**Eawag-Soil in enviPath: A new resource for exploring regulatory pesticide soil biodegradation pathways and half-life data**

Diogo A. R. S. Latino[a],*, Jörg Wicker[b], Martin Gütlein[b], Emanuel Schmid[c], Stefan Kramer[b], Kathrin Fenner[a,d]

[a] Eawag, Swiss Federal Institute of Aquatic Science and Technology, 8600 Dübendorf, Switzerland

[b] Institute of Computer Science, Johannes Gutenberg University Mainz, 55128 Mainz, Germany

[c] Scientific IT Services, ETH Zürich, 8092 Zürich, Switzerland

[d] Department of Chemistry, University of Zürich, 8057 Zürich, Switzerland

*Corresponding author: diogo.latino@eawag.ch, Tel.: +41 58 765 5737

**Abstract**

Developing models for the prediction of microbial biotransformation pathways and half-lives of trace organic contaminants in different environments requires as training data easily accessible and sufficiently large collections of respective biotransformation data that are annotated with metadata on study conditions. Here, we present the *Eawag-Soil* package, a public database that has been developed to contain all freely accessible regulatory data on pesticide degradation in laboratory soil simulation studies for pesticides registered in the EU (282 degradation pathways, 1535 reactions, 1619 compounds and 4716 biotransformation half-life values with corresponding metadata on study conditions). We provide a thorough description of this novel data resource, and discuss important features of the pesticide soil degradation data that are relevant for model development. Most notably, the variability of half-life values for individual compounds is large and only about one order of magnitude lower than the entire range of median half-life values spanned by all compounds, demonstrating the need to consider study conditions in the development of more accurate models for biotransformation prediction. We further show how the data can be used to find missing rules relevant for predicting soil biotransformation pathways. From this analysis, eight examples of reaction types were presented that should trigger the formulation of new biotransformation rules, e.g., Ar-OH methylation, or the extension of existing rules e.g., hydroxylation in aliphatic rings. The data were also used to exemplarily explore the dependence of half-lives of different amide pesticides on chemical class and experimental parameters. This analysis highlighted the value of considering initial transformation reactions for the development of meaningful quantitative-structure biotransformation relationships (QSBR), which is a novel opportunity offered by the simultaneous encoding of transformation reactions and corresponding half-lives in *Eawag-Soil*. Overall, *Eawag-Soil* provides an unprecedentedly rich collection of manually extracted and curated biotransformation data, which should be useful in a great variety of applications.

**Introduction**

When chemicals are released into the environment during or at the end of their product life cycle, their persistence in the environment is highly undesirable. Biotransformation by microbial communities in technical and environmental systems such as sewage treatment plants, aquatic sediments, and soils is a very efficient mechanism to reduce their environmental persistence, but might also lead to the formation of potentially hazardous transformation products [1-3]. Since the experimental assessment of chemical persistence and transformation product formation on a compound-by-compound basis is highly laborious and costly, so-called *in silico* or *non-testing* approaches that rely on computer-based algorithms to predict biotransformation have gained in importance for the evaluation of new and existing chemicals [4]. It has been suggested that such approaches would also be of use in the implementation of the "benign by design" concept where the environmental risk of a chemical is considered early in the development process or even before synthesis [5].

Quantitative structure-biodegradation relationships (QSBRs) predict chemical persistence, i.e., half-lives or readiness of biodegradation, based on chemical structure. They range from chemical class-specific to more broadly applicable models, and from simple regression models to models developed with machine learning methods [5-7]. Chemical class-specific models [8-10] typically yield reasonably accurate predictions of actual degradation half-lives, but are of limited use for risk assessment purposes due to their restricted applicability domain. In contrast, more widely applicable models are typically trained on a number of databases containing collections of ready biodegradability data from standardized tests carried out according to the OECD test guidelines for a wide variety of chemicals [11]. They usually show reasonable predictive power with approximately 80% correct binary classification as to whether a chemical is readily biodegradable or not (e.g., [12-14]). However, the accuracy of these more widely applicable models for quantitatively predicting biotransformation rates or half-lives under specific environmental conditions, which is what would actually be needed for risk assessment purposes, remains rather low [15-17].

Pathway prediction systems (e.g., PathPred [18], Catalogic [19], BNICE [20], and Eawag-PPS (former UM-PPS [21])) typically rely on dictionaries of biotransformation rules that recognize compound functional groups and transform them into product substructures. These biotransformation rules were designed to reflect known microbial transformation pathways of chemical contaminants. They are mostly based on the respective data collected in the Eawag Biodegradation/Biocatalysis Database (Eawag-BBD), formerly known as the University of Minnesota Biodegradation/Biocatalysis Database (UM-BBD) [22], which is considered the most extensive collection of manually curated biotransformation pathways of chemical contaminants [23] and, more recently, is available as Eawag-BBD from two online platforms (Eawag-BBD/PPS [24] and enviPath [25, 26]). Rule-based systems have been shown to predict transformation products observed in the environment fairly comprehensively (i.e., to display high sensitivity), but to notoriously predict many irrelevant products that are not likely to occur under specific environmental conditions (i.e., to display low selectivity) [27]. While application of machine

75    learning methods to improve relative reasoning between rules has increased the selectivity for the

76    training database (i.e., Eawag-BBD), selectivity on a set of pesticide soil degradation data used for

77    external validation remained low [28]. This poses a problem if the models were to be used in a chemical

78    risk assessment context where resources for assessing the risk associated with transformation products

79    are limited.

80    We argue that the low accuracy of QSBRs and the low selectivity of pathway prediction have at least

81    two common causes. First, almost all approaches to biotransformation prediction are based on

82    chemical structure only and have so far mostly ignored the fact that half-lives for the same compound

83    can vary strongly within the same type of environmental compartment [15]. This observed variance

84    stems from the fact that slightly different environmental conditions shape different microbial

85    communities that differ in their taxonomic composition and hence their pool of enzymes that catalyze

86    biotransformation reactions of chemical contaminants. Thus, different enzyme-catalyzed reactions

87    might occur at vastly different rates across different microbial communities [29, 30]. This suggests that

88    biotransformation prediction could be greatly improved by not only considering chemical structure,

89    but by also factoring in specific environmental conditions. Second, most of the data in Eawag-BBD

90    stems from studies with pure cultures of microorganisms or laboratory cultures with elongated

91    adaptation periods. Thus, while the organisms and degrading enzymes are typically well-characterized

92    in these studies and hence reported in the database, the current data in Eawag-BBD cannot be used for

93    understanding the influence of environmental factors on biotransformation pathways nor is the

94    relevance of the reported pathways under actual environmental conditions known. In pure and

95    enrichment culture systems, besides being known to be impacted by culturing artifacts [31], the

96    compound of concern serves as sole growth substrate. The latter is most likely also true for the ready

97    biodegradability tests that are run at high concentrations of the test chemicals as dominant carbon

98    source [32]. Under actual environmental conditions, contaminant trace concentrations are likely

99    transformed co-metabolically by mixed microbial communities alongside varying amounts of other,

100    natural organic material. The determinants of such co-metabolic transformations are typically not of

101    thermodynamic nature as in growth-related metabolism, but rather the available pool of catalytic

102    enzymes of the microbial community as shaped by the prevailing environmental conditions [23]. Finally,

103    it is worth noting that QSBRs and systems for the prediction of biotransformation pathways have so

104    far mostly been developed independently. However, given the fact that observed biotransformation

105    half-lives and transformation product spectra (i.e., the observed biotransformation pathways) both

106    depend on the rates of individual enzyme-catalyzed biotransformation reactions, treating these two

107    types of information separately may lead to a loss of information content.

108    In summary, we hypothesize that development of more accurate QSBRs and pathway prediction

109    models is impeded by a lack of biotransformation data (i.e., half-lives and pathway information) from

110    environmentally relevant mixed microbial communities and associated metadata on environmental

111    and/or experimental conditions. The latter are needed to account for their influence on the observed

112   biotransformation outcomes. Recently, we have introduced enviPath as a new database and pathways

113   prediction system that is suited to approach these information gaps [25]. enviPath offers a database

114   environment that, first, facilitates the annotation of biotransformation half-life and pathway

115   information, and, second, allows for supplementing the half-life and pathway information with

116   metadata, e.g., environmental and/or experimental conditions, through so-called scenarios. One fairly

117   consistent and large resource of chemical biotransformation data is data submitted for regulatory risk

118   assessments. These substance-specific dossiers typically contain information on biotransformation

119   half-lives and pathways from so-called simulation studies conducted for different relevant

120   environmental compartments (i.e., agricultural soil, aquatic sediments, activated sludge). Such data is

121   currently mostly available for pesticides [33], but upon implementation of REACH should increasingly

122   also become available for industrial chemicals [34]. However, these data are currently not readily

123   available in electronic format [35], and, if so (e.g., PPDB [36]), do not contain pathway information, lack

124   annotation with metadata on study conditions, or, to the best of our knowledge, are not publically

125   available (e.g., MetaPath [35, 37]).

126   Therefore, the objective of the work presented here was to electronically encode all freely accessible

127   regulatory data on pesticide degradation in laboratory soil simulation studies and to make these data

128   publically available for the development of improved QSBR models. Here, we present a thorough

129   description of this novel data resource, discuss characteristics of pesticide soil degradation data that

130   are relevant to model development, and give two examples of explorative analyses that should support

131   the further development of QSBRs and pathway prediction models.

132

133   **Materials and Methods**

134   *Pesticide soil degradation data*

135   We extracted pesticide soil degradation information from pesticide registration dossiers made

136   publically available through the European Food Safety Authority (EFSA) [33]. Specifically, only results

137   from laboratory studies conducted under aerobic conditions as reported in "Annex B.8: Fate and

138   behavior, B8.1: Route and rate of degradation in soil" in the respective dossiers were considered.

139   Initially, assessment reports, draft assessment reports (DARs) and additional reports available between

140   6/2015 and 6/2016 for 375 active substances were screened. Of these, dossiers for 93 active substances

141   were not considered further because the pesticides agents were not actual chemicals (e.g., bacteria) or

142   complex mixtures (e.g., clover oil), or because no degradation scheme was available or no aerobic

143   degradation studies in soil had been submitted at all. For the remaining 282 pesticides and agriculture-

144   related compounds, degradation information and accompanying metadata on study conditions were

145   encoded as separate data package, *Eawag-Soil*, in enviPath [38].

146   In *the Eawag-Soil* package, pathway information is stored in a biotransformation reaction scheme in

147   the entity *pathway* (see example in Figure 1). Compounds and reactions participating in a given

148   pathway are stored separately in the entities *compound* and *reaction* (see Figure 1). Metadata on the

149    experimental conditions (e.g., soil texture, soil moisture, pH, etc.) are stored in the entity *scenario* (see

150    Figure 1). A detailed list and explanation of all experimental conditions considered in the *Eawag-Soil*

151    package as well as the conventions used to store the data as standardized as possible are given in the

152    Supporting Information (SI) (Section S1 *Eawag-Soil* metadata and conventions) When available, one

153    or several biotransformation half-lives (in the form of dissipation half-lives, DT50) are additionally

154    associated with a given compound in the pathway and a specific scenario. *Pathways* depict all

155    reactions and compounds observed in aerobic soil experiments under any experimental condition.

156    Since not all transformation products in a pathway scheme are always observed, compounds in the

157    pathway are associated with a given scenario only when they have been experimentally observed

158    under the specific experimental conditions (see example in Figure 1). The associated scenarios are

159    listed on the compound page.
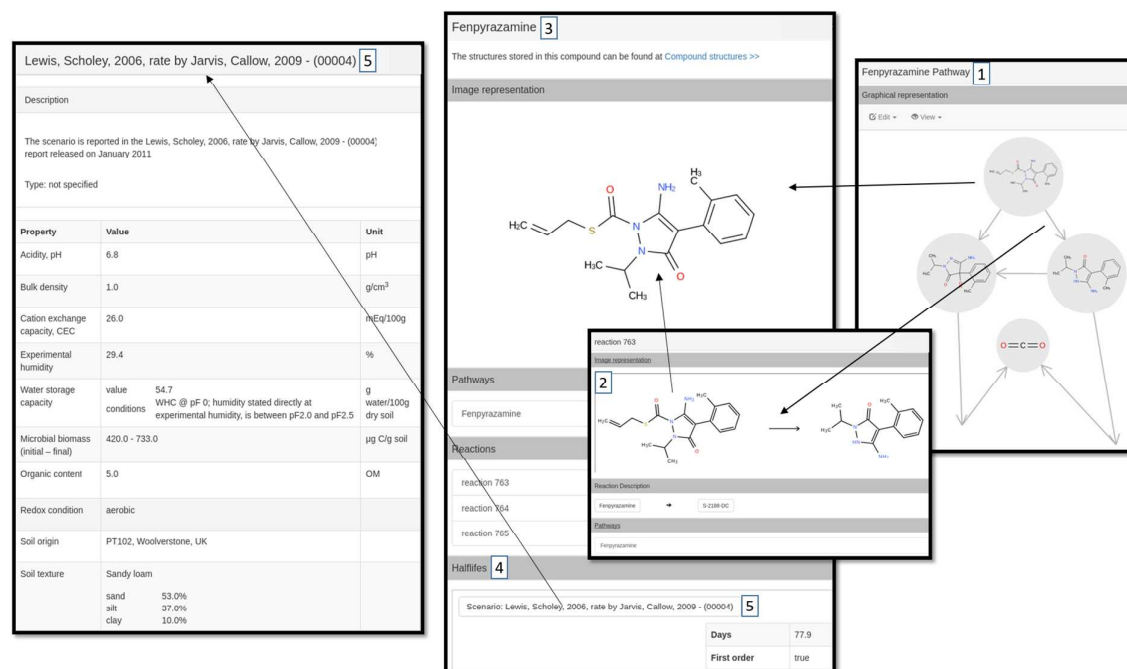
160



161

162    Figure 1: Scheme of assembled screenshots showing the most important elements of the *Eawag-Soil*

163    package. 1: *Pathway* page, 2: *Reaction* page, 3: *Compound* page, 4: List of half-lives determined for

164    the compounds and the associated scenario names, 5: *Scenario* page, containing the metadata on study

165    conditions (i.e., experimental parameters).

166

167    ***Chemical space analysis***

168    The chemical space covered by a set of compounds is defined by the multidimensional property space

169    of the compounds and is used to define the applicability domain of a model. We compared the

170    chemical spaces covered by *Eawag-BBD* and *Eawag-Soil*, and also compared them to the chemical

171    space covered by a third set of 1024 pharmaceutical compounds prevalently used in Switzerland and

172    the EU as extracted from Singer et al. [39]. This latter set of compounds was used to explore the

173    hypothesis that the addition of the *Eawag-Soil* package extends the chemical space of enviPath

174    towards more polar, multifunctional compounds such as pharmaceuticals.

175    The comparison was performed using a qualitative approach, i.e., the visualization of the top three

176    principal components of the compounds in the three datasets, and a quantitative approach using a one-

177    class support vector machine to identify objects that lie outside the chemical space. The top three

178    principal components were calculated using the DataWarrior software [40]

179    (http://www.openmolecules.org/datawarrior/) with the compounds represented using structural

180    fragment fingerprints (i.e., binary structural features (ECFP4) [41]) calculated with the CDK software [42].

181    One-class support vector machines (SVM) [43] is a machine learning technique that was used to

182    determine whether a compound belongs to the feature distribution space of an existing dataset or

183    rather has to be considered as an outlier or a novel compound. One-class SVM models were trained on

184    the *Eawag-BBD* dataset and the combined *Eawag-BBD* and *Eawag-Soil* datasets using the LIBSVM

185    implementation [44]) with the compounds represented using structural fragment fingerprints as explained

186    above. The ν-Parameter, which limits the number of predicted outliers in the training dataset, was set

187    to a value of 2%.

188

189    *Procedure for missing rule analysis (explorative analysis I)*

190    The Eawag-PPS system, which is hosted and further maintained in our research group, currently uses a

191    set of 249 biotransformation rules (btrules) that recognize specific functional groups in a molecule and

192    transform them according to the generalized biotransformation reaction encoded in the rule [22]. These

193    same rules also form the basis of the successor system enviPath [25]. When adding a new set of

194    biotransformation data such as the *Eawag-Soil* package to enviPath, a first pre-requisite to use its

195    information for improving pathway prediction models is to test the ability of the current rule set to

196    cover the reactions in the new database. As a first explorative analysis, we therefore conducted a

197    missing rule analysis.

198    Missing rule analysis was carried out in two steps: *i)* submission of the reactants of all chemical

199    reactions in *Eawag-Soil* to the complete set of btrules contained in Eawag-PPS, and *ii)* comparison of

200    the predicted reactions, i.e., reactant-product pairs, with the experimentally observed reactions. More

201    specifically, the reactants of all reactions in *Eawag-Soil* were submitted to the Eawag-PPS system and

202    three generations of transformation products were predicted through three times iterative application

203    of the prediction cycle. Then all the predicted first-generation reactant-product pairs were compared

204    with the experimentally observed reactant-product pairs and three different outcomes were noted: (*i*)

205    an experimentally observed reaction is matched by a predicted reaction indicating that there is a rule

206    that is correctly triggered by that reactant and that the system is therefore able to predict the

207    biotransformation observed in the soil degradation studies; (*ii*) a predicted reaction does not match any

208    of the experimentally observed reactions indicating that the system either predicts products that are not

209    actually relevant for a given reactant, or that the product, although plausible, escaped analytical

210    identification in the soil degradation studies; and (*iii*) an experimentally observed reaction is not

211    matched by any of the predicted reactions pointing towards a missing rule for that specific kind of

212    biotransformation reaction. While (*i*) gives the current sensitivity of the system towards pesticide

213    active ingredients, and (*ii*) is indicative of its current selectivity (prior to the addition of new rules),

214    reactions in (*iii*) were further explored to identify possible missing rules.

215    In a next step, for the set of reactions in (*iii*), the 2$^{nd}$ and 3$^{rd}$ generation products associated with their

216    respective substrates were explored to see whether any of them matched with the experimentally

217    observed product of the reaction. If so, this reaction is predicted in Eawag-PPS through a series of

218    multiple reactions where the intermediates might actually be readily transformed further and therefore

219    were not necessarily analytically observed and identified. These reactions were assigned as multi-step

220    reactions and added to the pool of reactions in (*i*).

221    To find a first set of missing rules, the remaining reactions in (*iii*) were, first, sorted according to mass

222    differences between the reactant and the product and, second, for a given mass difference, further

223    manually sorted into types of reactions based on our perceived similarity of the reaction center. This

224    approach is potentially limited because it will not group together reactions of the same type if

225    functional groups of different size are cleaved off. Therefore, in a next step, some groups of reactions

226    were joined together to make the reaction type more general and to more easily implement it as a new

227    rule later. For example, O-demethylations and O-deethylations were grouped together as O-

228    dealkylations, and, similarly, N-demethylations and N-deethylations were grouped together as N-

229    dealkylations. While fully manual, the approach was found appropriate for the identification of the

230    most populated reaction types. This approach also provides well-curated information that can be used

231    for validation of semi- and fully automated chemoinformatics methods for reaction classification,

232    which we plan to explore in a follow-up study.

233

234    *Univariate and multivariate analysis of half-life data (explorative analysis II)*

235    We performed an exploratory analysis of relationships between the experimental parameters and the

236    DT50 values for a group of six sulfonamide herbicides. We calculated Spearman rank correlation

237    coefficients and their related significance of being different from zero; significance was tested using a

238    two-tailed t-test on an approximation of the Student's t distribution equated by *abs(r((n-2)/(1-r$^2$))$^{0.5}$)*

239    where *r* is the Spearman rank correlation coefficient and *n* the number of half-lives per compounds.

240    Due to the multidimensional problem, the experimental parameters were also used to build multiple

241    linear regression models. Selection of the most relevant descriptors was performed with the

242    Correlation-based Feature Subset Selection (CFS) algorithm [45] implemented in Weka 3.8.1 [46]. The

243    algorithm takes into account the usefulness of the individual parameters for predicting the DT50

244    together with the level of intercorrelation among them. The experiments were carried out using the

245    AttributeSelectedClassifier routine of Weka with the CfsSubsetEval option for evaluator and BestFirst

246    or LinearForwardSelection options for search. The final set of parameters selected to build the model

247 will be in principle also the most relevant to explain the transformation of compounds across a

248 structural class of compounds and transformation reaction.

249

250 **Results and Discussion**

251 *Relevant characteristics of Eawag-Soil data*

252 In its current form, the *Eawag-Soil* package contains 282 degradation pathways, 1535 reactions

253 (excluding reactions leading to $CO_2$ through an unknown sequence of reactions) and 1619 compounds

254 (282 parent pesticides and agriculture-related compounds and 1337 biotransformation products). Of

255 these 1619 compounds, 777 (282 parent compounds and 495 biotransformation products) have at least

256 one associated half-life value. Since multiple half-lives may be available for individual compounds,

257 the *Eawag-Soil* package altogether contains 4716 biotransformation half-life values with

258 corresponding scenarios. These numbers will increase over time as the package is being further

259 developed. The size of the *Eawag-Soil* package in terms of numbers of pathways, reactions and

260 compounds lies in a similar range as the current size of the *Eawag-BBD* package (i.e., 219 pathways,

261 1503 reactions, 1396 compounds). Introducing it thus not only doubles the amount of

262 biotransformation pathway information to learn QSBRs and pathway prediction models from, but also

263 extends the chemical space covered from mostly legacy chemicals (i.e., persistent organic pollutants

264 and a few pesticides with long and extensive usage history) to modern, polar, and structurally more

265 complex pesticide active ingredients (see section on Chemical Space Analysis for a detailed

266 discussion).

267 A descriptive statistical analysis of the entire data set was performed. In the pre-processing of the data

268 set, values that seemed to be physically implausible based on the frequency distribution of the

269 parameters and on our knowledge about the different soil properties and their ranges were removed.

270 For example, values for the three *soil texture* parameters, i.e., % *sand*, *silt* and *clay*, were removed if

271 the sum of the three parameters was higher than 100%, for *humidity* parameter values higher than

272 100% were removed (5 values removed), or for soil *organic content* parameter values higher than 10 g

273 OC/100g soil were removed (12 values removed). In general, only a few values per parameter were

274 removed due to this analysis, corresponding to less than 1% of the values for most part of the

275 parameters.

276 For the experimental parameter *organic content* the soil organic content reported as soil organic matter

277 (OM) was transformed to organic carbon (OC) using the relationship OC = OM/1.724 [47]. For each

278 parameter, the number of missing values was determined and the distribution of values was

279 characterized by several statistical measures. A summary of results for DT50 and all experimental

280 parameters is given in Table 1 and additional statistical measures and histogram plots of all parameters

281 can be found in the SI (Table S1 and Figures S1 and S2).

282 The number of DT50s per compound in *Eawag-Soil* varies strongly. From the 777 compounds with

283 associated half-lives, 113 have more than ten DT50s (each one for a set of unique experimental

8

284  conditions) and more than half of the compounds, i.e., 419, have more than five DT50s (see SI Figures

285  S3 for a plot of the frequency distribution of half-lives per compounds). Figure 2 gives the maximum,

286  minimum and median DT50 values for all compounds in *Eawag-Soil* with at least 10 DT50 values per

287  compound (N=113). A corresponding graph for all 777 compounds with one or more associated half-

288  lives is given in the SI (Figure S4). The median half-lives across all of these compounds cover about

289  three orders of magnitude, suggesting that these data should provide a valuable resource to develop

290  QSBR models for soil half-life prediction. Also, the dataset separates out into 139 out of 777

291  compounds, i.e., 18%, having a median DT50 above 120 days, which is the persistence criterion for

292  pesticides and industrial chemicals in soil [48, 49]. Using appropriate sampling procedures the dataset can

293  be fairly balanced for the potential purpose of developing a persistence classification model from the

294  data. The presence of half-lives values close to zero is due to the presence of some rapidly degradable

295  or volatile compounds in the data set, e.g., metam-sodium or dazomet. On the opposite side, half-lives

296  above 1000 days indicate the presence of stable compounds in the data set, e.g., flutriafol and butralin,

297  where the DT50 values are extrapolated well beyond the study duration for most cases. For these cases

298  of extreme behavior, the DT50 values should be considered merely approximate values of the

299  behavior of the respective pesticides in soil.

300  Another important aspect that can be learned from the DT50 data shown in Figure 2 is the large

301  variability of DT50 values for individual compounds observed across different experimental

302  conditions. Considering the maximum and minimum half-lives, 244 out of 777 compounds, i.e., 31%,

303  show a variability in the DT50 values of two orders of magnitude or more. This shows that the

304  consideration of the experimental parameters is crucial for the development of QSBR models with

305  improved accuracy for the prediction of soil half-lives. The extensive collection of metadata made

306  available in *Eawag-Soil* should provide a useful resource for this purpose.

307 Table 1. Summary statistic of *DT50* and experimental parameters *% sand*, *% silt*, *% clay*, *pH*, *temperature*, *water storage capacity*, *% humidity*, *organic content*

308 (*OC*), *cation exchange capacity* (*CEC*), *bulk density*, *biomass start*, *biomass end* and *spike concentration*.

309

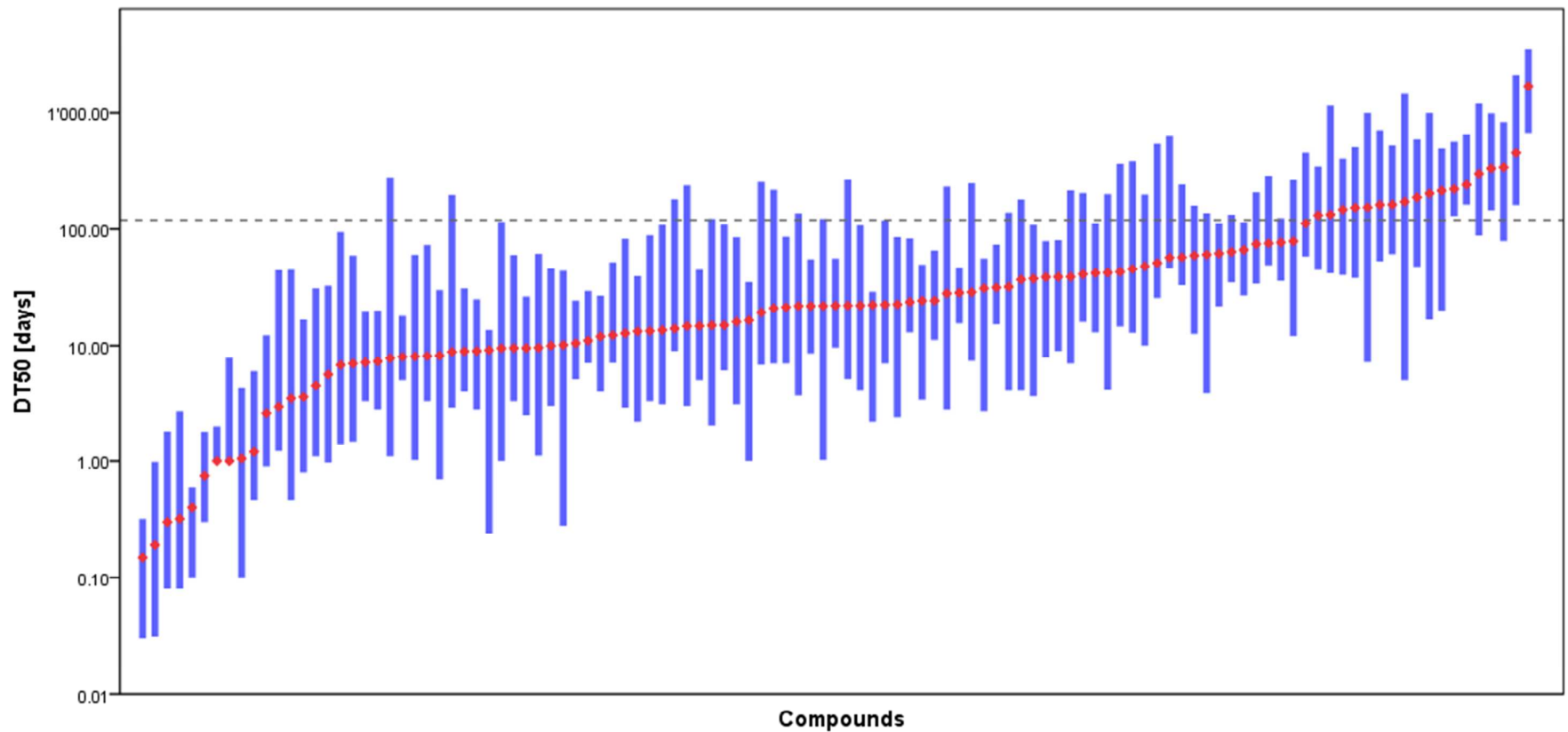| Parameter | Dimension | N | | Mean | Median | Std. Deviation | Minimum | Maximum |
|---|---|---|---|---|---|---|---|---|
| | | Valid | Missing | | | | | |
| DT50 | days | 4716 | 0 | 88.8 | 24.0 | 206 | 0.003 | 3690 |
| % Sand | – | 4108 | 608 | 51.9 | 55.0 | 24.6 | 0.00 | 99.0 |
| % Silt | – | 4114 | 602 | 32.0 | 28.0 | 18.6 | 0.00 | 88.6 |
| % Clay | – | 4132 | 584 | 16.0 | 12.6 | 11.0 | 0.10 | 94.3 |
| Acidity, pH | – | 4641 | 75 | 6.62 | 6.70 | 0.888 | 3.60 | 8.80 |
| Temperature | °C | 4659 | 57 | 20.1 | 20.0 | 3.50 | 1.00 | 49.0 |
| Water Storage Capacity | g water / 100 g dry soil | 3998 | 718 | 40.2 | 38.9 | 18.1 | 1.54 | 137 |
| % Humidity | – | 4350 | 366 | 55.4 | 45.0 | 21.6 | 5.00 | 100 |
| OC | g OC/100 g soil | 4540 | 176 | 1.81 | 1.62 | 1.08 | 0.02 | 10.0 |
| CEC | mEq/100 g soil | 3785 | 931 | 14.0 | 12.2 | 7.92 | 1.20 | 60.0 |
| Bulk Density | g/cm$^3$ | 1741 | 2975 | 1.34 | 1.40 | 0.282 | 0.00 | 2.66 |
| Biomass Start | µg C/g | 3147 | 1569 | 408 | 312 | 360 | 1.00 | 2445 |
| Biomass End | µg C/g | 2651 | 2065 | 345 | 257 | 339 | 0.05 | 2452 |
| Spike Concentration | mg/kg dry soil | 3917 | 799 | 1.28 | 0.400 | 2.55 | 0.003 | 25.0 |

310

311

312

10

Figure 2: Median DT50 values (red diamonds) and DT50 distributions (minimum to maximum) for 113 compounds with more than 10 associated DT50 values in *Eawag-Soil* (data used to build the Figure is available in Table S2 in SI). The dashed line indicates a persistence criterion in soil of 120 days.

11

316    Regarding the collection of metadata, missing values may pose a problem in data analysis. The

317    analysis of missing values in Table 1 shows that for 9 of the 13 numeric experimental parameters the

318    number of missing values is reasonably low and varies between 1.2% for *% humidity* and 17% for

319    *spike concentration*. The *soil texture*, a categorical parameter not shown in Table 1, is another

320    parameter with a small number of missing values, i.e., only 134 out of 4716, or 3% (see Figures S1 in

321    the SI for a frequency plot of the distribution of the 12 soil textural classes). The *water storage*

322    *capacity* only shows 8% of missing values, but will need some further pre-processing to be used for

323    modeling purposes. First, the *water storage capacity* is a property of the soil that could per se have an

324    influence on observed degradation. However, it has been measured under slightly varying water

325    tensions and therefore needs to be harmonized to one set of conditions. Second, the *water storage*

326    *capacity* does not directly describe the experimental moisture content of the soil, which also

327    potentially influences degradation. The latter could be obtained from multiplying the *water storage*

328    *capacity* with the *% humidity*.

329    Considerably more problematic in terms of missing values are the *cation exchange capacity*, *bulk*

330    *density*, *biomass start* and *biomass end* parameters, with 20%, 63%, 34% and 44% of missing values,

331    respectively. Bulk density shows the highest number of missing values, with 2975 scenarios out of

332    4716 not containing any information on it. Preliminary experiments showed that the missing bulk

333    densities could be imputed using the known relation between *bulk density*, soil texture and *organic*

334    *content*. A k-nearest neighbor model trained using the available *bulk density* values and *% clay* and

335    *organic content* as descriptors yielded results with a mean absolute error of 0.06 $g/cm^3$ in 10-fold cross

336    validation. Similarly, *cation exchange capacity* could be imputed from *% clay*, *organic content* and

337    *pH* with a mean absolute error of 1.12 mEq/100 g soil in 10-fold cross validation. While imputation of

338    some parameters based on known relationships between soil properties thus seems feasible and useful

339    for further model development, other missing values such as *biomass start* and *biomass end* will be

340    more difficult to address in this way because of the unknown relationship between these parameters

341    and other experimental parameters.

342    Overall, *Eawag-Soil* provides an unprecedentedly rich collection of half-lives and experimental

343    parameters manually extracted and curated, which should be useful in a great variety of applications,

344    some of which will be demonstrated in the following subsections.

345

346    *Chemical space analysis*

347    The projection of the top three principal components deduced from the structural fingerprints of the

348    compounds in the *Eawag-BBD* and *Eawag-Soil* packages is shown in Figure 3A, and with the

349    inclusion of the pharmaceuticals in Figure 3B. While *Eawag-Soil* overlaps with *Eawag-BBD* in some

350    regions of the space, it clearly extends it to a point where it also allows for a better coverage of other

351    relevant classes of compounds, e.g., the set of pharmaceuticals included for illustrative purposes

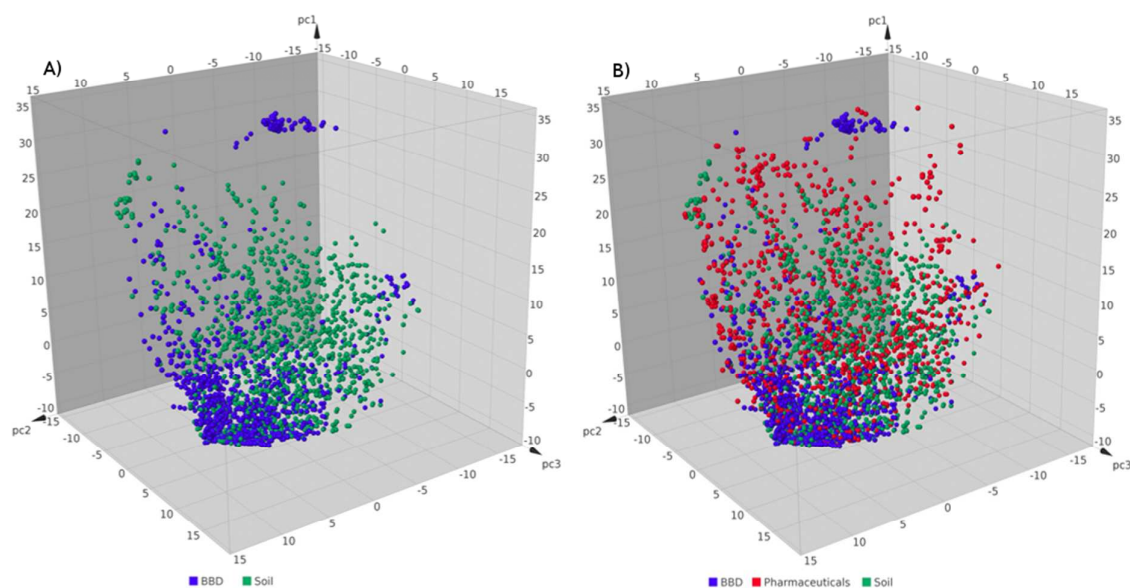352    (Figure 3B).

353



354

355 Figure 3. Projection of the top three principal components of A) *Eawag-BBD* and *Eawag-Soil*, and B)
356 *Eawag-BBD*, *Eawag-Soil* and pharmaceuticals set.

357

358 To quantitatively confirm the results shown in Figure 3, an analysis was also carried out based on one-
359 class SVM for outlier detection. The v-Parameter, which limits the number of outliers in the training
360 dataset, was initially set to a default value of 2%. To confirm that this value was reasonable, the 29
361 compounds identified as outliers in the *Eawag-BBD* package were visually checked and confirmed as
362 reasonable outliers. The outlier analysis indeed showed that the combined set of *Eawag-BBD* and
363 *Eawag-Soil* compounds covers a wider area of the relevant chemical space, compared to using only
364 the *Eawag-BBD* compounds. First, the outliers for the set of parent pesticides from *Eawag-Soil*
365 showed a reduction from 134 to 6 (Table 2). Similarly, the number of outliers for all compounds in
366 *Eawag-BBD* and *Eawag-Soil* combined showed a reduction from 590 to 61 (Table 2). These results
367 confirm that compounds such as pesticides would have been badly covered by the *Eawag-BBD* dataset
368 alone and, at the same time, can be considered as an internal validation of the outlier detection method.
369 More importantly, the number of outliers for the set of pharmaceuticals is reduced by 70% (152
370 instead of 515) when adding the *Eawag-Soil* dataset to the *Eawag-BBD* dataset. Three examples of
371 pharmaceuticals that were outliers in *Eawag-BBD* but are considered inside the combined chemical
372 space of *Eawag-BBD* and *Eawag-Soil* are shown, together with their corresponding three nearest
373 neighbors, in Table S3 of the SI. With the inclusion of the *Eawag-Soil* package, 85% of the
374 pharmaceuticals are now covered by the chemical space of the compounds in enviPath. Based on the
375 combination of the *Eawag-BBD* and *Eawag-Soil* packages, it should therefore become possible to train
376 models with a significantly enlarged application domain, and hence strongly increased prediction
377 accuracy and reliability for structurally complex and polar compounds, which are of particular concern
378 for water quality [50].

13

379

380 Table 2. Number of outliers detected in the different datasets.

|  | Number of compounds detected as outliers | | | |
|---|---|---|---|---|
|  | *Eawag-BBD* (N = 1399) | *Eawag-BBD* and *Eawag-Soil* (N = 3018) | Parent Pesticides in *Eawag-Soil* (N =280) | Pharmaceuticals (N = 1006) |
| *Eawag-BBD*[a] | 29 | 590 | 134 | 515 |
| *Eawag-BBD & Eawag-Soil*[a] | 21 | 61 | 6 | 152 |

381 [a]Compound sets used to define the chemical space

382

383 *Missing rules analysis*

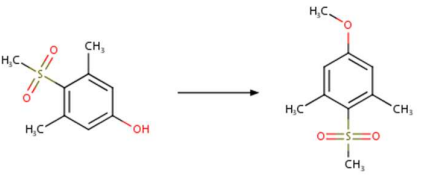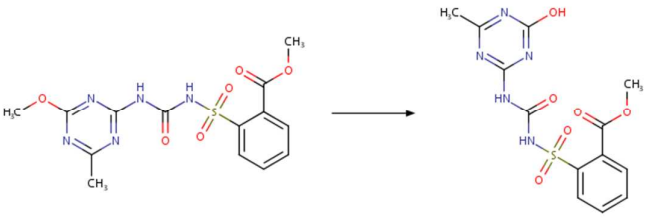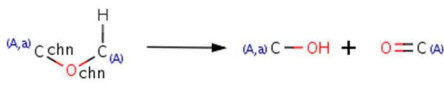384 From the set of 1535 experimental reactions, 711 reactions (i.e., 46%) are not predicted by the Eawag-

385 PPS system over a three-generation prediction cycle. Thus, the current sensitivity of Eawag-PPS to

386 predict transformation of the compounds contained in *Eawag-Soil* is only 54%. For the 711 reactions

387 not predicted, the respective type of transformation might be missing completely, or corresponding

388 rules exist but their specificity does not fully cover the substrate spectrum of compounds in *Eawag-*

389 *Soil*. Table 3 shows the eight reaction types that contain at least eight reactions that were not predicted

390 by Eawag-PPS. Together they cover 20% of the 711 reactions not predicted by Eawag-PPS. Another

391 253 reactions were classified into 48 more reaction types or combination of reaction types of smaller

392 size. 315 reactions could so far not be classified at all, amongst other reasons because of incompletely

393 documented biotransformation pathways with missing intermediates and improbable reported

394 structures of intermediates. These cases will need further attention, but might not be fully resolvable.

395     Table 3: Most populated reaction types not predicted by Eawag-PPS with example reactions.

| Reaction Type | Subclass of Reaction Type | N. # | Example reaction | Related *btrule* |
|---|---|---|---|---|
| Hydroxylation (52 reactions) | Hydroxylation in (hetero)aromatic ring | 33 |  | bt0013  |
| | Hydroxylation in aliphatic ring | 7 |  | |
| | Hydroxylation in alpha-position to allyl/aryl/carbonyl group | 6 |  | bt0242  |
| | Hydroxylation in (hetero)aromatic ring followed by keto enol tautomerism | 6 |  | |
| Scission of aryl-heteroaryl ether bond | | 19 |  | |

15

| | | | | |
|---|---|---|---|---|
| Amide hydrolysis (16 reactions) | (Cyclic) N-acylurea hydrolysis | 5 |  | |
| | (Cyclic) Sulfonyl urea hydrolysis | 7 |  | bt0067  |
| | Aliphatic amide hydrolysis | 4 |  | |
| Decarboxylation | | 13 |  | bt0051  |
| N-dealkylation | | 11 |  | bt0063  |

| | | | | |
|---|---|---|---|---|
| Ar-OH methylation | | 9 |  | |
| O-dealkylation | | 8 |  | bt0023  |
| Reduction of ketone to alcohol | | 8 |  | |

396

397   The most populated reaction types cover hydrolysis, oxidative N- and O-dealkylations, decarboxylations, reductions and also addition reactions. For most of these reaction types, i.e., hydroxylations, amide hydrolyses, decarboxylations, and oxidative N- and O-dealkylations, similar *btrules* already exist but are too specific to predict the observed reactions. For instance, they may be too specific in the definition of the neighborhood atoms of the reaction center. This is case for rules bt0011, bt0012, bt0013 and bt0014, which predict hydroxylations of monosubstituted benzene and pyridine rings in *o*-, *m*-, and *p*-position, but do not cover multiply substituted aromatic rings and N-heteroaromatic rings other than pyridine, which, however, are present as observed reactions in the *Eawag-Soil* dataset. Another case of existing rules being too specific are rules that have restrictions related to the presence of specific functional groups in the molecule. Here, hydroxylations in aliphatic rings and in alpha position to allyl, aryl and carbonyl groups are an interesting example since the existing rule bt0242 handles hydroxylation of secondary aliphatic carbon atoms in a ring, adjacent to a carbon that is sp2 hybridized, or bound to N or O, and should thus cover the observed reactions. However, rule bt0242 has been prevented from acting on these compounds due to "functional group restrictions" that state that the rule should not act on esters and amides, which, however, are present in the substrates of the respective, not predicted reactions. In these and similar cases, existing rules should be extended to cover the structural patterns observed in the *Eawag-Soil* reaction set, and functional group restrictions removed where they are in contradiction to evidence from *Eawag-Soil*. The latter point suggests itself even more since approaches have been developed to learn relative priorities between rules (e.g., between hydroxylation and amide or ester hydrolysis) from the data and implement them in terms of relative reasoning models, rather than hardcode them into the rules [25, 27, 28]. Of the most populated reactions in Table 3 only the scission of the aryl-heteroaryl ether bond, Ar-OH methylation and the reduction of the ketone group are not covered by any existing rule. The scission of the ether bond is an interesting case because the ether linkage is a common feature that, due to the high energy of the ether bond, confers stability and consequently renders these compounds typically rather resistant to microbial degradation. However, abiotic hydrolysis data [36, 51] demonstrate that at least some of these aryl-heteroaryl ether bonds can be hydrolyzed at acidic pH, resulting in the formation of a phenol and a 2-pyridone derivative. This must be due to the presence of the N-atom in the aromatic heterocycle, which renders it more e-deficient and draws electron density from the carbon atom that is attached to the ether bridge, making it more vulnerable for nucleophilic attack. These considerations suggest that the observed scission of aryl-heteroaryl ether bonds is due to hydrolysis, but based on the data available in *Eawag-Soil* it remains difficult to judge whether it is a purely abiotic or enzyme-catalyzed hydrolysis.

430   The Ar-OH methylation is another interesting case since it is an addition reaction. Methylation of phenol to anisol is a common transformation in many biosynthetic pathways (e.g., lignol, hormone, and flavonoid biosynthesis) and may be catalyzed by enzymes from the methyltransferase class (EC numbers of class 2.1.1.-). The substrates of six out of the nine Ar-OH methylations in *Eawag-Soil* are

434 halogenated substituted phenols. Therefore, likely candidate enzymes that could perform this type of

435 biotransformation are EC 2.1.1.136, a halogenated phenol O-methyltransferase that acts on mono-, bi-

436 and trichlorophenols, and EC 2.1.1.25, a phenol O-methyltransferase that acts on a wide variety of

437 simple alkyl-, methoxy- and halophenols. Interestingly, EC 2.1.1.136 has so far only been found to

438 occur in fungi, which suggests that the *Eawag-Soil* data set might also highlight some fungal

439 transformations that are only scarcely covered in *Eawag-BBD* and hence not only extend the coverage

440 of enviPath towards new types of compounds but also other types of catalyzing enzymes and

441 microbial organisms. There are also other addition reactions such as formylations, acetylations or

442 conjugations with more complex groups that are increasingly observed in microbial communities [52-54],

443 and, at least for the case of N-formylation (1 reaction) and N-acetylation (3 reactions), have also been

444 observed in *Eawag-Soil*. Addition reactions have typically not been implemented in pathway

445 prediction systems so far because their focus was on catabolic reactions. As a consequence, as is the

446 case for Ar-OH methylation, a *btrule* for the reverse reaction, i.e., the oxidative O-dealkylation, often

447 exists. Therefore, if addition reactions were to be implemented in the future, care has to be taken that

448 the system does not predict any products from a given substrate that are identical with any of its

449 precursor compounds in the pathway to avoid "dead cycles" in the prediction. The same is true if

450 rules for reductions such as the reduction of ketones to secondary alcohols (Table 3) were to be

451 implemented. Since these are more likely to proceed under anaerobic conditions, assigning them a low

452 aerobic likelihood within the Eawag-PPS system [55] could further restrict the application of such rules.

453

454 *Exploring half-life variability*

455 As discussed in a previous section, median half-lives across all pesticides cover about three orders of

456 magnitude, yet variability in half-lives for individual pesticides also spans about two orders of

457 magnitude. Thus, improved QSBR models need to account for both inter- and intra-compound

458 variability in half-lives. This should become more achievable if mechanistic understanding about the

459 fate of the pesticides in soil and about the ongoing transformation processes is included in the model

460 development as much as possible. Here, we therefore explore whether the consideration of initial

461 transformation reactions can support such an endeavor. The underlying hypothesis is that if a set of

462 structurally similar compounds, i.e., from the same pesticide class, showed the same initial

463 transformation reaction, this transformation is likely catalyzed by the same type of enzyme. The two

464 ensuing hypotheses, for which we did some initial exploration here, are *(i)* that inter-compound half-

465 life variability is considerably smaller within compounds that belong to the same pesticide class and

466 undergo the same transformation than across all compounds, and *(ii)* that intra-compound half-life

467 variability shows similar, characteristic dependencies on environmental conditions across compounds

468 that belong to the same pesticide class and undergo the same transformation. Hypothesis *(ii)* is

469 restricted to compounds belonging to the same pesticide class because we assume that direct effects of

19

470  the environmental conditions on their bioavailability and abiotic stability are consistent within a class

471  of pesticides but not necessarily across classes.

472  We explored these two hypotheses for all amide pesticides in the *Eawag-Soil* package. Amides were

473  selected because they constitute a large class in the *Eawag-Soil* package (i.e., 40 out of 282 parent

474  compounds are classified as amides according to ref [56]). The class also contains several compounds

475  that have particularly large numbers of half-lives and scenarios associated with them (i.e., 10 half-

476  lives/amide on average compared to 6 half-lives/compound on average across all of *Eawag-Soil*). The

477  amides were then further grouped into consistent sub- and subsubclasses, first according to ref [56] and

478  later through further manual curation. Finally, every amide was annotated manually with its initial

479  transformation reaction(s) according to the pathway maps in the *Eawag-Soil* package. This resulted in

480  three amide subsubclasses (sulfonamides, chloroacetanilides, anilides with N-substituted pyrazole

481  ring) that contained four or more structurally similar compounds undergoing the same type of initial

482  transformation reaction (Table 4). All other subsubclasses were either smaller or their members

483  underwent different initial transformation reactions. In the following, the three groups in Table 4 form

484  the basis for testing hypotheses (*i*) and (*ii*).

485

486  Table 4. Initial transformation reactions and half-lives (range and median) for sulfonamides,

487  chloroacetanilides, and anilides with N-substituted pyrazole ring.

| Pesticide class | Compound | Initial transformation reaction | Median DT50 [days] | DT50 range [days] | Number of DT50 values |
|---|---|---|---|---|---|
| Sulfonamides | Penoxsulam | O-demethylation | 24.5 | 15-137 | 7 |
| | Pyroxsulam | O-demethylation | 3.6 | 1-17 | 25 |
| | Florasulam | O-demethylation | 3.5 | 0-45 | 17 |
| | Metosulam | O-demethylation | 9.15 | 4-25 | 4 |
| | *Asulam* [a] | *Several, O-demethylation not possible* | *3.89* | *3-10* | *5* |
| | *Oryzalin* [a] | *Several, O-demethylation not possible* | *182* | *63-468* | *8* |
| Chloroacetanilides | Acetochlor | Substitution with GSH & hydrolysis; reductive dechlorination | 10.2 | 0.28-44 | 50 |
| | Dimethachlor | Substitution with GSH; reductive dechlorination; others | 7.15 | 3.31-19.8 | 12 |
| | Dimethenamide | Substitution with GSH & hydrolysis | 13 | 7.8-43.4 | 5 |

| | Metazachlor | Substitution with GSH; other | 13.6 | 3.1-109 | 49 |
|---|---|---|---|---|---|
| | Propisochlor | Substitution with GSH & hydrolysis | 11.4 | 8.4-40.2 | 5 |
| Anilides with N-substituted pyrazole ring | BAS 700 | N-dealkylation; amide hydrolysis; other | 326 | 72.7-810 | 6 |
| | Benzovindiflupyr | N-dealkylation; amide hydrolysis; hydroxylation | 550 | 349-1000 | 7 |
| | Bixafen | N-dealkylation; amide hydrolysis | >365 | n.a. | 4 [b] |
| | Isopyrazam | N-dealkylation; amide hydrolysis; hydroxylation | 231 | 29.8-976 | 9 |
| | Penflufen | Hydroxylation | 231 | 117-434 | 6 |
| | Penthiopyrad | Hydroxylation; N-dealkylation; amide hydrolysis | 146 | 60.5-413 | 6 |
| | Sedaxane | N-dealkylation; other | 74.2 | 57.6-138 | 8 |

488    [a] These two sulfonamides were included to demonstrate the effect of adding structurally similar
489    compounds that undergo different initial transformation reactions.

490    [b] For all four scenarios, half-lives were given as >365 d.

491

492    The median half-life data given in Table 4 and the distribution of individual half-lives of the three
493    groups as compared to the entirety of all amides shows that the distribution of half-lives for the
494    chloroacetanilides and the sulfonamides are overlapping, whereas the median half-life distribution for
495    the anilides is quite distinct. More importantly, however, the distributions for the three groups in Table
496    4 are more narrow than the half-life distribution for all amides (Figure 4), which is also demonstrated
497    by the coefficients of variation (CV) of the median half-lives, which are 0.97, 0.23, and 0.64 for the
498    sulfonamides, the chloroacetanilides and the anilides, respectively, as compared to 1.24 for all amides.
499    For demonstration purposes, we also considered half-life data for the remaining two sulfonamides in
500    the *Eawag-Soil* package (i.e., asulam and oryzalin), which, however, do not contain the
501    methoxypyrimidine moiety that is subject to O-demethylation in penoxsulam, pyroxsulam, florasulam
502    and metosulam and therefore undergo different initial transformation reactions. This resulted in a CV
503    across all six sulfonamides of 1.89. Altogether, these observations lend some support to hypothesis *(i)*
504    in that they demonstrate smaller inter-compound variability within groups of compounds undergoing
505    the same initial transformation reaction amongst the class of amide pesticides.
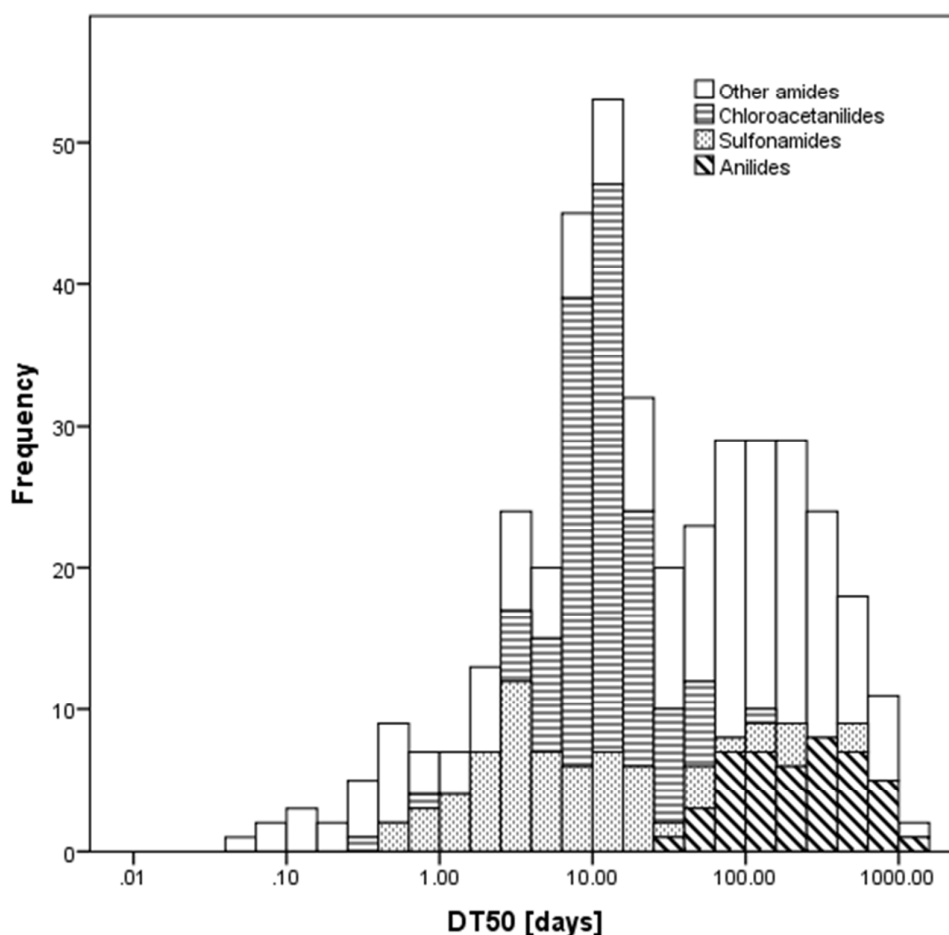
506

Figure 4. Half-life distribution (frequency plot) for chloroacetanilides, sulfonamides, anilides with N-substituted pyrazole ring, and remaining compounds classified as amides in the *Eawag-Soil* package.

To explore hypothesis (*ii*), first, Spearman rank correlation coefficients for the relationship between the DT50 values of the six sulfonamides and nine experimental parameters were calculated (Table 5). For some combinations of compounds and parameters the analysis was not possible because of missing values. A consistent negative relationship with temperature was observed (albeit only significant for one compound), which could be expected because of previous reports of Arrhenius-type dependence of pesticide soil degradation [57]. Only few additional univariate relationships were found that are significant at the 5% level. The case of OC seems most interesting because contradictory, yet significant dependencies are found for pyroxsulam (negative dependence, O-demethylation) and oryzalin (positive dependence, transformations other than O-demethylation). Additionally, a close to significant negative dependence on OC was also found for florasulam (O-demethylation). These results lend some support to the hypothesis that dependencies on environmental conditions are more similar if compounds undergo the same initial transformation reaction.

524 Table 5. Spearman rank correlation coefficients for relationships between DT50 values of six

525 sulfonamides and selected experimental parameters. Significance of rank correlation coefficients are

526 given in parenthesis and coefficients that are significant at the 5% significance level are highlighted in

527 bold.

| | Log(DT50) | | | | | |
|---|---|---|---|---|---|---|
| | Penoxsulam N=7 | Pyroxsulam N=25 | Florasulam N=17 | Metosulam N=4 | Asulam N=5 | Oryzalin N=8 |
| % Sand | 0.213 (0.685) | 0.095 (0.651) | 0.204 (0.661) | -0.738 (0.262) | 0.205 (0.741) | 0.627 (0.070) |
| % Silt | -0.213 (0.685) | 0.005 (0.981) | -0.204 (0.661) | 0.000 (1.00) | -0.205 (0.741) | **-0.814 (0.008)** |
| % Clay | -0.213 (0.685) | -0.172 (0.411) | 0.337 (0.460) | 0.800 (0.200) | -0.308 (0.614) | -0.627 (0.070) |
| pH | 0.273 (0.601) | **-0.503 (0.010)** | 0.163 (0.727) | -0.400 (0.600) | 0.718 (0.172) | -0.579 (0.102) |
| Temperature | **-0.845 (0.034)** | -0.283 (0.170) | -0.628 (0.131) | - | -0.707 (0.182) | -0.548 (0.127) |
| Organic Content | 0.395 (0.438) | **-0.523 (0.007)** | -0.738 (0.058) | -1.00 (-) | -0.205 (0.741) | **0.848 (0.004)** |
| CEC | 0.030 (0.955) | -0.312 (0.129) | 0.535 (0.216) | -0.80 (0.200) | 0.103 (0.869) | 0.034 (0.931) |
| Biomass Start | 0.516 (0.295) | **-0.611 (0.001)** | -0.553 (0.198) | -1.00 (-) | -0.205 (0.741) | - |
| Spike concentration | - | 0.248 (0.232) | - | - | - | -0.402 (0.283) |

528

529 Because univariate relationships are strongly confounded by the influence of all other experimental

530 parameters on the observed half-lives, multiple linear regression models were developed to further

531 explore the validity of hypothesis *(ii)* for the example of the sulfonamide herbicides. Considering the

532 fact that the DT50 ranges of all the sulfonamide herbicides are in the same range with exception of the

533 oryzalin, multiple linear regression were developed using as training set a combination of all DT50

534 values and corresponding scenarios for only those four sulfonamides undergoing O-demethylation

535 (experiment 1) or for all six sulfonamides (experiment 2). In Table 6, a summary of the resulting

536 models is given.

537

538 Table 6. Multiple linear regression models developed for DT50 values of sulfonamides.

| Compounds | N desc | Desc. Selected | Training | | | Test (10-fold cross validation) | | |
|---|---|---|---|---|---|---|---|---|
| | | | R | MAE | RMSE | R | MAE | RMSE |
| Penoxsulam Florasulam Metosulam Pyroxsulam (N=53) | 5 | pH; T; OC; CEC; Biomass | 0.813 | 0.252 | 0.315 | 0.729 | 0.312 | 0.377 |
| Penoxsulam Florasulam Metosulam Pyroxsulam Oryzalin Asulam (N=66) | 4 | T; CEC; Biomass; Spike Concentration | 0.743 | 0.371 | 0.476 | 0.627 | 0.426 | 0.563 |

539

540  After parameter selection, the final MLR model for experiment 1 yielded the following equation:

$$Log(DT50) = -0.166 * pH - 0.0467 * T - 0.166 * OC + 0.0249 * CEC - 0.0005$$
$$* Biomass\ Start + 2.75$$

541  The final MLR model had a mean absolute error of 0.312 (corresponding to roughly a factor of two) in
542  10-fold cross validation and showed an only minor decrease of R between training and cross-
543  validation, indicating that the data were not over fitted by the model. Also, at least two of the observed
544  dependencies are plausible based on our understanding of the fate of pesticides in soil: The
545  temperature-dependence again follows the logic of an Arrhenius relationship, and the negative
546  dependence on biomass follows the logic of a second-order rate constant that depends on biomass and
547  compound concentration. The fact that a model could be built that encompassed all four sulfonamides
548  in experiment 1 and that remained robust in rigorous 10-fold cross-validation, without the need to
549  include information on the structure of the compound or other molecular descriptors, clearly support
550  both hypotheses *(i)* and *(ii)*. Hypothesis *(ii)* is further supported by the fact that the regression model
551  from experiment 2 performed worse than for experiment 1, suggesting that adding data for structurally
552  similar compounds that, however, undergo a different type of transformation weakens the observed
553  dependences on experimental parameters. The finding that experiment 2 performs worse than
554  experiment 1 still holds true when the half-lives for each compound are z-normalized to account for
555  the different half-life range of oryzalin compared to the other sulfonamides (i.e., R values of 0.680 and
556  0.648 are obtained for training the model on only those sulfonamides containing the
557  methoxypyrimidine moiety and on all six sulfonamides, respectively). Overall, the most encouraging
558  outcome from this exploration of half-life variability is that for groups of structurally similar
559  compounds undergoing the same transformation models can be built that capture a relevant part of the
560  observed variability (i.e., >60%).

561

562  **Conclusions & outlook**

563  In this article we presented the *Eawag-Soil* package as a novel biotransformation data package made
564  available through the enviPath environment. *Eawag-Soil* contains a comprehensive collection of all
565  freely accessible regulatory data on pesticide degradation in laboratory soil simulation studies under
566  aerobic conditions for pesticides registered in the EU. This data resource has been developed in order
567  to respond to the need for more environmentally relevant training data sets to develop models for the
568  microbial biotransformation of polar, structurally complex trace organic contaminants such as
569  pesticides and pharmaceuticals. An analysis of the chemical spaces covered by the existing *Eawag-*
570  *BBD* dataset and *Eawag-Soil* confirmed a strongly improved coverage of these types of chemicals,
571  suggesting that through the combination of the *Eawag-BBD* and *Eawag-Soil* packages it should
572  become possible to train models with an increased prediction accuracy and reliability for structurally
573  complex and polar compounds.

574 We have further explored two lines of research that can greatly profit from the data in *Eawag-Soil*: *(i)*

575 the formulation of new rules or adaptation of existing rules to obtain a better coverage for the

576 prediction of soil biotransformation of structurally complex trace organic contaminants, and *(ii)* the

577 elucidation of the dependency of observed half-life variability on the study conditions as expressed by

578 the experimental parameters. Based on the analysis of missing rules, eight examples of reaction types

579 were presented that should trigger the formulation of new biotransformation rules, e.g., Ar-OH

580 methylation, or the extension of existing rules e.g., hydroxylation in aliphatic rings. The exploration of

581 the half-lives of different amide pesticides not only showed that different subsubclasses of structurally

582 similar amides have significantly different median half-lives, but also yielded some first evidence that

583 the consideration of initial transformation reactions within groups of structurally similar amides seems

584 to support a more accurate description of how half-lives depend on environmental conditions. Based

585 on these results, we argue that the consideration of the type of initial transformation reactions in the

586 development of QSBRs should greatly facilitate the consideration of the influence of experimental

587 parameters on half-lives in such models. Doing so is a novel opportunity offered by the simultaneous

588 encoding of transformation reactions and corresponding half-lives in *Eawag-Soil*. Ultimately, a

589 combined pathway prediction system could be developed where the reactivity pattern of the compound

590 (as encoded by an extended set of btrules) is used as one type of descriptor in combination with

591 molecular descriptors and experimental conditions to predict half-lives. To work towards this end, a

592 more complete analysis of the reactions not predicted by Eawag-PPS will be sought, and an automated

593 procedure for defining new rules based on a chemoinformatics approach for semi-automatic analysis

594 and assignment of reaction types will be implemented.

595 Overall, *Eawag-Soil* makes an unprecedented amount of manually curated soil biotransformation

596 information available to the public in an easily accessible manner. This should not only be of high

597 interest for researchers developing QSBR-type models and pathway prediction systems, but also for

598 regulators and the general public as an information resource.

599

609

610 **References**

611 1.    D. M. Cwiertny, S. A. Snyder, D. Schlenk and E. P. Kolodziej, *Environmental Science &*
612    *Technology* **2014,** *48*, 11737-11745.

613 2.    B. Escher and K. Fenner, *Environmental Science & Technology* **2011,** *45*, 3835 - 3847.

614 3.    A. B. A. Boxall, C. J. Sinclair, K. Fenner, D. Kolpin and S. J. Maud, *Environmental Science &*
615    *Technology* **2004,** *38*, 368A-375A.

616 4.    ECHA *Guidance on information requirements and chemical safety assessment, Chapter R.7b:*
617    *Endpoint specific guidance*; European Chemicals Agency, Helsinki, Finland, 2016.

618 5.    C. Rücker and K. Kümmerer, *Green Chemistry* **2012,** *14*, 875-887.

619 6.    L. Mamy, D. Patureau, E. Barriuso, C. Bedos, F. Bessac, X. Louchart, F. Martin-Laurent, C.
620    Miege and P. Benoit, *Critical Reviews in Environmental Science and Technology* **2015,** *45*, 1277-
621    1377.

622 7.    M. Pavan and A. P. Worth, *QSAR Comb. Sci.* **2008,** *27*, 32-40.

623 8.    S. Banerjee, P. H. Howard, A. M. Rosenberg, A. E. Dombrowski, H. Sikka and D. L. Tullis,
624    *Environmental Science & Technology* **1984,** *18*, 416-422.

625 9.    D. F. Paris and N. L. Wolfe, *Appl. Environ. Microbiol.* **1987,** *53*, 911-916.

626 10.    Y. Urushigawa, S. Masunaga and Y. Yonezawa, *Water Sci. Technol.* **1988,** *20*, 459-461.

627 11.    OECD *guidelines for testing of chemicals; 301, Ready Biodegradability* OECD: Paris, 1992.

628 12.    R. S. Boethling, P. H. Howard, W. Meylan, W. Stiteler, J. Beauman and N. Tirado,
629    *Environmental Science & Technology* **1994,** *28*, 459-465.

630 13.    S. Vorberg and I. V. Tetko, *Mol. Inf.* **2014,** *33*, 73-85.

631 14.    K. Mansouri, T. Ringsted, D. Ballabio, R. Todeschini and V. Consonni, *J. Chem Inf. Model.*
632    **2013,** *53*, 867-878.

633 15.    D. Aronson, R. S. Boethling, P. H. Howard and W. Stiteler, *Chemosphere* **2006,** *63*, 1953-
634    1960.

635 16.    K. Fenner, S. Canonica, B. I. Escher, L. Gasser, S. Spycher and H. C. Tulp, *Chimia* **2006,** *60*,
636    683-690.

637 17.    R. Kuhne, R. U. Ebert and G. Schuurmann, *QSAR Comb. Sci.* **2007,** *26*, 542-549.

638 18.    Y. Moriya, D. Shigemizu, M. Hattori, T. Tokimatsu, M. Kotera, S. Goto and M. Kanehisa,
639    *Nucleic Acids Res.* **2010,** *38*, W138-W143.

640 19.    S. Dimitrov, T. Pavlov, N. Dimitrova, D. Georgieva, D. Nedelcheva, A. Kesova, R. Vasilev
641    and O. Mekenyan, *SAR QSAR Environ. Res.* **2011,** *22*, 719-755.

642 20.    S. D. Finley, L. J. Broadbelt and V. Hatzimanikatis, *Biotechnology and Bioengineering* **2009,**
643    *104*, 1086-1097.

644 21.    L. B. M. Ellis, J. Gao, K. Fenner and L. P. Wackett, *Nucleic Acids Res.* **2008,** *36*, W427-
645    W432.

646 22.    J. Gao, L. B. M. Ellis and L. P. Wackett, *Nucleic Acids Res.* **2010,** *38*, D488-D491.

647 23.    V. de Lorenzo, *Current Opinion in Biotechnology* **2008,** *19*, 579-589.

648    24.    Eawag Biocatalysis/Biodegradation Database. http://eawag-bbd.ethz.ch (accessed November
649    2016).

650    25.    J. Wicker, T. Lorsbach, M. Gütlein, E. Schmid, D. Latino, S. Kramer and K. Fenner, *Nucleic*
651    *Acids Res.* **2016,** *44*, D502-D508.

652    26.    enviPath – The environmental contaminant biotransformation pathway resource.
653    https://envipath.org (accessed November 2016).

654    27.    K. Fenner, J. F. Gao, S. Kramer, L. Ellis and L. Wackett, *Bioinformatics* **2008,** *24*, 2079-2085.

655    28.    J. Wicker, K. Fenner, L. Ellis, L. Wackett and S. Kramer, *Bioinformatics* **2010,** *26*, 814-821.

656    29.    G. D. Bending, S. D. Lincoln and R. N. Edmondson, *Environmental Pollution* **2006,** *139*, 279-
657    287.

658    30.    D. E. Helbling, D. R. Johnson, M. Honti and K. Fenner, *Environmental Science & Technology*
659    **2012,** *46*, 10579-10588.

660    31.    T. Kaeberlein, K. Lewis and S. S. Epstein, *Science* **2002,** *296*, 1127-1129.

661    32.    A. Kowalczyk, T. J. Martin, O. R. Price, J. R. Snape, R. A. van Egmond, C. J. Finnegan, H.
662    Schafer, R. J. Davenport and G. D. Bending, *Ecotox. Environ. Safe.* **2015,** *111*, 9-22.

663    33.    EFSA Draft Assessment Reports (DAR). dar.efsa.europe/dar-web/provision (accessed
664    November 2016)

665    34.    EU, Regulation (EC) No 1907/2006 of the European Parliament and of the Council of 18
666    December 2006 concerning the Registration, Evaluation, Authorisation and Restriction of Chemicals
667    (REACH), establishing a European Chemicals Agency, amending Directive 1999/45/EC and repealing
668    Council Regulation (EEC) No 793/93 and Commission Regulation (EC) No 1488/94 as well as
669    Council Directive 76/769/EEC and Commission Directives 91/155/EEC, 93/67/EEC, 93/105/EC and
670    2000/21/EC. *Official Journal of the European Union* **2006,** *L 396*, 1-849.

671    35.    R. C. Kolanczyk, P. Schmieder, W. J. Jones, O. G. Mekenyan, A. Chapkanov, S. Temelkov, S.
672    Kotov, M. Velikova, V. Kamenska, K. Vasilev and G. D. Veith, *Regul. Toxicol. Pharmacol.* **2012,** *63*,
673    84-96.

674    36.    PPDB: Pesticide Properties DataBase. http://sitem.herts.ac.uk/aeru/ppdb/en/index.htm
675    (accessed November 2016).

676    37.    S. Dimitrov, D. Nedelcheva, N. Dimitrova and O. Mekenyan, *Science of the Total*
677    *Environment* **2010,** *408*, 3811-3816.

678    38.    Eawag-Soil package. https://envipath.org/package/5882df9c-dae1-4d80-a40e-db4724271456
679    (accessed December 2016).

680    39.    H. P. Singer, A. E. Wössner, C. S. McArdell and K. Fenner, *Environmental Science &*
681    *Technology* **2016,** *50*, 6698-6707.

682    40.    T. Sander, J. Freyss, M. von Korff and C. Rufener, *J. Chem Inf. Model.* **2015,** *55*, 460-473.

683    41.    D. Rogers and M. Hahn, *J. Chem Inf. Model.* **2010,** *50*, 742-754.