

1 enviLink: A database linking contaminant biotransformation rules to enzyme classes in  
2 support of functional association mining

3

4 Emanuel Schmid<sup>1</sup>, Kathrin Fenner<sup>2,3,4</sup>

5

6 <sup>1</sup>Scientific IT Support, ETH Zürich, 8092 Zürich, Switzerland; <sup>2</sup>Eawag, Swiss Federal

7 Institute of Aquatic Science and Technology, 8600 Dübendorf, Switzerland; <sup>3</sup>Institute of

8 Biogeochemistry and Pollutant Dynamics, ETH Zürich, 8092 Zürich, Switzerland;

9 <sup>4</sup>Department of Chemistry, University of Zürich, 8057 Zürich, Switzerland.

10

11

## 12 **Abstract**

13 **Motivation:** The ability to assess and engineer biotransformation of chemical contaminants  
14 present in the environment requires knowledge on which enzymes can catalyze specific  
15 contaminant biotransformation reactions. For the majority of over 100'000 chemicals in  
16 commerce such knowledge is not available. Enumeration of enzyme classes potentially  
17 catalyzing observed or *de novo* predicted contaminant biotransformation reactions can  
18 support research that aims at experimentally uncovering enzymes involved in contaminant  
19 biotransformation in complex natural microbial communities.

20 **Database:** enviLink is a new data module integrated into the enviPath database and contains  
21 316 theoretically derived linkages between generalized biotransformation rules used for  
22 contaminant biotransformation prediction in enviPath and 3<sup>rd</sup> level EC classes. Rule-EC  
23 linkages have been derived using two reaction databases, i.e., Eawag-BBD in enviPath,  
24 focused on contaminant biotransformation reactions, and KEGG. 32.6% of identified rule-EC  
25 linkages overlap between the two databases, whereas 40.2% and 27.2%, respectively, are  
26 originating from Eawag-BBD and KEGG only.

27 **Implementation and availability:** enviLink is encoded in RDF triples as part of the enviPath  
28 RDF database. enviPath is hosted on a public webserver ([envipath.org](http://envipath.org)) and all data is freely  
29 available for non-commercial use. enviLink can be searched online for individual  
30 transformation rules of interest (<https://tinyurl.com/y63ath3k>) and is also fully downloadable  
31 from the supporting materials (i.e., Jupyter notebook “enviLink” and tsv files provided  
32 through GitHub at <https://github.com/emanuel-schmid/enviLink>).

33

## 34 **Introduction**

35 Refined understanding of contaminant degradation in environmental microbial  
36 communities depends on knowledge about catalyzing enzymes. For co-metabolic  
37 transformation at low substance concentrations that knowledge is hardly available. Available  
38 experimental approaches (gene knock outs or overexpression) are very costly and labor-  
39 intensive and therefore rely on strong hypotheses about potential enzyme candidates. To  
40 identify such potential enzyme candidates, functional association mining between  
41 metatranscriptomic or -genomic profiles and contaminant biotransformation information (i.e.,  
42 rate constants and reaction pathway) has been suggested as a promising way forward <sup>1-3</sup>.  
43 However, association mining suffers from low significance due to massive multiple  
44 hypothesis testing unless the range of enzymes plausibly catalyzing a given, observed  
45 transformation reaction can be restricted.

46 Currently available tools that, given a transformation reaction, allow predicting potentially  
47 catalyzing enzymes (or enzyme-encoding genes) are E-zyme/E-zyme2 <sup>4,5</sup> and BridgIT <sup>6</sup>. One  
48 obvious drawback for their application to contaminant biotransformation reactions is that both  
49 tools are trained on KEGG data only. KEGG very extensively covers reactions associated  
50 with primary metabolism and secondary metabolism of natural products, but only contains  
51 limited information on contaminants.

52 Eawag-BBD instead exclusively contains information on experimentally observed  
53 contaminant biotransformation reactions <sup>7</sup>. These have served as a basis for deriving a set of  
54 manually curated generalized biotransformation rules (btrules) which are used for *de novo*  
55 contaminant pathway prediction <sup>8</sup>. Most contaminant biotransformation reactions in Eawag-  
56 BBD are annotated with an EC number, which has been manually extracted by a data curator  
57 from the original publication reporting the experimental evidence. Most reactions are  
58 annotated with a 4<sup>th</sup> or 3<sup>rd</sup> level EC number (44.2% and 43.3%, respectively). The remaining  
59 2<sup>nd</sup> and 1<sup>st</sup> level annotations are based on educated guesses of the data curators rather than  
60 actual experimentally proven linkages (personal communication, Prof. Lynda Ellis). Both,  
61 Eawag-BBD and Eawag-PPS have recently been implemented in a more flexible and state-of-  
62 the-art successor system called enviPath <sup>9</sup>.

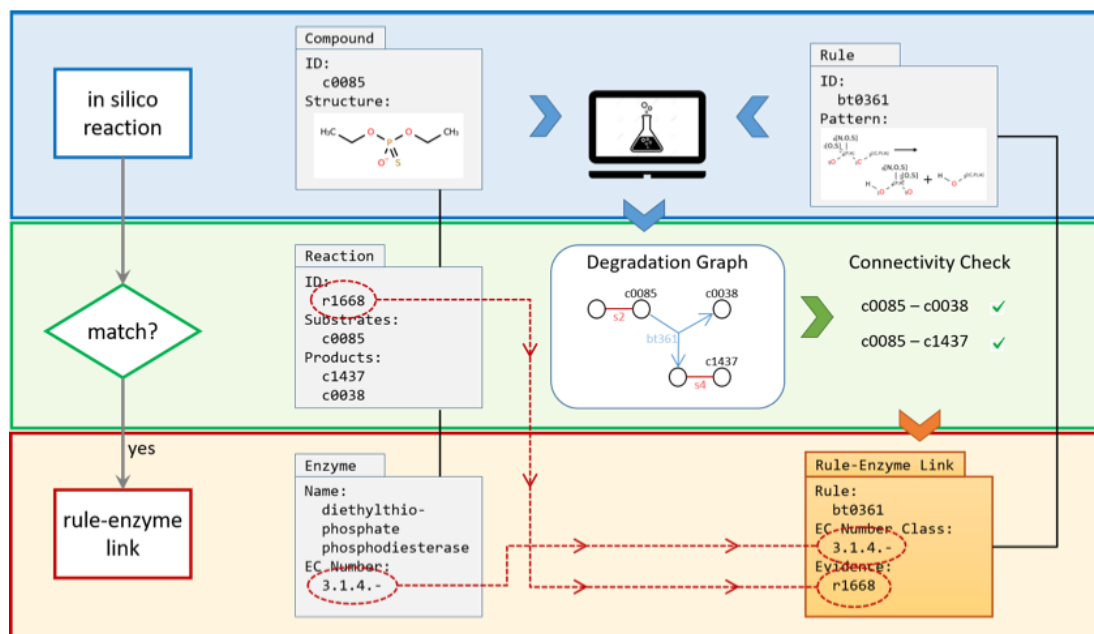
63 In developing enviLink, the database presented here, we therefore used reactions and their  
64 experimentally associated enzymes from both Eawag-BBD and, for completion, KEGG. We  
65 derived linkages between generalized biotransformation rules and 3<sup>rd</sup>-level EC classes rather  
66 than between actual reactions and 4<sup>th</sup> level EC classes (as in BridgIT or E-zyme) for two  
67 reasons. First, given the enormous structural diversity of synthetic chemicals, the number of  
68 experimentally validated enzyme-reaction associations for contaminants is simply too low to  
69 derive a finer linkage scheme and validate it. Second, for the purpose of functional

70 association mining, there is no need to target one enzyme only, but rather the goal is to  
71 produce a reasonably restricted, yet comprehensive list of suspect enzymes.

72

## 73 Methods

74 The workflow for creating enviLink included three major steps (see Figure 1): (i) “in silico”  
75 reaction of Eawag-BBD and KEGG substrates against all Eawag-BBD biotransformation  
76 rules (btrules); (ii) Comparison of “in silico” generated reaction pairs (i.e., substrate(s) and  
77 product(s)) with Eawag-BBD or KEGG database reactions to find matching reactions; and  
78 (iii) generation of rule-enzyme links by extracting enzyme class of matching reactions and  
79 associating them with the btrule that predicted this reaction. Finally, to derive linkages  
80 between generalized biotransformation rules and 3<sup>rd</sup>-level EC classes, 4<sup>th</sup>-level EC numbers  
81 were summarized into the corresponding 3<sup>rd</sup>-level EC classes. All analyses were carried out  
82 separately for Eawag-BBD (1479 contaminant biotransformation reactions with 1301  
83 associated EC classes) and KEGG (9952 reactions with 7007 associated EC classes, as of  
84 June 5<sup>th</sup> 2020), and resulting links were compared as discussed below (note that for BBD 3<sup>rd</sup>  
85 and 4<sup>th</sup> level ECs were extracted, whereas for KEGG only the 4<sup>th</sup> level ECs were considered).  
86 Details on each step of the workflow are given as Supporting Information in the form of  
87 interlinked Jupyter notebooks, which are available through GitHub  
88 (<https://github.com/emanuel-schmid/enviLink>). All data required to run the notebooks are  
89 available at this repository in the form of tsv files, but can alternatively also be downloaded  
90 following the code provided in the Jupyter notebooks.



91

92 **Figure 1:** Overview of workflow to produce rule-enzyme linkages demonstrated for the example of  
93 Eawag-BBD and including three major steps: (i) Enumeration of “in silico” reactions by running all  
94 btrules against all BBD compound structures to produce predicted degradation graph for each BBD

95 compound (blue upper panel, blue arrows), (ii) comparison of “in silico” prediction reactions from  
96 degradation graph with database reactions to check for matching reactions (green middle panel, green  
97 arrows), and (iii) generation of rule-enzyme links for matching reactions (red lower panel, orange  
98 arrow). Data entities are named in accordance with Eawag-BBD, and the example shown is taken from  
99 Eawag-BBD. In the degradation graph, blue lines stand for predicted reactions and red lines for  
100 standardizations. s2 and s4 represent standardizations and, in the specific case, stand for  
101 protonation/deprotonation reactions at differently substituted phosphate groups. Black connectors  
102 between data entities represent database relations, and red dashed connectors visualize the information  
103 flow from “Reaction”, “Enzyme” and “Rule” to yield entries for the new data entity “Rule-Enzyme  
104 Link” in enviLink.

105

## 106 **Results**

107 Resulting linkages from Eawag-BBD and KEGG are accessible through the Jupyter notebook  
108 “enviLink results” at the GitHub repository. Alternatively, enviLink can be searched online at  
109 [envipath.org](http://envipath.org) for individual transformation rules of interest (see information given under “EC  
110 numbers” on the rule pages of the EAWAG-BBD package (<https://tinyurl.com/y63ath3k>)).  
111 Altogether 316 linkages between 169 btrules and 107 3<sup>rd</sup> level EC classes were found and  
112 compiled in enviLink. For 39 btrules, no corresponding 3<sup>rd</sup> level EC class could be identified.  
113 32.6% of the identified rule-EC linkages overlap between the two databases, whereas 40.2%  
114 and 27.2%, respectively, are originating from either Eawag-BBD or KEGG only. The fact  
115 that more than one third of the linkages originate from Eawag-BBD exclusively demonstrates  
116 its unique information content with respect to contaminant biotransformation. One example of  
117 such an Eawag-BBD-exclusive linkage is the link between bt0241 and bt0242, two rules for  
118 hydroxylation of secondary and tertiary aliphatic groups, and 1.14.15, which contains  
119 monooxygenases using a reduced iron-sulfur protein as additional electron donor. Eawag-  
120 BBD contains literature entries reporting hydroxylating activity of camphor 5-  
121 monooxygenase (EC 1.14.15.1) on specific contaminants (e.g., adamantanone, tetralin) that  
122 are obviously not in the scope of KEGG and hence not reported therein.

123 In the “enviLink results” notebook, a histogram is provided showing how the linkages cover  
124 the space of btrules and 3<sup>rd</sup> level EC classes. It can be observed that several 3<sup>rd</sup> level EC are  
125 linked to multiple btrules (e.g., EC 1.14.12 is linked to bt0042, bt0072, bt0216 etc., which all  
126 encode for *vic*-dihydroxylation reactions at differently substituted aromatic rings). This  
127 illustrates that btrules in enviPath are divergent from the EC classification system in that they  
128 were optimized for specificity in contaminant biotransformation prediction<sup>8, 10</sup>.

129 Finally, to illustrate application of enviLink, consider the neonicotinoide acetamiprid, for  
130 which we observed enzymatic hydrolysis to the corresponding amide in activated sludge<sup>11</sup>.  
131 This reaction is predicted by bt0028 in enviPath, which in turn is linked to EC 4.2.1.- (hydro-  
132 lases) in enviLink. When screening for associations between abundance of gene transcripts

133 annotated to 4<sup>th</sup> level EC classes belonging to 4.2.1.- and rate constants of acetamiprid  
134 biotransformation in activated sludge, nitrile hydratase transcript abundances (EC 4.2.1.84)  
135 showed significant correlations<sup>1</sup>. Indeed, own and literature evidence later confirmed that  
136 different nitrile hydratase homologs can turn over acetamiprid<sup>1,12</sup>.

137

### 138 **Acknowledgements**

139 We acknowledge financial support from the European Research Council under the European  
140 Union's Seventh Framework Programme (ERC grant agreement 614768, PROduCTS).

141

### 142 **References**

- 143 1. Achermann, S.; Mansfeldt, C. B.; Müller, M.; Johnson, D. R.; Fenner, K., Relating  
144 Metatranscriptomic Profiles to the Micropollutant Biotransformation Potential of Complex  
145 Microbial Communities. *Environmental Science & Technology* **2020**, *54*, (1), 235-244.
- 146 2. Johnson, D. R.; Helbling, D. E.; Men, Y.; Fenner, K., Can meta-omics help to  
147 establish causality between contaminant biotransformations and genes or gene products?  
148 *Environmental Science: Water Research & Technology* **2015**, *1*, (3), 272-278.
- 149 3. Stadler, L. B.; Love, N. G., Oxygen Half-Saturation Constants for Pharmaceuticals in  
150 Activated Sludge and Microbial Community Activity under Varied Oxygen Levels.  
151 *Environmental Science & Technology* **2019**, *53*, (4), 1918-1927.
- 152 4. Tabei, Y.; Yamanishi, Y.; Kotera, M., Simultaneous prediction of enzyme orthologs  
153 from chemical transformation patterns for de novo metabolic pathway reconstruction.  
154 *Bioinformatics* **2016**, *32*, (12), i278-i287.
- 155 5. Yamanishi, Y.; Hattori, M.; Kotera, M.; Goto, S.; Kanehisa, M., E-zyme: predicting  
156 potential EC numbers from the chemical transformation pattern of substrate-product pairs.  
157 *Bioinformatics* **2009**, *25*, (12), I179-I186.
- 158 6. Hadadi, N.; MohammadiPeyhani, H.; Miskovic, L.; Seijo, M.; Hatzimanikatis, V.,  
159 Enzyme annotation for orphan and novel reactions using knowledge of substrate reactive  
160 sites. *Proceedings of the National Academy of Sciences* **2019**, *116*, (15), 7298.
- 161 7. Gao, J.; Ellis, L. B. M.; Wackett, L. P., The University of Minnesota  
162 Biocatalysis/Biodegradation Database: Improving Public Access. *Nucleic Acids Res.* **2010**,  
163 *38*, D488-D491.
- 164 8. Gao, J. F.; Ellis, L. B. M.; Wackett, L. P., The University of Minnesota Pathway  
165 Prediction System: multi-level prediction and visualization. *Nucleic Acids Res.* **2011**, *39*,  
166 W406-W411.
- 167 9. Wicker, J.; Lorsbach, T.; Gütlein, M.; Schmid, E.; Latino, D.; Kramer, S.; Fenner, K.,  
168 enviPath – The environmental contaminant biotransformation pathway resource. *Nucleic*  
169 *Acids Res.* **2016**, *44*, (Database issue), D502-D508.

- 170 10. Fenner, K.; Gao, J. F.; Kramer, S.; Ellis, L.; Wackett, L., Data-driven extraction of  
171 relative reasoning rules to limit combinatorial explosion in biodegradation pathway  
172 prediction. *Bioinformatics* **2008**, *24*, (18), 2079-2085.
- 173 11. Achermann, S.; Falås, P.; Joss, A.; Mansfeldt, C. B.; Men, Y.; Vogler, B.; Fenner, K.,  
174 Trends in Micropollutant Biotransformation along a Solids Retention Time Gradient.  
175 *Environmental Science & Technology* **2018**, *52*, (20), 11601-11611.
- 176 12. Guo, L.; Fang, W.-W.; Guo, L.-L.; Yao, C.-F.; Zhao, Y.-X.; Ge, F.; Dai, Y.-J.,  
177 Biodegradation of the Neonicotinoid Insecticide Acetamiprid by Actinomycetes *Streptomyces*  
178 *canus* CGMCC 13662 and Characterization of the Novel Nitrile Hydratase Involved. *Journal*  
179 *of Agricultural and Food Chemistry* **2019**, *67*, (21), 5922-5931.
- 180