# Letter to the Editor:

# FAIR-ifying the Exposome Journal: Templates for Chemical Structures and Transformations

Emma L. Schymanski[1]* and Evan E. Bolton[2]*

[1]Luxembourg Centre for Systems Biomedicine (LCSB), University of Luxembourg, 6 avenue du Swing, 4367, Belvaux, Luxembourg. ORCID 0000-0001-6868-8145

[2]National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, 20894, USA. ORCID 0000-0002-5959-6190

* Correspondence: ELS: emma.schymanski@uni.lu and EEB: bolton@ncbi.nlm.nih.gov

***Running title***: FAIRifying Chemical Structures and Transformations

## Abstract

The exposome, the totality of lifetime exposures, is a new and highly complex paradigm for health and disease. Tackling this challenge requires an effort well beyond single individuals or laboratories, where every piece of the puzzle will be vital. The launch of this new Exposome journal coincides with the evolution of the exposome through its teenage years and into a growing maturity in an increasingly open and FAIR (***findable, accessible, interoperable, reusable***) world. This letter discusses how both authors and the Exposome journal alike can help increase the ***FAIR***ness of the chemical structural information and the associated metadata in the journal, aiming to capture more details about the chemistry of exposomics. The proposed chemical structure template can serve as an ***interoperable*** supplementary format that is made ***accessible*** through the website and more ***findable*** by linking the DOI of this data file to the article DOI metadata, supporting further ***reuse***. An additional Transformations template provides authors with a means to connect predecessor (parent, substrate) molecules to successor (transformation product, metabolite) molecules and thus provide ***FAIR*** connections between observed (*i.e.,* experimental) chemical exposures and biological responses, to help improve the public knowledgebase on exposome-related transformations. These connections are vital to extend current biochemical knowledge and to fulfil the current Exposome definition of "the cumulative measure of environmental influences and associated biological responses throughout the lifespan including exposures from the environment, diet, behaviour, and endogenous processes".

## Keywords

34 Open science, chemical information, FAIR, transformation products, data workflows, data sharing

## Main Text

36

### Motivation

38 The "exposome" is a concept first mentioned in 2005 by Wild[1] to offer an environmental complement to
39 the genome[2] in considering health and disease. Now that the exposome is in its adolescence and
40 "emerging from the primordial swamp" sufficiently to warrant its own journal[2], it is a good time to reflect
41 on what steps are required to enable exposomics to mirror the achievements of genomics. A quick search
42 reveals, for instance, that global investment in genomics is projected into the tens of billions in the coming
43 years[3,4], while the global investment in the exposome or exposomics is rather of the order of tens of
44 millions. Yet, exposomics is an extraordinarily complex paradigm that will certainly require concerted
45 global effort comparable to that of the human genome[5]. Although capturing "the cumulative measure of
46 environmental influences and associated biological responses throughout the lifespan including
47 exposures from the environment, diet, behaviour, and endogenous processes"[6] may seem unachievable
48 for some, sequencing the human genome was also considered an almost impossible task only a few
49 decades ago. While the success of genomics is arguably due to many factors (including extensive
50 investment), one very significant factor in its success is the open exchange of genomics data and the
51 ecosystem of open resources that has been built around genomics, enabling scientists around the world
52 to achieve extraordinary progress in a relatively short time. Can exposomics achieve the same?

53 With this letter, we provide some perspectives and guidance on how both authors of articles in Exposome
54 and the Exposome journal itself can contribute to the cumulative efforts needed to tackle the exposomics
55 challenge from a chemical information and chemical informatics standpoint. Exposomics is inherently a
56 data-driven discipline. The interlinking of chemical, disease and reference information is already providing
57 support to exposomics efforts, as shown in Figure 1 using an examples from PubChem[7] and the
58 Comparative Toxicogenomics Database (CTD)[8], as well as from the CompTox Chemicals Dashboard[9,10].
59 Such information gathering and cross-resource integration efforts are much easier if data is both open
60 and **FAIR** (*findable, accessible, interoperable, reusable*). Providing guidance and coordinating at a journal
61 level is one way to enable such information gathering; genomics data deposition is mandated in most
62 major journals and this has been key to building the open genomics data resources that are so critical for
63 food-based pathogen surveillance, COVID-19 disease variant tracking, and so much more. If sufficient
64 information for exposomics was available, what can we as a community achieve?

65 Authors need guidance to properly and uniformly capture and report chemical structure information and
66 transformations, *i.e.,* connecting either endogenous or exogenous chemicals with their metabolites – thus
67 helping capture the associated biological responses. The flexible templates provided here (see sections
68 "Chemical Structure Data" and "Transformations Data") show how authors can consistently submit this
69 information to the Exposome journal as supplementary materials with their articles. These templates are
70 designed such that authors can include as much or as little information as is available, yet still contribute
71 their knowledge and outcomes to the exposomics "pool" (and beyond) in an open and FAIR manner. The
72 "Chemical Structure Data" template is identical to the template introduced recently in the Journal of
73 Cheminformatics[11].

**A**

# PubChem 1-Chloro-2,4-dinitrobenzene (Compound)

## 13 Associated Disorders and Diseases

Page 3 of 14 items · View More Rows & Details · Download

| Disease | Evidence Type | Evidence PMID |
|---|---|---|
| Inflammation | marker/mechanism | 19647056 20096324 25449201 |
| Melanoma | therapeutic | 12202904 17334785 |
| Necrosis | marker/mechanism | 28826779 |
| Respiratory Hypersensitivity | marker/mechanism | 17693426 |

‹ Previous   1   2   3

▶ Comparative Toxicogenomics Database (CTD)

**B**

| Chemical | CAS RN | DSSToxID | PMID Ct | Seizures | Nervous System Diseases | Peripheral Nervous System Diseases | Brain Diseases | Muscular Diseases | Basal Ganglia Diseases | Parkinson Disease, Secondary | Coma | Hallucinations | Tremor | Memory Disorders | Central Nervous |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Cisplatin | 15663-27-1 | DTXSID4024983 | 1032 | 20 | 47 | 140 | 13 | 0 | 4 | 1 | 1 | 0 | 1 | 2 | 4 |
| Ethanol | 64-17-5 | DTXSID9020584 | 768 | 100 | 23 | 11 | 18 | 26 | 1 | 3 | 20 | 6 | 17 | 54 | 2 |
| Lead | 7439-92-1 | DTXSID2024161 | 740 | 28 | 107 | 68 | 102 | 4 | 2 | 2 | 1 | 3 | 4 | 19 | 30 |
| Lithium | 7439-93-2 | DTXSID5036761 | 689 | 30 | 50 | 9 | 22 | 5 | 36 | 13 | 25 | 6 | 93 | 12 | 15 |
| Valproic Acid | 76584-70-8 | DTXSID70227388 | 666 | 32 | 10 | 3 | 65 | 6 | 10 | 18 | 45 | 5 | 18 | 4 | 2 |
| 1-Methyl-4-phe | 28289-54-5 | DTXSID8040933 | 638 | 1 | 24 | 0 | 11 | 0 | 6 | 289 | 0 | 0 | 5 | 0 | 1 |
| Vincristine | 2068-78-2 | DTXSID8044331 | 567 | 17 | 59 | 125 | 15 | 5 | 1 | 1 | 5 | 3 | 2 | 1 | 8 |
| Phenytoin | 57-41-0 | DTXSID8020541 | 560 | 37 | 24 | 25 | 16 | 9 | 3 | 1 | 9 | 3 | 8 | 4 | 6 |
| Haloperidol | 52-86-8 | DTXSID4034150 | 555 | 6 | 6 | 1 | 10 | 6 | 153 | 51 | 4 | 4 | 11 | 1 | 0 |
| Cocaine | 50-36-2 | DTXSID2038443 | 530 | 151 | 16 | 0 | 8 | 0 | 2 | 3 | 3 | 8 | 6 | 12 | 11 |
| Aspirin | 50-78-2 | DTXSID5020108 | 489 | 8 | 3 | 0 | 3 | 2 | 2 | 0 | 9 | 4 | 1 | 0 | 5 |
| Paclitaxel | 33069-62-4 | DTXSID9023413 | 485 | 4 | 43 | 217 | 9 | 14 | 0 | 0 | 0 | 0 | 0 | 1 | 2 |
| Aluminum | 7429-90-5 | DTXSID3040273 | 477 | 13 | 41 | 1 | 105 | 4 | 0 | 0 | 1 | 0 | 1 | 13 | 12 |
| Lidocaine | 6108-05-0 | DTXSID80209953 | 464 | 150 | 26 | 15 | 3 | 2 | 0 | 0 | 8 | 4 | 6 | 2 | 10 |
| Methotrexate | 59-05-2 | DTXSID4020822 | 451 | 17 | 25 | 1 | 79 | 4 | 0 | 1 | 5 | 0 | 1 | 9 | 18 |
| Mercury | 7439-97-6 | DTXSID1024172 | 450 | 6 | 79 | 22 | 23 | 2 | 3 | 5 | 2 | 2 | 38 | 7 | 25 |

*Figure 1: FAIRifying and opening up exposomics information is critical to "big data" exposomics, empowering information discovery and cross-resource integration. Top (A): associated disorders and diseases (and references) for a single chemical, 1-chloro-2,4-dinitrobenzene in PubChem[7], with information sourced from the Comparative Toxicogenomics Database (CTD)[8]. Source: https://pubchem.ncbi.nlm.nih.gov/compound/6#section=Associated-Disorders-and-Diseases. Bottom (B): Individual chemical – disease endpoint mappings via Name, Chemical Abstract Services Registry Numbers (CAS RN), CompTox Chemicals Dashboard identifiers (DSSToxID or DTXSIDs), plus total and endpoint-specific reference counts in the context of neurotoxicity, embedded in an excel macro[10,12].*

An incredible amount of knowledge relevant for exposomics has already been gathered, yet current studies are based primarily on using public resources to find existing information. To extend exposomics into the future, we need to enable the discovery and reporting of new findings via rapid integration into

84  public resources. Thus, author contributions, no matter how small, will gradually help build the bigger
85  picture needed to unravel and comprehend the exposome. Before we launch into the template
86  descriptions, a few definitions are covered in the next section.

87

## Definitions

89  While "FAIR" and "Open" are used somewhat interchangeably in this article as we strongly believe that
90  chemical data should be both where possible, there is a distinction that is particularly relevant for
91  exposomics, as sensitive human data cannot necessarily be made open. Data can be "open" but not
92  "FAIR", and vice versa. Open science has many facets; of most relevance to this article is open access.
93  Open access (OA) is a set of principles and a range of practices through which research outputs are
94  distributed online, free of cost or other access barriers[13]. The FAIR principles for digital assets, on the other
95  hand, include guidance on how to make data more Findable, Accessible, Interoperable and Reusable[14,15].
96  For example, if you have open data that is not findable, no one can use it; whereas if you have "FAIR" data
97  that is not "open", it is not available for integration into open community resources.  Thus, the most
98  powerful data is both open and FAIR

99  In Table 1, we provide some definitions of chemical and transformation terms used later in this article.

100  *Table 1: Definition of chemical and transformation terms used in this article and/or templates.*

| Concept | Definition |
|---|---|
| Biosystem | The medium in which the predecessor is transformed into the successor (e.g., environment, human liver, etc.) |
| Identifier | An identifier or name that you (the author) have for a chemical structure |
| InChI | IUPAC International Chemical Identifier is a descriptor of a chemical structure[16] |
| InChIKey | A 27-character long, layered "hash" of an InChI[16] |
| PubChem CID | PubChem Compound Identifier |
| Predecessor | Substrate/parent that is transformed (somehow) into a successor product |
| SMILES | Chemical structure notation expressed as a string |
| Successor | Transformation product/metabolite resulting from transformation (somehow) of a substrate/parent |

101

102

## Templates for FAIR Exposomics Chemical Data

104

### Chemical Structure Data

106  Better consideration of chemical factors in the exposome requires high-quality chemical information in
107  research articles. Many exposomics resources are based (mostly) on literature mining using name and
108  synonym matching, which can be notoriously prone to errors. In this section, we provide some guidance
109  on what information authors should consider providing, as well as the pros and cons of various choices.
110  Since this Chemical Structure Data template was presented recently to the Journal of Cheminformatics[11],
111  some of the material in this section overlaps with the previous article.

112 Authors should consider submitting their chemical structure information with their manuscript as
113 Supplementary Material using the suggested template as comma separated value (CSV; *.csv); or,
114 alternatively, as tab-separated value (TSV; *.tsv) or structure data file (SDF; *.sdf) formats. These formats
115 ensure maximum interoperability between resources and operating systems. The popular XLS(X) format
116 is not truly interoperable (options to save as CSV or TSV are offered), while the extraction of information
117 from PDF format is difficult without introducing errors. The content below describes the CSV/TSV formats,
118 SDF instructions are available elsewhere[17] (however, the SD fields should match the CSV/TSV headers). In
119 our experience, so far CSV often proves most interoperable for the widest audience, although the other
120 formats also have certain advantages.

121 For CSV/TSV files, the header (first row) indicates the data content of each column; each subsequent row
122 corresponds to a complete chemical record description: chemical structure, chemical names, identifiers,
123 comments, and any other data the authors wish to provide (as additional columns). The interoperable
124 case-insensitive template CSV/TSV column headers (or SDF SD fields) are: **SMILES, InChI**, and **InChIKey** for
125 chemical structure; **Name** and **Synonym** for chemical names; and **Comment** for textual comments. Any
126 additional columns headers (*e.g.*, for data, additional identifiers, or desired metadata) are up to the author
127 (*e.g.,* the **PubChem_CID** identifier header in Figure 2). Note that there may be many **Synonym** and
128 **Comment** columns in the file to provide space for more chemical names and metadata, respectively.

129 The author-submitted template file[18] should contain **at least one** of the following columns: **SMILES, InChI,**
130 **Name** or **InChIKey**. The **Name** column corresponds to a single primary name for the chemical structure.
131 Each **Synonym** column corresponds to an additional chemical name (one name entry per column). Each
132 **Comment** column can be added to provide additional text that may be important to the downstream user.
133 Authors can also provide additional CSV/TSV columns (or SDF SD fields) containing information about their
134 chemical substances (with unique, descriptive headers) for additional context. Chemical database
135 identifiers or registry numbers could be included in this manner (as additional columns or fields), or as a
136 **Synonym**. Note that chemical records indicating chemical structure with only **InChIKey** or **Name** will not
137 contain sufficient information to describe a chemical structure; and *can only be mapped to existing entries*
138 in destination resources. Batch services are available (*e.g.* from PubChem[7,19] or CompTox[9,20]) for authors
139 to add, *e.g., SMILES* and/or **InChI** to their records, based upon the **Name** or other identifiers.

140 Figure 1 in Schymanski & Bolton 2021[11] shows the template file, which is available for download[18] and as
141 Supporting Information with this article. Figure 2 below shows an example submission according to the
142 proposed template, created by sub-setting the "HSDBTPS" dataset of literature-mined and curated
143 transformation products from the Hazardous Substance Data Bank (HSDB) in PubChem[21,22]. This example
144 provides the **Name**, **SMILES** and **InChIKey** fields as suggested, and an identifier (the PubChem Compound
145 Identifier, CID) as an additional (optional) column (**PubChem_CID**) with a unique and easily recognizable
146 header that can be processed by other resources as they choose, helping with interoperability.

| PubChem_CID | Name | SMILES | InChIKey |
|---|---|---|---|
| 2256 | Atrazine | CCNC1=NC(=NC(=N1)Cl)NC(C)C | MXWJVTOOROXGIU-UHFFFAOYSA-N |
| 2328 | Bentazone | CC(C)N1C(=O)C2=CC=CC=C2NS1(=O)=O | ZOMSMJKLGFBRBS-UHFFFAOYSA-N |
| 3030 | Dicamba | COC1=C(C=CC(=C1C(=O)O)Cl)Cl | IWEDIXLBFLAXBO-UHFFFAOYSA-N |
| 3120 | Diuron | CN(C)C(=O)NC1=CC(=C(C=C1)Cl)Cl | XMTQQYYKAHVGBJ-UHFFFAOYSA-N |
| 4169 | Metolachlor | CCC1=CC=CC(=C1N(C(C)COC)C(=O)CCl)C | WVQBLGZPHOPPFO-UHFFFAOYSA-N |
| 7257 | 3,4-Dichloroaniline | C1=CC(=C(C=C1N)Cl)Cl | SDYWXFYBZPNOFX-UHFFFAOYSA-N |
| 12584 | Ammelide | C1(=NC(=O)NC(=O)N1)N | YSKUZVBSHIWEFK-UHFFFAOYSA-N |

148  *Figure 2: An example chemical structure data file constructed according to the proposed template[18] by taking a subset of the*
149  *HSDBTPS structure data[21]. Image created in RStudio (Version 1.2.5042). The HSDBTPS efforts resulted in the deposition of 5 new*
150  *structures to PubChem all documented in HSDB text snippets, CIDs 146035700, 146035701, 146035702, 146035703 and*
151  *146037633.*

## Transformations Data

153  The advancement of modern science is data driven[23,24]. Providing key data in a ready to use format helps
154  to assist in its reuse in research articles, regulatory reports, or machine learning data models. Exposomics
155  especially needs access to ready-to-use, high-quality chemical information from individual research
156  articles (*e.g.,* such as the connection of detected chemicals with the disease endpoint investigated or the
157  aggregation of known metabolites of thousands of common chemicals). For instance, HSDB contains
158  metabolites and metabolism information for 3220 chemicals gathered over 40 years, but these are only
159  available as text snippets that need to be matched to chemical structures by synonyms followed by
160  manual curation (initial efforts have covered only 1/100th of this dataset[22]). However, as mentioned above,
161  a key challenge in exposomics is to connect chemicals (*e.g.,* of anthropogenic origin, but also endogenous
162  or exogenous chemicals) that are associated with exposures with their biological response. Since
163  metabolism is the most dynamic of the biological responses, and metabolites per definition fall into the
164  same molecular mass category as many anthropogenic chemicals of concern, a key gap in exposomics
165  knowledge is the connection between chemicals and their metabolites. The efforts of many will be needed
166  to help fill this knowledge gap, and the timing could not be better for exposomics with several recent
167  studies emerging using *in vitro* enzymes to investigate parent-metabolite relationships of drugs and other
168  relevant chemicals[25,26].

169  The Transformations template provided here has been designed on the basis of recent efforts to fill the
170  gaps of transformation products in PubChem using literature data[27], in collaboration with the NORMAN
171  Suspect List Exchange (NORMAN-SLE)[28–30]. Several datasets from a variety of sources have now been
172  processed. Transformations from the NORMAN-SLE, where S## refers to the list number, followed by the
173  list code, include: S60 SWISSPEST19[31,32], S66 EAWAGTPS[33,34], S68 HSDBTPS[21,22], S73 METXBIODB[35,36], S74
174  REFTPS[37], S78 SLUPESTTPS[38,39], S79 UACCSCEC[40,41] and S81 THSTPS[42] (list available from https://git-
175  r3lab.uni.lu/eci/pubchem/-/raw/master/annotations/tps/Transformation_Datasets.txt).    Of    these,
176  MetXBioDB also contains enzyme information, while the rest are primarily environmental data. Figure 3
177  shows an example "environmental" dataset compiled from several of these lists, using the proposed
178  template. In addition to the NORMAN-SLE datasets, a dataset of more than 1200 transformations from

179 ChEMBL[43] has also been added, including enzyme, gene and protein information (where available). An
180 example of Transformations with more biological information available is given in Figure 4.

181 Information about both the predecessor (parent/precursor) and successor (transformation
182 product/metabolite) must be given for a valid transformation. The template can accept *at least one* of
183 **Name**, **SMILES** or **PubChem CID** for each, where **SMILES** or **CID** is preferred, and **SMILES** will be the most
184 interoperable. Note that these need not be consistent – for instance, it is possible to provide **SMILES** of
185 the successor and a **CID** of the predecessor if a **Name** or **CID** is not available for the successor. It is
186 preferable to give two fields, Figure 3 shows the example of **Name** and **CID**, while Figure 4 an example of
187 **SMILES** and **Name** (top panel on each figure).

| Predecessor_CID | Predecessor_Name | Transformation | Successor_CID | Successor_Name |
|---|---|---|---|---|
| 13101 | 6PPD | Ozone | 154926030 | 6PPD-quinone |
| 2256 | Atrazine | Environmental | 13878 | Deisopropyl-atrazine |
| 2256 | Atrazine | Mammalian metabolism | 135408770 | Ammeline |
| 2256 | Atrazine | Fungal metabolism | 22563 | Desethyl-atrazine |
| 2256 | Atrazine | Dehalogenation | 135398733 | Atrazine-2-hydroxy |
| 13450 | Terbutryn | Mammalian metabolism | 13019211 | Desethyl-terbutryn |
| 5216 | Simazine | Plant metabolism | 12584 | Ammelide |

| Biosystem | Reference_ID | Reference_Description |
|---|---|---|
| Environment | DOI:10.1126/science.abd6951 | Tian, Z. et al. (2020) A ubiquitous tire rubber-derived chemic... |
| Soil | DOI:10.5281/zenodo.4687924 | S78 \| SLUPESTTPS \| Pesticides and TPs from SLU, Sweden |
| Mammal | DOI:10.5281/zenodo.3827487 | Kearney, P.C., and D. D. Kaufman (eds.) Herbicides: Chemistr... |
| Fungus | PMID:8967773 | S68 \| HSDBTPS \| Transformation Products Extracted from HS... |
| Environment | DOI:10.1007/s13361-017-1797-6 | Schollee et al, Similarity of High-Resolution Tandem Mass S... |
| Mammal | DOI:10.1002/bms.1200050604 | S68 \| HSDBTPS \| Transformation Products Extracted from HS... |
| Plant | DOI:10.5281/zenodo.3827487 | USEPA/Office of Pesticides and Toxic Substances; Simazine: ... |

189 *Figure 3: An example of various environmental transformations constructed according to the proposed Transformations*
190 *template[44] (using Name and PubChem CID), taking a subset of transformations from NORMAN-SLE datasets (REFTPS[37], HSDBTPS[21],*
191 *SLUPESTTPS[38], EAWAGTPS[33] and SWISSPEST19[31]). Image created in RStudio (Version 1.2.5042).*

192

| Predecessor_Name | Predecessor_SMILES | Successor_Name | Successor_SMILES | Transformation |
|---|---|---|---|---|
| Carbamazepine | C1=CC=C2C(=C1)C=CC3=CC=CC=C3N2C(=O)N | Carbamazepine-10,11-epoxide | C1=CC=C2C(=C1)C3C(O3)C4=CC=CC=C4N2C(=O)N | Epoxidation of 1,2-disubstituted alkene / Human Phase I |
| Acrolein | C=CC=O | Acrylic Acid | C=CC(=O)O | |
| Furan | C1=COC=C1 | (E)-2-Butenedial | C(=C\C=O)\C=O | Oxidation / Human Phase I |
| Benzene | C1=CC=CC=C1 | Phenol | C1=CC=C(C=C1)O | Hydroxylation of aromatic carbon / Human Phase I |
| Nicotinamide | C1=CC(=CN=C1)C(=O)N | MNAM | C[N+]1=CC=CC(=C1)C(=O)N | |

| Biosystem | Enzyme | Gene_ID | Protein_ID | Reference_ID | Reference_Description |
|---|---|---|---|---|---|
| Human | CYP3A4\|CYP2C8 | | | DOI:10.1186/s13321-018-0324-5 | Brown, C.M. et al. (2008) Cytochromes P450: A Structure-Bas... |
| | Aldehyde dehydrogenase 1A1 | 216 | P00352 | DOI:10.1111/j.1365-2125.2006.02690.x | Data from ChEMBL - IDs (pred,succ,enzyme): CHEMBL721\|C... |
| Human | CYP2E1 | | | PMID:20043645 | S73 \| METXBIODB \| Metabolite Reaction Database from BioT... |
| Human | CYP2E1 | | | DOI:10.1186/s13321-018-0324-5 | Brown, C.M. et al. (2008); Cytochromes P450: A Structure-Ba... |
| | Nicotinamide N-methyltransferase | 4837 | P40261 | DOI:10.1124/dmd.112.049734 | Data from ChEMBL - ChEMBL IDs (pred,succ,enzyme): CHEM... |

194 *Figure 4: An example of biological transformations constructed according to the proposed Transformations template[44] (using*
195 *Name and SMILES), taking a subset of transformations from NORMAN-SLE dataset MetXBioDB[35] (from BioTransformer[36]) and*
196 *the ChEMBL[43] datasets on PubChem; both datasets have some degree of enzyme, gene and/or protein information available.*

197  If available, a brief description of the transformation is useful and can be provided in the "**Transformation**"
198  field (top panel, Figure 3 and Figure 4). Short, informative descriptions are preferred; the current entries
199  have been either extracted automatically from existing datasets or entered manually. In the future, it may
200  be possible to provide some guidance via an ontology as the public dataset grows to improve the machine
201  readability. Similarly, if information on the biosystem is available (*i.e.,* where the transformation takes
202  place), this can be included in the **Biosystem** column (see Figure 3 and Figure 4 for examples).

203  For datasets with biological information, this can be provided (optionally) in the **Enzyme**, **Gene_ID** and
204  **Protein_ID** columns. At this stage the template allows flexible input (see Figure 4 for examples) but
205  recommend **Enzyme** are provided as either: Enzyme Commission (EC) number,[45–47] such as "EC 2.3.2.23";
206  gene symbol, such as "CYP1A1"; or as enzyme names, such as "Aryl hydrocarbon hydroxylase".  The
207  **Gene_ID** is expected to be an NCBI Gene[48] ID, such as "1543".  The **Protein_ID** is expected to be either an
208  NCBI Protein[49] accession, such as "NP_059488.2" or an UniProt identifier,[50] such as "P08684".  If multiple
209  entries for **Enzyme**, **Gene_ID** and **Protein_ID** are provided, they should be separated by a "pipe" symbol
210  ("|") or provided as new rows.

211  Finally, the **Reference_ID** and **Reference_Description** columns provide the opportunity to credit the
212  original sources of the information. **Reference_ID** entries should be either PubMed identifiers[51] (PMIDs)
213  or Digital Object Identifiers[52] (DOIs), preceded with "PMID:" or "DOI:", respectively, for easy recognition,
214  and separated by a "pipe" ("|") if multiple IDs exist (they can be mixed – for example,
215  "PMID:33929905|DOI:10.1186/s13321-018-0324-5"). The **Reference_Description** can be used to provide
216  a free text form of the reference, to describe the data source (if no PMID / DOI available) or to describe
217  evidence of the transformation. Only **Reference_ID** can be processed automatically. Again, see Figure 3
218  and Figure 4 and the Transformations template[44] for examples.
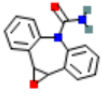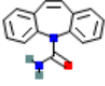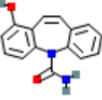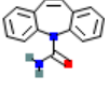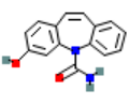
219  So far, about 6000 Transformations have been processed using these templates, from nine different
220  sources (many of these being composite data from several sources themselves, including ChEMBL[43],
221  MetXBioDB[35] and REFTPS[37]). The Transformations are being integrated into current computational mass
222  spectrometry workflows (such as patRoon[53] and as documented in Krier *et al.*[22]) and are openly available
223  for all. The summarized files are likewise available for comprehensive efforts such as BioTransformer[36] to
224  add this new data to their training set (MetXBioDB[35] is the library behind BioTransformer) and likewise
225  improve predictions. Overall, FAIR transformations data will greatly support exposomics, and discussions
226  to extend these templates into fields with formal ontologies and/or other formats such as mzTab[54,55] in
227  the future are welcomed.  As demonstrated in Figure 5 and Figure 6, one can see the benefits of arranging
228  data in FAIR templates. Figure 5 is an example of a resulting Transformation entry in PubChem, while
229  Figure 6 can be created automatically in CDK Depict using simple code in R to create annotated reaction
230  SMILES from the fields shown in Figure 3 only.

Figure 5 content: PubChem Carbamazepine (Compound) Transformations table

**Pub⬡hem** Carbamazepine (Compound)

## 8.10  Transformations

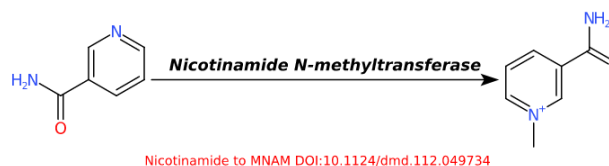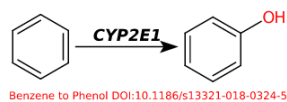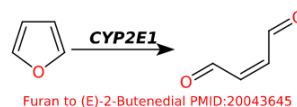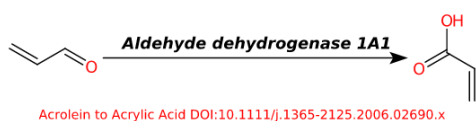7 items  View More Rows & Details ⬈                                    ⬇ Download

SORT BY  ⬍  Please Choose One                                    ⌄

| Predecessor | Predecessor Name | Successor | Successor Name | Transformation | Enzyme |
|---|---|---|---|---|---|
| | carbamazepine | | Carbamazepine 10,11-epoxide | Epoxidation of 1,2-disubstituted alkene / Human Phase I | CYP3A4 CYP2C8 |
| | carbamazepine | | 9-Hydroxycarbamazepine | Aromatic hydroxylation of fused benzene ring / Human Phase I | CYP3A4 CYP2C8 |
| | carbamazepine | | 3-Hydroxycarbamazepine | Aromatic hydroxylation of fused benzene ring / Human Phase I | CYP3A4 CYP2C8 |

232 *Figure 5: Example "Transformations" table in PubChem for Carbamazepine, demonstrating possible display options (including*
233 *hyperlinking) for FAIR Transformations. Source: https://pubchem.ncbi.nlm.nih.gov/compound/2554#section=Transformations.*

234



Acrolein to Acrylic Acid DOI:10.1111/j.1365-2125.2006.02690.x — *Aldehyde dehydrogenase 1A1*

Furan to (E)-2-Butenedial PMID:20043645 — *CYP2E1*

Benzene to Phenol DOI:10.1186/s13321-018-0324-5 — *CYP2E1*

Nicotinamide to MNAM DOI:10.1124/dmd.112.049734 — *Nicotinamide N-methyltransferase*

236 *Figure 6: Example reactions corresponding with the last four rows of Figure 3, automatically created and depicted with CDK*
237 *Depict[56] (https://www.simolecule.com/cdkdepict/depict.html) directly from template content shown in Figure 3 (SMILES, Name,*
238 *Enzyme and Reference_ID fields).*

239

## Closing

Exposomics is a data-driven science, and vast quantities of information will be needed for it to be successful. By making the output of exposomics research available in a more machine-readable way, we can accelerate our progress and rise to the challenge. The templates provided here are a means to make primary outputs FAIR (***Findable, Accessible, Interoperable***, ***Reusable***). When authors provide this content as Supplementary Information, it can be readily accessed and utilized, ideally without human intervention. When the journal interlinks these Supplementary Material files with the article DOI and associated metadata, other resources can rapidly find and integrate this content and provide enhanced services for the entire community. Improving the ***FAIR***ness of Supplementary Material greatly decreases the effort to combine and aggregate information between papers and improves the correctness of the information over text-mining based approaches. It also greatly enhances the visibility of the individual works and research outputs. As a young scientific discipline, the exposome should learn from its closely related 'elder' disciplines. Genomic approaches gained incredible traction due to the widely encouraged and eventually mandated sharing of information. Let us take these lessons to heart and advance together as a field. We need to share information – and lots of it – to help make sense of the exposome. The use of these facile, ready-to-use templates will help advance exposomics by contributing vital information to complete the exposomics "puzzle".

257

272

## Supplementary Materials

The chemical structure data submission template and transformations template are provided as Supplementary Material and are also available online[18,44,57].

All Transformations mentioned in this article are openly available on the NORMAN-SLE and PubChem.

277

# References

1. Wild CP. Complementing the genome with an "exposome": the outstanding challenge of environmental exposure measurement in molecular epidemiology. *Cancer Epidemiol Biomarkers Prev*. 2005;14(8):1847-1850. doi:10.1158/1055-9965.EPI-05-0456

2. Miller GW. Exposome: a new field, a new journal. *Exposome*. 2021;1(1). doi:10.1093/exposome/osab001

3. GlobeNewswire, Inc. Genomics Market to Reach USD 94.66 Billion by 2028; Increasing Genomics' Application & Rising Government Investments to Amplify Market Growth: Says Fortune Business Insights™. Accessed September 5, 2021. https://www.globenewswire.com/news-release/2021/05/20/2233128/0/en/Genomics-Market-to-Reach-USD-94-66-Billion-by-2028-Increasing-Genomics-Application-Rising-Government-Investments-to-Amplify-Market-Growth-Says-Fortune-Business-Insights.html

4. P&S Intelligence. Global Genomics Market to Reach $68 Billion by 2030: P&S Intelligence. Accessed September 5, 2021. https://www.prnewswire.com/news-releases/global-genomics-market-to-reach-68-billion-by-2030-ps-intelligence-301125318.html

5. Vermeulen R, Schymanski EL, Barabási AL, Miller GW. The exposome and health: Where chemistry meets biology. *Science*. 2020;367(6476):392. doi:10.1126/science.aay3164

6. Miller GW, Jones DP. The nature of nurture: refining the definition of the exposome. *Toxicol Sci*. 2014;137(1):1-2. doi:10.1093/toxsci/kft251

7. Kim S, Chen J, Cheng T, et al. PubChem 2019 update: improved access to chemical data. *Nucleic Acids Research*. 2019;47(D1):D1102-D1109. doi:10.1093/nar/gky1033

8. Davis AP, Grondin CJ, Johnson RJ, et al. Comparative Toxicogenomics Database (CTD): update 2021. *Nucleic Acids Research*. 2021;49(D1):D1138-D1143. doi:10.1093/nar/gkaa891

9. Williams AJ, Grulke CM, Edwards J, et al. The CompTox Chemistry Dashboard: a community data resource for environmental chemistry. *Journal of Cheminformatics*. 2017;9(1):61. doi:10.1186/s13321-017-0247-6

10. Schymanski EL, Baker NC, Williams AJ, et al. Connecting environmental exposure and neurodegeneration using cheminformatics and high resolution mass spectrometry: potential and challenges. *Environ Sci: Processes Impacts*. 2019;21(9):1426-1445. doi:10.1039/C9EM00068B

11. Schymanski EL, Bolton EE. FAIR chemical structures in the Journal of Cheminformatics. *J Cheminform*. 2021;13(1):50. doi:10.1186/s13321-021-00520-4

12. Baker NC, Schymanski EL, Williams AJ. Literature Neurotoxicants: Excel Macro File. *FigShare*. doi:10.23645/epacomptox.7334603

13. Peter Suber. Open Access Overview (definition, introduction). Accessed July 3, 2021. http://legacy.earlham.edu/~peters/fos/overview.htm

314     14. GO FAIR. FAIR Principles. Published 2021. Accessed March 23, 2021. https://www.go-fair.org/fair-
315         principles/

316     15. Wilkinson MD, Dumontier M, Aalbersberg IjJ, et al. Comment: The FAIR Guiding Principles for
317         scientific data management and stewardship. *Scientific Data*. 2016;3(1):1-9.
318         doi:10.1038/sdata.2016.18

319     16. Heller S, McNaught A, Stein S, Tchekhovskoi D, Pletnev I. InChI - the worldwide chemical structure
320         identifier standard. *Journal of Cheminformatics*. 2013;5(1):7. doi:10.1186/1758-2946-5-7

321     17. NCBI/NLM/NIH. PubChem Documentation: Substance SDF Submission. Published 2021. Accessed
322         March 23, 2021.
323         https://pubchem.ncbi.nlm.nih.gov/upload/docs/examples/substance_submission.sdf

324     18. NCBI/NLM/NIH. Chemical Structure Data Template (CSV). Published 2021. Accessed May 9, 2021.
325         https://ftp.ncbi.nlm.nih.gov/pubchem/Other/Submissions/Chemical_Structure_Data_Template.csv

326     19. NCBI/NLM/NIH. PubChem Identifier Exchange. Published 2021. Accessed March 23, 2021.
327         https://pubchem.ncbi.nlm.nih.gov/idexchange/idexchange.cgi

328     20. United States Environmental Protection Agency. CompTox Batch Search. Published 2021. Accessed
329         March 23, 2021. https://comptox.epa.gov/dashboard/dsstoxdb/batch_search

330     21. LCSB-ECI, Krier, Jessy, Schymanski, Emma, et al. S68 | HSDBTPS | Transformation Products Extracted
331         from HSDB Content in PubChem. Published online June 11, 2020. doi:10.5281/ZENODO.3827487

332     22. Krier J, Singh RR, Kondić T, et al. Discovering pesticides and their TPs in Luxembourg waters using
333         open cheminformatics approaches. *Environment International*. 2022;158:106885.
334         doi:10.1016/j.envint.2021.106885

335     23. Montáns FJ, Chinesta F, Gómez-Bombarelli R, Kutz JN. Data-driven modeling and learning in science
336         and engineering. *Comptes Rendus Mécanique*. 2019;347(11):845-855.
337         doi:10.1016/j.crme.2019.11.009

338     24. Clauset A, Larremore DB, Sinatra R. Data-driven predictions in the science of science. *Science*.
339         2017;355(6324):477-480. doi:10.1126/science.aal4217

340     25. Liu K, Lee C, Singer G, et al. Enzyme-Based Chemical Identification for Metabolomics. *FASEB j*.
341         2021;35(S1):fasebj.2021.35.S1.04277. doi:10.1096/fasebj.2021.35.S1.04277

342     26. Ross DH, Seguin RP, Krinsky AM, Xu L. *High-Throughput Measurement and Machine Learning-Based
343         Prediction of Collision Cross Sections for Drugs and Drug Metabolites*. Bioinformatics; 2021.
344         doi:10.1101/2021.05.13.443945

345     27. Schymanski EL, Kondić T, Neumann S, Thiessen PA, Zhang J, Bolton EE. Empowering large chemical
346         knowledge bases for exposomics: PubChemLite meets MetFrag. *J Cheminform*. 2021;13(1):19.
347         doi:10.1186/s13321-021-00489-0

348    28. NORMAN Network. NORMAN Suspect List Exchange. NORMAN Suspect List Exchange. Accessed
349         June 9, 2019. https://www.norman-network.com/nds/SLE/

350    29. NORMAN Network. NORMAN Suspect List Exchange on Zenodo. NORMAN Suspect List Exchange:
351         Zenodo Community. Accessed June 9, 2019. https://zenodo.org/communities/norman-sle/

352    30. NORMAN Network, NCBI/NLM/NIH. NORMAN SLE Classification Browser. Accessed May 7, 2020.
353         https://pubchem.ncbi.nlm.nih.gov/classification/#hid=101

354    31. Kiefer, Karin, Müller, Adrian, Singer, Heinz, Hollender, Juliane. S60 | SWISSPEST19 | Swiss Pesticides
355         and Metabolites from Kiefer et al 2019. Published online November 17, 2019.
356         http://doi.org/10.5281/zenodo.3544760

357    32. Kiefer K, Müller A, Singer H, Hollender J. New relevant pesticide transformation products in
358         groundwater detected using target and suspect screening for agricultural and urban micropollutants
359         with LC-HRMS. *Water Research*. 2019;165:114972. doi:10.1016/j.watres.2019.114972

360    33. Schollee, Jennifer, Schymanski, Emma. S66 | EAWAGTPS | Parent-Transformation Product Pairs
361         from Eawag. Published online April 23, 2020. doi:10.5281/ZENODO.3754448

362    34. Schollée JE, Schymanski EL, Stravs MA, Gulde R, Thomaidis NS, Hollender J. Similarity of High-
363         Resolution Tandem Mass Spectrometry Spectra of Structurally Related Micropollutants and
364         Transformation Products. *J Am Soc Mass Spectrom*. 2017;28(12):2692-2704. doi:10.1007/s13361-
365         017-1797-6

366    35. Djoumbou-Feunang, Yannick, Schymanski, Emma, Zhang, Jeff, Wishart, David S. S73 | METXBIODB |
367         Metabolite Reaction Database from BioTransformer. Published online November 5, 2020.
368         doi:10.5281/ZENODO.4056560

369    36. Djoumbou-Feunang Y, Fiamoncini J, Gil-de-la-Fuente A, Greiner R, Manach C, Wishart DS.
370         BioTransformer: a comprehensive computational tool for small molecule metabolism prediction and
371         metabolite identification. *J Cheminform*. 2019;11(1):2. doi:10.1186/s13321-018-0324-5

372    37. Schymanski, Emma. S74 | REFTPS | Transformation Products and Reactions from Literature.
373         Published online December 12, 2020. doi:10.5281/ZENODO.4318838

374    38. Menger, Frank, Boström, Gustaf. S78 | SLUPESTTPS | Pesticides and TPs from SLU, Sweden.
375         Published online May 10, 2021. doi:10.5281/ZENODO.4687924

376    39. Menger F, Boström G, Jonsson O, et al. Identification of Pesticide Transformation Products in
377         Surface Water Using Suspect Screening Combined with National Monitoring Data. *Environ Sci
378         Technol*. 2021;55(15):10343-10353. doi:10.1021/acs.est.1c00466

379    40. Belova L, Caballero-Casero N, van Nuijs ALN, Covaci A. Ion Mobility-High-Resolution Mass
380         Spectrometry (IM-HRMS) for the Analysis of Contaminants of Emerging Concern (CECs): Database
381         Compilation and Application to Urine Samples. *Anal Chem*. 2021;93(16):6428-6436.
382         doi:10.1021/acs.analchem.1c00142

383   41. Belova, Lidia, Caballero-Casero, Noelia, van Nuijs, Alexander L. N., Covaci, Adrian. S79 | UACCSCEC |
384       Collision Cross Section (CCS) Library from UAntwerp. Published online May 10, 2021.
385       doi:10.5281/ZENODO.4704648

386   42. Merino C, Vinaixa M, Ramirez N. S81 | THSTPS | Thirdhand Smoke Specific Metabolites. Published
387       online September 2, 2021. doi:10.5281/ZENODO.5394629

388   43. Gaulton A, Hersey A, Nowotka M, et al. The ChEMBL database in 2017. *Nucleic Acids Res*.
389       2017;45(D1):D945-D954. doi:10.1093/nar/gkw1074

390   44. NCBI/NLM/NIH. Transformations Data Template (CSV). Published 2021. Accessed May 25, 2021.
391       https://ftp.ncbi.nlm.nih.gov/pubchem/Other/Submissions/Transformations_Template.csv

392   45. McDonald AG, Boyce S, Tipton KF. ExplorEnz: the primary source of the IUBMB enzyme list. *Nucleic*
393       *Acids Research*. 2009;37(Database):D593-D597. doi:10.1093/nar/gkn582

394   46. Bairoch A. The ENZYME database in 2000. *Nucleic Acids Res*. 2000;28(1):304-305.
395       doi:10.1093/nar/28.1.304

396   47. Chang A, Jeske L, Ulbrich S, et al. BRENDA, the ELIXIR core data resource in 2021: new developments
397       and updates. *Nucleic Acids Research*. 2021;49(D1):D498-D508. doi:10.1093/nar/gkaa1025

398   48. Brown GR, Hem V, Katz KS, et al. Gene: a gene-centered information resource at NCBI. *Nucleic Acids*
399       *Res*. 2015;43(Database issue):D36-42. doi:10.1093/nar/gku1055

400   49. Benson DA, Cavanaugh M, Clark K, et al. GenBank. *Nucleic Acids Res*. 2018;46(D1):D41-D47.
401       doi:10.1093/nar/gkx1094

402   50. The UniProt Consortium. UniProt: the universal protein knowledgebase. *Nucleic Acids Res*.
403       2017;45(D1):D158-D169. doi:10.1093/nar/gkw1099

404   51. Sayers EW, Beck J, Bolton EE, et al. Database resources of the National Center for Biotechnology
405       Information. *Nucleic Acids Res*. 2021;49(D1):D10-D17. doi:10.1093/nar/gkaa892

406   52. International DOI Foundation. Frequently Asked Questions about the DOI® System. Accessed
407       September 7, 2021. https://www.doi.org/faq.html

408   53. Helmus R, ter Laak TL, van Wezel AP, de Voogt P, Schymanski EL. patRoon: open source software
409       platform for environmental mass spectrometry based non-target screening. *J Cheminform*.
410       2021;13(1):1. doi:10.1186/s13321-020-00477-w

411   54. Griss J, Jones AR, Sachsenberg T, et al. The mzTab Data Exchange Format: Communicating Mass-
412       spectrometry-based Proteomics and Metabolomics Experimental Results to a Wider Audience.
413       *Molecular & Cellular Proteomics*. 2014;13(10):2765-2775. doi:10.1074/mcp.O113.036681

414   55. Hoffmann N, Rein J, Sachsenberg T, et al. mzTab-M: A Data Standard for Sharing Quantitative
415       Results in Mass Spectrometry Metabolomics. *Anal Chem*. 2019;91(5):3302-3310.
416       doi:10.1021/acs.analchem.8b04310

417    56. Mayfield J. CDK Depict Web Interface. Accessed October 30, 2018.
418        http://simolecule.com/cdkdepict/depict.html

419    57. NCBI/NLM/NIH. PubChem Submissions Template Folder. Published 2021. Accessed May 25, 2021.
420        https://ftp.ncbi.nlm.nih.gov/pubchem/Other/Submissions/

421