



HAL
open science

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-03109818v1

HAL Id: hal-03109818

<https://hal.inrae.fr/hal-03109818v1>

Preprint submitted on 14 Jan 2021 (v1), last revised 7 Feb 2022 (v3)

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.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

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

Rémi Servien^{a,b,*}, Eric Latrille^{a,b}, Dominique Patureau^a, Arnaud Hélias^{c,d}

^aINRAE, Univ. Montpellier, LBE, 102 Avenue des étangs, F-11000 Narbonne, France

^bChemHouse Research Group, Montpellier, France

^cITAP, Univ Montpellier, INRAE, Institut Agro, Montpellier, France

^dELSA, Research group for environmental life cycle sustainability assessment and ELSA-Pact industrial chair, Montpellier, France

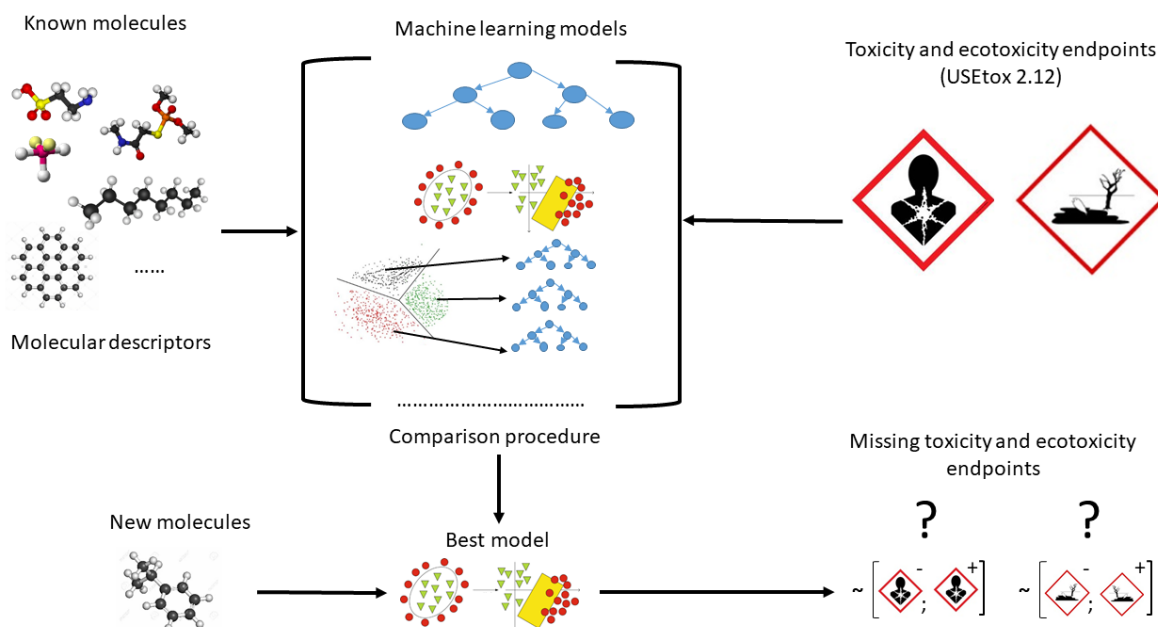
*corresponding author : remi.servien@inrae.fr

Highlights:

- Characterization factors (for human health and ecotoxicological impacts) were predicted using molecular descriptors.
- Several linear or non-linear machine learning methods were compared.
- A train and test procedure was applied to assess the performances of the methods.
- Predictions using machine learning were good.
- This methodology was then used to derive tens of characterization factors for USEtox.

Abstract: It is a real challenge for life cycle assessment practitioners to identify all relevant substances contributing to the ecotoxicity. Once this identification has been made, the lack of corresponding ecotoxicity factors can make the results partial and difficult to interpret. So, it is a real and important challenge to provide ecotoxicity factors for a wide range of compounds. Nevertheless, obtaining such factors using experiments is tedious, time-consuming, and made at a high cost. A modeling method that could predict these factors from easy-to-obtain information on each chemical would be of great value. Here, we present such a method, based on machine learning algorithms, that used molecular descriptors to predict two specific endpoints in continental freshwater for ecotoxicological and human impacts. The method shows good performances on a learning database. Then, predictions were derived from the validated model for compounds with missing toxicity/ecotoxicity factors.

Graphical abstract:



1
2

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

5
6
7

1. Introduction

8 Recent legislations such as the Registration, Evaluation, Authorization and restriction of
9 Chemicals (REACH) regulation in the EU requires that manufacturers of substances and
10 formulators register to provide eco/toxicological data for substances with volume higher than
11 one metric ton per year. As an example, the U.S. Environmental Protection Agency (EPA) has
12 more than 85,000 chemicals listed under the Toxic Substances Control Act (Hinds and Weller,
13 2016). The needed information has to be equivalent to the standard information requirement
14 and adequate to draw overall conclusions with respect to the regulatory endpoints
15 classification and labeling. Beyond specific regulatory needs, the same questions concern
16 chemical substances that came from various sources and are potentially present in the
17 environment.

18

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

23

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

1 factors (CF) that multiply the emitted masses to determine the impacts. They are not
2 recalculated for each study but provided within a Life Cycle Impact Assessment (LCIA)
3 method.

4
5 For a given impact, the LCIA method designer refers to the knowledge of the scientific
6 community to model the mechanisms involved. For human toxicity and freshwater ecotoxicity,
7 USEtox (Rosenbaum et al., 2008), was developed by life cycle initiative under the United
8 Nations Environmental Programme (UNEP) and the Society for Environmental Toxicology and
9 Chemistry (SETAC) (Henderson et al. 2011) to produce a transparent and consensus
10 characterization model. USEtox is also used for the European Product Environmental
11 Footprint (PEF) (Saouter et al., 2020). This model gathers in one single characterization factor
12 the chemical fate, the exposure, and the effect for each of the several thousands of organic
13 and inorganic compounds. If the structure of this multimedia model is always the same, to
14 determine the CF of a molecule, numerous physico-chemical parameters (such as solubility,
15 hydrophobicity, degradability) and detailed toxicological and ecotoxicological data must be
16 provided. For example, EC50 values for at least three species from three different trophic
17 levels are required for the ecotoxicological effect factor.

18
19 Over the past few decades, thousands of tests (in laboratory and field) have been carried out
20 to evaluate the potential hazard effects of chemicals (He et al., 2017). Usually, toxicity testing
21 has relied on *in vivo* animal models, which is extremely costly and time-consuming (Xia et al.,
22 2008). In recent years, under societal pressures, there has been a significant paradigm shift
23 in toxicity testing of chemicals from traditional *in vivo* tests to less expensive and higher
24 throughput *in vitro* methods (National Research Council, 2007). However, it is still extremely
25 hard to test the number of existing and ever-increasing numbers of new chemicals, which
26 leaves their impacts largely unknown. That's why more computational models are needed to
27 complement experimental approaches to decrease the experimental cost and determine the
28 prioritization for those chemicals which may need further *in vivo* studies. Such models already
29 exist, like QSAR models that are mostly linear models based on the chemical structure of
30 compounds (Danish QSAR database (DTU, 2015), ECOSAR (Mayo-Bean et al., 2011), VEGA
31 (Benfenati et al., 2013)) and are used to predict ecotoxicological data (LC50) needed for
32 REACH for example. Recently, machine learning algorithms have been used to predict
33 hazardous concentration 50% (HC50) based on 14 physico-chemical characteristics (Hou et
34 al., 2020a) or on 691 more various variables (Hou et al., 2020b). In the case of USEtox, despite
35 its wide use in LCA, it only offers characterization factors for approximately 3000 chemicals
36 and even for this limited number of compounds, 19% of ecotoxicity CFs and 67% of human
37 toxicity CFs are missing. The objective of this article is thus to propose a new way of
38 calculating CFs using machine learning approaches to solve the problem of nonlinearity that
39 could affect a linear QSAR method. This makes it possible, when the CFs are not determined
40 due to lack of time or lack of data, to propose values based solely on easily identifiable
41 molecular descriptors. Here, the main differences with the above-cited methods are twofold:
42 first, our input variables are only molecular descriptors that could be easily collected for any
43 newly available compounds; second, our output variables are directly the CFs that are closer
44 to the endpoints than the HC/LC50.

45
46 Indeed, the USEtox model results can be extended to determine endpoint effects expressed
47 as disability-adjusted life years (DALY) for human health impacts and potentially disappeared
48 fraction of species (PDF) for ecotoxicological impacts. The PDF represents an increase in the

1 fraction of species potentially disappearing as a consequence of emission in a compartment
 2 while the DALY represents an increase in adversely affected life years. These endpoints are
 3 now consensual at an international level (Verones et al., 2017). These two specific endpoints
 4 will be studied in the present paper through the emission of compounds in continental
 5 freshwater and will be named CF_{ET} for ecotoxicological impacts and CF_{HT} for human ones. For
 6 this aim, we rely on the TyPol tool with associated molecular descriptors and classification tool
 7 (Servien et al., 2014).

10 2. Materials & Methods

11 2.1. USEtox database

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

16 2.2. TyPol database

18 We recently developed TyPol (Typology of Pollutants), a classification method based on
 19 statistical analyses combining several environmental parameters (i.e., sorption coefficient,
 20 degradation half-life, Henry constant) and an ecotoxicological parameter (bioconcentration
 21 factor BCF), and structural molecular descriptors (i.e., number of atoms in the molecule,
 22 molecular surface, dipole moment, energy of orbitals). Molecular descriptors are calculated
 23 using an *in silico* approach (combining Austin Model1 and Dragon software). In the present
 24 paper, we only extract and use the molecular descriptors from the TyPol database, as this
 25 information could be easily collected for any new compound. The 40 descriptors included in
 26 the TyPol database have been selected based on a literature review on QSAR equations used
 27 to predict the main environmental processes as degradation, sorption, volatilization. These 40
 28 descriptors were the ones most frequently used in the equations, meaning describing the best
 29 the behaviour of organic compounds in the environment. They are constitutional, geometric,
 30 topological, and quantum-chemical descriptors (see Table 1). For more details, we refer the
 31 interested reader to Servien et al. 2014. Now, TyPol includes 549 compounds, which are
 32 mainly pesticides and their transformation products (Benoit et al. 2017, Traoré et al. 2018).

34 **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
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36

2.3. Machine learning methods

To predict the CFs using the molecular descriptors we use three modelling methods combined. The first method is a linear well-known prediction method namely the Partial Least Squares (PLS) (Wold, 1985). It finds the multidimensional directions in the observable variable (molecular descriptor) space that explains the maximum multidimensional variance direction in the predicted variable (CF) space. That provides a linear regression model based on the observable variables to predict the predicted variable. We also choose to compare two non-linear machine learning methods: the random forest (Breiman 2001) and the support vector machines (SVM) (Drucker et al. 1996). Random forests are a machine learning method, for classification or, in our case, regression, that operate by constructing a multitude of decision trees that uses a random subset of the training data and limits the number of variables used at each split and outputting the mean prediction (regression) of the individual trees. SVM constructs a hyperplane or set of hyperplanes in a high- or infinite-dimensional space in which the problem is linearly separable.

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

All the models were computed in the freeware R (R core team, 2019). The PLS has been computed using the package mixOmics (Rohart et al., 2017), the random forests using the package randomForest (Liaw et al., 2002), and the SVM using the package e1071 (Meyer et al., 2019). These 3 modelling methods have some parameters that needed to be fixed: the number of latent components for the PLS (fixed using the tune.pls function), the number of variables randomly sampled as candidates at each split for the random forests (selected using the tune.randomForest function) and, for the SVM, the gamma parameter of the radial kernel and the cost of constraints violation (using the tune.svm function). All these different tune functions are based on cross-validation.

2.4. Clustering-based model

1 A recent popular way to make predictions is to use a cluster-then-predict approach. That is,
2 clustering is used for pre-classification which is to arrange a given collection of input patterns
3 into natural meaningful clusters. Then, the clustering results are used to construct a predictor
4 in each cluster. The main idea of the cluster-then-predict approach is that if the clustering
5 performs well the prediction will be easier by modeling only similar compounds. If a new
6 compound with no CF_{ET} and/or CF_{HT} is investigated, the clustering can easily be applied to it
7 before the prediction model itself. The cluster-then-predict approach has already been applied
8 with success in various domains such as sentiment prediction (Sony et al., 2015), finance
9 (Tsai et al., 2014), chemometrics (Minh Mai Le et al., 2018). So we decided to use the
10 clustering given by the TyPol application (more details in Servien et al., 2014) based on the
11 whole database and the molecular descriptors. Note that the TyPol clustering has already
12 been shown relevant on various occasion: in combination with mass spectrometry to
13 categorize tebuconazole products in soil (Storck et al., 2016), to explore the potential
14 environmental behaviour of putative chlordecone transformation products (Benoit et al., 2017)
15 or to classify pesticides with similar environmental behaviors (Traore et al., 2018). This
16 clustering is given in Supplementary Figure S1.

17

18 **2.5. Comparison procedure**

19

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

- 21 1. Split each cluster between a training set (85% of the dataset) and a test set (15%).
22 The test set is not used for any step of the procedure (such as the imputation of the
23 missing data, the calibration of the parameters ...).
- 24 2. Imputation of the NA values (less than 1%) in the descriptor matrix using the NIPALS
25 algorithm (Wold, 1985).
- 26 3. Tune the parameters and train the specific models on the training set. We have 3 global
27 models to train (PLS, random forest, and SVM) and the cluster-then-test models (PLS,
28 random forest and SVM for each cluster).
- 29 4. Test the different models on the test set. Compute the absolute error.
- 30 5. Back to step 1.

31

32 For cluster 5, the 3 global models are the only ones available as we can't define a cluster-
33 then-test model due to a lack of data. The whole algorithm is repeated 200 times. All the
34 performances are compared in terms of absolute error. The absolute error is the absolute
35 difference between the prediction and the true value. It has been shown to be the most natural
36 and unambiguous measure of error (Willmott et Matsuura, 2005). For each cluster, we chose
37 the model with the lowest median absolute error.

38

39 Then, the best model is calibrated and computed on the whole cluster. Finally, it is applied to
40 the compounds, according to their clusters, with a CF_{ET} (or a CF_{HT}) equals to NA to provide a
41 prediction. For the compounds in cluster 5, this best model cannot be a cluster-then-predict
42 one and, by consequence, is a global one. A 95% prediction interval is also derived for each
43 prediction. The type of model and its corresponding parameters are fixed during this process,
44 according to the best model of the cluster. For example, if the best model of cluster 1 was the
45 random forest approach, random forest models are used with the parameters optimized during
46 the previous step. Then, we perform a leave-one-out bootstrap on the dataset that was used
47 to compute the model (the whole dataset if the model is global, only the data lying in the
48 dedicated cluster if that is a cluster-then-predict model) and a new model is computed on this

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

4

5 The five more important descriptors are then derived for each chosen model. For a random
6 forest model, these descriptors are calculated using variable permutations (Breiman, 2001),
7 for the SVM they are the descriptors with the higher coefficients in absolute value.

8

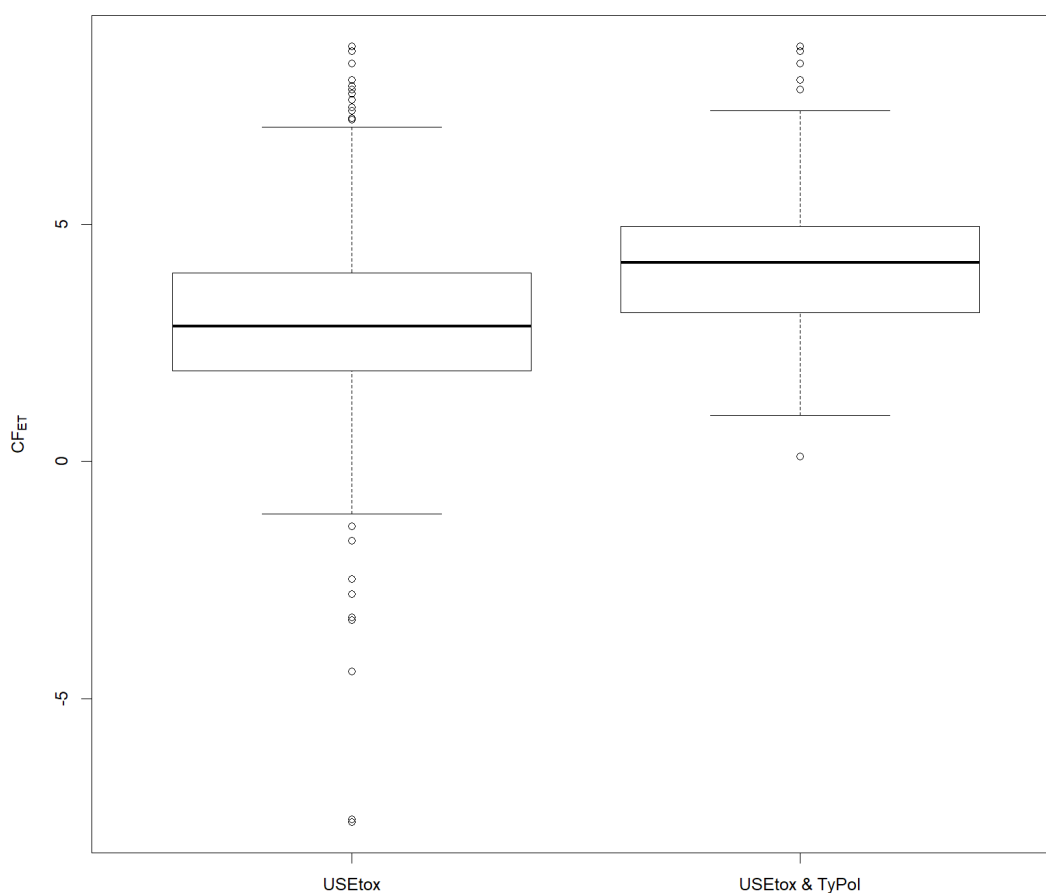
9 3. Results

10

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

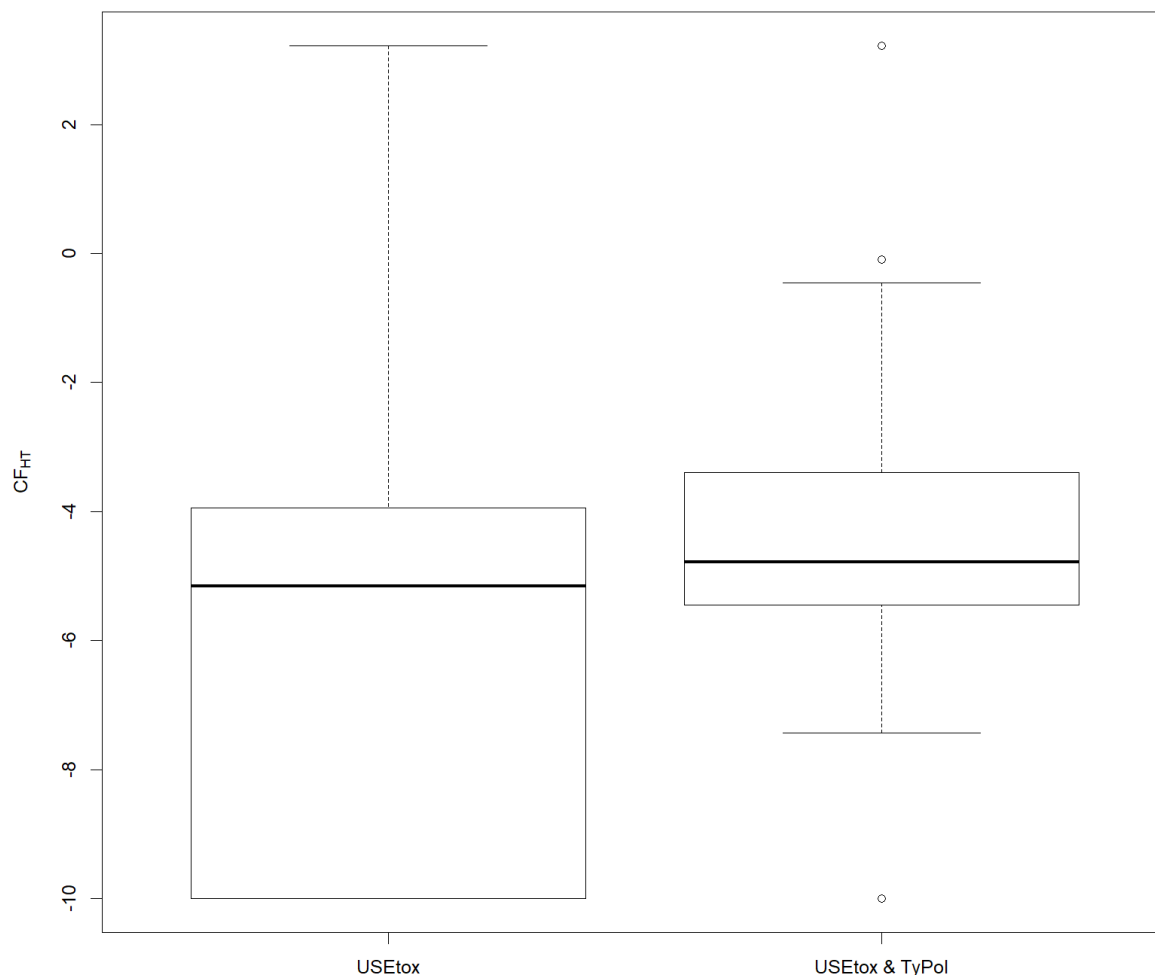
13

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} .



20

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



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 ϵ is
 4 needed as some values of the DALY are exactly equal to zero. ϵ 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.

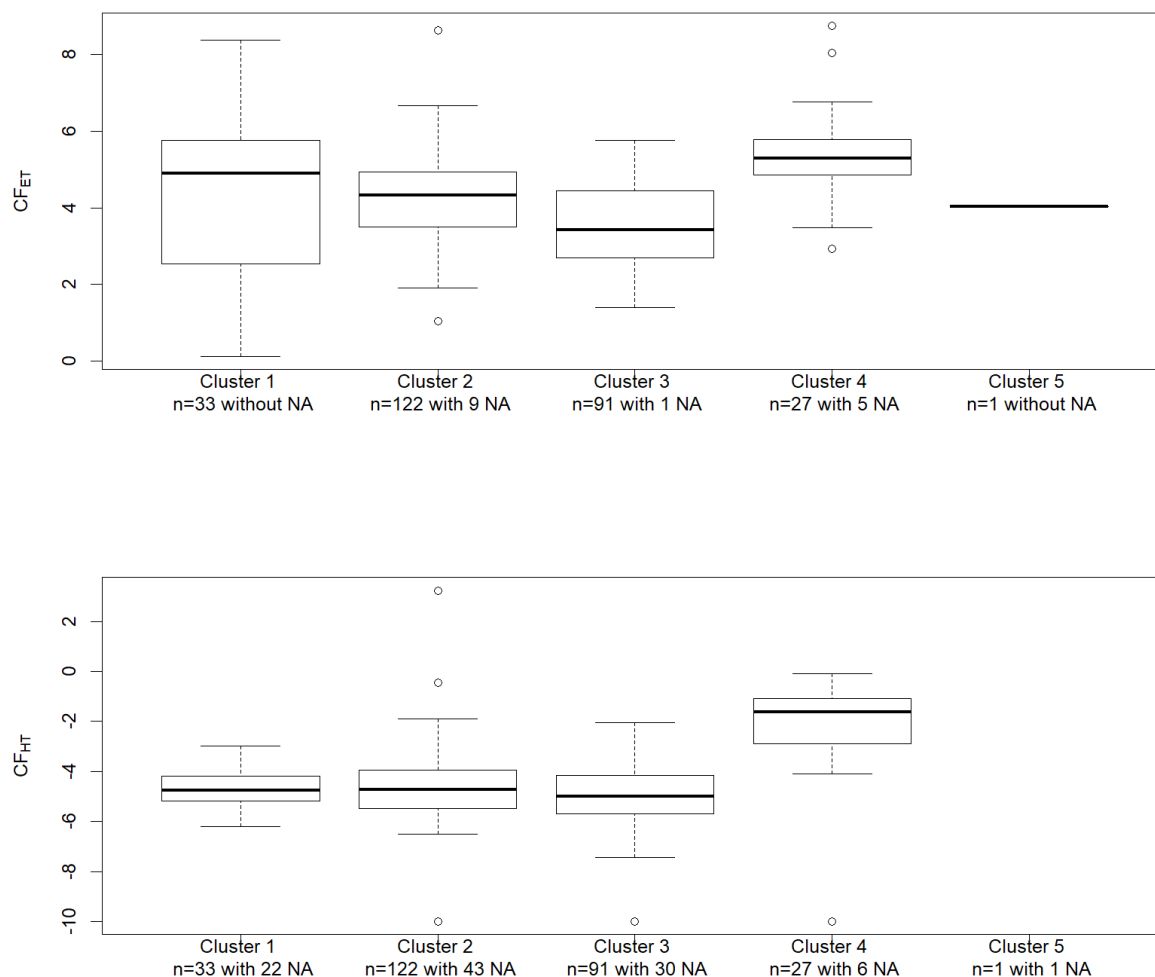
10
 11 The TyPol clustering focused on the common compounds is plotted in Supplementary Figure
 12 S2 and the boxplots of each molecular descriptor per cluster are given in Supplementary
 13 Figure S3 with different indicators in Table S2. We could see that they are clustered in 5 groups
 14 with different sizes (respectively 33 compounds in the first black cluster, 122 compounds in
 15 the second red cluster, 91 compounds in the third green cluster, 27 compounds in the fourth
 16 blue cluster, and one compound in the fifth brown cluster). Cluster 1 grouped compounds with
 17 a high number of aromatic bonds, double bonds, rotatable bonds, and multiple bonds. Cluster
 18 2 is an intermediate one between clusters 1 and 3, with less extreme values. Cluster 3 is made
 19 of compounds with the lowest molecular mass. Cluster 4 gathered compounds presenting a

1 high number of halogens, rings, and circuits. The unique compound in the fifth cluster is
2 erythromycin (highest molecular mass and number of H and C, lowest number of rings) and,
3 obviously, no cluster-then-predict model could be built for this cluster

4

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

7



8

9 **Figure 3-** Boxplot by cluster for the CF_{ET} and CF_{HT} values. Note that the unique compound of
10 Cluster 5 has no CF_{HT} value. The size of the clusters and the numbers of NA are gathered in
11 the legend.

12

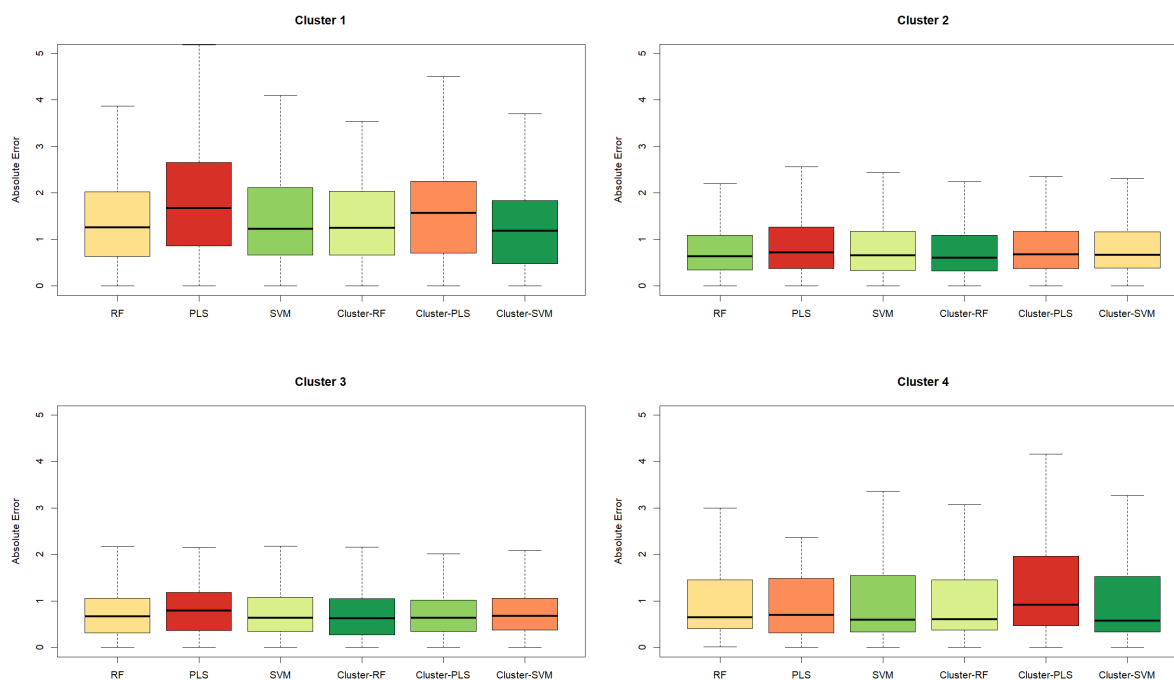
13 The predictions will be made difficult for the CF_{ET} of cluster 1 as it covers a wide range whereas
14 it includes a relatively small number of compounds. On the contrary, cluster 3 covers a small
15 range with no extreme values and includes a high number of compounds, for this cluster the
16 cluster-then-predict approach could produce interesting results.

17

18 **3.2. Models and prediction of the CF_{ET}**

19 **3.2.1. Performances of the machine learning methods**

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



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

9
10 The performances are not similar in each cluster. For example, performances of all methods
11 for cluster 1 are very poor (median absolute error above 1) whereas performances for cluster
12 4 seem good despite its smallest size (median absolute error around 0.6). So, a future
13 prediction of an unknown compound which lies in cluster 1 will be less reliable than in other
14 clusters. Note that we could not test this in the next section as no NA value is present in this
15 cluster 1.

16
17 The cluster-then-predict methods seem more appropriate in each cluster. The cluster-then-RF
18 approach has the best performances (with a global median absolute error equals to 0.64 and
19 the best performances on clusters 2 and 3), even if there is not a big difference between the
20 different methods. The cluster-then-SVM is also the best method for the two clusters 1 and 4.
21 The linear methods (PLS and cluster-then-PLS) have higher absolute errors but are
22 competitive. The individual predictions of the best method in each cluster are reported in
23 Figure S5.

24 3.2.2. Prediction with the best model

25
26
27 Then we apply the best model in each cluster: a cluster-then-predict approach using SVM for
28 clusters 1 and 4 and using random forest for clusters 2 and 3. To compare the different models
29 in each cluster and give an idea of what are the important molecular descriptors we provide
30 the five most important molecular descriptors for each cluster in the following table.

1
2
3

Table 1- The five most important molecular descriptors for each best model for each cluster. 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

4
5
6
7

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

8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

Then the models were used to predict the missing CF_{ET} of the common compounds between USEtox and TyPol databases. These values are by consequence new estimations of the CF_{ET} for compounds on which we have no information. The prediction intervals are relatively small: less than $0.5 \log_{10}$ in a log scale which highlights the robustness of the estimation. They are given in Table S3. No NA value was present in cluster 1 with no prediction for this cluster. For cluster 2 gathering molecules with intermediate molecular mass, 9 CF_{ET} values were predicted for various kinds of compounds. One value concerns the antibiotic sulfamethazine and its value is quite near to the one of sulfamethoxazole and sulfadiazine of the same sulphonamide antibiotic family constituted of the sulphonamide group ($-S(=O)_2-NR_2R_3$). Cluster 3 grouped compounds with the lowest molecular mass and the lowest median CF_{ET} like ibuprofen, phthalates, cresol constituted of monoaromatic ring substituted with methyl, carboxylic groups. The CF_{ET} prediction for acetylsalicylic acid seemed coherent with the value of the nearest compounds (herbicides mecoprop) of this group. Cluster 4 gathered compounds with the highest median CF_{ET} and that presented a high number of rings halogenated or not, like PAH and hormones. The 5 CF_{ET} predicted concerned 4 PAHs and 1 hormone. By comparison to the 2 other PAHs present in this cluster, the 4 predicted CF_{ET} are quite similar and higher. Concerning the prediction for the hormone, the CF_{ET} is intermediate between the CF_{ET} of the 3 other hormones in the cluster. It seems that all these 5 predicted values are very closed, falling near the median value of this cluster.

27
28
29
30
31

3.3. Models and prediction of the CF_{HT}

3.3.1. Performances of the methods

32
33

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.

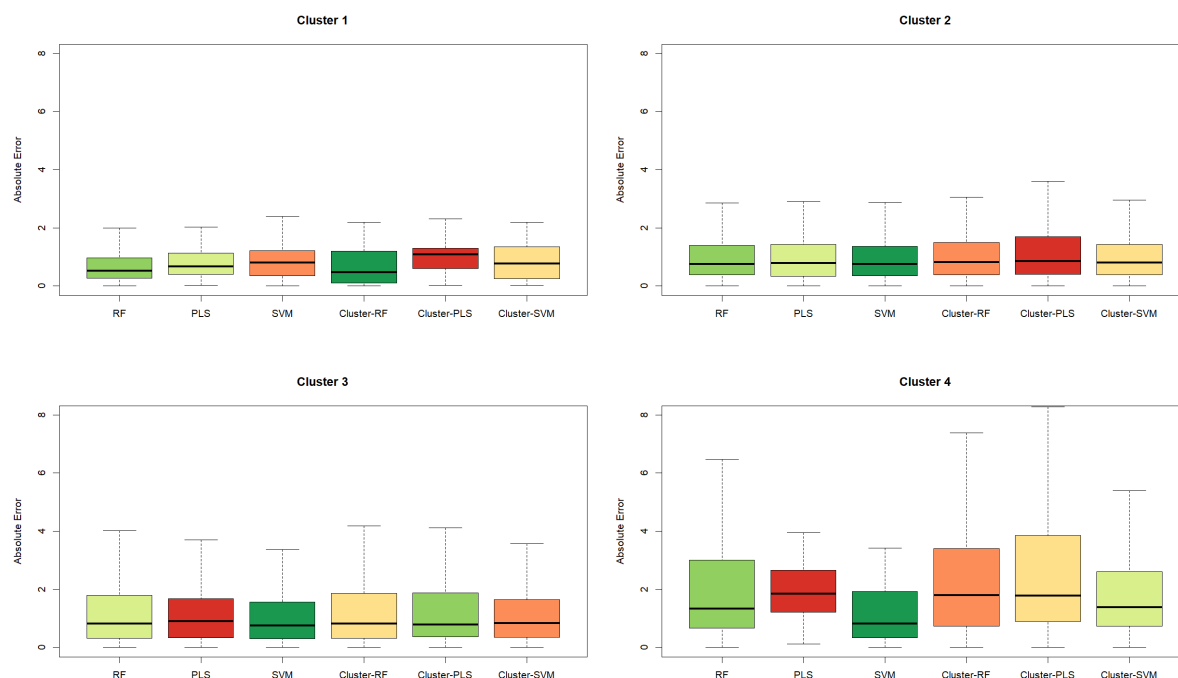


Figure 5- Performances of the different methods in terms of the log of the absolute error of the CF_{HT} with respect to the different clusters. In each cluster, the models are coloured from green (best) to red (worst) according to their median of the absolute error.

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 0.82, all the medians of the other methods are above 1.3. Global performances of the different methods are given in Supplementary Figure S6.

3.3.2. Prediction with the best model

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

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

1
2 Then, this model was used to predict the CF_{HT} value for the 102 common compounds without
3 a CF_{HT} value. These predictions are reported in Supplementary Table S4. As for the CF_{ET} , the
4 small width of the prediction interval (less than a \log_{10} in a log scale) highlights the robustness
5 of the approach even with a relatively small number like estimations made for compounds that
6 lie in cluster 1. In this cluster 1, CF_{HT} for a phthalate (DEHP) is already known, but the one for
7 diisodecyl and diisononyl phthalate was predicted with value in the same range. The 3 cyclines
8 (tetracycline, aureomycin, and oxytetracycline) present in cluster 1, presented also similar
9 predicted CF_{HT} . This was also the case for triclosan and triclocarban in cluster 2. Similar
10 predicted and known CF_{HT} were found for four herbicides from the substituted urea family
11 (linuron, diuron, monolinuron, isoproturon) in cluster 3. Cluster 4 gathered a small number of
12 molecules but with the highest median CF_{HT} , the predicted CF_{HT} of the organochlorine
13 insecticide isodrin was similar to another congener of the same family, aldrin.

14 15 **4. Discussion**

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

29
30 The idea of predicting ecotoxicity characterization factors for chemicals using machine
31 learning algorithms has already been used (Hou et al., 2020a and 2020b). But, here, our
32 findings go further. Indeed, we show that we could directly obtain accurate estimations of
33 endpoint values from easy-to-obtain molecular descriptors. This will open the door to the fast
34 characterization of each new unknown compound that appears, including transformation
35 products. We also show that the cluster-then-predict approach can give better performances
36 than the usual ones. This local approach confirms that local models could be an efficient
37 prediction method when heterogeneity of data generates nonlinear relations between the
38 response and the explicative variables (Lesnoff et al., 2020).

39 40 **5. Conclusion**

41
42 In a recent study, Aemig et al. (2021) studied the potential impacts on Human health and
43 aquatic environment of the release of 286 micropollutants (organic and inorganic) at the scale

1 of France. One of their conclusion was that, due to a lack of characterization factors, these
2 impacts could be assessed only for 1/3 of these molecules. This paper fills this gap by
3 providing a new modeling method to derive characterization factors from easily obtainable
4 molecular descriptors. By consequence, these missing characterization factors, as well as
5 those of new molecules, could now be quickly estimated with an overall good precision. More
6 generally, one of the key factors in the evaluation of toxicity and ecotoxicity in LCA lies in the
7 construction of the characterization factors: a task requiring a large amount of data and a
8 consequent investment of time. The use of machine learning allows us to go beyond these
9 constraints. This makes it possible to obtain characterization factor values in a fast and simple
10 way, which can be used as long as conventionally established CFs are not available.

11

12 **Declaration of Competing Interest**

13

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

16

17 **Acknowledgments**

18

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

21

22 **Funding**

23

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

26

27 **Supplementary materials**

28

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

30

31 **References**

32

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

36

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

39

40 Benoit, P., Mamy, L., Servien, R., Li, Z., Latrille, E., Rossard, V., Bessac, F., Patureau, D.,
41 Martin-Laurent, F., 2017. Categorizing chlordecone potential degradation products to explore
42 their environmental fate, *Science of the Total Environment*, 574, 781–795.
43 <https://doi.org/10.1016/j.scitotenv.2016.09.094>.

44

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

47

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

1
2 Douziech, M., Oldenkamp, R., van Zelm, R., King, H., Hendriks, A.J., Ficheux, A.-S.,
3 Huijbregts, M.A.J., 2019. Confronting variability with uncertainty in the ecotoxicological impact
4 assessment of down-the-drain products, *Environment International*, 126, 37-45,
5 <https://doi.org/10.1016/j.envint.2019.01.080>.
6
7 Drucker, H., Burges, C.C., Kaufman, L., Smola, A.J., Vapnik, V., 1997, Support Vector
8 Regression Machines, *Advances in Neural Information Processing Systems* 9, NIPS, 155–
9 161, MIT Press. <https://dl.acm.org/doi/10.5555/2998981.2999003>.
10
11 DTU, 2015. Danish QSAR database. Danish QSAR group, National Food Institute, Technical
12 University of Denmark.

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

17 He, J., Tang, Z., Zhao, Y., Fan, M., Dyer, S. D., Belanger, S. E., Wu, F., 2017. The combined
18 QSAR-ICE models: practical application in ecological risk assessment and water quality
19 criteria, *Environmental Science & Technology*, 51, 8877.
20 <https://doi.org/10.1021/acs.est.7b02736>.
21
22 Henderson, A.D., Hauschild, M.Z., Van De Meent, D., Huijbregts, M.A.J., Larsen, H.F., Margni,
23 M., McKone, T.E., Payet, J., Rosenbaum, R.K., Jolliet O., 2011. USEtox fate and ecotoxicity
24 factors for comparative assessment of toxic emissions in life cycle analysis: sensitivity to key
25 chemical properties, *The International Journal of Life Cycle Assessment*, 16, pp. 701-709
26 <https://doi.org/10.1007/s11367-011-0294-6>.
27
28 Hinds, R.d.C., Weller, J.L., 2016. Toxic Substances Control Act. *Environmental Law Practice*
29 *Guide*, vol. 4.
30
31 Hou, P., Jolliet, O., Zhu, J., Xu, M., 2020a. Estimate ecotoxicity characterization factors for
32 chemicals in life cycle assessment using machine learning models. *Environment International*,
33 135, 105393. <https://doi.org/10.1016/j.envint.2019.105393>.
34
35 Hou, P., Zhao, B., Jolliet, O., Zhu, J., Wang, P., Xu, M., 2020b. Rapid Prediction of Chemical
36 Ecotoxicity Through Genetic Algorithm Optimized Neural Network Models, *ACS Sustainable*
37 *Chemistry & Engineering*, 8 (32), 12168-12176.
38 <https://dx.doi.org/10.1021/acssuschemeng.0c03660>.
39
40 Lesnoff, M., Metz, M., Roger, JM., 2020. Comparison of locally weighted PLS strategies for
41 regression and discrimination on agronomic NIR data, *Journal of Chemometrics*, 34(5), e3209,
42 <https://doi.org/10.1002/cem.3209>.
43
44 Liaw, A., Wiener, M., 2002. Classification and Regression by randomForest, *R News*, 2(3),
45 18-22. <http://CRAN.R-project.org/doc/Rnews/>.
46

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 Maï Le, L., Kégl, B., Gramfort, A., Marini, C., Nguyen, D., Cherti, M., Tfaili, 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

28 Rosenbaum, R.K., Margni, M., Jolliet, O., 2007. A flexible matrix algebra framework for the
29 multimedia multipathway modeling of emission to impacts, *Environment International*,
30 33(5),624-634. <https://doi.org/10.1016/j.envint.2007.01.004>.

31

32 Rosenbaum, R.K., Bachmann, T. M., Gold, L. S., Huijbregts, M.A.J., Jolliet, O., Juraske, R.,
33 Koehler, A., Larsen, H.F., MacLeod, M., Margni, M., McKone, T.E., Payet, J., Schuhmacher,
34 M., van de Meent, D., Hauschild, M.Z., 2008. USEtox—the UNEP-SETAC Toxicity Model:
35 Recommended Characterisation Factors for Human Toxicity and Freshwater Ecotoxicity in
36 Life Cycle Impact Assessment, *The International Journal of Life Cycle Assessment*, 13 (7),
37 532–546. <https://doi.org/10.1007/s11367-008-0038-4>.

38

39 Saouter, E., Biganzoli, F., Ceriani, L., Versteeg, D., Crenna, E., Zampori, L., Sala, S., Pant,
40 R., 2020. Environmental Footprint: Update of Life Cycle Impact Assessment Methods –
41 Ecotoxicity freshwater, human toxicity cancer, and non-cancer, Publications Office of the
42 European Union, Luxembourg, <https://doi.org/10.2760/300987>.

43

44 Servien, R., Mamy, L., Li, Z., Rossard, V., Latrille, E., Bessac, F., Patureau, D., Benoit, P.,
45 2014. TyPol - a new methodology for organic compounds clustering based on their molecular
46 characteristics and environmental behaviour, *Chemosphere*, 111, 613–622.
47 <https://doi.org/10.1016/j.chemosphere.2014.05.020>.

48

1 Soni, R., Mathai, K.J., 2016. An Innovative 'Cluster-then-Predict' Approach for Improved
2 Sentiment Prediction. In: Choudhary R., Mandal J., Auluck N., Nagarajaram H. (eds)
3 Advanced Computing and Communication Technologies. Advances in Intelligent Systems and
4 Computing, vol 452. Springer, Singapore. https://doi.org/10.1007/978-981-10-1023-1_13.
5
6 Storck, V., Lucini, L., Mamy, L., Ferrari, F., Papadopoulou, E.S., Nikolaki, S., Karas, P.A.,
7 Servien, R., Karpouzas, D.G., Trevisan, M., Benoit, P., and Martin-Laurent, F., 2016.
8 Identification and characterization of tebuconazole transformation products in soil by
9 combining suspect screening and molecular typology, Environmental Pollution, 208 B, 537-
10 545. <https://doi.org/10.1016/j.envpol.2015.10.027>.
11
12 Traore, H., Crouzet, O., Mamy, L., Sireyjol, C., Rossard, V., Servien, R., Latrille, E., Martin-
13 Laurent, F., Patureau, D., Benoit, P., 2018. Clustering pesticides according to their molecular
14 properties, fate and effects by considering additional ecotoxicological parameters in the TyPol
15 method, Environmental Science and Pollution Research, 25(5), 4728-4738.
16 <https://doi.org/10.1007/s11356-017-0758-8>.
17
18 Tsai, C.-F., 2014. Combining cluster analysis with classifier ensembles to predict financial
19 distress, Information Fusion, 16, 46-58. <https://doi.org/10.1016/j.inffus.2011.12.001>.
20
21 UNEP-SETAC, 2019. Global Guidance for Life Cycle Impact Assessment Indicators: Volume
22 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/)
23 [assessment-indicators-volume-2/](https://www.lifecycleinitiative.org/training-resources/global-guidance-for-life-cycle-impact-assessment-indicators-volume-2/) (accessed Nov 22, 2020).
24
25 USEtox 2020: USEtox database system, <https://usetox.org/model/download>.

26 Verones, F., Bare, J., Bulle, C., Frischknecht, R., Hauschild, M., Hellweg, S., Henderson, A.,
27 Jolliet, O., Laurent, A., Liao, X., et al., 2017. LCIA Framework and Cross-Cutting Issues
28 Guidance within the UNEP-SETAC Life Cycle Initiative, Journal of Cleaner Production, 161,
29 957–967. <https://doi.org/10.1016/j.jclepro.2017.05.206>.

30 Willmott, C., Matsuura, K., 2005. Advantages of the Mean Absolute Error (MAE) over the Root
31 Mean Square Error (RMSE) in Assessing Average Model Performance, Climate Research,
32 30, 79. <https://doi.org/10.3354/cr030079>.
33
34 Wold, H., 1985. Partial least squares, In Kotz, Samuel; Johnson, Norman L. (eds.),
35 Encyclopedia of statistical sciences, vol 6, New York, Wiley.
36
37 Xia, M., Huang, R., Witt, K.L., Southall, N., Fostel, J., Cho, M.-H., Jadhav, A., Smith, C.S.,
38 Inglese, J., Portier, C.J., Tice, R.R., Austin, C.P., 2008. Compound cytotoxicity profiling using
39 quantitative high-throughput screening, Environmental Health Perspectives, 116 (3), 284–
40 291, <https://doi.org/10.1289/ehp.10727>.
41

Supplementary Material for 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

1. Supplemental Figures

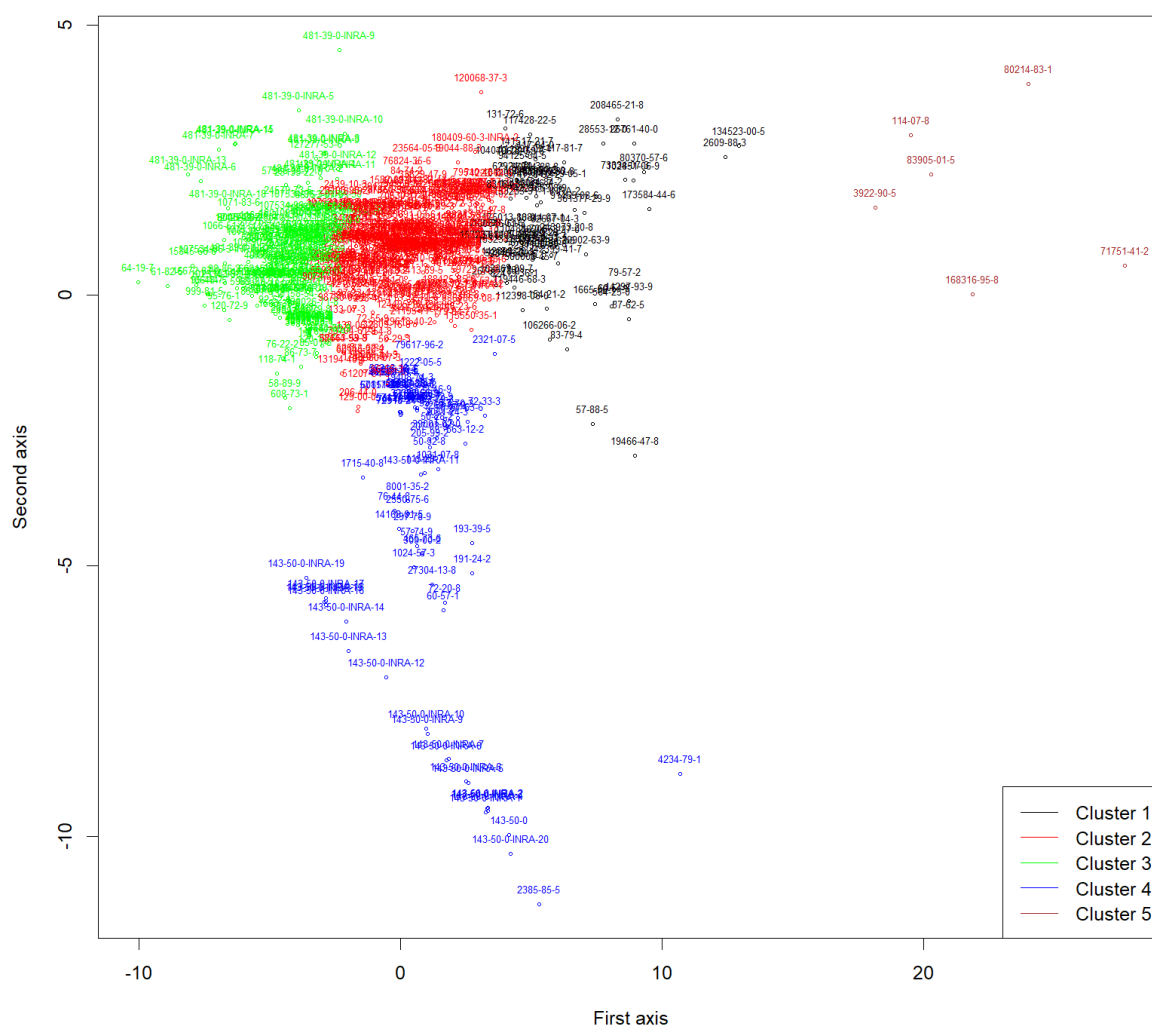


Figure S1- Clustering based on molecular descriptors produced by TyPol on the 526 molecules of the database. We represent here the two first axes of the PLS and the five different clusters in different colours.

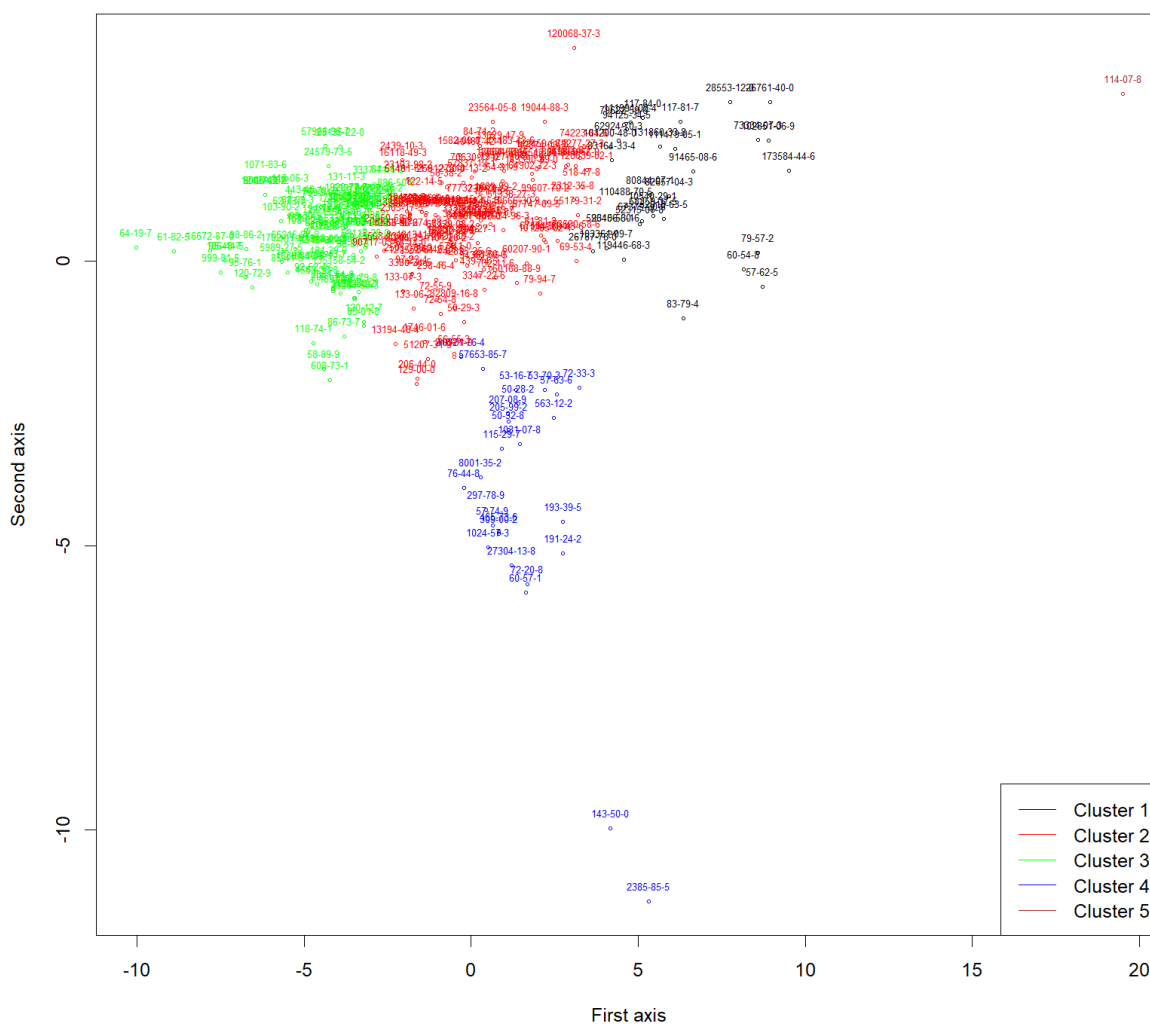


Figure S2- Focus on the 274 common molecules of TyPol & USEtox. The cluster 5 in brown is reduced to a single molecule so the cluster-then-predict methodology cannot be applied for it.

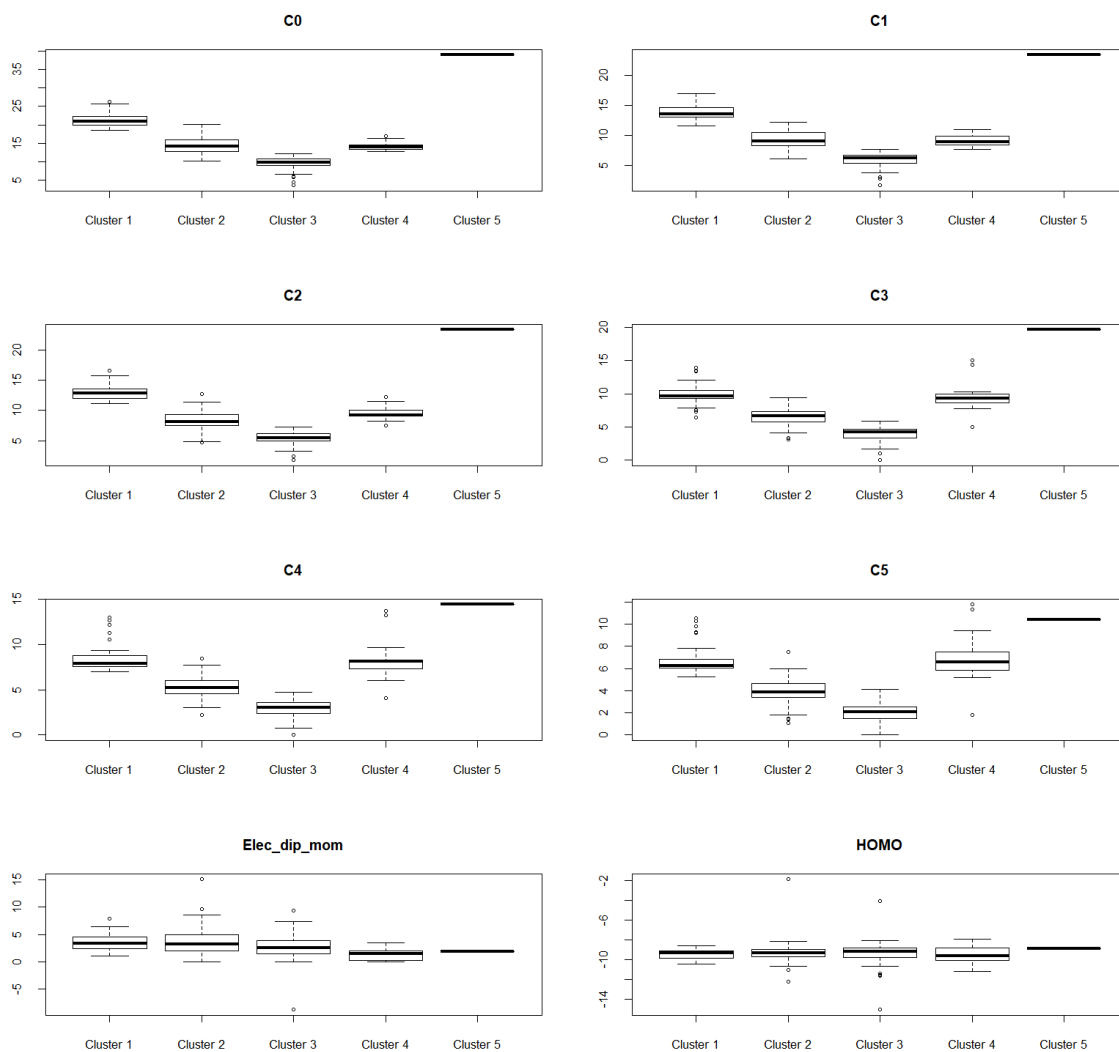


Figure S3 – Boxplots of the 40 molecular descriptors for the clustering given by TyPol on the common compounds of TyPol & USEtox.

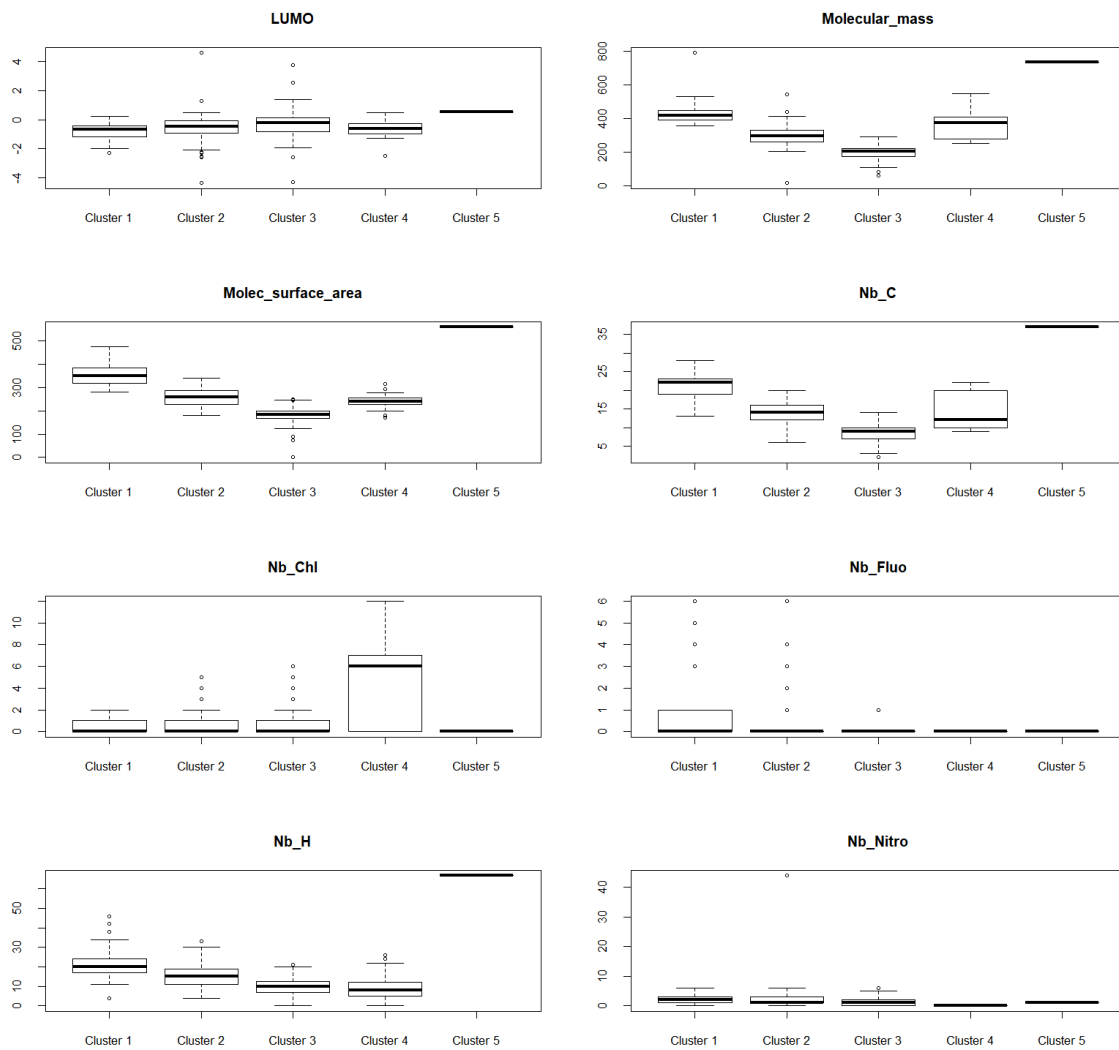


Figure S3 (continued)

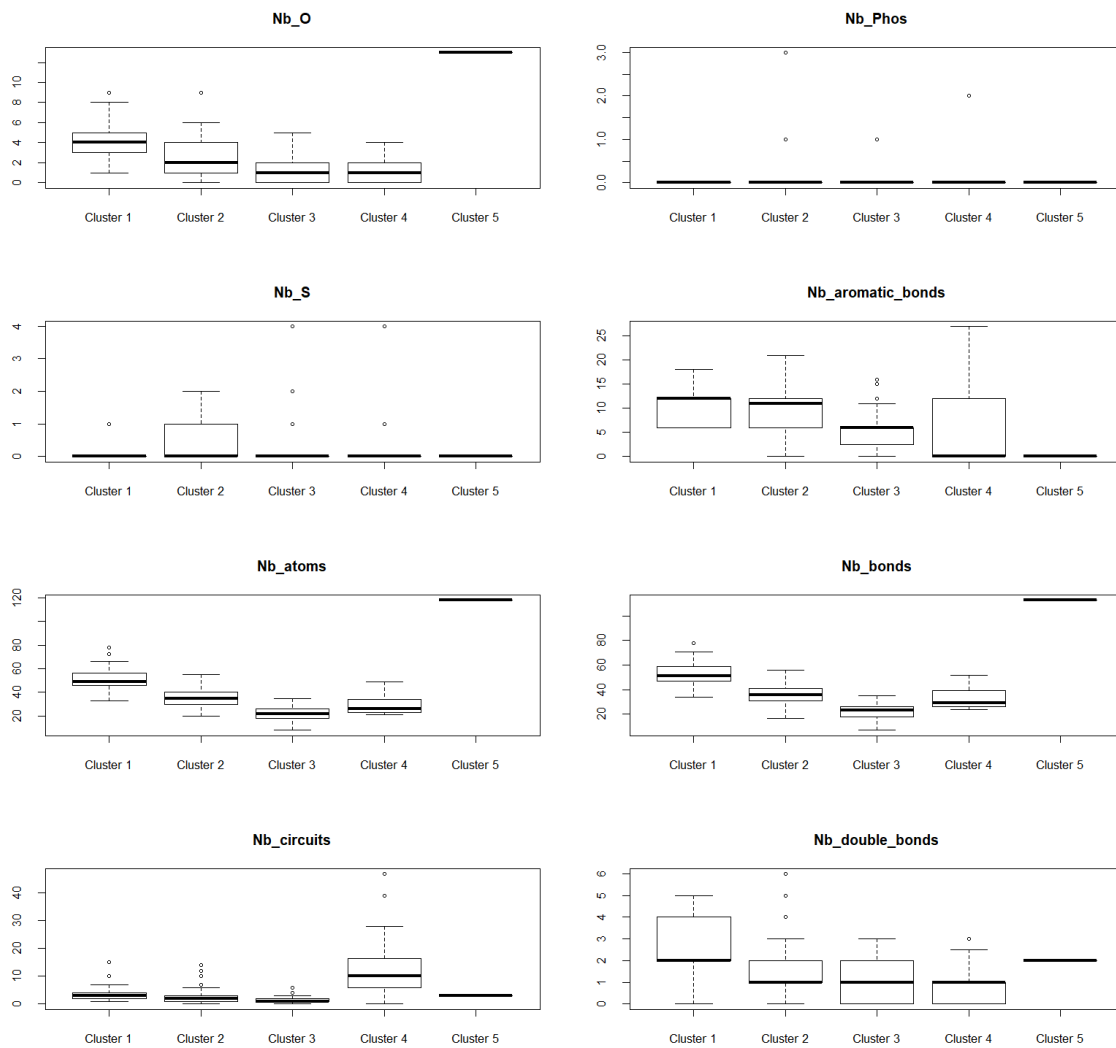


Figure S3 (continued)

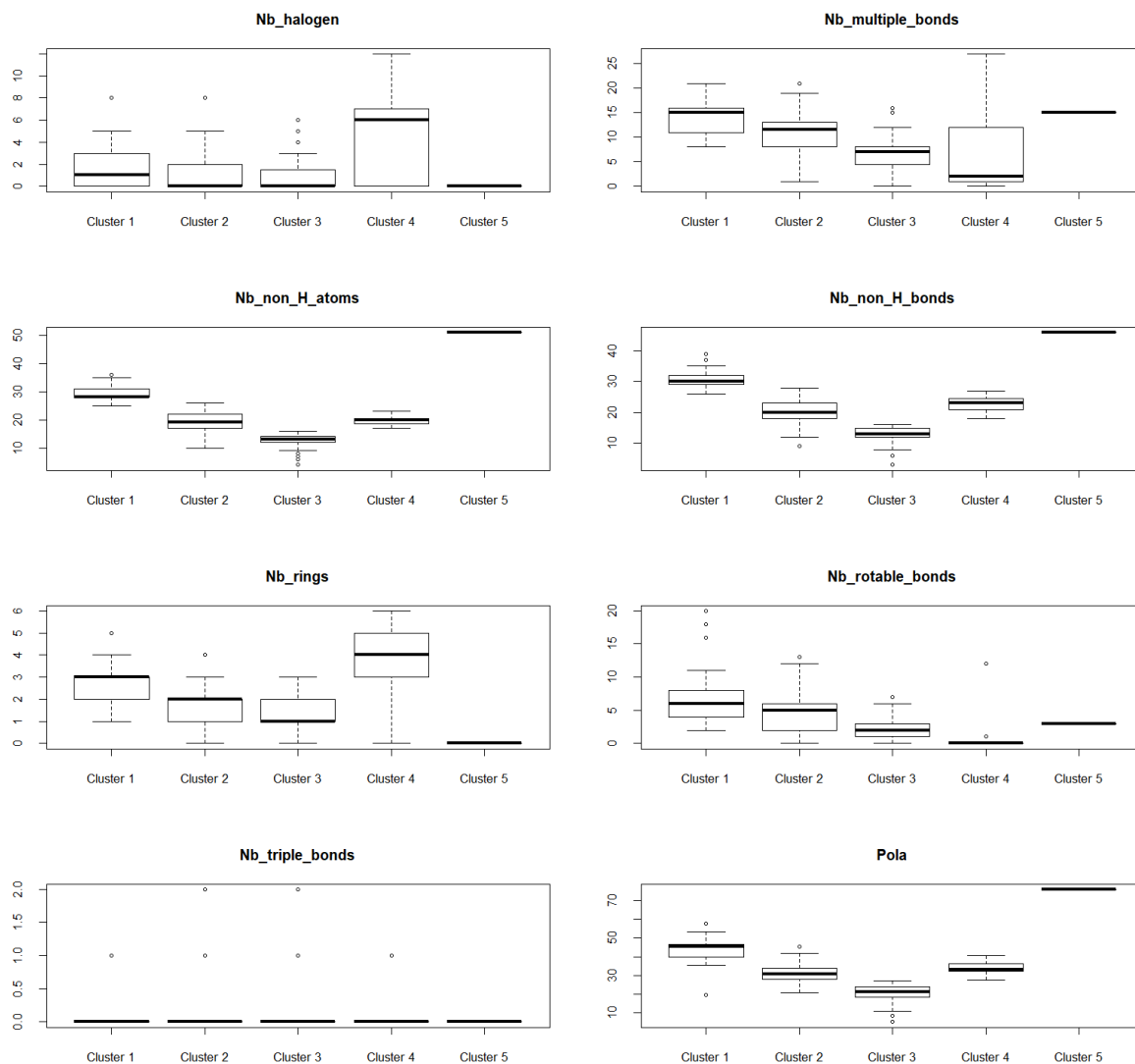


Figure S3 (continued)

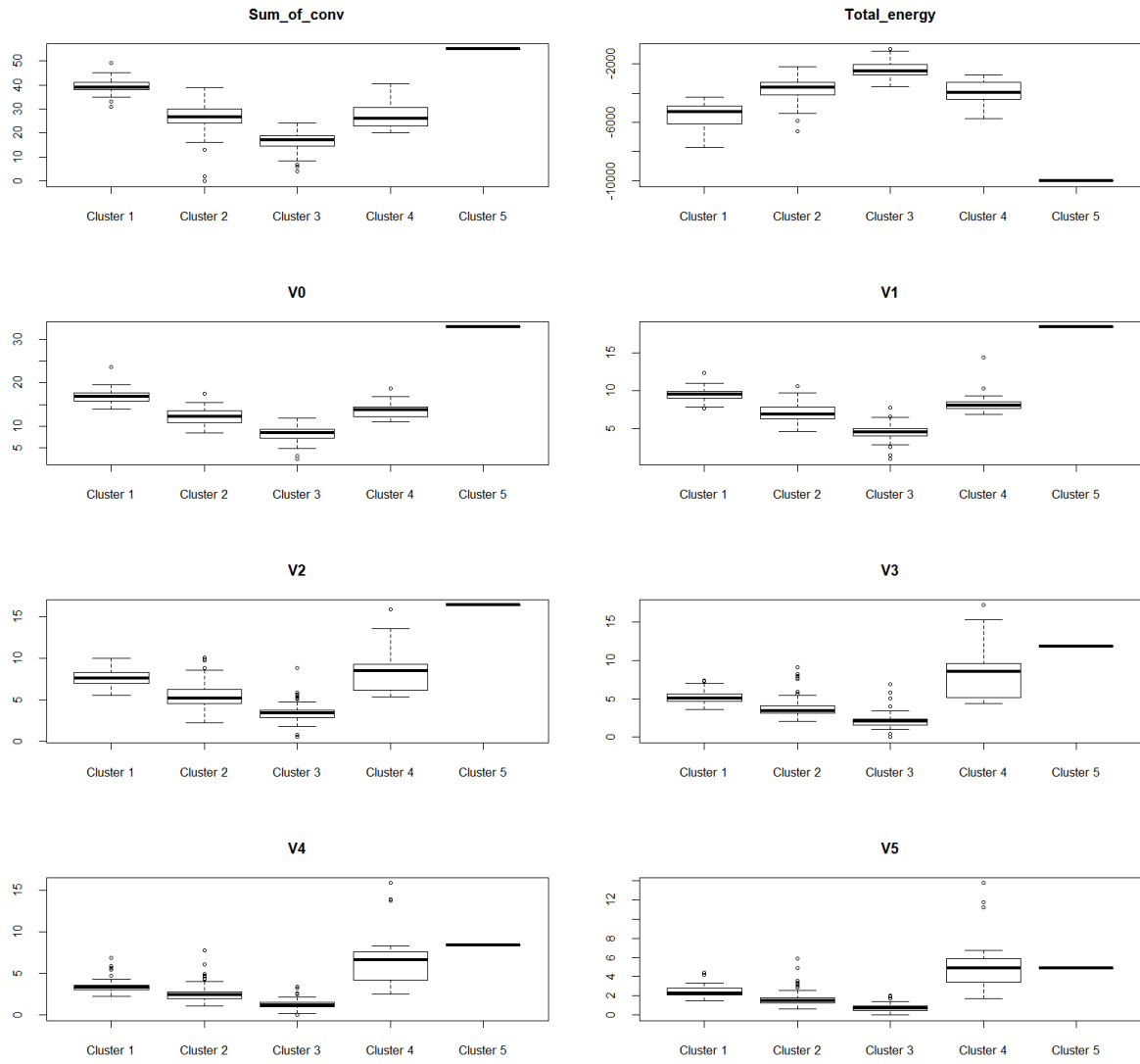


Figure S3 (continued)

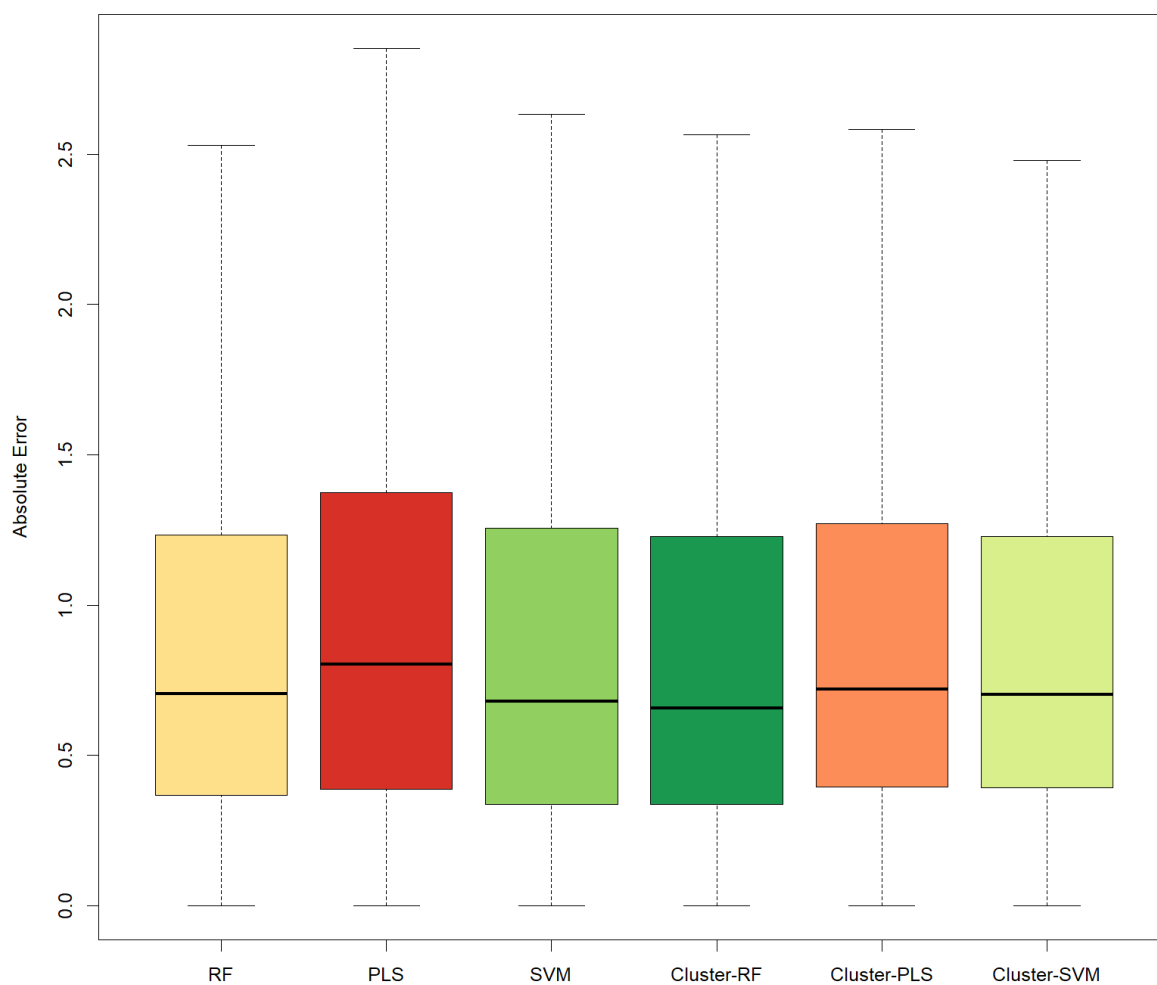


Figure S4- Performances of the different methods in terms of absolute error of the CF_{ET} . The models are coloured from green (best) to red (worst) according to their median of the absolute error.

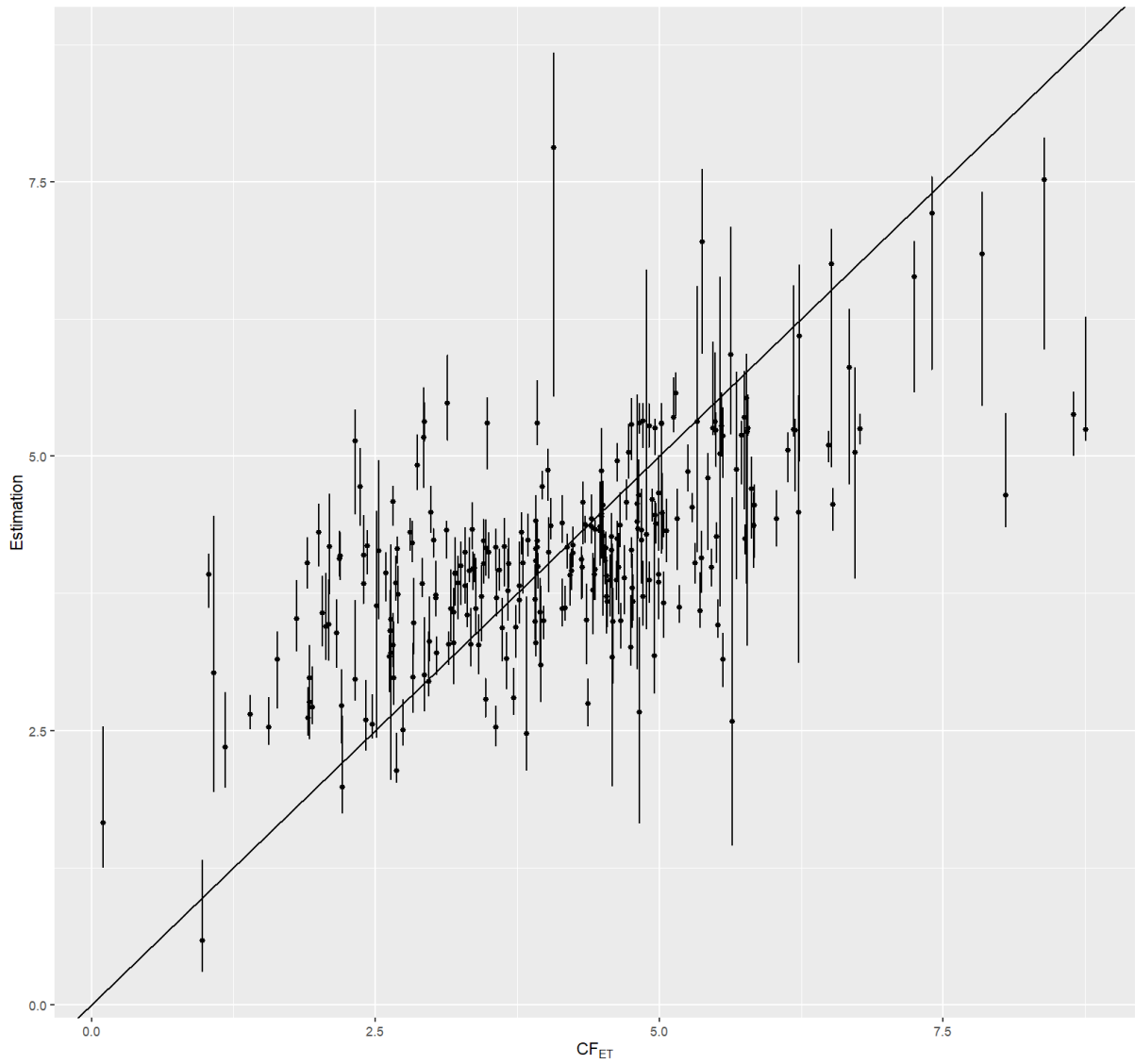


Figure S5- Estimation of CF_{ET} according to the value in *Usetox* . The estimation is the median of the estimation made using the best method of the cluster during the comparison procedure. The bar represents the 5% and the 95% quantiles of these individual estimations.

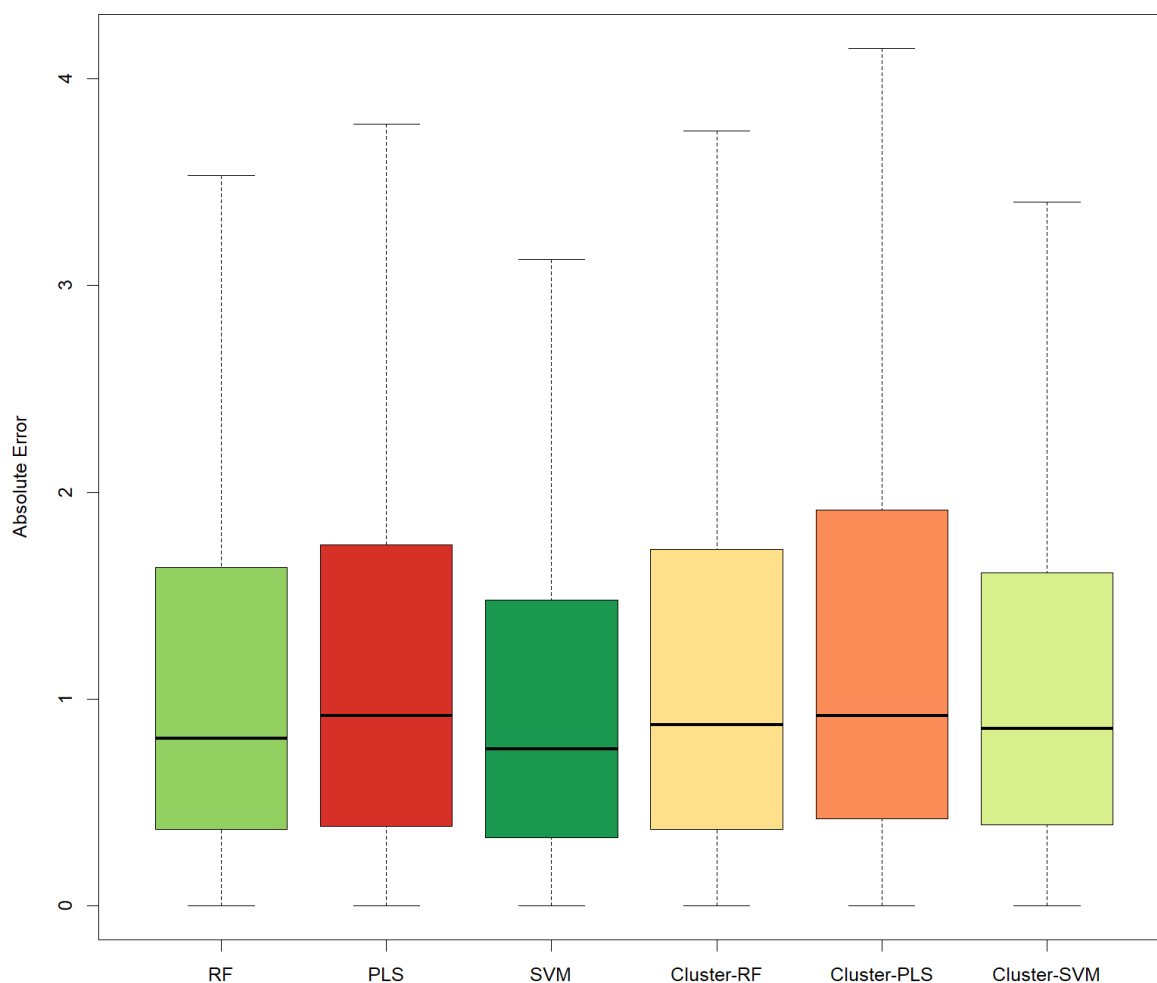


Figure S6- Boxplots of the log of the absolute error for the CF_{HT} estimation for the 6 different methods. The models are coloured from green (best) to red (worst) according to their median of the absolute error.

2. Supplemental Tables

Table S1- CAS number and name of the 274 common compounds between TyPol and USEtox databases and their associated CF_{ET} and CF_{HT} values. NA means that there is no value in USEtox for this compound.

CAS	Name	CF_{HT}	CF_{ET}	Cluster
101-20-2	Triclocarban	NA	6.79E+05	2
101-21-3	Chlorpropham	9.60E-06	2.74E+03	3
101-42-8	Fenuron	NA	1.39E+03	3
101200-48-0	Tribenuron-methyl	1.80E-05	3.39E+02	1
101205-02-1	Cycloxydim	NA	1.56E+02	2

1024-57-3	Heptachlor epoxide	0.81	3.17E+05	4
102851-06-9	tau-Fluvalinate	NA	4.28E+05	1
103-90-2	Acetamide, n-(4-hydroxyphenyl)	1.00E-06	4.33E+01	3
1031-07-08	Endosulfan sulfate	NA	1.05E+05	4
103361-09-7	Flumioxazin	NA	2.38E+05	1
104-40-5	P-nonylphenol	NA	3.24E+04	2
10540-29-1	Tamoxifen	NA	4.40E+05	1
105512-06-9	Clodinafop-propargyl	NA	1.39E+04	2
106-44-5	P-cresol	NA	5.51E+02	3
1071-83-6	Glyphosate	4.30E-07	1.60E+02	3
107534-96-3	Tebuconazole	2.00E-05	3.43E+04	2
108-62-3	Metaldehyde (tetramer)	NA	1.23E+02	3
110488-70-5	Dimethomorph	NA	1.37E+03	1
111479-05-1	Propaquizafop	NA	6.71E+04	1
111991-09-4	Nicosulfuron	NA	3.25E+02	1
114-07-08	Erythromycin	NA	1.07E+04	5
114369-43-6	Fenbuconazole	2.80E-05	5.87E+04	2
115-29-7	Endosulfan	8.10E-05	2.97E+05	4
116-06-03	Aldicarb	0.00028	2.35E+04	3
117-81-7	Di-(2-ethylhexyl)-phthalate (DEHP)	4.10E-06	1.61E+02	1
117-84-0	Di(n-octyl) phthalate	NA	1.51E+01	1
118-74-1	Hexachlorobenzene	0.0091	5.13E+04	3
119446-68-3	Difenoconazole	NA	6.43E+04	1
0120-12-7	Anthracene	0.0029	1.51E+05	3
120-72-9	Indole	0	2.95E+03	3
120068-37-3	Fipronil	0.00089	1.08E+06	2
121-75-5	Malathion	5.80E-07	3.11E+04	2
1214-39-7	1h-purin-6-amine, n-(phenylmethyl)-	NA	5.01E+02	2
121552-61-2	Cga 219417 (cyprodinil)	NA	1.40E+04	2
122-14-5	Fenitrothion	8.80E-05	9.87E+04	2
122-34-9	Simazine	7.50E-05	3.89E+04	3

12427-38-2	Maneb	1.20E-05	3.44E+04	3
128639-02-1	Carfentrazone-ethyl	NA	1.17E+05	2
129-00-0	Pyrene	0.00047	6.47E+05	2
131-11-3	Dimethylphthalate (DMP)	NA	8.35E+01	3
131341-86-1	Fludioxonil	NA	4.94E+04	2
131860-33-8	Azoxystrobin	NA	3.85E+04	1
13194-48-4	O-ethyl s,s-dipropyl phosphorodithioate	0.00049	1.06E+05	2
133-06-02	Captan	7.60E-06	4.24E+04	2
133-07-03	Folpet	4.80E-06	5.58E+05	2
135158-54-2	Cga 245704	NA	9.02E+03	3
13684-56-5	Desmedipham	NA	4.23E+04	2
13684-63-4	Phenmedipham	1.30E-06	2.10E+04	2
137-26-8	Thiram	1.20E-05	2.90E+05	3
138261-41-3	Imidacloprid	6.80E-06	1.60E+03	2
140-66-9	P-(1,1,3,3-tetramethylbutyl)phenol	NA	1.74E+04	2
142459-58-3	Fluthiamide	NA	8.71E+04	2
143-50-0	Kepone	0.042	5.95E+05	4
143390-89-0	Bas 490f	1.20E-06	8.18E+04	2
14698-29-4	Oxolinic acid	6.50E-06	1.09E+05	2
148-79-8	Thiabendazole	3.70E-06	1.70E+04	3
15299-99-7	N,n-diethyl-2-(1-naphthalenyloxy)propanamide	1.90E-06	1.96E+03	2
15307-86-5	Diclofenac	0.00043	9.72E+02	2
15545-48-9	Chlortoluron	NA	1.34E+03	3
1563-38-8	Carbofuran phenol	NA	2.57E+03	3
1563-66-2	Carbofuran	1.00E-04	5.61E+04	2
15687-27-1	Ibuprofen	0	1.17E+02	3
1570-64-5	2-methyl-4-chlorophenol	NA	3.64E+03	3
1582-09-08	Trifluralin	9.30E-05	5.38E+04	2
15972-60-8	Alachlor	NA	3.81E+04	2
16118-49-3	Carbetamide	NA	1.08E+03	2

16672-87-0	Ethephon	1.40E-05	6.80E+02	3
1689-84-5	Bromoxynil	8.80E-06	8.23E+03	3
1689-99-2	Bromoxynil octanoate	5.10E-06	9.27E+04	2
1698-60-8	Chloridazon	NA	4.65E+03	3
1702-17-6	3,6-dichloropicolinic acid	NA	4.55E+02	3
173584-44-6	Dpx-mp062	NA	7.78E+04	1
1746-01-06	2,3,7,8-TetraCDD	1.70E+03	4.72E+06	2
1746-81-2	Monolinuron	NA	9.65E+03	3
1897-45-6	Chlorothalonil	1.00E-05	5.72E+05	3
19044-88-3	Oryzalin	3.10E-06	1.10E+05	2
191-24-2	Benzo[g,h,i]perylene	0.00073	NA	4
1912-24-9	Atrazine	5.40E-05	4.37E+04	3
1918-00-9	Dicamba	6.30E-06	9.43E+02	3
1918-02-1	Picloram	2.00E-06	1.59E+03	3
1918-16-7	Propachlor	4.40E-06	3.72E+04	3
1929-77-7	Vernolate	1.50E-05	2.20E+03	3
193-39-5	Indeno[1,2,3-cd]-pyrene	0.019	NA	4
19666-30-9	Oxadiazon	0.00075	3.20E+05	2
205-99-2	Benzo[b]fluoranthene	0.081	NA	4
2050-68-2	PCB-15	NA	2.74E+04	3
2051-60-7	PCB-1	NA	2.05E+03	3
2051-61-8	PCB-2	NA	1.55E+03	3
206-44-0	Fluoranthene	0.001	5.70E+04	2
207-08-09	Benzo[k]fluoranthene	0.035	NA	4
21087-64-9	Metribuzin	4.20E-06	4.73E+03	3
21725-46-2	Cyanazine	0.00043	4.28E+04	3
218-01-09	Chrysene	0.013	NA	2
22071-15-4	Ketoprofen	0	NA	2
2303-16-4	Diallate	0.00021	2.25E+03	3
2303-17-5	Triallate	4.30E-05	9.34E+03	2
23103-98-2	Pirimicarb	6.90E-06	8.24E+02	2

2312-35-8	Propargite	1.00E-04	7.21E+04	2
23135-22-0	Oxamyl	1.10E-05	8.09E+03	3
23564-05-08	Thiophanate-methyl	4.70E-06	3.64E+03	2
2385-85-5	Mirex	0.024	8.59E+02	4
23950-58-5	Pronamide	3.70E-05	2.15E+03	2
197143	Dodine	4.40E-07	8.51E+03	2
24579-73-5	Propamocarb	1.40E-06	8.27E+01	3
25057-89-0	Bentazone	3.30E-06	1.00E+02	2
25812-30-0	Gemfibrozil	3.60E-05	NA	2
26225-79-6	Ethofumesate	NA	1.96E+03	2
26761-40-0	Diisodecyl phthalate	NA	1.30E+00	1
26787-78-0	Amoxicillin	NA	5.28E+06	1
27304-13-8	Oxychlorane	NA	7.16E+04	4
27314-13-2	Norflurazon	4.10E-06	2.54E+04	2
28553-12-0	Diisononyl phthalate	NA	9.50E+00	1
2921-88-2	Chloropyrifos	0.0012	3.12E+06	2
297-78-9	Isobenzan	0	8.14E+04	4
298-46-4	Carbamazepine	6.30E-06	3.90E+02	2
3060-89-7	Metobromuron	NA	6.72E+02	3
309-00-2	Aldrin	0.033	1.34E+05	4
32809-16-8	Procymidone	3.30E-06	4.51E+02	2
330-54-1	Diuron	1.80E-05	3.00E+04	3
330-55-2	Linuron	9.90E-05	9.93E+04	3
33284-50-3	PCB-7	NA	2.21E+04	3
333-41-5	Diazinon	0.00042	9.26E+04	2
3337-71-1	Asulam	2.20E-06	1.08E+02	3
3347-22-6	Dithianone	1.40E-05	2.12E+04	2
33629-47-9	Butralin	NA	9.85E+04	2
3380-34-5	5-chloro-2-(2,4-dichlorophenoxy)phenol	NA	6.60E+04	2
34014-18-1	Tebuthiuron	4.10E-06	6.35E+03	3
34123-59-6	Isoproturon	NA	5.78E+04	3

34256-82-1	Acetochlor	NA	3.38E+04	2
34883-43-7	2,4'-dichlorobiphenyl	NA	2.52E+04	3
35554-44-0	Imazalil base	2.50E-05	8.14E+03	2
36734-19-7	Rovral (Iprodione)	2.30E-05	3.11E+04	2
3739-38-6	M-phenoxybenzoic acid	NA	2.31E+02	2
39148-24-8	Fosetyl-aluminium	3.30E-07	7.45E+02	2
40321-76-4	1,2,3,7,8-pentachlorodibenzo-p-dioxin	NA	5.71E+08	4
40487-42-1	Pendimethalin	1.60E-06	2.29E+05	2
41394-05-02	Metamitron	NA	2.49E+02	3
41483-43-6	Bupirimate	NA	8.41E+03	2
41859-67-0	Bezafibrate	3.00E-05	6.43E+02	2
42835-25-6	Flumequine	NA	4.33E+03	2
42874-03-03	Oxyfluorfen	0.002	3.19E+04	2
439-14-5	Diazepam	0	NA	2
443-48-1	Metronidazole	3.80E-06	8.07E+01	3
465-73-6	Isodrin	NA	6.08E+05	4
481-39-0	5-hydroxy-1,4-naphthoquinone	NA	4.60E+04	3
50-28-2	Estradiol	0	1.12E+08	4
50-29-3	p,p'-DDT	0.0065	1.39E+05	2
50-32-8	Benzo[a]pyrene	0.032	8.44E+03	4
50-78-2	Acetylsalicylic acid	0	NA	3
51-03-6	Piperonyl butoxide	1.80E-05	2.06E+04	2
51207-31-9	2,3,7,8-TetraCDF	NA	4.45E+08	2
51218-45-2	Metolachlor	3.30E-06	3.35E+04	2
51338-27-3	Diclofop-methyl	NA	6.48E+04	2
51481-61-9	Cimetidine	0	NA	2
518-47-8	Fluorescein sodium	NA	1.09E+01	2
52315-07-08	Cypermethrin	1.10E-05	2.51E+07	1
52645-53-1	Permethrin	4.10E-06	5.88E+05	1
52888-80-9	Prosulfocarb	NA	1.55E+04	2
52918-63-5	Deltamethrin	2.00E-05	1.72E+06	1

53-16-7	Estrone	NA	1.18E+04	4
53-70-3	Dibenz(a,h)anthracene	0.14	3.05E+03	4
53112-28-0	Pyrimethanil	NA	1.70E+03	3
54-31-9	Furosemide	3.70E-06	NA	2
55179-31-2	Bitertanol	9.30E-05	8.11E+03	2
55219-65-3	Triadimenol	1.50E-05	2.85E+03	2
55335-06-03	Triclopyr	NA	2.43E+03	3
555-37-3	Neburon	NA	2.68E+04	2
5598-13-0	Chlorpyrifos methyl	0.0012	3.64E+05	2
56-38-2	Parathion	0.00011	3.40E+06	2
56-55-3	Benz[a]anthracene	0.0086	6.77E+05	2
563-12-2	Ethion	0.0013	1.05E+05	4
57-41-0	Phenytoin	3.30E-05	NA	2
57-62-5	Aureomycin	NA	4.33E+02	1
57-63-6	Ethinyl estradiol	0.0079	1.57E+06	4
57-68-1	Sulfamethazine	1.20E-06	NA	2
57-74-9	Chlordane	0.12	9.17E+04	4
57653-85-7	1,2,3,6,7,8-hexachlorodibenzo-p-dioxin	NA	1.52E+06	4
57837-19-1	Metalaxyl	1.60E-06	4.78E+02	2
57966-95-7	Cymoxanil	NA	5.45E+03	3
58-08-2	Caffeine	0	3.49E+04	3
58-14-0	Pyrimethamine	0	2.98E+03	2
58-89-9	Gamma-HCH (lindane)	0.0012	1.44E+05	3
5915-41-3	Terbuthylazine	NA	2.36E+05	3
5989-27-5	D-limonene	4.80E-06	1.45E+02	3
60-51-5	Dimethoate	1.10E-05	8.95E+03	3
60-54-8	Tetracycline	NA	1.25E+02	1
60-57-1	Dieldrin	0.15	3.10E+05	4
60168-88-9	Fenarimol	0.00012	1.73E+04	2
60207-90-1	Propiconazole	4.10E-05	1.11E+04	2
608-73-1	1,2,3,4,5,6-hexachlorocyclohexane	0.00077	6.99E+04	3

61-82-5	Amitrole	7.00E-05	4.90E+02	3
61213-25-0	Flurochloridone	NA	1.05E+04	2
62-73-7	Dichlorvos	0.00041	3.62E+05	3
62924-70-3	Flumetralin	NA	4.81E+05	1
63-25-2	Carbaryl	9.50E-05	2.29E+04	3
64-19-7	Acetic acid	NA	2.50E+01	3
64902-72-3	Chlorsulfuron	7.80E-06	6.12E+03	2
66215-27-8	Cyromazine	2.10E-05	1.56E+03	3
66246-88-6	Penconazole	0.00013	8.39E+03	2
67129-08-02	Metazachlor	NA	3.72E+03	2
67375-30-8	alpha-Cypermethrin	1.40E-05	1.75E+07	1
67564-91-4	Fenpropimorph	NA	5.89E+03	2
67747-09-05	Prochloraz	0.0027	1.96E+05	2
68-35-9	Sulfadiazine	NA	5.87E+03	2
68359-37-5	Cyfluthrin	3.80E-05	2.44E+08	1
69-53-4	Ampicillin	NA	1.53E+02	2
69377-81-7	Fluroxypyr	NA	1.46E+03	3
70630-17-0	Metalaxyl-M	NA	1.08E+03	2
7085-19-0	Mecoprop	3.80E-05	4.31E+02	3
709-98-8	Propanil	1.20E-05	2.07E+05	3
72-20-8	Endrin	0.019	5.90E+06	4
72-33-3	Mestranol	0	NA	4
72-54-8	DDD	0.35	1.36E+06	2
72-55-9	p,p'-DDE	0.0042	3.51E+05	2
723-46-6	Sulfamethoxazole	1.30E-06	2.35E+03	2
731-27-1	Tolyfluanide	NA	1.80E+05	2
732-11-6	Phosmet	2.20E-05	6.91E+05	2
73334-07-03	Iopromide	6.40E-07	1.20E+01	1
73590-58-6	Omeprazole	1.30E-05	NA	2
738-70-5	Trimethoprim	7.50E-06	4.98E+02	2
74070-46-5	Aclonifen	NA	3.31E+05	2

74223-64-6	Metsulfuron-methyl	1.60E-06	1.07E+04	2
759-94-4	Eptc	6.20E-06	8.54E+02	3
76-44-8	Heptachlor	0.021	6.73E+04	4
77732-09-03	Oxadixyl	NA	7.93E+01	2
79-57-2	Oxytetracycline	NA	6.81E+03	1
79-94-7	2,2-bis(4-hydroxy-3,5-dibromophenyl)propane	NA	3.09E+04	2
79127-80-3	Fenoxycarb	NA	1.65E+04	2
79277-27-3	Harmony	3.10E-05	6.43E+04	2
79622-59-6	Fluazinam	NA	3.45E+05	1
80-05-7	4,4'-Isopropylidenediphenol	3.00E-06	4.18E+03	2
8001-35-2	Toxaphene	0.23	5.27E+05	4
8018-01-7	Mancozeb	5.80E-06	2.63E+04	3
80844-07-01	Etofenprox	0.0011	2.11E+02	1
81-81-2	Warfarin	0.0011	2.70E+02	2
81777-89-1	Clomazone	NA	3.89E+03	2
82558-50-7	Isoxaben	1.80E-05	2.72E+04	2
82657-4-3	Bifenthrin	0.00034	3.29E+06	1
83-79-4	Rotenone	0.00012	2.16E+05	1
83164-33-4	Diflufenican	NA	8.48E+02	1
0834-12-8	Ametryne	NA	3.80E+04	3
84-66-2	Diethylphthalate (DEP)	3.70E-08	2.11E+02	3
84-74-2	Dibutylphthalate (DBP)	3.20E-07	3.16E+03	2
85-01-8	Phenanthrene	0.00039	8.21E+03	3
85-41-6	Phthalimide	NA	4.21E+02	3
85-68-7	Butyl benzyl phthalate	7.70E-07	2.83E+03	2
86-50-0	Methyl azinphos	8.40E-05	2.69E+05	2
86-73-7	Fluorene	7.70E-05	1.80E+03	3
86-87-3	Naphthaleneacetic acid	0	6.43E+01	3
87-51-4	Indole-3-acetic acid	0	4.60E+02	3
87-86-5	Pentachlorophenol	0.00038	4.53E+04	3
87392-12-9	S-Metolachlor	NA	5.72E+04	2

87674-68-8	Dimethenamid	NA	7.02E+04	2
88-99-3	O-phthalic acid	NA	2.60E+02	3
886-50-0	Terbutryn	0.00063	3.22E+04	3
88671-89-0	Myclobutanil	6.30E-06	1.49E+04	2
90-43-7	2-Phenylphenol	4.10E-06	4.55E+03	3
9006-42-2	Metiram	4.90E-07	1.03E+03	3
90717-03-06	Quinmerac	NA	2.49E+02	2
91465-08-06	Lambda-cyhalothrin	NA	6.93E+07	1
92-52-4	Biphenyl	6.80E-07	1.10E+03	3
93106-60-6	Enrofloxacin	NA	1.69E+06	1
94-74-6	2-Methyl-4-chlorophenoxyacetic acid	6.80E-05	9.40E+02	3
94-75-7	2-(2,4-dichlorophenoxy)acetic acid	1.60E-05	4.30E+02	3
94-82-6	2,4-DB	9.50E-06	6.92E+02	3
94125-34-5	Prosulfuron	NA	9.07E+04	1
94361-06-05	Cyproconazole	NA	2.30E+03	2
95-48-7	o-cresol	5.40E-07	2.96E+02	3
95-76-1	3,4-Dichloroaniline	NA	5.24E+03	3
97-23-4	Phenol,2,2'-methylenebis 4-chloro	NA	3.02E+04	2
98-86-2	Acetophenone	4.50E-08	3.63E+01	3
99-30-9	2,6-dichloro-4-nitroaniline	1.00E-05	8.25E+03	3
99607-70-2	Cloquintocet-mexyl	NA	7.00E+03	2
999-81-5	Chlormequat chloride	0	8.83E+01	3

Table S2 - Summary of the descriptors included in the whole TyPol database (in the first three columns) and for the 274 compounds common between TyPol and UseTox databases (in the last three columns)

Descriptors	TyPol			TyPol & UseTox		
	Min global	Max global	Nb NA (%)	Min commun	Max commun	Nb NA (%)
Connectivity index chi-0	3.58	44.67	1.09	3.58	38.96	1.46

Connectivity index chi-1	1.73	29.5	0.18	1.73	23.43	0
Connectivity index chi-2	1.73	27.87	1.09	1.73	23.46	1.46
Connectivity index chi-3	0	24.49	0.18	0	19.66	0
Connectivity index chi-4	0	20.34	1.09	0	14.47	1.46
Connectivity index chi-5	0	16.52	1.09	0	11.81	1.46
Electric dipole moment	-8.8	24.14	0.18	-8.8	15.19	0
HOMO energy	-15.04	-0.26	0.18	-15.04	-1.81	0
LUMO energy	-9.96	8.47	0.18	-4.38	4.63	0
Molecular mass	16	873.2	0	16	791.12	0
Molecular surface area (Connolly)	0	698.85	0	0	560.27	0
Number of Carbon atoms	2	48	0	2	37	0
Number of Chlorine atoms	0	12	0	0	12	0
Number of Fluorine atoms	0	6	0	0	6	0
Number of Hydrogen atoms	0	116	0	0	67	0
Number of Nitrogen atoms	0	44	0	0	44	0
Number of Oxygen atoms	0	15	0	0	13	0
Number of Phosphorus atoms	0	3	0	0	3	0
Number of Sulfur atoms	0	4	0	0	4	0
Number of aromatic bonds	0	27	0.18	0	27	0

Number of atoms	8	134	0	8	118	0
Number of bonds	4	140	0.18	7	113	0
Number of circuits	0	47	0.18	0	47	0
Number of double bonds	0	10	0.18	0	6	0
Number of halogen atoms	0	12	0	0	12	0
Number of multiple bonds	0	27	0.18	0	27	0
Number of non-H atoms	4	62	0	4	51	0
Number of non-H bonds	2	68	0.18	3	46	0
Number of rings	0	7	0.18	0	6	0
Number of rotatable bonds	0	28	0.18	0	20	0
Number of triple bonds	0	3	0.18	0	2	0
Polarizability	5.13	94.85	0.18	5.13	75.99	0
Sum of conventional bond order	0	74	0.18	0	55	0
Total energy	-11625.4	5030.03	0.18	-10037.4	-952.91	0
Valence connectivity index chi-0	2.36	38.22	1.09	2.36	32.94	1.46
Valence connectivity index chi-1	0.93	22.86	1.09	0.93	18.49	1.46
Valence connectivity index chi-2	0.52	18.96	1.09	0.52	16.47	1.46
Valence connectivity index chi-3	0	17.47	1.09	0	17.29	1.46

Valence connectivity index chi-4	0	15.94	1.09	0	15.9	1.46
-------------------------------------	---	-------	------	---	------	------

Table S3- Predicted CF_{ET} for the common compounds of the two databases with NA CF_{ET} in USEtox. The unit is the USEtox one.

CAS	Name	Cluster	Predicted CF _{ET}	Lower bound of the prediction intervals	Upper bound of the prediction intervals
191-24-2	Benzo[g,h,i]perylene	4	176978	164318	187562
193-39-5	Indeno[1,2,3-cd]-pyrene	4	176978	164318	187562
205-99-2	Benzo[b]fluoranthene	4	176846	164198	187499
207-08-9	Benzo[k]fluoranthene	4	176896	164244	187523
218-01-9	Chrysene	2	25996	14315	28110
22071-15-4	Ketoprofen	2	5318	4395	6023
25812-30-0	Gemfibrozil	2	13174	12189	16126
439-14-5	Diazepam	2	4687	4545	7272
50-78-2	Acetylsalicylic acid	3	451	399	542
51481-61-9	Cimetidine	2	5371	4304	5863
54-31-9	Furosemide	2	23463	20978	30042
57-41-0	Phenytoin	2	4109	2941	4368
57-68-1	Sulfamethazine	2	5177	4826	6983
72-33-3	Mestranol	4	178155	165342	188398
73590-58-6	Omeprazole	2	6781	4992	7587

Table S4- Predicted CF_{HT} for the common compounds without a CF_{HT} value. The predicted CF_{HT} are rounded at two decimal digits (in USEtox unit).

CAS	Name	Cluster	Predicted CF _{HT}	Lower bound for the prediction intervals	Upper bound for the prediction intervals
101-20-2	Triclocarban	2	2.3E-04	2.0E-04	2.4E-04
101-42-8	Fenuron	3	2.5E-05	1.9E-05	3.2E-05
101205-02-1	Cycloxydim	2	8.6E-06	6.9E-06	1.2E-05
102851-06-9	tau-Fluvalinate	1	4.7E-05	3.1E-05	1.1E-04
1031-07-08	Endosulfan sulfate	4	1.4E-03	1.1E-03	1.6E-03

103361-09-7	Flumioxazin	1	2.1E-05	1.7E-05	2.7E-05
104-40-5	P-nonylphenol	2	3.7E-06	3.2E-06	5.3E-06
10540-29-1	Tamoxifen	1	1.0E-04	2.3E-05	1.1E-04
105512-06-9	Clodinafop-propargyl	2	4.9E-05	4.3E-05	5.6E-05
106-44-5	P-cresol	3	1.7E-06	1.2E-06	2.1E-06
108-62-3	Metaldehyde (tetramer)	3	5.5E-06	4.8E-06	6.7E-06
110488-70-5	Dimethomorph	1	2.0E-05	1.5E-05	2.8E-05
111479-05-1	Propaquizafop	1	2.9E-05	1.8E-05	6.2E-05
111991-09-4	Nicosulfuron	1	2.7E-05	1.9E-05	2.9E-05
114-07-08	Erythromycin	5	1.8E-04	1.5E-04	2.2E-04
117-84-0	Di(n-octyl) phthalate	1	9.9E-06	7.7E-06	2.4E-05
119446-68-3	Difenoconazole	1	3.3E-05	2.0E-05	4.3E-05
1214-39-7	1h-purin-6-amine, n- (phenylmethyl)	2	6.2E-06	5.5E-06	7.7E-06
121552-61-2	Cga 219417 (Cyprodinil)	2	2.0E-05	1.8E-05	2.4E-05
128639-02-1	Carfentrazone-ethyl	2	8.3E-05	6.6E-05	9.9E-05
131-11-3	Dimethylphthalate (DMP)	3	2.0E-06	1.9E-06	2.2E-06
131341-86-1	Fludioxonil	2	1.7E-05	1.5E-05	2.0E-05
131860-33-8	Azoxystrobin	1	7.0E-05	3.6E-05	8.2E-05
135158-54-2	Cga 245704	3	2.1E-06	2.0E-06	2.5E-06
13684-56-5	Desmedipham	2	1.0E-05	9.9E-06	1.2E-05
140-66-9	P-(1,1,3,3- tetramethylbutyl)phe nol	2	3.9E-06	3.2E-06	5.4E-06
142459-58-3	Fluthiamide	2	2.9E-05	2.3E-05	3.3E-05
15545-48-9	Chlortoluron	3	6.4E-06	5.7E-06	7.2E-06
1563-38-8	carbofuran phenol	3	3.0E-06	2.5E-06	3.8E-06
1570-64-5	2-methyl-4- chlorophenol	3	9.8E-07	7.9E-07	1.3E-06
15972-60-8	Alachlor	2	7.0E-06	6.4E-06	9.3E-06

16118-49-3	Carbetamide	2	2.5E-06	2.3E-06	2.9E-06
1698-60-8	Chloridazon	3	6.4E-06	5.9E-06	7.3E-06
1702-17-6	3,6-dichloropicolinic acid	3	5.4E-06	4.9E-06	6.1E-06
173584-44-6	Dpx-mp062	1	4.6E-05	2.7E-05	1.0E-04
1746-81-2	Monolinuron	3	4.8E-06	4.4E-06	5.4E-06
2050-68-2	PCB-15	3	7.5E-05	5.8E-05	8.6E-05
2051-60-7	PCB-1	3	1.4E-05	1.2E-05	1.7E-05
2051-61-8	PCB-2	3	1.4E-05	1.1E-05	1.6E-05
26225-79-6	Ethofumesate	2	6.1E-06	5.5E-06	7.4E-06
26761-40-0	Diisodecyl phthalate	1	1.2E-05	8.1E-06	4.0E-05
26787-78-0	Amoxicillin	1	1.5E-05	1.2E-05	2.3E-05
27304-13-8	Oxychlorane	4	4.6E-02	4.1E-02	4.9E-02
28553-12-0	Diisononyl phthalate	1	1.5E-05	1.0E-05	4.5E-05
3060-89-7	Metobromuron	3	8.5E-06	7.9E-06	9.2E-06
33284-50-3	PCB-7	3	5.2E-05	4.1E-05	6.0E-05
33629-47-9	Butralin	2	2.3E-06	2.2E-06	2.8E-06
3380-34-5	5-chloro-2-(2,4-dichlorophenoxy)phenol	2	2.2E-04	1.8E-04	2.4E-04
34123-59-6	Isoproturon	3	3.2E-06	2.8E-06	3.9E-06
34256-82-1	Acetochlor	2	6.8E-06	6.1E-06	9.1E-06
34883-43-7	2,4'-dichlorobiphenyl	3	4.8E-05	3.8E-05	5.5E-05
3739-38-6	M-phenoxybenzoic acid	2	6.5E-04	5.2E-04	7.3E-04
40321-76-4	1,2,3,7,8-pentachlorodibenzo-p-dioxin	4	1.3E-02	1.1E-02	1.4E-02
41394-05-02	Metamitron	3	4.0E-06	3.6E-06	4.7E-06
41483-43-6	Bupirimate	2	3.6E-06	3.4E-06	4.4E-06
42835-25-6	Flumequine	2	1.1E-05	9.9E-06	1.3E-05
465-73-6	Isodrin	4	1.7E-02	1.4E-02	1.8E-02

481-39-0	5-hydroxy-1,4-naphthoquinone	3	2.8E-06	2.4E-06	3.2E-06
51207-31-9	2,3,7,8-TetraCDF	2	3.4E-03	2.8E-03	3.8E-03
51338-27-3	Diclofop-methyl	2	6.7E-05	5.9E-05	7.4E-05
518-47-8	Fluorescein sodium	2	2.3E-04	2.0E-04	2.6E-04
52888-80-9	Prosulfocarb	2	4.6E-06	4.1E-06	6.1E-06
53-16-7	Estrone	4	7.1E-05	5.9E-05	9.1E-05
53112-28-0	Pyrimethanil	3	8.3E-06	6.9E-06	1.0E-05
55335-06-03	Triclopyr	3	2.4E-05	2.2E-05	2.8E-05
555-37-3	Neburon	2	1.1E-05	9.8E-06	1.2E-05
57-62-5	Aureomycin	1	2.9E-05	2.0E-05	8.7E-05
57653-85-7	1,2,3,6,7,8-hexachlorodibenzo-p-dioxin	4	3.0E-02	2.3E-02	3.2E-02
57966-95-7	Cymoxanil	3	4.2E-06	3.9E-06	4.5E-06
5915-41-3	Terbutylazine	3	5.9E-06	5.3E-06	6.8E-06
60-54-8	Tetracycline	1	2.9E-05	2.1E-05	8.3E-05
61213-25-0	Flurochloridone	2	7.4E-05	6.5E-05	8.6E-05
62924-70-3	Flumetralin	1	3.4E-05	2.0E-05	3.9E-05
64-19-7	Acetic acid	3	4.7E-06	3.2E-06	5.4E-06
67129-08-02	Metazachlor	2	1.4E-05	1.3E-05	1.6E-05
67564-91-4	Fenpropimorph	2	1.1E-05	9.4E-06	1.6E-05
68-35-9	Sulfadiazine	2	2.6E-06	2.4E-06	3.0E-06
69-53-4	Ampicillin	2	1.2E-05	1.1E-05	1.4E-05
69377-81-7	Fluroxypyr	3	1.8E-05	1.6E-05	2.1E-05
70630-17-0	Metalaxyl-M	2	2.9E-06	2.6E-06	3.7E-06
731-27-1	Tolyfluanide	2	2.9E-05	2.5E-05	3.1E-05
74070-46-5	Aclonifen	2	7.6E-06	7.1E-06	8.7E-06
77732-09-03	Oxadixyl	2	2.6E-06	2.4E-06	3.0E-06
79-57-2	Oxytetracycline	1	2.5E-05	2.0E-05	8.0E-05

79-94-7	2,2-bis(4-hydroxy-3,5-Dibromophenyl)propane	2	8.0E-04	6.0E-04	8.9E-04
79127-80-3	Fenoxycarb	2	8.1E-06	7.6E-06	9.8E-06
79622-59-6	Fluazinam	1	4.7E-05	2.8E-05	6.0E-05
81777-89-1	Clomazone	2	1.4E-06	1.2E-06	1.8E-06
83164-33-4	Diflufenican	1	5.7E-05	2.9E-05	6.4E-05
834-12-8	Ametryne	3	6.7E-06	6.3E-06	7.9E-06
85-41-6	Phthalimide	3	1.6E-06	1.4E-06	1.9E-06
87392-12-9	S-Metolachlor	2	8.6E-06	7.7E-06	1.2E-05
87674-68-8	Dimethenamid	2	9.8E-06	8.9E-06	1.2E-05
88-99-3	O-phthalic acid	3	1.9E-06	1.8E-06	2.1E-06
90717-03-06	Quinmerac	2	4.5E-06	4.3E-06	5.2E-06
91465-08-06	Lambda-cyhalothrin	1	6.8E-05	2.9E-05	9.0E-05
93106-60-6	Enrofloxacin	1	1.9E-05	1.4E-05	2.6E-05
94125-34-5	Prosulfuron	1	2.9E-05	1.9E-05	3.3E-05
94361-06-05	Cyproconazole	2	1.4E-05	1.4E-05	1.7E-05
95-76-1	3,4-dichloroaniline	3	5.0E-06	4.0E-06	5.9E-06
97-23-4	Phenol,2,2'-methylenebis 4-chloro	2	1.1E-04	9.0E-05	1.2E-04
99607-70-2	Cloquintocet-mexyl	2	2.0E-05	1.8E-05	2.4E-05