



King's Research Portal

DOI:

[10.1016/j.scitotenv.2018.08.122](https://doi.org/10.1016/j.scitotenv.2018.08.122)

Document Version

Peer reviewed version

[Link to publication record in King's Research Portal](#)

Citation for published version (APA):

Miller, T. H., Gallidabino, M. D., MacRae, J. R., Owen, S. F., Bury, N. R., & Barron, L. P. (2019). Prediction of bioconcentration factors in fish and invertebrates using machine learning. *Science of the Total Environment*, 648, 80-89. <https://doi.org/10.1016/j.scitotenv.2018.08.122>

Citing this paper

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

General rights

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

Take down policy

If you believe that this document breaches copyright please contact librarypure@kcl.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.

1 **PREDICTION OF BIOCONCENTRATION FACTORS IN FISH AND**
2 **INVERTEBRATES USING MACHINE LEARNING**

3 Thomas H. Miller^{a*}, Matteo D. Gallidabino^b, James R. MacRae^c, Stewart F. Owen^d,
4 Nicolas R. Bury^{ef}, Leon P. Barron^{a*}

5

6 *^aDepartment of Analytical, Environmental & Forensic Sciences, School of Population*
7 *Health & Environmental Sciences, Faculty of Life Sciences and Medicine, King's*
8 *College London, 150 Stamford Street, London, SE1 9NH, UK.*

9 *^bDepartment of Applied Sciences, Northumbria University, Newcastle Upon Tyne, NE1*
10 *8ST, UK.*

11 *^cMetabolomics Laboratory, The Francis Crick Institute, 1 Midland Road, London, NW1*
12 *1AT, UK.*

13 *^dAstraZeneca, Global Environment, Alderley Park, Macclesfield, Cheshire SK10 4TF,*
14 *UK.*

15 *^eDivision of Diabetes and Nutritional Sciences, Faculty of Life Sciences and Medicine,*
16 *King's College London, Franklin Wilkins Building, 150 Stamford Street, London, SE1*
17 *9NH, UK.*

18 *^fFaculty of Science, Health and Technology, University of Suffolk, James Hehir*
19 *Building, University Avenue, Ipswich, Suffolk, IP3 0FS, UK.*

20

21 *Corresponding authors

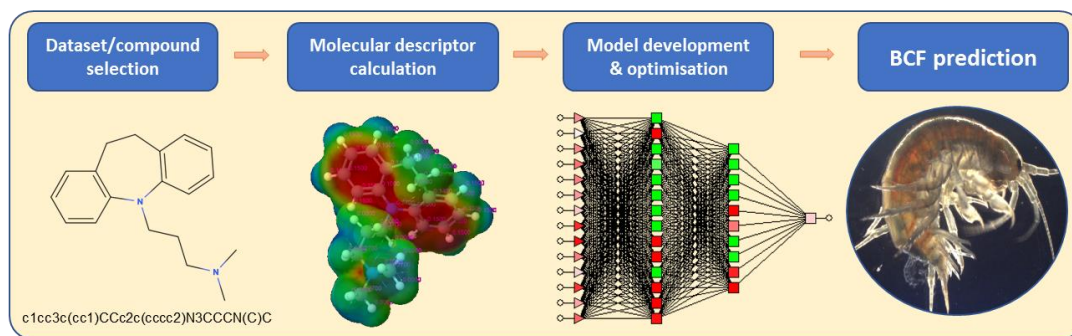
22 Email: thomas.miller@kcl.ac.uk (Tel: +44 20 7848 4978) or leon.barron@kcl.ac.uk;
23 (Tel.: +44 20 7848 3842)

24

25

26 **GRAPHICAL ABSTRACT**

27
28
29



30 **Abstract**

31 The application of machine learning has recently gained interest from ecotoxicological
32 fields for its ability to model and predict chemical and/or biological processes, such as
33 the prediction of bioconcentration. However, comparison of different models and the
34 prediction of bioconcentration in invertebrates has not been previously evaluated. A
35 comparison of 24 linear and machine learning models is presented herein for the
36 prediction of bioconcentration in fish and important factors that influenced
37 accumulation identified. R^2 and root mean square error (RMSE) for the test data (n =
38 110 cases) ranged from 0.23 – 0.73 and 0.34 – 1.20, respectively. Model performance
39 was critically assessed with neural networks and tree-based learners showing the best
40 performance. An optimised 4-layer multi-layer perceptron (14 descriptors) was
41 selected for further testing. The model was applied for cross-species prediction of
42 bioconcentration in a freshwater invertebrate, *Gammarus pulex*. The model for *G.*
43 *pulex* showed good performance with R^2 of 0.99 and 0.93 for the verification and test
44 data, respectively. Important molecular descriptors determined to influence
45 bioconcentration were molecular mass (MW), octanol-water distribution coefficient
46 (logD), topological polar surface area (TPSA) and number of nitrogen atoms (nN)
47 among others. Modelling of hazard criteria such as PBT, showed potential to replace
48 the need for animal testing. However, the use of machine learning models in the
49 regulatory context has been minimal to date and is critically discussed herein. The
50 movement away from experimental estimations of accumulation to *in silico* modelling
51 would enable rapid prioritisation of contaminants that may pose a risk to environmental
52 health and the food chain.

53 **Keywords** modelling, PBT, pharmaceutical, bioconcentration, BCF, machine
54 learning

55 **Introduction**

56 Both terrestrial and aquatic environments experience pollution from a wide
57 range of chemical contaminants. The presence of these contaminants is a cause for
58 concern as they may elicit adverse effects to environmental and public health.
59 Bioaccumulation of chemicals is critically important for understanding the risk of
60 chemicals in the environment. The complexity of confounding factors that affect uptake
61 make simple relationships that can confidently predict the accumulation elusive; but it
62 may not have to be that way.

63 Live animal exposure studies are currently the norm, using many hundreds of
64 fish for each assessment [1]. Across the European Union (EU), various guidelines
65 have been established for industry to minimise the risk posed by their chemical
66 products. For pharmaceuticals in the EU this is regulated by the European Medicines
67 Agency (EMA) and for other chemicals substances the regulations are outlined by the
68 Registration, Evaluation, Authorisation and restriction of CHemicals (REACH) [2, 3].
69 According to REACH, any manufacturer of a chemical that exceeds quantities of 10
70 tonnes per annum must submit a chemical safety assessment (CSA). For
71 environmental risk assessment, part of the CSA includes persistence,
72 bioaccumulation and toxicity (PBT) assessments. Alternatively, for pharmaceuticals
73 environmental risk assessment (ERA) follows an initial screening (Phase I) where
74 physico-chemical properties of the compound are determined (e.g. logP) and the
75 expected exposure is estimated. The Phase I exposure estimation is calculated as the
76 predicted environmental concentration (PEC). If the PEC is $>0.01 \mu\text{g L}^{-1}$ then the
77 pharmaceutical must undergo further testing to assess environmental fate and toxicity.
78 However, it should be noted that substances with a logP >4.5 , will trigger a PBT
79 assessment (following REACH guidelines) regardless of the Phase I PEC.

80 For PBT assessments, existing available screening data and prior assessment
81 information are used to determine whether a chemical is bioaccumulative (B) or very
82 bioaccumulative (vB) by estimation of a bioconcentration factor (BCF) or
83 bioaccumulation factor (BAF). Currently, pharmaceuticals are not restricted or
84 replaced as would normally be defined under REACH. Furthermore, whilst PBT
85 assessments are implemented, the persistence and bioaccumulation outcome of
86 these assessments are not taken into consideration for authorisation purposes, as no
87 legal provisions specifically cover persistent, bioaccumulative and toxic substances
88 for pharmaceuticals [4].

89 Laboratory testing for PBT brings with it a significant level of planning, quality
90 control and cost [1]. Therefore, *in silico* methodologies to predict BCF or BAF offers a
91 potential advantage to more intelligently use data to characterise potential exposure
92 and risk. Quantitative Structure Activity Relationships (QSARs) are becoming
93 increasingly popular within ecotoxicological fields as they represent, perhaps, the only
94 realistically feasible scenario to assess the environmental risk of the several thousand
95 chemicals that are available on the market [5]. In addition, such models can be used
96 to ethically reduce or replace animal testing and falls under the replacement, reduction
97 and refinement (3Rs) framework [6]. Further, effective *in silico* models could also be
98 utilised to help shape future drugs in terms of 'green by design' ambitions [7].

99 More recently, more complex machine learning-based QSAR models involving
100 artificial neural networks (ANNs), tree-based learners or support vector machines
101 (SVMs) have been used to model BCF in fish [8-11]. However, several variations of
102 machine learning-type models exist and wider applications of such models for
103 bioaccumulation prediction have not yet been evaluated to identify any added benefits.
104 Furthermore, current QSAR models have only been applied to modelling fish

105 bioaccumulation data and do not incorporate pharmaceutical data. The potential for
106 application to other taxa such as invertebrates is also non-existent, mainly due to a
107 shortage of available data.

108 The aim of this work was to develop and critically evaluate several machine
109 learning-based modelling tools for prediction of bioconcentration factor (BCF) in both
110 a fish (*Cyprinus carpio*) and an invertebrate species (*Gammarus pulex*) for the first
111 time. An open access fish BCF dataset was used in the first instance to build and
112 compare 24 different models for 352 different compounds. Subsequently, the best
113 model was applied to both a set of fish and invertebrate BCF data to assess its
114 potential for cross-species prediction. The invertebrate dataset also contained mainly
115 pharmaceuticals. In parallel, independent models were developed *ab initio* on a
116 smaller set of invertebrate BCF data alone to assess the degree of commonality with
117 the model developed on fish BCF data. Finally, the importance of molecular
118 descriptors to understand the potential for a chemical to accumulate in biota was
119 assessed. The use of such rapid and flexible modelling approaches is now critical to
120 support the 3Rs, aid greener design and to help meet the demand for PBT
121 assessments of potentially large numbers of compounds, which could be expanded to
122 new and emerging environmental contaminants across different species.

123

124 **Materials and Methods**

125 *Dataset generation and pre-processing*

126 Bioconcentration factors were collated from the European Chemical Industry
127 Council Long-range Research Initiative (Cefic LRI) project EC07 in collaboration with
128 European Academy for Standardisation e.V (EURAS) which established the BCF gold
129 standard database across multiple fish species and is freely available at

130 <http://ambit.sourceforge.net/euras/>. BCFs were down-selected to reduce variability
131 between different species and experimental conditions within the database. The BCF
132 data used herein were specific to *C. carpio* and were included by the Chemicals
133 Inspection and Testing Institute [12]. Out of all BCF data, this sub-selection resulted
134 in the largest dataset with a single fish species ($n=352$) for modelling purposes. The
135 reported BCFs represented whole-body values only and included pigments,
136 pesticides, fungicides, herbicides, insecticides, polyaromatic hydrocarbons (PAHs)
137 and polychlorinated biphenyls (PCBs), organochlorines, nitroaromatics, alkylphenols,
138 aromatic hydrocarbons, organosulfurs and organotins. Approximately 36 % of the
139 dataset contained ionisable compounds (estimated from ACD labs, Percepta
140 software). The invertebrate BCF dataset ($n=34$) was collated from literature reported
141 data [13-17] for the benthic freshwater organism, *G. pulex*. This species was selected
142 as there was a relatively large amount of BCF data available when compared with
143 other invertebrate species. For these, BCF data were only available for
144 pharmaceuticals and pesticides and, again, represented whole-body values.

145 Simplified molecular input line entry system (SMILES) strings were generated
146 for each compound using Chemspider (Royal Society of Chemistry, UK). Molecular
147 descriptors were generated from SMILES strings using Parameter Client (Virtual
148 Computational Chemistry Laboratory, Munich, Germany), and ACD Labs Percepta
149 (Advanced Chemistry Development Laboratories, ON, Canada). Approximately 450
150 descriptors were initially generated covering constitutional, topological, geometrical
151 and physico-chemical properties. The fish and invertebrate datasets were pre-
152 processed to remove any zero variance descriptors or descriptors that were
153 erroneous. All BCF data used for modelling was log transformed for improved
154 predictive accuracy.

155

156 *Feature selection*

157 Descriptors were down-selected using three different feature selection
158 algorithms, the first of which was a genetic algorithm (GA). The GA parameters were
159 set to population = 500, generations = 250, mutation rate = 0.1 and cross-over rate =
160 1. The remaining two selection methods were part of stepwise regression which
161 included a forward selection algorithm (FA) and backwards selection algorithm (BA).
162 The feature selection algorithms used a generalised regression neural networks
163 (GRNN) to monitor the error associated with the selected descriptors, where descriptor
164 sets were optimised when the error showed no improvement. The use of GRNN for
165 descriptor selection is very fast and requires minimal processing power. The
166 performance of each feature selection algorithm was characterised by then testing
167 several thousand neural networks and evaluating the predictive performance of the
168 models based on the error of the predictions. The best feature selection method was
169 the GA, which resulted in the down-selection of descriptors to a total of 14 that included
170 6 topological descriptors; radial centric information index (ICR), Narumi harmonic
171 topological function (Hnar), ramification index (Ram), superpendentic index (SPI),
172 spanning tree number (STN), topological polar surface area (TPSA), 4 constitutional
173 descriptors; number of hydrogens (nH), number of carbons (nC), number of nitrogens
174 (nN), molecular weight (MW), 3 electrotopological descriptors; maximal
175 electrotopological negative variation (MAXDN), maximal electrotopological positive
176 variation (MAXDP), mean atomic Sanderson electronegativity (Me) and 1 physico-
177 chemical property; the octanol-water distribution coefficient (logD) (See SI, Table S3).

178

179 *Modelling approaches*

180 Two different software packages were used to assess the applicability of
181 several *in silico* models in predicting bioconcentration. Trajan 6.0 (Trajan Software
182 Ltd., Lincolnshire, UK) was used to build and evaluate artificial neural networks. In
183 addition, this software was also used for the feature selection and the same
184 descriptors were used in both modelling software packages. Models developed and
185 optimised in Trajan included generalised regression neural networks (GRNN), radial
186 basis function networks (RBF) and 3-/4-layer multilayer perceptrons (MLP). Training
187 of the MLPs used two training algorithms referred to as back propagation (BP) and
188 conjugate gradient descent (CGD), models were trained for 100 iterations. The
189 optimised model was a four-layer MLP. The first and fourth layers were the inputs
190 (molecular descriptors) and outputs (logBCF), respectively. The second and third
191 layers (hidden layers) contained 14 and 10 nodes, respectively. Regularisation was
192 performed with the use of early stopping to prevent over-training of the dataset.
193 Parameter tuning was performed by changing the number of hidden layers and nodes
194 and assessing the model performance on the verification and test subsets. The
195 subsets of cases presented to the neural networks were split so that 242 compounds
196 (70 %) were used for training, 55 compounds (15 %) for verification and 55 compounds
197 (15 %) for testing the networks. Normalisation of the input features showed no
198 improvement in performance of the networks and training was performed without
199 centred or scaled descriptors.

200 In the second software package, modelling was performed using the R
201 statistical computing language (freely available from <https://www.r-project.org>). Here,
202 19 predictive models from different kinds of learner categories including both linear
203 and non-linear models were trained and tested. These included, ordinary least-
204 squares regression (OLM, package: *stats*), partial least-squares (PLS, package: *pls*),

205 ridge regression (RR, package: *elasticnet*), elastic net (EN, package: *elasticnet*),
206 quantile regression with LASSO penalty (QRL, package: *rqPen*) multivariate adaptive
207 regression splines (MARS & B-MARS, package: *earth*), k-nearest neighbours
208 regression (KNN, package: *caret*), extreme learning machines (ELM, package:
209 *elmNN*), support vector machines with radial basis function (SVM-R, package:
210 *kernlab*) and polynomial (SVM-P, package: *kernlab*) kernels, random forest exploiting
211 classification and regression trees (RF-CART, package: *randomForest*) and
212 conditional inference trees (RF-CIT, package: *party*) algorithms as base learners,
213 boosted trees (BT, package: *gbm*) and Cubist regression (CR, package: *Cubist*). MLPs
214 (3-5 layers) with 1 hidden layer (ANN-1HL, package: *nnet*), averaged 1 hidden layer
215 (ANN-a1HL, package: *nnet*), 2 hidden layers (ANN-2HL, package: *RSNNS*) and 3
216 hidden layers (ANN-3HL, package: *RSNNS*) were also tested. For this modelling
217 approach, the same molecular descriptors and logBCF were used again as input and
218 output variables. The dataset was split into two subsets, training data (70 %) and test
219 data (30 %). Normalisation of the data was required for the modelling application and
220 the dataset was both centred and scaled. Parameter tuning was performed by
221 resampling of the training subset following a 10-fold cross-validation scheme repeated
222 five times and implemented through the *caret* package. Performance of each model
223 was assessed from the root-mean square error (RMSE) and the correlation coefficient
224 (R^2). The best model for each regression method was then selected, retrained on the
225 entire training dataset and used to predict cases in the test dataset. Final datasets
226 used for modelling the optimised models are given in the SI (Table S1 & S2). The
227 finalised models were all tested according to OECD guidelines [18] for QSAR model
228 validation.

229

230 **Results and Discussion**

231

232 *Down-selection of input features for modelling BCFs in fish*

233 The down-selection of the input features was assessed using three different
234 feature-selection algorithms. Stepwise methods that included forwards or backwards
235 selection (FA/BA) reduced the number of descriptors from 180 down to 72, whilst the
236 GA reduced the number of descriptors to 66. The GA showed better correlation
237 between selected descriptors with logBCF compared to stepwise algorithms (Figure
238 S1). For both BA and FA, the selection process converged to the same local minima
239 indicating that there was no difference in using either algorithm. The improved
240 performance of the GA is due to selection of descriptors from multiple points in the
241 descriptor space, as opposed to FA or BA that start selection from a single point. Thus,
242 approaching global minima is more likely to arise when using the GA over stepwise
243 selection methods.

244 From the 66 descriptors selected by the GA, the top 22 descriptors plus an
245 additional two user curated descriptors were selected for further modelling (See SI,
246 Table S3). These additional descriptors were logD and number of hydrogen acceptor
247 groups (nHAcc) and were chosen for their previously demonstrated influence on
248 accumulation in biota [19, 20]. All descriptors were then tested across several
249 thousand MLPs (three and four-layer) where the Trajan software sub-selected the best
250 from the group of 24 descriptors based on model performance (MLPs yielded the best
251 performance over other model types in terms of R^2 and RMSE). The descriptors were
252 down-selected to a total of 14 that showed relatively good performance across MLPs
253 tested and were subsequently used in both modelling approaches discussed herein
254 (Table S3). Given the scale of BCF data used for training ($n=242$), the 5:1 Topliss

255 threshold set out by the OECD guidelines [18] for the ratio of numbers of cases to
256 descriptors was acceptable at 17:1.

257

258 *Comparison of model performances for prediction of fish BCFs*

259 The results of both modelling approaches are shown in Table 1. For models
260 trained in R, the highest RMSE values were observed for OLM (1.203), followed by
261 PLS (1.164) and then QRL (1.112). The relatively poor performance of such linear
262 models may be expected as modelling such a biologically complex process is not likely
263 to follow linear relationships using simple molecular descriptors. Even with well-
264 studied descriptors, such as logP, there is a non-linear trend with accumulation over
265 a specific threshold (generally, logP >6) [21]. However, when used as a sole
266 descriptor, logP may exclude processes that are also important for accumulation. For
267 example, elimination and metabolism rates may impact net accumulation as well as
268 more specific physiology such as carrier mediated transport and protein binding [22]
269 will also influence accumulation, especially for emerging contaminant classes such as
270 pharmaceuticals. By comparison, better performance was achieved using higher
271 complexity models. The lowest RMSEs were observed for RF-CART (0.771), followed
272 by BT (0.789) and RF-CIT (0.821), i.e. three tree-based machine learners. Next, ANNs
273 and SVMs performed very similarly to tree learners, e.g. SVM-R (0.841), ANN-a1HL
274 (0.859) and ANN-3HL (0.880).

275 Models tested in Trajan showed particularly good performance, in comparison
276 to those built in R. The lowest RMSE value was observed for a 4-layer MLP (0.524),
277 followed by 3-layer MLP (0.538), RBF (0.689), GRNN (0.893) and Linear (1.052). In
278 absolute terms, definitive conclusions cannot be drawn from direct comparison of
279 modelling approaches (i.e., Trajan vs. R), as tuning and training methods between

280 modelling software packages are slightly different. However, overall results converged
281 to support the higher reliability of non-linear approaches for modelling logBCF from
282 molecular descriptors.

283 Model complexity does not necessarily mean better predictive performance by
284 default, as several non-linear machine learners did not perform well at all. These
285 included ELM and SVM-P, where the RMSE values observed on the test set were >1.
286 Although ELM is a feedforward neural network, the weights associated with the
287 neurons in the network are not updated and thus the initialisation of the network is a
288 random selection of weights that may not model the output reliably. The EN
289 outperformed QRL and RR models, where the EN is a combination of the penalties
290 (L1 and L2 regularisation) used by both models that usually leads to better predictive
291 performance. The RR model RMSE for the test set data was also lower than the RMSE
292 for the QRL model. This can be observed when comparing RR and QRL methods, as
293 the penalty associated with LASSO can lead to the omission of highly correlated
294 covariables and thus lead to lower model robustness.

295 Limitations of predictive performance may also stem from the raw data. For
296 example, the dataset used herein did not report individual experimental pH, but instead
297 reported a range from 6.0 to 8.5. Therefore, descriptors such as logD that require pH
298 data may become limited and especially where molecular pK_a lies within this 2.5 pH
299 unit range. LogD has been shown in several works to influence uptake and
300 accumulation [23-25]. As a compromise, we calculated logD at pH 7, but this may have
301 been different to the exact experimental pH and may have added to predictive
302 inaccuracy across the whole analyte set. Lastly, it is also likely that BCF/BAF
303 prediction will be influenced by variance in biotic factors such as ventilation rates, age,

304 genetic factors and metabolism and lay beyond our ability to determine in more detail
305 [26, 27].

306 MLP models trained in Trajan offered the best performance. Consequently, this
307 model was chosen for further investigation in line with the OECD validation guidelines
308 to assess validity of QSAR modelling. The mean absolute error (MAE) corresponded
309 to 0.38 logBCF units for the verification subset (internal validation set) and 0.53
310 logBCF for the test subset (external validation set), as shown in Table 1. The RMSE
311 for verification and test subsets were 0.524 and 0.644, respectively. The predictive
312 performance of this model was better or comparable to all models in the literature that
313 have attempted to model accumulation processes. Dearden and Shinnawei [28] used
314 a linear QSAR approach to predict BCFs for 135 chemicals with an R^2 of 0.637 and
315 RMSE of 0.661 logBCF units. Another QSAR model by Sahu and Singh [29] used
316 multiple linear regression to predict BCFs for 131 organic compounds with a RMSE of
317 0.556 log units. However, this model was not validated against a test subset and
318 therefore generalised applicability of the model performance is arguably limited.

319 In alternative approaches to linear QSAR models, other machine learning
320 approaches have also been reported [8-10]. A MLP predicted BCFs for 9 test
321 compounds with an average absolute error of 0.33 ± 0.22 log units [8]. Whilst the errors
322 were low, too few compounds were tested to provide a reliable assessment of its
323 generalisability. In another approach, Zhao et al., [10] used SVM, RBF and MLR
324 models individually. Better performance was observed when two RBF models (using
325 different descriptors) were combined into a 'hybrid' model to predict logBCF. The
326 developed model showed an R^2 of 0.6917 for an external test set with a reported
327 RMSE of 0.69 logBCF units for 119 compounds showing similar performance to the
328 fish-based MLP presented here, using a single MLP. The hybrid model also showed

329 a limitation in the training set, where several cases were not modelled correctly
330 between the ranges of logBCF 4 to 5 and was observed by a plateau in the regression
331 analysis.

332

333 *A remark on outliers and the applicability domain*

334 Training and testing of all models led to the observation of several common
335 outliers. The reason for poor prediction for such cases may stem from under
336 representation in the dataset used for modelling. The spread of input and output data
337 between training and validation subsets showed that there was no significant
338 difference between the spread or skew of the data (Figure S2). However, using PCA
339 analysis and distances between the descriptor spaces there were several cases that
340 did not cluster well with the remaining data (Figure 1a). For example, logBCF for
341 perfluorotributylamine was predicted poorly across the majority of trained models. The
342 use of PCA and descriptor data spacing in this way enabled characterisation of the
343 applicability domain (AD) for a given model. A threshold may then be used to
344 determine cases that fall outside the domain and are likely to have higher predictive
345 error (Figure 1b) [30, 31].

346 According to the OECD QSAR model validation guidance [18], consideration of
347 models for regulatory purposes must be associated with a defined domain of
348 applicability under Principle 3. However, one key consideration in the use of distance-
349 based ADs is that input descriptors are not used equally by the model [32]. Therefore,
350 such ADs may not accurately identify those cases having a greater predictive error in
351 every case. This was observed for outliers in the PCA analysis, but where logBCF was
352 predicted relatively well and *vice versa*. For example, di-2-naphthyldisulfide was not
353 an outlier in the AD but was poorly predicted across all models. On the other hand,

354 pigment yellow-12 was an AD outlier, but logBCF was predicted well by the majority
355 of models.

356 Poor predictive accuracy for molecularly similar compounds could be also
357 caused by other factors such as poor quality raw data or too few representative training
358 cases for the model to learn from. It has been shown previously that experimental BCF
359 data can vary from 0.42 to 0.75 log units [9, 33, 34]. Nevertheless, even with the
360 limitations associated with defining an AD, it is useful and important to identify any
361 cases that might not be reliably predicted so that rapid prioritisation of compounds can
362 begin. Only for these cases, may it then be appropriate to revert to experimental
363 testing.

364

365 *Machine learning in a regulatory context*

366 Several of the developed machine learning tools in Table 1 showed potential
367 for the replacement and reduction in animal use. However, it is important to recognise
368 the complexities of machine learning approaches from the outset, especially where
369 they are intended for use in regulation. Under Principle 2 of the OECD guidelines,
370 models used in this way must be based on “unambiguous algorithms”. In particular, it
371 is highlighted that two significant limitations exist regarding artificial neural networks,
372 for example. These are: (a) the necessity for large (BCF) datasets to develop suitable
373 models (which do not exist for some classes of compounds, like pharmaceuticals) and
374 also (b) that these types of machine learning tools are more ambiguous than other
375 types of model, especially those that are linear in nature. For the latter, the guidance
376 is vague concerning appropriateness of ANNs for use under this specific principle but
377 infers that it is an acceptable limitation. Furthermore, the definition of an unambiguous
378 algorithm is in fact ambiguous and should be further refined to prevent confusion to

379 the reader. This principle could be applied in different ways to different models and
380 may cover the generation of molecular descriptors, the feature selection algorithms
381 used, the learning process (for machine learners where the ambiguity lies) and the
382 final model [35]. The majority of the literature seems to have focused on linear models
383 perhaps as a result, mainly to aid in mechanistic understanding and to allow expert
384 interpretation of individual chemicals to provide extra assurance in predicted data
385 (linked to Principle 5).

386 Principle 5 of the OECD guidelines relates to mechanistic interpretability of
387 QSAR models (if possible). This can be considered a limitation for machine learning
388 algorithms if the aim is to achieve an interpretable model, such as would normally be
389 expected of linear models such as OLS or PLS regression. The OECD guidelines also
390 remain vague regarding mechanistic interpretation of machine learners. However,
391 whilst linear relationships may not be apparent, descriptor sensitivity analyses can
392 indicate the importance of individual descriptors and thus enables interpretation of
393 factors that influence the modelled process. Bioconcentration processes are not
394 simple and extensive datasets are extremely impractical to curate experimentally.
395 Therefore, complex non-linear models may provide a more rapid solution to regulatory
396 decision-making meantime. Therefore, we suggest that guidelines for QSAR model
397 validation need to be expanded to better define the scope of applicability of all the
398 different types of machine learning tools and their fitness for purpose in a regulatory
399 context.

400 For PBT testing, the same regulations are triggered when a threshold for
401 bioaccumulation is reached, regardless of the extent to which the threshold is
402 exceeded. Thus, if the value is classified within the correct category of non-
403 bioaccumulative (nB), bioaccumulative (B) or very bioaccumulative (vB), the model will

404 be useful in the context of PBT assessments. Variability in measurement can arise
405 from kinetic modelling approaches [17], biological/physiological variability (age, health,
406 lipid content etc.) [27, 36-39] and experimental conditions (pH, temperature, etc.) [23,
407 40]. As such, reported BCFs have been shown to differ by 1-2 orders of magnitude
408 even within the same species [27].

409 The 4-layer MLP here showed a correct classification rate of 90 % across the
410 verification and test subsets. The 10 % misclassification of cases was split to 6 % of
411 cases predicted as false negatives and 4 % of cases predicted as false positives (See
412 SI, Figure S3). This is consistent with the hybrid model developed by Zhao et al. which
413 has shown classification accuracies ranging from 91 % to 98 % [9, 10]. It is possible
414 that using QSARs for classification instead of regression analysis may improve the
415 accuracy and without the need for the application of a bias. This would be particularly
416 suitable for bioaccumulation assessments where only a threshold value determines
417 the level of regulation enforced.

418 Some studies have reported the application of models for classification of
419 bioaccumulation thresholds, with accuracies ranging from 84.5 – 91.1 % (depending
420 on model type) [41] and 91.7 % [11]. The authors that used tree-based learners also
421 used these models for quantitative prediction achieving RMSE of 0.554 and R^2 of
422 0.836 on the test set data [11]. The models tested across the literature have tended to
423 achieve similar performance for both classification and prediction. The agreement in
424 performance between different works and the comprehensive model evaluation here,
425 support that *in silico* methods should be adopted for chemicals where environmental
426 uptake data are limited to enable flexible, cheap and rapid PBT assessment for
427 compound prioritisation. Furthermore, it suggests that the use of chemical descriptors
428 may only be able to achieve a certain level of predictive or classification performance

429 for modelling approaches where other variables become important as mentioned
430 above.

431

432 *Can the developed model be used for cross phylum prediction?*

433 There is little understanding of whether accumulation will be similar across the
434 invertebrate phylum. The dominant site of uptake for waterborne micropollutants in
435 fish is across the gills and therefore accumulation across taxa may be significantly
436 different for differing modes of respiration. Other factors such as size, enzyme
437 speciation and lipid content may also influence the accumulation potential [27]. The
438 optimised model for fish was applied to the prediction of logBCF in a freshwater
439 invertebrate, *Gammarus pulex* (Figure 3a). The accumulation data in *G. pulex*
440 predominantly covered pharmaceuticals and pesticides. The fish-based MLP showed
441 relatively low predictive performance for the invertebrate accumulation factors. The
442 correlation between observed and predicted BCF was R^2 0.3295 with a MAE of 0.80
443 \pm 0.65 log units, which indicated that the model generalisations between species were
444 limited. The largest predictive error was for the compound imipramine that was
445 overestimated by 2.7 logBCF units. This compound in a previous study had
446 considerable variation in the estimated BCF (212 – 4533) depending on the method
447 of estimation used [17].

448 A significant difference in BCFs between trophic levels has been shown with
449 higher trophic levels displaying increased BCFs [42]. This trend would suggest that
450 the BCF predictions of the invertebrates might be overestimated but the opposite was
451 observed (62 % of cases were underestimated). In addition to the biological complexity
452 between species, another confounding factor to affect the predictive accuracy and
453 generalisability is the compound class. The fish model included no pharmaceutical

454 compounds whereas the invertebrate BCF data contained 18 cases (~53%).
455 Inspection of the molecular similarity between the datasets indicated that the
456 invertebrate and fish datasets were dissimilar (Figure S4). Thus, the bioconcentration
457 potential may not follow the same relationships with neutral hydrophobic organic
458 contaminants.

459 The fish-based model was subsequently reinitialised and trained on the
460 invertebrate dataset only (using the same descriptors) (Figure 3b). The invertebrate
461 model showed good correlation with R^2 of 0.9605 with 0.972 for the training set, 0.9932
462 for the verification set and 0.9323 for the test set. The model demonstrated good
463 accuracy across the verification and test subset with a MAE of 0.07 ± 0.08 logBCF
464 units for the verification set and 0.29 ± 0.27 logBCF units for the test set. The
465 successful retraining of the model to invertebrate data suggests that case
466 representation (i.e. compound class) is likely to limit models that are applied across
467 taxa. An alternative approach to overcome this could involve development of a model
468 with two or more outputs to represent different species, but commonality in BCF cases
469 would be required for both species. Whilst the predictive accuracy of the retrained
470 model was very good, it is also limited by the small number of cases used.
471 Generalisability is also likely to be limited given the ratio of cases to descriptors
472 (Topliss ratio of ~2.5:1) Nevertheless, and as new BCF data emerges, this approach
473 holds excellent potential by using the same molecular descriptors for BCF predictions
474 in two very different species. In addition, to using the fish-based model to predict
475 invertebrate BCFs we also used the invertebrate-based model to predict fish BCFs of
476 pharmaceuticals reported in the literature (Figure S5). The invertebrate model was
477 able to predict BCFs within the reported range for 45 % of the compounds selected (n
478 = 11). The remaining compounds, with the exception of sertraline and gemfibrozil,

479 were predicted relatively well even though they were not within the reported ranges.
480 Sertraline is an interesting case as although it has not shown very high
481 bioconcentration in fish (BCFs: <1 – 626) [43-47] there have been reported BCF values
482 of up to 32,022 in invertebrates (namely, *Lasmigona costata* [48] and 990 in *Planorbis*
483 *sp.* [49]). As the model used here was trained on BCFs from an invertebrate species,
484 it may not correlate well with fish BCF data, suggesting that cross-phylum predictive
485 modelling may be limited by both case representation and biological variation.
486 However, as the models here used the same descriptors this enables flexibility in
487 retraining optimised models and inevitably as more BCF data is generated for the
488 same compounds in different species, this technology could be used to map
489 accumulation across taxa more effectively. It is critically important to understand
490 uptake (internal concentration) across taxa as the conservation of pharmaceutical
491 targets extends widely [50].

492

493 *Model sensitivity to descriptors: interpreting accumulation through chemistry*

494 Whilst machine learning models are more difficult to interpret due to the non-
495 linear functionality, collinearity and/or curvilinearity; the importance of the 14
496 descriptors described here still offered some mechanistic understanding of the
497 processes involved (Figure 4). For the fish-based model, the most important descriptor
498 was TPSA with an error ratio of 2.08. Higher error ratios correspond to increased
499 predictive error for all compounds upon removal of this descriptor from the dataset.
500 Previous investigations have demonstrated that descriptors related to polarisability,
501 hydrophobicity and hydrogen bonding of the molecule is important to modelling BCFs
502 [10, 28, 51]. TPSA is defined as the surface area occupied by nitrogen and oxygen
503 atoms including connected hydrogen atoms [52]. Polar surface area has also been

504 shown to influence drug absorption in humans, where increasing polar surface area
505 decreases the drug fraction absorbed [19, 53]. The relationship between
506 bioconcentration and TPSA may be dependent on several factors such as permeation
507 through the lipid bilayer, binding of polar functional groups to epithelial membranes
508 and the size of hydration shell around a molecule [54].

509 Permeation through cellular membranes was further supported by the
510 importance of MW to the model. The size of a molecule also affects permeation and
511 diffusion through membranes (Lipinski's rule of five [55]). It has previously been
512 demonstrated that dye pigments did not show bioaccumulation in fish due to their large
513 molecular size [56]. In another study, it was suggested that there is a threshold
514 diameter value of 1.5 nm which governed bioconcentration in addition to
515 hydrophobicity [57]. Stempel et al., [11] also found that molecular weight, molecular
516 diameter, TPSA and logD were important for classification and prediction of
517 bioaccumulation.

518 Topological descriptors such as STN, Hnar, Ram, SPI and ICR were also found
519 to be important. These indices are useful especially for differentiating constitutional
520 isomers (except enantiomers) [58]. Error ratios for STN, Hnar, ICR, SPI and Ram
521 spanned from 1.31 – 1.72. These indices are related to molecular branching/shape
522 and the importance of these descriptors relate to molecular size which can influence
523 bioconcentration [59, 60]. MAXDN and MAXDP relate to the partial charges on atoms
524 relative to their topological position within the molecule and therefore relate to the
525 nucleophilicity and electrophilicity of a molecule [61]. Aside from polarity-related
526 accumulation across cellular membranes, it is also possible that these are associated
527 with metabolic activity (from nucleophilic or electrophilic attack). The importance of

528 other electrotopological descriptors (along with molecular flexibility) has been
529 previously shown for modelling bioconcentration [62].

530 Interpretation of the relative importance of descriptors is affected by collinearity
531 or multicollinearity (See SI, Table S4 & S5). The collinearity of the descriptors showed
532 that molecular weight was collinear with SPI ($R=0.794$) and Ram ($R=0.696$). The
533 descriptor Ram was also collinear with SPI ($R=0.787$) and STN was collinear with
534 HNar ($R=0.748$). The relation between these topological descriptors and molecular
535 weight is that they all describe molecular size (shape, volume, weight) to some extent.
536 Therefore, the rank importance of these particular descriptors should be approached
537 with some caution. Whilst the error ratio is higher for certain descriptors that are
538 collinear, their removal from the network model may not correctly determine the ratio
539 value due to redundant information. Nevertheless, the descriptor sensitivity can still be
540 useful for directing mechanistic and experimental studies. This was shown recently in
541 a neural network application to passive sampling [63] which was later followed by a
542 mechanistic study [64], that supported the interpretation of the model.

543 The invertebrate-based MLP used the same descriptors as the fish-based
544 model, but the network was reinitialised and retrained. The retraining of the network
545 also showed that the importance of the descriptors changed from the fish-based
546 model. The most important descriptor was HNar (error ratio = 5.75) followed by nN
547 (error ratio = 5.09) and logD (error ratio = 4.71). The increased importance of the
548 number of nitrogen atoms likely reflected the number of pharmaceutical compounds
549 in the dataset. In addition, logD increased in rank to the top three descriptors in the
550 invertebrate model. The increased sensitivity of the model to logD also relates to
551 training of the model with ionisable pharmaceuticals and is in agreement with other
552 studies showing logD to be important in accumulative processes [11, 64]. Whilst

553 hydrophobicity may be a principal factor of bioconcentration, it is possible that carrier-
554 mediated transport may also play an important role. Both models here demonstrated
555 that other variables also strongly influence BCF prediction. Thus, QSAR models that
556 rely solely on logP or logD in our opinion are limited in their application.

557 It is important to consider that descriptors not used in this work may also have
558 a potential for BCF modelling. For example, the major mechanism of transport across
559 epithelia tissue is passive diffusion and so it is also possible that diffusion coefficients
560 could potentially be an important descriptor for consideration among others, however
561 these descriptors are difficult to acquire and therefore reduce the practicability of a
562 model based on these.

563 **Conclusions**

564 The work presented herein has shown that *in silico* modelling approaches are
565 a powerful approach to predict bioconcentration of environmental contaminants,
566 enabling rapid prioritisation of compounds during ERA. The approach could be used
567 to better understand bioaccumulation, and the molecular descriptors that drive it;
568 moving the science beyond simple hydrophobicity models that poorly account for the
569 complexity of pharmaceuticals. Cross-species prediction of accumulation warrants
570 further investigation as the results indicate both case representation and biological
571 variability might limit prediction of accumulation between different taxonomic groups.
572 Nevertheless, the use of machine learning has been increasing within the field and is
573 necessary to improve our understanding of biological processes that affect
574 environmental health. The interpretation of descriptors here is critical as it
575 demonstrates that, in addition to rapid prediction of bioconcentration factors, *in silico*
576 models are useful for mechanistic understanding which in turn can be used to direct
577 further work. This is particularly true for pharmaceutical uptake in biota, where the

578 mechanisms that govern uptake, elimination and accumulation processes are still not
579 fully understood. Excellent potential exists for rapid screening using machine learning
580 technology in future ERA, without the need for costly and ethically challenging animal
581 experiments. Finally, the OECD QSAR validation guidelines for machine learners are
582 inexplicit and we suggest these guidelines should be expanded with more focus on
583 this type of modelling approach. This will begin to address the applicability and
584 usefulness of these models for regulatory schemes such as REACH where PBT
585 assessments are required for several thousand chemicals.

586 **Acknowledgements**

587 This work was conducted under funding from the Biotechnology and Biological
588 Sciences Research Council (BBSRC) CASE industrial scholarship scheme
589 (Reference BB/K501177/1), iNVERTOX project (Reference BB/P005187/1) and
590 AstraZeneca Global SHE research programme. AstraZeneca is a biopharmaceutical
591 company specialising in the discovery, development, manufacturing and marketing of
592 prescription medicines, including some products reported here. SFO is an employee
593 of AstraZeneca and a partner of the Innovative Medicines Initiative Joint Undertaking
594 under iPiE grant agreement n° 115735, resources of which are composed of financial
595 contribution from the European Union's Seventh Framework Programme (FP7/2007-
596 2013) and EFPIA companies' in-kind contribution. Funding bodies played no role in
597 the design of the study or decision to publish. The authors thank Jason Snape
598 (AstraZeneca) for critical review and declare no financial conflict of interest.

599 **References**

- 600 1. Rovida, C. and T. Hartung, *Re-evaluation of animal numbers and costs for in vivo tests to*
601 *accomplish REACH legislation requirements for chemicals-a report by the Transatlantic Think*
602 *Tank for Toxicology (t4)*. ALTEX-Alternatives to animal experimentation, 2009. **26**(3): p. 187-
603 208.
- 604 2. Commission, E., *Regulation (EC) No 1907/2006 of the European Parliament and of the*
605 *Council of 18 December 2006 concerning the Registration, Evaluation, Authorisation and*

- 606 *Restriction of Chemicals (REACH), establishing a European Chemicals Agency, amending*
607 *Directive 1999/45/EC and repealing Council Regulation (EEC) No 793/93 and Commission*
608 *Regulation (EC) No 1488/94 as well as Council Directive 76/769/EEC and Commission*
609 *Directives 91/155/EEC, 93/67/EEC, 93/105/EC and 2000/21/EC. 2006: Official Journal of the*
610 *European Union. p. 1 - 849.*
- 611 3. Agency, E.M., *Guideline on the environmental risk assessment of medicinal products for*
612 *human use. 2006, European Medicines Agency*
 - 613 4. Agency, E.M., *Reflection paper on the authorisation of veterinary medicinal products*
614 *containing (potential) persistent, bioaccumulative and toxic (PBT) or very persistent and very*
615 *bioaccumulative (vPvB) substances. 2016.*
 - 616 5. Gissi, A., et al., *Integration of QSAR models for bioconcentration suitable for REACH. Science*
617 *of The Total Environment, 2013. 456–457: p. 325-332.*
 - 618 6. de Wolf, W., et al., *Animal use replacement, reduction, and refinement: Development of an*
619 *integrated testing strategy for bioconcentration of chemicals in fish. Integrated*
620 *Environmental Assessment and Management, 2007. 3(1): p. 3-17.*
 - 621 7. Lockwood, S. and N. Saïdi, *Background document for public consultation on pharmaceuticals*
622 *in the environment. 2017.*
 - 623 8. Fatemi, M.H., M. Jalali-Heravi, and E. Konuze, *Prediction of bioconcentration factor using*
624 *genetic algorithm and artificial neural network. Analytica Chimica Acta, 2003. 486(1): p. 101-*
625 *108.*
 - 626 9. Lombardo, A., et al., *Assessment and validation of the CAESAR predictive model for*
627 *bioconcentration factor (BCF) in fish. Chemistry Central Journal, 2010. 4(1): p. 1.*
 - 628 10. Zhao, C., et al., *A new hybrid system of QSAR models for predicting bioconcentration factors*
629 *(BCF). Chemosphere, 2008. 73(11): p. 1701-1707.*
 - 630 11. Stempel, S., et al., *Using conditional inference trees and random forests to predict the*
631 *bioaccumulation potential of organic chemicals. Environmental Toxicology and Chemistry,*
632 *2013. 32(5): p. 1187-1195.*
 - 633 12. Institute, C.I.a.T., *Biodegradation and Bioaccumulation data of existing chemicals based on*
634 *the CSCL Japan. 1992, Japan: Chemical Industry Ecology-Toxicology & Information Center.*
 - 635 13. Ashauer, R., A. Boxall, and C. Brown, *Uptake and Elimination of Chlorpyrifos and*
636 *Pentachlorophenol into the Freshwater Amphipod Gammarus pulex. Archives of*
637 *Environmental Contamination and Toxicology, 2006. 51(4): p. 542-548.*
 - 638 14. Ashauer, R., et al., *Bioaccumulation kinetics of organic xenobiotic pollutants in the*
639 *freshwater invertebrate Gammarus pulex modeled with prediction intervals. Environmental*
640 *Toxicology and Chemistry, 2010. 29(7): p. 1625-1636.*
 - 641 15. Meredith-Williams, M., et al., *Uptake and depuration of pharmaceuticals in aquatic*
642 *invertebrates. Environmental Pollution, 2012. 165(Supplement C): p. 250-258.*
 - 643 16. Miller, T.H., et al., *Uptake, biotransformation and elimination of selected pharmaceuticals in*
644 *a freshwater invertebrate measured using liquid chromatography tandem mass*
645 *spectrometry. Chemosphere, 2017. 183(Supplement C): p. 389-400.*
 - 646 17. Miller, T.H., et al., *Assessing the reliability of uptake and elimination kinetics modelling*
647 *approaches for estimating bioconcentration factors in the freshwater invertebrate,*
648 *Gammarus pulex. Science of The Total Environment, 2016. 547(Supplement C): p. 396-404.*
 - 649 18. OECD, *Guidance Document on the Validation of (Q)SAR Models. 2007.*
 - 650 19. Palm, K., et al., *Polar Molecular Surface Properties Predict the Intestinal Absorption of Drugs*
651 *in Humans. Pharmaceutical Research, 1997. 14(5): p. 568-571.*
 - 652 20. Kah, M. and C.D. Brown, *LogD: Lipophilicity for ionisable compounds. Chemosphere, 2008.*
653 *72(10): p. 1401-1408.*
 - 654 21. Devillers, J., et al., *Fish Bioconcentration Modelling With LogP. Toxicology Methods, 1998.*
655 *8(1): p. 1-10.*

- 656 22. Dobson, P.D. and D.B. Kell, *Carrier-mediated cellular uptake of pharmaceutical drugs: an*
657 *exception or the rule?* Nature Reviews Drug Discovery, 2008. **7**: p. 205.
- 658 23. Nakamura, Y., et al., *The effects of pH on fluoxetine in Japanese medaka (Oryzias latipes):*
659 *Acute toxicity in fish larvae and bioaccumulation in juvenile fish.* Chemosphere, 2008. **70**(5):
660 p. 865-873.
- 661 24. Rendal, C., K.O. Kusk, and S. Trapp, *Optimal choice of pH for toxicity and bioaccumulation*
662 *studies of ionizing organic chemicals.* Environmental Toxicology and Chemistry, 2011. **30**(11):
663 p. 2395-2406.
- 664 25. Karlsson, M.V., et al., *Novel Approach for Characterizing pH-Dependent Uptake of Ionizable*
665 *Chemicals in Aquatic Organisms.* Environmental Science & Technology, 2017. **51**(12): p.
666 6965-6971.
- 667 26. Mackay, D. and A. Fraser, *Bioaccumulation of persistent organic chemicals: mechanisms and*
668 *models.* Environmental Pollution, 2000. **110**(3): p. 375-391.
- 669 27. Rubach, M.N., et al., *Toxicokinetic variation in 15 freshwater arthropod species exposed to*
670 *the insecticide chlorpyrifos.* Environmental Toxicology and Chemistry, 2010. **29**(10): p. 2225-
671 2234.
- 672 28. Dearden, J.C. and N.M. Shinnawei, *Improved prediction of fish bioconcentration factor of*
673 *Hydrophobic Chemicals.* SAR and QSAR in Environmental Research, 2004. **15**(5-6): p. 449-
674 455.
- 675 29. Sahu, V.K. and R.K. Singh, *Prediction of the Bioconcentration Factor of Organic Compounds in*
676 *Fish.* CLEAN – Soil, Air, Water, 2009. **37**(11): p. 850-857.
- 677 30. Aalizadeh, R., et al., *Quantitative Structure–Retention Relationship Models To Support*
678 *Nontarget High-Resolution Mass Spectrometric Screening of Emerging Contaminants in*
679 *Environmental Samples.* Journal of Chemical Information and Modeling, 2016. **56**(7): p.
680 1384-1398.
- 681 31. Weaver, S. and M.P. Gleeson, *The importance of the domain of applicability in QSAR*
682 *modeling.* Journal of Molecular Graphics and Modelling, 2008. **26**(8): p. 1315-1326.
- 683 32. Netzeva, T.I., et al., *Current status of methods for defining the applicability domain of*
684 *(quantitative) structure-activity relationships.* ATLA, 2005. **33**: p. 155-173.
- 685 33. Dimitrov, S., et al., *Base-line model for identifying the bioaccumulation potential of*
686 *chemicals.* SAR and QSAR in Environmental Research, 2005. **16**(6): p. 531-554.
- 687 34. Arnot, J.A. and F.A.P.C. Gobas, *A review of bioconcentration factor (BCF) and*
688 *bioaccumulation factor (BAF) assessments for organic chemicals in aquatic organisms.*
689 Environmental Reviews, 2006. **14**(4): p. 257-297.
- 690 35. Gramatica, P., *Principles of QSAR models validation: internal and external.* QSAR &
691 Combinatorial Science, 2007. **26**(5): p. 694-701.
- 692 36. Verhaar, H.J.M., J. de Jongh, and J.L.M. Hermens, *Modeling the Bioconcentration of Organic*
693 *Compounds by Fish: A Novel Approach.* Environmental Science & Technology, 1999. **33**(22):
694 p. 4069-4072.
- 695 37. Hendriks, A.J., et al., *The power of size. 1. Rate constants and equilibrium ratios for*
696 *accumulation of organic substances related to octanol-water partition ratio and species*
697 *weight.* Environmental Toxicology and Chemistry, 2001. **20**(7): p. 1399-1420.
- 698 38. Buchwalter, D.B., J.J. Jenkins, and L.R. Curtis, *Respiratory strategy is a major determinant of*
699 *[3H]water and [14C]chlorpyrifos uptake in aquatic insects.* Canadian Journal of Fisheries and
700 Aquatic Sciences, 2002. **59**(8): p. 1315-1322.
- 701 39. Rubach, M.N., D.J. Baird, and P.J. Van den Brink, *A new method for ranking mode-specific*
702 *sensitivity of freshwater arthropods to insecticides and its relationship to biological traits.*
703 Environmental Toxicology and Chemistry, 2010. **29**(2): p. 476-487.
- 704 40. Karara, A.H. and W.L. Hayton, *A pharmacokinetic analysis of the effect of temperature on the*
705 *accumulation of di-2-ethylhexyl phthalate (DEHP) in sheepshead minnow.* Aquatic
706 Toxicology, 1989. **15**(1): p. 27-36.

- 707 41. Sun, X., et al., *Classification of bioaccumulative and non-bioaccumulative chemicals using*
708 *statistical learning approaches*. *Molecular Diversity*, 2008. **12**(3): p. 157.
- 709 42. LeBlanc, G.A., *Trophic-Level Differences in the Bioconcentration of Chemicals: Implications in*
710 *Assessing Environmental Biomagnification*. *Environmental Science & Technology*, 1995.
711 **29**(1): p. 154-160.
- 712 43. Grabicova, K., et al., *Tissue-specific bioconcentration of antidepressants in fish exposed to*
713 *effluent from a municipal sewage treatment plant*. *Science of The Total Environment*, 2014.
714 **488-489**: p. 46-50.
- 715 44. Lajeunesse, A., et al., *Distribution of antidepressants and their metabolites in brook trout*
716 *exposed to municipal wastewaters before and after ozone treatment – Evidence of biological*
717 *effects*. *Chemosphere*, 2011. **83**(4): p. 564-571.
- 718 45. Tanoue, R., et al., *Simultaneous determination of polar pharmaceuticals and personal care*
719 *products in biological organs and tissues*. *Journal of Chromatography A*, 2014. **1355**: p. 193-
720 205.
- 721 46. Togunde, O.P., et al., *Determination of Pharmaceutical Residues in Fish Bile by Solid-Phase*
722 *Microextraction Couple with Liquid Chromatography-Tandem Mass Spectrometry*
723 *(LC/MS/MS)*. *Environmental Science & Technology*, 2012. **46**(10): p. 5302-5309.
- 724 47. Xie, Z., et al., *Occurrence, bioaccumulation, and trophic magnification of pharmaceutically*
725 *active compounds in Taihu Lake, China*. *Chemosphere*, 2015. **138**: p. 140-147.
- 726 48. de Solla, S.R., et al., *Bioaccumulation of pharmaceuticals and personal care products in the*
727 *unionid mussel *Lasmigona costata* in a river receiving wastewater effluent*. *Chemosphere*,
728 2016. **146**: p. 486-496.
- 729 49. Du, B., et al., *Pharmaceutical bioaccumulation by periphyton and snails in an effluent-*
730 *dependent stream during an extreme drought*. *Chemosphere*, 2015. **119**: p. 927-934.
- 731 50. Verbruggen, B., et al., *ECODrug: a database connecting drugs and conservation of their*
732 *targets across species*. *Nucleic Acids Research*, 2018. **46**(D1): p. D930-D936.
- 733 51. Gramatica, P. and E. Papa, *QSAR Modeling of Bioconcentration Factor by theoretical*
734 *molecular descriptors*. *QSAR & Combinatorial Science*, 2003. **22**(3): p. 374-385.
- 735 52. Pajouhesh, H. and G.R. Lenz, *Medicinal Chemical Properties of Successful Central Nervous*
736 *System Drugs*. *NeuroRX*, 2005. **2**(4): p. 541-553.
- 737 53. Kelder, J., et al., *Polar Molecular Surface as a Dominating Determinant for Oral Absorption*
738 *and Brain Penetration of Drugs*. *Pharmaceutical Research*, 1999. **16**(10): p. 1514-1519.
- 739 54. Skyner, R., et al., *A review of methods for the calculation of solution free energies and the*
740 *modelling of systems in solution*. *Physical Chemistry Chemical Physics*, 2015. **17**(9): p. 6174-
741 6191.
- 742 55. Tice, C.M., *Selecting the right compounds for screening: does Lipinski's Rule of 5 for*
743 *pharmaceuticals apply to agrochemicals?* *Pest Management Science*, 2001. **57**(1): p. 3-16.
- 744 56. Anliker, R., P. Moser, and D. Poppinger, *Advances in Environmental Hazard and Risk*
745 *Assessment 1987 Bioaccumulation of dyestuffs and organic pigments in fish. Relationships to*
746 *hydrophobicity and steric factors*. *Chemosphere*, 1988. **17**(8): p. 1631-1644.
- 747 57. Dimitrov, S.D., et al., *Predicting bioconcentration factors of highly hydrophobic chemicals.*
748 *Effects of molecular size, in Pure and Applied Chemistry*. 2002. p. 1823.
- 749 58. Randić, M., et al., *A rational selection of graph-theoretical indices in the QSAR*. *International*
750 *Journal of Quantum Chemistry*, 1988. **34**(S15): p. 267-285.
- 751 59. Anliker, R., P. Moser, and D. Poppinger, *Bioaccumulation of dyestuffs and organic pigments*
752 *in fish. Relationships to hydrophobicity and steric factors*. *Chemosphere*, 1988. **17**(8): p.
753 1631-1644.
- 754 60. Opperhulzen, A., et al., *Relationship between bioconcentration in fish and steric factors of*
755 *hydrophobic chemicals*. *Chemosphere*, 1985. **14**(11): p. 1871-1896.

- 756 61. Gramatica, P., M. Corradi, and V. Consonni, *Modelling and prediction of soil sorption*
757 *coefficients of non-ionic organic pesticides by molecular descriptors*. Chemosphere, 2000.
758 **41**(5): p. 763-777.
- 759 62. Wang, Y., et al., *Estimation of bioconcentration factors using molecular electro-topological*
760 *state and flexibility*. SAR and QSAR in Environmental Research, 2008. **19**(3-4): p. 375-395.
- 761 63. Miller, T.H., et al., *The First Attempt at Non-Linear in Silico Prediction of Sampling Rates for*
762 *Polar Organic Chemical Integrative Samplers (POCIS)*. Environmental Science & Technology,
763 2016. **50**(15): p. 7973-7981.
- 764 64. Morin, N.A.O., et al., *Kinetic accumulation processes and models for 43 micropollutants in*
765 *“pharmaceutical” POCIS*. Science of The Total Environment, 2018. **615**(Supplement C): p.
766 197-207.

Table 1: Comparison of model performance for the prediction of BCF in *Cyprinus carpio*. MAE is the mean absolute error and NA indicates the metric was not applicable.

Model		RMSE			R ²			MAE		
		Training	Verification	Test	Training	Verification	Test	Training	Verification	Test
Trajan	<i>Linear</i>	0.785	1.052	0.832	0.532	0.390	0.521	0.619	0.835	0.608
	<i>GRNN</i>	0.830	0.893	0.873	0.673	0.400	0.569	0.664	0.893	0.718
	<i>RBF</i>	0.723	0.689	0.584	0.651	0.635	0.725	0.565	1.600	0.450
	<i>3-MLP</i>	0.689	0.538	0.337	0.675	0.770	0.659	0.548	1.608	0.553
	<i>4-MLP</i>	0.403	0.524	0.644	0.887	0.819	0.702	0.313	0.380	0.530
Model		Training	Cross-Validation	Test	Training	Cross-Validation	Test	Training	Cross-Validation	Test
R	<i>OLM</i>	0.719	0.771	1.203	0.621	0.570	0.234	0.560	NA	0.778
	<i>PLS</i>	0.722	0.769	1.164	0.618	0.571	0.254	0.564	NA	0.765
	<i>RR</i>	0.725	0.766	1.083	0.614	0.576	0.304	0.568	NA	0.753
	<i>EN</i>	0.729	0.760	1.054	0.612	0.582	0.314	0.577	NA	0.754
	<i>QRL</i>	0.733	0.757	1.112	0.607	0.585	0.284	0.562	NA	0.770
	<i>KNN</i>	0.517	0.683	0.902	0.807	0.665	0.468	0.404	NA	0.648
	<i>ELM</i>	0.673	0.756	1.014	0.668	0.593	0.346	0.529	NA	0.768
	<i>ANN-1HL</i>	0.596	0.751	0.877	0.739	0.597	0.505	0.462	NA	0.620
	<i>ANN-a1HL</i>	0.395	0.672	0.859	0.888	0.678	0.518	0.319	NA	0.612
	<i>ANN-2HL</i>	0.232	0.834	1.022	0.962	0.560	0.370	0.174	NA	0.680
	<i>ANN-3HL</i>	0.454	0.795	0.880	0.860	0.582	0.520	0.345	NA	0.624
	<i>MARS</i>	0.539	0.730	1.014	0.787	0.632	0.390	0.425	NA	0.696
	<i>B-MARS</i>	0.500	0.681	0.899	0.819	0.673	0.479	0.395	NA	0.633
	<i>SVM-R</i>	0.383	0.644	0.841	0.893	0.704	0.537	0.261	NA	0.590
	<i>SVM-P</i>	0.699	0.747	1.029	0.643	0.594	0.340	0.539	NA	0.729
	<i>RF-CART</i>	0.292	0.675	0.771	0.956	0.688	0.633	0.231	NA	0.589
	<i>RF-CIT</i>	0.605	0.739	0.821	0.762	0.630	0.586	0.485	NA	0.652
	<i>BT</i>	0.249	0.660	0.789	0.957	0.687	0.593	0.187	NA	0.587
<i>CR</i>	0.353	0.678	0.973	0.910	0.673	0.431	0.282	NA	0.628	

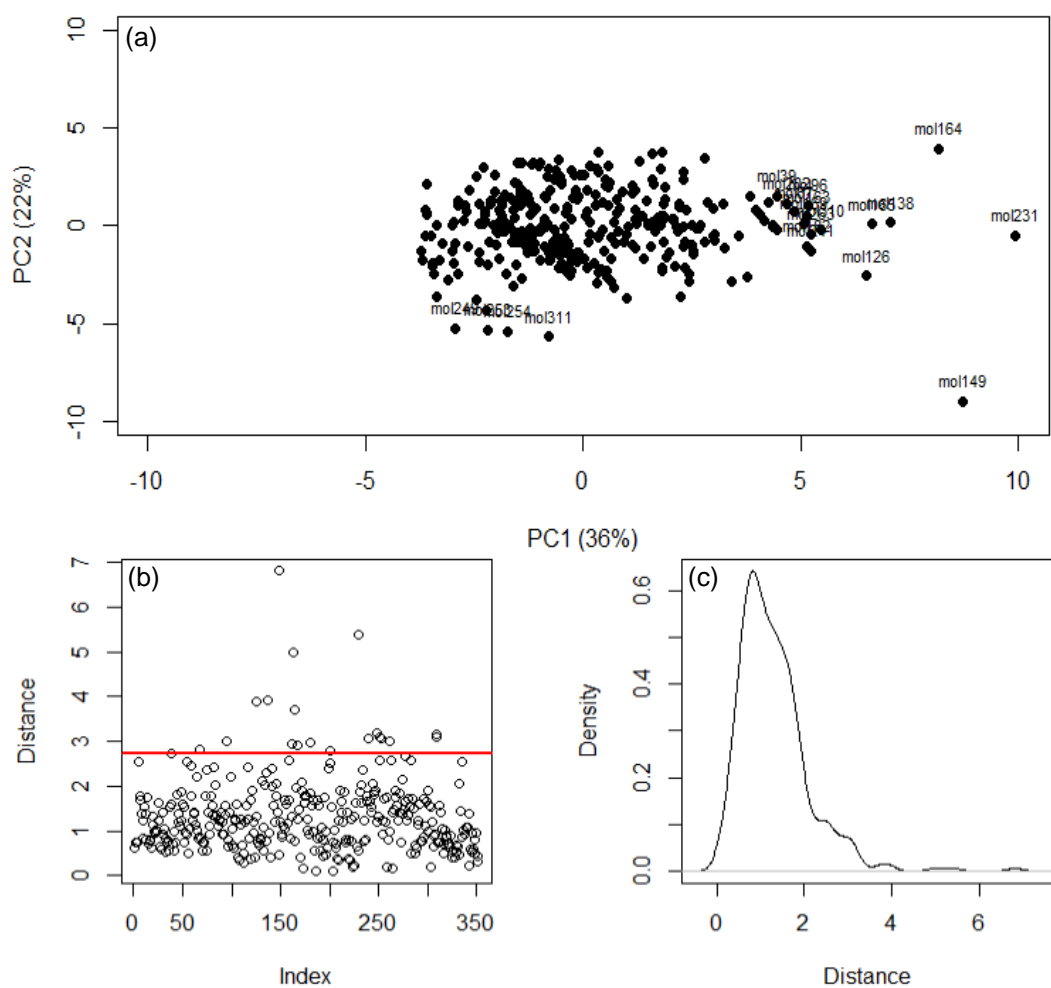


Figure 1: (a) Principal component analysis used for visualisation of the case similarity based on the 14 modelled descriptors (i.e. applicability domain). (b) Distances between cases in the PCA space with a threshold applied (0.975 quantile of χ^2 distribution) designated by the red line (c) the distribution of cases based on distance in the PCA space.

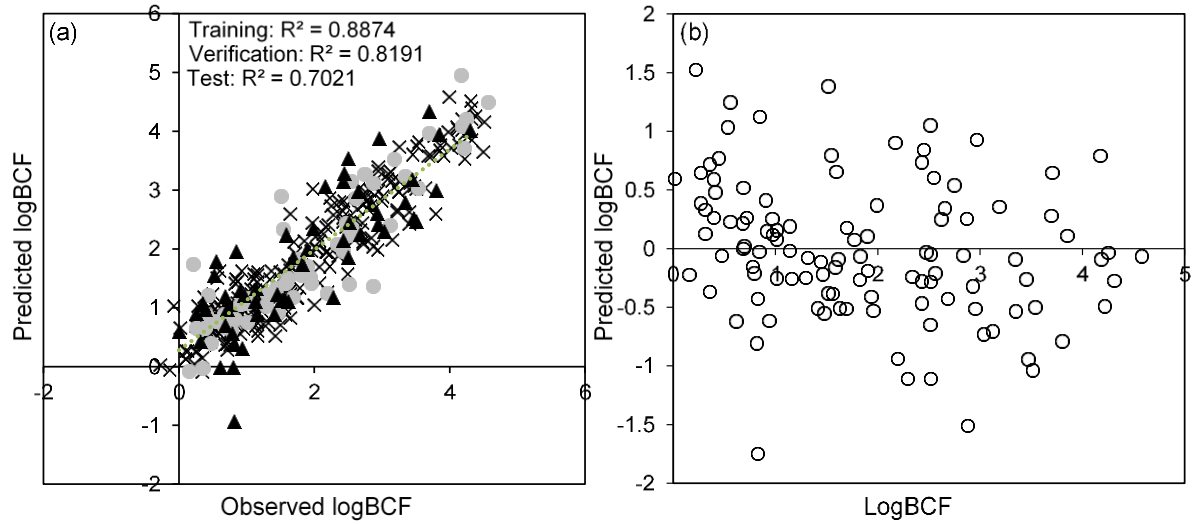


Figure 2: (a) linear regression of the predicted logBCF values versus the observed logBCF values in fish using the 4-MLP developed in approach 1, training data (crosses, $n = 242$), verification data (circles, $n = 55$) and test data (triangles, $n = 55$). (b) Raw residuals of the predicted logBCF data in fish for the verification and test data only.

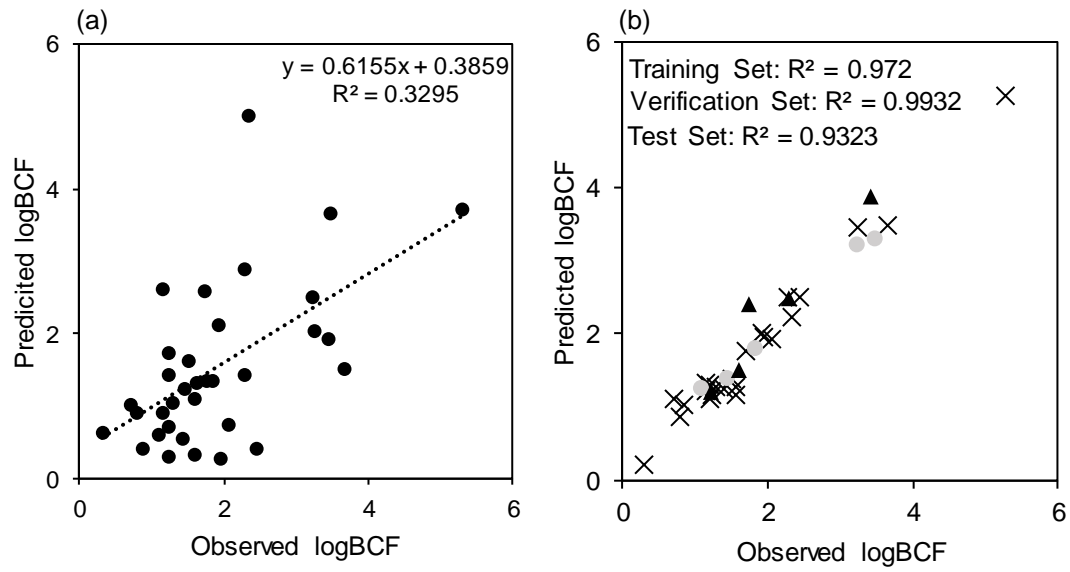


Figure 3: (a) Comparison of the predicted logBCF data versus the observed logBCF in invertebrates using the fish-based 4-layer MLP. (b) Regression of a separately developed and optimised model trained with the invertebrate BCF data (*Gammarus pulex*), training set (crosses, $n = 24$), verification set (circles, $n = 5$) and test set (triangles, $n = 5$)

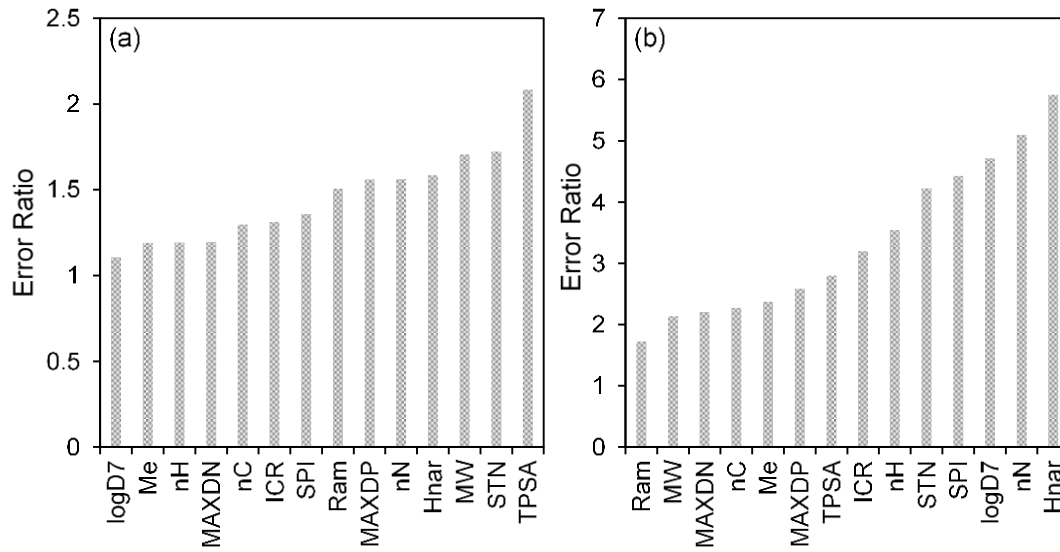


Figure 4: Descriptors sensitivity analysis performed by removing a descriptor from the model and assessing the affected performance. Increased error ratios indicate more important descriptors. (a) descriptor sensitivity for the fish-based model and (b) for the invertebrate-based model.