., Roger, JM., 2020. Comparison of locally weighted PLS strategies for  
23 regression and discrimination on agronomic NIR data, Journal of Chemometrics, 34(5), e3209,  
24 <https://doi.org/10.1002/cem.3209>.

25

26 Liaw, A., Wiener, M., 2002. Classification and Regression by randomForest, R News, 2(3),  
27 18-22. <http://CRAN.R-project.org/doc/Rnews/>.

28

1 Mayo-Bean, K., Nabholz, J., Clements, R., Zeeman, M., Henry, T., Rodier, D., Moran, K.,  
2 Meylan, B., Ranslow, P., 2011. Methodology document for the ECOlogical Structure-Activity  
3 Relationship Model (ECOSAR) class program: estimating toxicity of industrial chemicals to  
4 aquatic organisms using ECOSAR class program (Ver. 1.1). In: US Environmental Protection  
5 Agency, Office of Chemical Safety and Pollution Prevention, Office of Pollution Prevention and  
6 Toxics, Washington, DC.

7

8 Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A., Leisch, F., 2019. e1071: Misc Functions  
9 of the Department of Statistics, Probability Theory Group (Formerly: E1071), TU Wien. R  
10 package version 1.7-2. <https://CRAN.R-project.org/package=e1071>.

11

12 Minh Mai Le, L., Kégl, B., Gramfort, A., Marini, C., Nguyen, D., Cherti, M., Tfaily, S., Tfayli, A.,  
13 Baillet-Guffroy, A., Prognon, P., Chaminade, P., Caudron, E., 2018. Optimization of  
14 classification and regression analysis of four monoclonal antibodies from Raman spectra using  
15 collaborative machine learning approach, *Talanta*, 184, 260-265,  
16 <https://doi.org/10.1016/j.talanta.2018.02.109>.

17

18 National Research Council, 2007. Toxicity Testing in the 21st Century: A Vision and a  
19 Strategy; National Academies Press, <https://doi.org/10.17226/11970>.

20

21 R Core Team, 2019. R: A Language and Environment for Statistical Computing. R Foundation  
22 for Statistical Computing, Vienna, Austria. <http://www.r-project.org/index.html>.

23

24 Rohart, F., Gautier, B., Singh, A., Le Cao, K.-A., 2017. mixOmics: An R package for omics  
25 feature selection and multiple data integration, *PLoS computational biology*, 13(11),  
26 e1005752. <https://doi.org/10.1371/journal.pcbi.1005752>.

27

1 Rosenbaum, R.K., Margni, M., Jolliet, O., 2007. A flexible matrix algebra framework for the  
2 multimedia multipathway modelling of emission to impacts, *Environment International*,  
3 33(5),624-634. <https://doi.org/10.1016/j.envint.2007.01.004>.  
4  
5 Rosenbaum, R.K., Bachmann, T. M., Gold, L. S., Huijbregts, M.A.J., Jolliet, O., Juraske, R.,  
6 Koehler, A., Larsen, H.F., MacLeod, M., Margni, M., McKone, T.E., Payet, J., Schuhmacher,  
7 M., van de Meent, D., Hauschild, M.Z., 2008. USEtox®—the UNEP-SETAC Toxicity Model:  
8 Recommended Characterisation Factors for Human Toxicity and Freshwater Ecotoxicity in  
9 Life Cycle Impact Assessment, *The International Journal of Life Cycle Assessment*, 13 (7),  
10 532–546. <https://doi.org/10.1007/s11367-008-0038-4>.  
11  
12 Saouter, E., Biganzoli, F., Ceriani, L., Versteeg, D., Crenna, E., Zampori, L., Sala, S., Pant,  
13 R., 2020. Environmental Footprint: Update of Life Cycle Impact Assessment Methods –  
14 Ecotoxicity freshwater, human toxicity cancer, and non-cancer, Publications Office of the  
15 European Union, Luxembourg, <https://doi.org/10.2760/300987>.  
16  
17 Servien, R., Mamy, L., Li, Z., Rossard, V., Latrille, E., Bessac, F., Patureau, D., Benoit, P.,  
18 2014. TyPol - a new methodology for organic compounds clustering based on their molecular  
19 characteristics and environmental behaviour, *Chemosphere*, 111, 613–622.  
20 <https://doi.org/10.1016/j.chemosphere.2014.05.020>.  
21  
22 Soni, R., Mathai, K.J., 2016. An Innovative ‘Cluster-then-Predict’ Approach for Improved  
23 Sentiment Prediction. In: Choudhary R., Mandal J., Auluck N., Nagarajaram H. (eds)  
24 Advanced Computing and Communication Technologies. *Advances in Intelligent Systems and*  
25 *Computing*, vol 452. Springer, Singapore. [https://doi.org/10.1007/978-981-10-1023-1\\_13](https://doi.org/10.1007/978-981-10-1023-1_13).  
26  
27 Storck, V., Lucini, L., Mamy, L., Ferrari, F., Papadopoulou, E.S., Nikolaki, S., Karas, P.A.,  
28 Servien, R., Karpouzas, D.G., Trevisan, M., Benoit, P., and Martin-Laurent, F, 2016.

1 Identification and characterization of tebuconazole transformation products in soil by  
2 combining suspect screening and molecular typology, *Environmental Pollution*, 208 B, 537-  
3 545. <https://doi.org/10.1016/j.envpol.2015.10.027>.  
4

5 Traore, H., Crouzet, O., Mamy, L., Sireyjol, C., Rossard, V., Servien, R., Latrille, E., Martin-  
6 Laurent, F., Patureau, D., Benoit, P., 2018. Clustering pesticides according to their molecular  
7 properties, fate and effects by considering additional ecotoxicological parameters in the TyPol  
8 method, *Environmental Science and Pollution Research*, 25(5), 4728-4738.  
9 <https://doi.org/10.1007/s11356-017-0758-8>.  
10

11 Tsai, C.-F., 2014. Combining cluster analysis with classifier ensembles to predict financial  
12 distress, *Information Fusion*, 16, 46-58. <https://doi.org/10.1016/j.inffus.2011.12.001>.  
13

14 UNEP-SETAC, 2019. Global Guidance for Life Cycle Impact Assessment Indicators: Volume  
15 2. [https://www.lifecycleinitiative.org/training-resources/global-guidance-for-life-cycle-impact-](https://www.lifecycleinitiative.org/training-resources/global-guidance-for-life-cycle-impact-assessment-indicators-volume-2/)  
16 [assessment-indicators-volume-2/](https://www.lifecycleinitiative.org/training-resources/global-guidance-for-life-cycle-impact-assessment-indicators-volume-2/) (accessed Nov 22, 2020).  
17

18 USEtox® 2020: USEtox® database system, <https://USEtox.org/model/download>.  
19

20 Verones, F., Bare, J., Bulle, C., Frischknecht, R., Hauschild, M., Hellweg, S., Henderson, A.,  
21 Jolliet, O., Laurent, A., Liao, X., et al., 2017. LCIA Framework and Cross-Cutting Issues  
22 Guidance within the UNEP-SETAC Life Cycle Initiative, *Journal of Cleaner Production*, 161,  
23 957–967. <https://doi.org/10.1016/j.jclepro.2017.05.206>.  
24

25 Willmott, C., Matsuura, K., 2005. Advantages of the Mean Absolute Error (MAE) over the Root  
26 Mean Square Error (RMSE) in Assessing Average Model Performance, *Climate Research*,  
30, 79. <https://doi.org/10.3354/cr030079>.

1 Wold, H., 1985. Partial least squares, In Kotz, Samuel; Johnson, Norman L. (eds.),  
2 Encyclopedia of statistical sciences, vol 6, New York, Wiley.

3

4 Xia, M., Huang, R., Witt, K.L., Southall, N., Fostel, J., Cho, M.-H., Jadhav, A., Smith, C.S.,  
5 Inglese, J., Portier, C.J., Tice, R.R., Austin, C.P., 2008. Compound cytotoxicity profiling using  
6 quantitative high-throughput screening, Environmental Health Perspectives, 116 (3), 284–  
7 291, <https://doi.org/10.1289/ehp.10727>.

8