

	<b>QMRF identifier (JRC Inventory): To be entered by JRC</b>
	<b>QMRF Title: Persistence (sediment) model (IRFMN) - v. 1.0.1</b>
	<b>Printing Date: Mar 31, 2022</b>

## 1. QSAR identifier

### 1.1. QSAR identifier (title):

Persistence (sediment) quantitative model (IRFMN) - v. 1.0.1

### 1.2. Other related models:

No

### 1.3. Software coding the model:

VEGA (<https://www.vegahub.eu/>)

The VEGA software provides QSAR models to predict tox, ecotox, environ, phys-chem and toxicokinetic properties of chemical substances.

emilio.benfenati@marionegri.it

## 2. General information

### 2.1. Date of QMRF:

31-03-2022

### 2.2. QMRF author(s) and contact details:

Emilio Benfenati Istituto di Ricerche Farmacologiche Mario Negri - IRCSS Via Mario Negri 2, 20156 Milano, Italy [emilio.benfenati@marionegri.it](mailto:emilio.benfenati@marionegri.it) <https://www.marionegri.it/>

### 2.3. Date of QMRF update(s):

No update

### 2.4. QMRF update(s):

No update

### 2.5. Model developer(s) and contact details:

Alberto Manganaro Istituto di Ricerche Farmacologiche Mario Negri - IRCSS Via Mario Negri 2, 20156 Milano, Italy [alberto.manganaro@marionegri.it](mailto:alberto.manganaro@marionegri.it) <https://www.marionegri.it/>

### 2.6. Date of model development and/or publication:

2015

### 2.7. Reference(s) to main scientific papers and/or software package:

[1] Benfenati E, Roncaglioni A, Lombardo A, Manganaro A. Integrating QSAR, Read-Across, and Screening Tools: The VEGAHUB Platform as an Example. *Advances in Computational Toxicology*; Springer; 2019. p. 365-81.

[2] A. Manganaro, F. Pizzo, A. Lombardo, A. Pogliaghi, E. Benfenati, "Predicting persistence in the sediment compartment with a new automatic software based on the k-Nearest Neighbor (k-NN) algorithm", *Chemosphere* (2015)

[3] Pizzo, Fabiola, Anna Lombardo, Marc Brandt, Alberto Manganaro, and Emilio Benfenati. 'A New Integrated in Silico Strategy for the Assessment and Prioritization of Persistence of Chemicals under REACH'. *Environment International* 88 (1 March 2016): 250–60. <https://doi.org/10.1016/j.envint.2015.12.019>.

### 2.8. Availability of information about the model:

The model is non-proprietary and the training set is available.

### 2.9. Availability of another QMRF for exactly the same model:

Another QMRF is not available.

### 3. Defining the endpoint - OECD Principle 1

#### 3.1. Species:

Persistence in sediment compartment measured as half-life in days.

#### 3.2. Endpoint:

ENV FATE 5.2.2. Biodegradation in water and sediment: simulation tests. OECD. Test No. 308: Aerobic and Anaerobic Transformation in Aquatic Sediment Systems (2002) [4]

#### 3.3. Comment on endpoint:

The model is based on half-life ultimate biodegradation test data and provide qualitative evaluation (four classes) of persistence property in the sediment compartment.

#### 3.4. Endpoint units:

Days.

#### 3.5. Dependent variable:

Classification in four class, P, P/vP, nP/P, nP, (see further explanation under point 3.7 below)

#### 3.6. Experimental protocol:

OECD. Test No. 308: Aerobic and Anaerobic Transformation in Aquatic Sediment Systems (2002) [4]

#### 3.7. Endpoint data quality and variability:

All employed HL ultimate degradation data referred to data on mono-constituent organic substances. HL data were collected from different sources. The first was Gouin et al., 2004 [5], which contains HL expressed as "hours" (h) for 233 organic compounds in nine classes (on a semi-decade log scale basis). This set of compounds covers four environmental compartments: water (not specified whether marine or freshwater), sediment (not specified whether marine or freshwater, but we consider that as freshwater sediment), soil and air. Mean HL is assigned to each class and is the only figure available for each chemical. As mentioned in Gouin et al 2004: "*half-lives were assigned on a semi-decade logarithmic scale to one of nine classes as follows: (1) 5 h (range: 0– 10 h), (2) 17 h (10–30), (3) 55 h (30–100), (4) 170 h (100– 300), (5) 550 h (300–1000), (6) 1700 h (1000–3000), (7) 5500 h (3000–10 000), (8) 17 000 h (10 000–30 000) and (9) 55 000 h (30 000–100 000)... To the extent possible, these allocations were made by careful analysis of experimental degradation rate data, but inevitably a high degree of scientific judgment was involved. It is recognized that by allocating a chemical to a half-life class, there is likely to be an estimation error of  $\pm 1$  to 2 classes (Mackay et al., 1999)*". The original data source is a handbook reporting an assessment based on scientific judgment of the available experimental and estimated data [9]. Another source was Gramatica and Papa, 2007 [6] which contains data for the same compartments and categories as in Gouin et al., 2004 on 250 organic compounds. Indeed, most of these figures came from Gouin et al., 2004.

For both datasets we only retained data on the sediment compartment. For this compartment, the threshold HL value for P compounds is 2880 h (120 days), and the threshold for vP compounds is 4320 h (180 days). These values have been chosen according to Annex XIII of REACH regulation (ECHA, 2014) [7]. We double-checked the correspondence between CAS number and chemical structure for all the compounds using the freely available databases ChemIDplus (ChemIDplus, 2015) and Pubchem (Pubchem, 2015). Salts, mixtures, doubtful compounds, duplicates and compounds present in both datasets but with different values were eliminated. We finally obtained a dataset with 297 mono-constituent organic compounds containing HL data for the sediment compartment.

It is important to note that both Gouin et al. (2004) and Gramatica and Papa (2007) bin HL data into one of a series of categories from which only models for nP-, P- and vP-classification purposes can be obtained, even though there are numerous degradation HL categories in Gouin et al 2004 and Gramatica & Papa 2007. Based on the threshold criteria defined above, the thresholds for identifying P and vP chemicals are representative of class 6 and class 7, respectively. Therefore, chemicals categorized as P could be found in both classes, making the selection of these substances impossible (Table 1). In an effort to provide greater

discrimination between the HL classes, we have added two classes (, and we grouped them to one of four categories: nP compounds (classes 1 to 5, i.e. compounds with HL below the P threshold), nP/P compounds (class 6, i.e. compounds in the class which includes both nP and P chemicals with HL values near the P threshold), P/vP compounds (class 7, i.e. compounds in the class which includes both P and vP chemicals with HL values near the vP threshold) and vP compounds (classes from 8 to 11, i.e. compounds above the vP threshold).

The distribution of the sediment HL data among the four experimental categories was balanced: 75 nP compounds, 69 nP/P compounds, 62 P/vP compounds, 91 vP compounds.

Differently from the model described in the cited paper [2], the k-NN has been rebuilt using all the available experimental data, indeed all compounds used as test set have been added to the training set, while the optimal configuration parameters for the k-NN model has been preserved. With this modification, the k-NN model exploits all the available information and performs slightly better than the statistics reported in the original paper.

## 4. Defining the algorithm - OECD Principle 2

### 4.1. Type of model:

The model is based on the half-lives test data and provides a qualitative evaluation (four persistency categories: nP, nP/P, P/vP and vP c.f. point 3.7 above) of persistence property in the sediment compartment. It has been developed using an ensemble of k-NN modelling and a set of alerts extracted with Sarpy software, by Istituto di Ricerche Farmacologiche Mario Negri.

### 4.2. Explicit algorithm:

This model has been built with the istKNN application (developed by Kode srl, <http://chm.kode-solutions.net>) and it is based on the similarity index developed inside the VEGA platform; the index takes into account several structural aspects of the compounds, such as their fingerprint, the number of atoms, of cycles, of heteroatoms, of halogen atoms, and of particular fragments (such as nitro groups). The index value ranges from 1 (maximum similarity) to 0. On the basis of this structural similarity index, the four compounds from the dataset resulting most similar to the chemical to be predicted are taken into account; compounds with a similarity value lower than 0.75 are discarded, and if only one compound remains available for prediction, it is kept only if it has a similarity value higher than 0.8. If no compounds fall under these conditions, no prediction is provided. The prediction is calculated as the most representative class in the compound selected with the above mentioned procedure, considering their similarity index values as a weight so that the most similar compounds have an higher influence in the prediction.

Differently from the model described in the cited paper ([2] in sec. 2.7), the k-NN has been rebuilt using all the available experimental data, indeed all compounds used as test set have been added to the training set, while the optimal configuration parameters for the k-NN model has been preserved. With this modification, the k-NN model exploits all the available information and performs slightly better than the statistics reported in the original paper.

The overall architecture provides firstly the prediction as calculated by the k-NN model. If some structural alerts are found, they do not change the prediction but modify the applicability domain value: if the alerts confirm the k-NN prediction, the applicability domain index (ADI) value increases, if the alert are in disagreement with the prediction the ADI value decreases. The alerts are anyway used to provide a prediction if the k-NN model is not able to predict the compound.

The predicted value is provided as one of the following four classes based on the standard labelling of non persistent (nP), persistent (P) and very persistent (vP) compounds: nP, nP/P, P/vP, vP. The nP class defines compound with HL certainly below the threshold of 120 days. The vP class defines compound with HL certainly higher the threshold of 180 days. The remaining two classes indicate a not entirely certain situation: nP/P class represent compounds with HL values falling near the P threshold of 120 days; the P/vP class represent compounds with HL values falling near the vP threshold of 180 days.

### 4.3. Descriptors in the model:

Similarity index

Alerts by SarPy manually modified by experts: fingerprints, the number of atoms, of cycles, of heteroatoms, of halogen atoms, and of particular fragments (such as nitro groups)

The following SAs have been extracted from the original dataset and are related to compounds which are nP in the sediment compartment:

- nP alert no. 1, defined by the SMARTS: O=CC
- nP alert no. 2, defined by the SMARTS: O=Cc1ccccc1
- nP alert no. 3, defined by the SMARTS: C(N)C
- nP alert no. 4, defined by the SMARTS: OC(C)
- nP alert no. 5, esters (aliphatic), defined by the SMARTS: O=[C;D2],[C;D3]C][O;D2][C,c]
- nP alert no. 6, aldehydes (aliphatic), defined by the SMARTS: [\$([C;D2])(=O)C]
- nP alert no. 7, ketones (aliphatic), defined by the SMARTS: [C;D3](=O)([C])[C]
- nP alert no. 8, primary amines, defined by the SMARTS: [N;D1][\$(C,c)];\$(C=[O,S])]
- nP alert no. 9, hydroxyl groups (multiple), defined by the SMARTS: [O;D1;!-]A

The following SAs have been extracted from the original dataset and are related to compounds which are vP in the sediment compartment:

- vP alert no. 1, defined by the SMARTS: c1(c2ccccc2)ccccc1
- vP alert no. 2, defined by the SMARTS: c1cc(c(cc1)Cl)Cl
- vP alert no. 3, defined by the SMARTS: CICC
- vP alert no. 4, defined by the SMARTS: c12c(cccc1)Oc3c(cccc3)O2
- vP alert no. 5, ethers (aromatic, multiple), defined by the SMARTS: [#6;!\$(C=O);!(C#N)]O[a]
- vP alert no. 6, halogen on ring C (sp<sup>3</sup>) (multiple), defined by the SMARTS: [Cl,Br,F,I][\$(C@[\*])@[\*]);!\$(C=\*)]

#### 4.4.Descriptor selection:

NA

#### 4.5.Algorithm and descriptor generation:

istKNN application

#### 4.6.Software name and version for descriptor generation:

istKNN application (developed by Kode srl, <http://chm.kodesolutions.net> )

alerts by SarPy

#### 4.7.Chemicals/Descriptors ratio:

297/15 = 20

### 5.Defining the applicability domain - OECD Principle 3

#### 5.1.Description of the applicability domain of the model:

The Applicability Domain (AD) is assessed using the original algorithm implemented within VEGA. An overall AD index is calculated, based on a number of parameters, which relate to the results obtained on similar substances within the training.

ADI is defined in this way for this QSAR model's predictions:

If  $1 \geq \text{AD index} \geq 0.85$ , the predicted substance is regarded in the Applicability Domain of the model. It corresponds to "good reliability" of prediction.

If  $0.85 > \text{AD index} \geq 0.65$ , the predicted substance could be out of the Applicability Domain of the model. It corresponds to "moderate reliability" of prediction.

If  $\text{AD index} < 0.65$ , the predicted substance is regarded out of the Applicability Domain of the model. It corresponds to "low reliability" of prediction.

The overall architecture provides firstly the prediction as calculated by the k-NN model. If some structural alerts are found, they do not change the prediction but modify the applicability domain value: if the alerts

confirm the k-NN prediction, the applicability domain index (ADI) value increases, if the alert are in disagreement with the prediction the ADI value decreases. The alerts are anyway used to provide a prediction if the k-NN model is not able to predict the sediment HL for the compound. (C.f. also point 5.2 below)

Indices are calculated on the first  $k$  = the number of  $k$  compounds used in the KNN model for the prediction most similar molecules, each having  $S_k$  similarity value with the target molecule.

**Similarity index** (*IdxSimilarity*) is calculated as:

$$\frac{\sum_k S_k}{k} \times (1 - Diam^2)$$

where *Diam* is the difference in similarity values between the most similar molecule and the  $k$ -th molecule.

**Accuracy index** (*IdxAccuracy*) is calculated as:

$$\frac{\sum_c \log(1 + S_c)}{\sum_k \log(1 + S_k)}$$

where the molecules with  $c$  index are the subset of the  $k$  molecules where the prediction of the model matches with the experimental value of the molecule.

**Concordance index** (*IdxConcordance*) is calculated as:

$$\frac{\sum_c \log(1 + S_c)}{\sum_k \log(1 + S_k)}$$

where the molecules with  $c$  index are the subset of the  $k$  molecules where the experimental value of the molecule matches with the prediction made for the target molecule.

**ACF contribution** (*IdxACF*) index is calculated as

$$ACF = rare \times missing$$

where: *rare* is calculated on the number of fragments found in the molecule and found in the training set in less than 3 occurrences as following: if the number is 0, *rare* is set to 1.0; if the number is 1, *rare* is set to 0.6; otherwise *rare* is set to 0.4

*missing* is calculated on the number of fragments found in the molecule and never found in the training set as following: if the number is 0, *missing* is set to 1.0; if the number is 1, *missing* is set to 0.6; otherwise *missing* is set to 0.4

**AD final index** is calculated as following:

$$ADI = (IdxSimilarity)^{0.5} \times IdxAccuracy^{0.25} \times IdxConcordance^{0.25} \times IdxACF$$

## 5.2.Method used to assess the applicability domain:

The Applicability Domain and the chemical similarity is measured with the algorithm developed for VEGA. Full details are in the VEGA website ([www.vegahub.eu](http://www.vegahub.eu)), including the open access paper describing it [8]. The VEGAAD also evaluates the correctness of the prediction on similar compounds (accuracy), the consistency between the predicted value for the target compound and the experimental values of the similar compounds, the range of the descriptors, and the presence of unusual fragments, using atom centred fragments.

**These indices are defined in this way for this QSAR model:**

**Similar molecules with known experimental value:**

This index takes into account how similar are the first two most similar compounds found. Values near 1 mean that the predicted compound is well represented in the dataset used to build the model, otherwise the prediction could be an extrapolation. Defined intervals are:

If  $1 \geq \text{index} > 0.8$ , strongly similar compounds with known experimental value in the training set have been found

If  $0.8 \geq \text{index} > 0.6$ , only moderately similar compounds with known experimental value in the training set have been found

If  $\text{index} \leq 0.6$ , no similar compounds with known experimental value in the training set have been found

**Accuracy (average error) of prediction for similar molecules:**

This index takes into account the classification accuracy in prediction for the two most similar compounds found. Values near 1 mean that the predicted compounds fall in an area of the model's space where the model gives reliable predictions (no misclassifications), otherwise the lower is the value, the worse the model behaves. Defined intervals are:

If  $1 \geq \text{index} > 0.9$ , accuracy of prediction for similar molecules found in the training set is good

If  $0.9 \geq \text{index} > 0.5$ , accuracy of prediction for similar molecules found in the training set is not optimal

If  $\text{index} \leq 0.5$ , accuracy of prediction for similar molecules found in the training set is not adequate

**Concordance for similar molecules:**

This index takes into account the difference between the predicted value and the experimental values of the two most similar compounds. Values near 0 mean that the prediction made disagrees with the values found in the model's space, thus the prediction could be unreliable. Defined intervals are:

If  $1 \geq \text{index} > 0.9$ , molecules found in the training set have experimental values that agree with the target compound predicted value

If  $0.9 \geq \text{index} > 0.5$ , similar molecules found in the training set have experimental values that slightly disagree with the target compound predicted value

If  $\text{index} \leq 0.5$ , similar molecules found in the training set have experimental values that completely disagree with the target compound predicted value

**Structural Alerts Concordance:**

This index takes into account the concordance between the prediction provided by the k-NN model and the alerts found. Defined values are:

If  $\text{index} = 1$ , all alerts are related to experimental values in agreement with the prediction, thus confirming the k-NN output

If  $\text{index} = 0.9$ , no alerts have been found, thus it is not possible to confirm the k-NN output

If  $\text{index} = 0.85$ , no k-NN prediction is available and the final prediction is based only on the found alerts

If  $\text{index} = 0.7$ , one or more alerts are related to experimental values not in agreement with the prediction, thus conflicting with the k-NN output

### **Atom Centered Fragments similarity check:**

This index takes into account the presence of one or more fragments that aren't found in the training set, or that are rare fragments. First order atom centered fragments from all molecules in the training set are calculated, then compared with the first order atom centered fragments from the predicted compound; then the index is calculated as following: a first index RARE takes into account rare fragments (those who occur less than three times in the training set), having value of 1 if no such fragments are found, 0.85 if up to 2 fragments are found, 0.7 if more than 2 fragments are found; a second index NOTFOUND takes into account not found fragments, having value of 1 if no such fragments are found, 0.6 if a fragments is found, 0.4 if more than 1 fragment is found. Then, the final index is given as the product RARE \* NOTFOUND. Defined intervals are:

If index = 1, all atom centered fragment of the compound have been found in the compounds of the training set

If  $1 > \text{index} \geq 0.7$ , some atom centered fragment of the compound have not been found in the compounds of the training set or are rare fragments

If index < 0.7, a prominent number of atom centered fragments of the compound have not been found in the compounds of the training set or are rare fragments

### **5.3. Software name and version for applicability domain assessment:**

VEGA ([www.vegahub.eu](http://www.vegahub.eu))

### **5.4. Limits of applicability:**

The model is not applicable to inorganic chemicals and substances containing unusual elements (i.e., different from C, O, N, S, P, Cl, Br, F, I). Salts can be predicted only if converted to the neutralized form.

## **6. Internal validation - OECD Principle 4**

### **6.1. Availability of the training set:**

Yes

### **6.2. Available information for the training set:**

CAS RN: Yes

Chemical Name: Yes

Smiles: Yes

Formula: Yes

INChI: Yes

MOL file: Yes

NanoMaterial: No

### **6.3. Data for each descriptor variable for the training set:**

All

### **6.4. Data for the dependent variable for the training set:**

All

### **6.5. Other information about the training set:**

NA

### **6.6. Pre-processing of data before modelling:**

NA

### **6.7. Statistics for goodness-of-fit:**

Leave-one-out approach (k-NN for each compound has been performed on the whole dataset without the compound itself): n = 297; Accuracy = 85%; Non predicted compounds: n = 11

Statistics are also calculated for three cases:

1. Class vP vs All other classes (nP, nP/P, vP/P)  
TP = 88, TN = 185; FP=11; FN =2, not assigned = 11, Accuracy = 95%, Specificity = 94%, Sensitivity = 98%
2. Class nP vs All other classes (vP, vP/P, nP/P)  
TP = 66, TN = 206; FP=8; FN =6, not assigned = 11, Accuracy = 95%, Specificity = 96%, Sensitivity = 92%
3. Classes vP/P + vP vs nP/P + nP  
TP = 138, TN = 124; FP=15, FN =9, not assigned = 11, Accuracy = 92%, Specificity = 89%, Sensitivity = 94%

### **6.8. Robustness - Statistics obtained by leave-one-out cross-validation:**

NA

### **6.9. Robustness - Statistics obtained by leave-many-out cross-validation:**

NA

### **6.10. Robustness - Statistics obtained by Y-scrambling:**

NA

### **6.11. Robustness - Statistics obtained by bootstrap:**

NA

### **6.12. Robustness - Statistics obtained by other methods:**

NA

## **7. External validation - OECD Principle 4**

### **7.1. Availability of the external validation set:**

NA

### **7.2. Available information for the external validation set:**

NA

### **7.3. Data for each descriptor variable for the external validation set:**

NA

### **7.4. Data for the dependent variable for the external validation set:**

NA

### **7.5. Other information about the external validation set:**

NA

### **7.6. Experimental design of test set:**

NA

### **7.7. Predictivity - Statistics obtained by external validation:**

### **7.8. Predictivity - Assessment of the external validation set:**

NA

### **7.9. Comments on the external validation of the model:**

NA



## 8. Providing a mechanistic interpretation - OECD Principle 5

### 8.1. Mechanistic basis of the model:

No assumption on the mechanism is done.

### 8.2. A priori or a posteriori mechanistic interpretation:

The model has a "a posteriori" mechanistic interpretation based on the descriptors choose

### 8.3. Other information about the mechanistic interpretation:

NA

## 9. Miscellaneous information

### 9.1. Comments:

NA

### 9.2. Bibliography:

[1] Benfenati E, Roncaglioni A, Lombardo A, Manganaro A. Integrating QSAR, Read-Across, and Screening Tools: The VEGAHUB Platform as an Example. *Advances in Computational Toxicology*; Springer; 2019. p. 365-81.

[2] A. Manganaro, F. Pizzo, A. Lombardo, A. Pogliaghi, E. Benfenati, "Predicting persistence in the sediment compartment with a new automatic software based on the k-Nearest Neighbor (k-NN) algorithm", *Chemosphere* (2015)

[3] Pizzo, Fabiola, Anna Lombardo, Marc Brandt, Alberto Manganaro, and Emilio Benfenati. 'A New Integrated in Silico Strategy for the Assessment and Prioritization of Persistence of Chemicals under REACH'. *Environment International* 88 (1 March 2016): 250–60.

<https://doi.org/10.1016/j.envint.2015.12.019>.

[4] OECD. Test No. 308: Aerobic and Anaerobic Transformation in Aquatic Sediment Systems. Paris: Organisation for Economic Co-operation and Development, 2002. [https://www.oecd-ilibrary.org/environment/test-no-308-aerobic-and-anaerobic-transformation-in-aquatic-sediment-systems\\_9789264070523-en](https://www.oecd-ilibrary.org/environment/test-no-308-aerobic-and-anaerobic-transformation-in-aquatic-sediment-systems_9789264070523-en).

[5] Gouin, T., Cousins, I., Mackay, D., 2004. Comparison of two methods for obtaining degradation half-lives. *Chemosphere* 56, 531–535.

[6] Gramatica, Paola, and Ester Papa. 'Screening and Ranking of POPs for Global Half-Life: QSAR Approaches for Prioritization Based on Molecular Structure'. *Environmental Science & Technology* 41, no. 8 (1 April 2007): 2833–39. <https://doi.org/10.1021/es061773b>.

[7] European Chemical Agency (ECHA) Guidance on Information Requirements and Chemical Safety Assessment. Chapter R.7b: Endpoint Specific Guidance Version 2.0 (2014)

[8] Floris, Matteo, Alberto Manganaro, Orazio Nicolotti, Ricardo Medda, Giuseppe Felice Mangiatordi, e Emilio Benfenati. «A generalizable definition of chemical similarity for read-across». *Journal of Cheminformatics* 6, n. 1 (18 october 2014): 39. <https://doi.org/10.1186/s13321-014-0039-1>

[9] Mackay, D., Shiu, W.Y., Ma, K.C., 1999. *Physical–Chemical Properties and Environmental Fate Handbook*; CRC netBASE CD-ROM. Chapman and Hall/CRC Press, Boca Raton, FL

### 9.3. Supporting information:

#### Training set(s) Test set(s) Supporting information:

All available dataset are present in the model inside the VEGA software.

## 10. Summary (JRC QSAR Model Database)

### 10.1. QMRF number:

To be entered by JRC

### 10.2. Publication date:

To be entered by JRC

**10.3.Keywords:**

To be entered by JRC

**10.4.Comments:**

To be entered by JRC