





FAIRifying the exposome journal: Templates for chemical structures and transformations

Emma L. Schymanski ^{1,*}, and Evan E. Bolton ^{2,*}

¹Luxembourg Centre for Systems Biomedicine (LCSB), University of Luxembourg, 6 Avenue du Swing, Belvaux, Luxembourg

²National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA

*To whom correspondence should be addressed: emma.schymanski@uni.lu; bolton@ncbi.nlm.nih.gov

Editor's Note

The following Letter to the Editor is rather unique. Based on a dialog between the authors and the Editor, it was decided that the journal would grant an exception and allow an extended Letter to the Editor to facilitate discussion on this topic.

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, and reusable*) world. This letter discusses how both authors and the Exposome journal alike can help increase the FAIRness 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 and substrate) molecules to successor (transformation product and 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, behavior, and endogenous processes”.

Keywords: Open science; chemical information; FAIR; transformation products; data workflows; data sharing

Motivation

The “exposome” is a concept first mentioned in 2005 by Wild¹ to offer an environmental complement to the genome² in considering health and disease. Now that the exposome is in its adolescence and “emerging from the primordial swamp” sufficiently to warrant its own journal,² it is a good time to reflect on what steps are required to enable exposomics to mirror the achievements of genomics. A quick search reveals, for instance, that global investment in genomics is projected into the tens of billions in the coming years,^{3,4} while the global investment in the exposome or exposomics is rather of the order of tens of millions. Yet, exposomics is an extraordinarily complex paradigm that will certainly require concerted global effort comparable to that of the human genome.⁵ Although capturing “the cumulative measure of environmental influences and associated biological responses throughout the lifespan including exposures from the environment, diet, behaviour, and endogenous processes”⁶ may seem unachievable for some, sequencing the human genome was also considered an almost impossible task only a few decades ago. While the success of genomics is arguably due to many factors (including extensive investment), one very significant factor in its success is the open exchange of genomics data and the

ecosystem of open resources that has been built around genomics, enabling scientists around the world to achieve extraordinary progress in a relatively short time. Can exposomics achieve the same?

With this letter, we provide some perspectives and guidance on how both authors of articles in Exposome and the Exposome journal itself can contribute to the cumulative efforts needed to tackle the exposomics challenge from a chemical information and chemical informatics standpoint. Exposomics is inherently a data-driven discipline. The interlinking of chemical, disease and reference information is already providing support to exposomics efforts, as shown in Figure 1 using examples from PubChem⁷ and the Comparative Toxicogenomics Database (CTD),⁸ as well as from the CompTox Chemicals Dashboard.^{9,10} Such information gathering and cross-resource integration efforts are much easier if data are both open and FAIR (*findable, accessible, interoperable, reusable*). Providing guidance and coordinating at a journal level is one way to enable such information gathering; genomics data deposition is mandated in most major journals and this has been key to building the open genomics data resources that are so critical for food-based pathogen surveillance, COVID-19 disease variant tracking, and so much more. If sufficient

Received: September 08, 2021. Revised: December 06, 2021. Accepted: December 14, 2021

© The Author(s) 2021. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited. Please note, elements of this work were written by an employee of the US Government.



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.^{10,11}

information for exposomics was available, what can we as a community achieve?

Authors need guidance to properly and uniformly capture and report chemical structure information and transformations, that is, connecting either endogenous or exogenous chemicals with their metabolites—thus helping capture the associated biological responses. The flexible templates provided here (see Sections “Chemical Structure Data” and “Transformations Data”) show how authors can consistently submit this information to the Exposome journal as [supplementary materials](#) with their articles. These templates are designed such that authors can include as much or as little information as is available, yet still contribute their knowledge and outcomes to the exposomics “pool” (and beyond) in an open and FAIR manner. The “Chemical Structure Data” template is identical to the template introduced recently in the Journal of Cheminformatics.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 public resources. Thus, author contributions, no matter how small, will gradually help build the bigger picture needed to unravel and comprehend the exposome. Before we launch into the template descriptions, a few definitions are covered in the next section.

Definitions

While “FAIR” and “Open” are used somewhat interchangeably in this article as we strongly believe that chemical data should be both where possible, there is a distinction that is particularly relevant for exposomics, as sensitive human data cannot necessarily be made open. Data can be “open” but not “FAIR,” and vice versa. Open science has many facets; of most relevance to this article is open access (OA). OA is a set of principles and a range of practices through which research outputs are distributed online, free of cost or other access barriers.¹³ The FAIR principles for digital assets, on the other hand, include guidance on how to make data

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

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

more Findable, Accessible, Interoperable, and Reusable.^{14,15} For example, if you have open data that is not findable, no one can use it; whereas if you have “FAIR” data that is not “open”, it is not available for integration into open community resources. Thus, the most powerful data are both open and FAIR.

In [Table 1](#), we provide some definitions of chemical and transformation terms used later in this article.

Templates for FAIR exposomics chemical data

Chemical structure data

Better consideration of chemical factors in the exposome requires high-quality chemical information in research articles. Many exposomics resources are based (mostly) on literature mining using name and synonym matching, which can be notoriously prone to errors. In this section, we provide some guidance on what information authors should consider providing, as well as the pros and cons of various choices. Since this Chemical Structure Data template was presented recently to the *Journal of Cheminformatics*,¹² some of the material in this section overlaps with the previous article.

Authors should consider submitting their chemical structure information with their manuscript as [Supplementary Material](#) using the suggested template as comma separated value (CSV; *.csv); or, alternatively, as tab-separated value (TSV; *.tsv) or structure data file (SDF; *.sdf) formats. These formats ensure maximum interoperability between resources and operating systems. The popular XLS(X) format is not truly interoperable (options to save as CSV or TSV are offered), while the extraction of information from PDF format is difficult without introducing errors. The content below describes the CSV/TSV formats, SDF instructions are available elsewhere¹⁷ (however, the SD fields should match the CSV/TSV headers). In our experience, so far CSV often proves most interoperable for the widest audience, although the other formats also have certain advantages.

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

The author-submitted template file¹⁸ should contain at least one of the following columns: *SMILES*, *InChI*, *Name*, or *InChIKey*. The *Name* column corresponds to a single primary name for the chemical structure. Each *Synonym* column corresponds to an additional chemical name (one name entry per column). Each *Comment* column can be added to provide additional text that may be important to the downstream user. Authors can also provide additional CSV/TSV columns (or SDF SD fields) containing information about their chemical substances (with unique, descriptive headers) for additional context. Chemical database identifiers or registry numbers could be included in this manner (as additional columns or fields), or as a *Synonym*. Note that chemical records indicating chemical structure with only *InChIKey* or *Name* will not contain sufficient information to describe a chemical structure; and *can only be mapped to existing entries* in destination resources. Batch services are available (e.g., from PubChem^{7,20} or CompTox^{9,21}) for authors to add, e.g., *SMILES* and/or *InChI* to their records, based upon the *Name* or other identifiers.

[Figure 1](#) in Schymanski and Bolton¹¹ shows the template file, which is available for download¹⁸ and as [Supplementary Material](#) with this article. [Figure 2](#) shows an example submission according to the proposed template, created by sub-setting the “HSDBTPS” dataset of literature-mined and curated transformation products from the Hazardous Substance Data Bank (HSDB) in PubChem.^{19,22} This example provides the *Name*, *SMILES*, and *InChIKey* fields as suggested, and an identifier (the PubChem Compound Identifier, CID) as an additional (optional) column (*PubChem_CID*) with a unique and easily recognizable header that can be processed by other resources as they choose, helping with interoperability.

Transformations data

The advancement of modern science is data driven.^{23,24} Providing key data in a ready to use format helps to assist in its reuse in research articles, regulatory reports, or machine-learning data models. Exposomics especially needs access to ready-to-use, high-quality chemical information from individual research articles (e.g., such as the connection of detected chemicals with the disease endpoint investigated or the aggregation of known metabolites of thousands of common chemicals). For instance, HSDB contains metabolites and metabolism information for 3220 chemicals gathered over 40 years, but these are only available as text snippets that need to be matched to chemical structures by synonyms followed by manual curation (initial efforts have covered only 1/100th of this dataset²²). However, as mentioned above, a key challenge in exposomics is to connect chemicals (e.g., of anthropogenic origin, but also endogenous or exogenous chemicals) that are associated with exposures with their biological response. Since metabolism is the most dynamic of the biological responses, and metabolites per definition fall into the

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	IWEDIXLBFLEXBO-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)CC)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

Figure 2. An example chemical structure data file constructed according to the proposed template¹⁸ by taking a subset of the HSDBTPS structure data.¹⁹ Image created in RStudio (Version 1.2.5042). The HSDBTPS efforts resulted in the deposition of five new structures to PubChem all documented in HSDB text snippets, CIDs 146035700, 146035701, 146035702, 146035703, and 146037633.

same molecular mass category as many anthropogenic chemicals of concern, a key gap in exposomics knowledge is the connection between chemicals and their metabolites. The efforts of many will be needed to help fill this knowledge gap and the timing could not be better for exposomics with several recent studies emerging using *in vitro* enzymes to investigate parent–metabolite relationships of drugs and other relevant chemicals.^{25,26}

The Transformations template provided here has been designed on the basis of recent efforts to fill the gaps of transformation products in PubChem using literature data,²⁷ in collaboration with the NORMAN Suspect List Exchange (NORMAN-SLE).^{28–30} Several datasets from a variety of sources have now been processed. Transformations from the NORMAN-SLE, where S## refers to the list number, followed by the list code, include: S60 SWISSPEST19,^{31,32} S66 EAWAGTPS,^{33,34} S68 HSDBTPS,^{21,22} S73 METXBIODB,^{35,36} S74 REFTPS,³⁷ S78 SLUPESTTPS,^{38,39} S79 UACCSCEC,^{40,41} and S81 THSTPS⁴² (list available from https://git-r3lab.uni.lu/eci/pubchem/-/raw/master/annotations/tps/Transformation_Datasets.txt). Of these, MetXBioDB also contains enzyme information, while the rest are primarily environmental data. Figure 3 shows an example “environmental” dataset compiled from several of these lists, using the proposed template. In addition to the NORMAN-SLE datasets, a dataset of more than 1200 transformations from ChEMBL⁴³ has also been added, including enzyme, gene, and protein information (where available). An example of Transformations with more biological information available is given in Figure 4.

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

If available, a brief description of the transformation is useful and can be provided in the “Transformation” field (top panel, Figures 3 and 4). Short, informative descriptions are preferred; the current entries have been either extracted automatically from existing datasets or entered manually. In the future, it may

be possible to provide some guidance via an ontology as the public dataset grows to improve the machine readability. Similarly, if information on the biosystem is available (i.e., where the transformation takes place), this can be included in the Biosystem column (see Figures 3 and 4 for examples).

For datasets with biological information, this can be provided (optionally) in the Enzyme, Gene_ID, and Protein_ID columns. At this stage, the template allows flexible input (see Figure 4 for examples) but recommend Enzyme are provided as either: Enzyme Commission (EC) number,^{45–47} such as “EC 2.3.2.23”; gene symbol, such as “CYP1A1”; or as enzyme names, such as “Aryl hydrocarbon hydroxylase.” The Gene_ID is expected to be an NCBI Gene⁴⁸ ID, such as “1543.” The Protein_ID is expected to be either an NCBI Protein⁴⁹ accession, such as “NP_059488.2” or an UniProt identifier,⁵⁰ such as “P08684.” If multiple entries for Enzyme, Gene_ID, and Protein_ID are provided, they should be separated by a “pipe” symbol (“|”) or provided as new rows.

Finally, the Reference_ID and Reference_Description columns provide the opportunity to credit the original sources of the information. Reference_ID entries should be either PubMed identifiers⁵¹ (PMIDs) or Digital Object Identifiers⁵² (DOIs), preceded with “PMID:” or “DOI:”, respectively, for easy recognition, and separated by a “pipe” (“|”) if multiple IDs exist (they can be mixed—for example, “PMID:33929905|DOI:10.1186/s13321-018-0324-5”). The Reference_Description can be used to provide a free text form of the reference, to describe the data source (if no PMID/DOI available) or to describe evidence of the transformation. Only Reference_ID can be processed automatically. Again, see Figures 3 and 4 and the Transformations template⁴⁴ for examples.

So far, about 6000 Transformations have been processed using these templates, from nine different sources (many of these being composite data from several sources themselves, including ChEMBL,⁴³ MetXBioDB,³⁵ and REFTPS³⁷). The Transformations are being integrated into current computational mass spectrometry workflows (such as patRoos⁵³ and as documented in Krier et al.²²) and are openly available for all. The summarized files are likewise available for comprehensive efforts such as BioTransformer³⁶ to add this new data to their training set (MetXBioDB³⁵ is the library behind BioTransformer) and likewise improve predictions. Overall, FAIR transformations data will greatly support exposomics, and discussions to extend these templates into fields with formal ontologies and/or other formats such as mzTab^{54,55} in the future are welcomed. As demonstrated

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: ...

Figure 3. An example of various environmental transformations constructed according to the proposed Transformations template⁴⁴ (using Name and PubChem CID), taking a subset of transformations from NORMAN-SLE datasets (REFTPS,³⁷ HSDBTPS,²⁰ SLUPESTTPS,³⁸ EAWAGTPS,³³ and SWISSPEST19³¹). Image created in RStudio (Version 1.2.5042).

Predecessor_Name	Predecessor_SMILES	Successor_Name	Successor_SMILES	Transformation
Carbamazepine	<chem>C1=CC=C2C(=C1)C=CC3=CC=CC=C3N2C(=O)N</chem>	Carbamazepine-10,11-epoxide	<chem>C1=CC=C2C(=C1)C3C(O3)C4=CC=CC=C4N2C(=O)N</chem>	Epoxidation of 1,2-disubstituted alkene / Human Phase I
Acetaminophen	<chem>CC(=O)NC1=CC=C(C=C1)O</chem>	Acetaminophen glucuronide	<chem>CC(=O)NC1=CC=C(C=C1)O[C@@H]2[C@@H]([C@@H]([C@@H]2)C(=O)O)C(=O)O</chem>	Glucuronide conjugation at the hydroxyl forming an ether
Acrolein	<chem>C=CC=O</chem>	Acrylic Acid	<chem>C=CC(=O)O</chem>	
Furan	<chem>C1=COC=C1</chem>	(E)-2-Butenedial	<chem>C=C\C=C=O</chem>	Oxidation / Human Phase I
Benzene	<chem>C1=CC=CC=C1</chem>	Phenol	<chem>C1=CC=C(C=C1)O</chem>	Hydroxylation of aromatic carbon / Human Phase I
Nicotinamide	<chem>C1=CC(=CN=C1)C(=O)N</chem>	MNAM	<chem>C[N+](=O)C=CC(=O)N</chem>	

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-Ba...
	UDP-glucuronosyltransferase 1-6 UDP-glucuronosyltransfer...	UGT1A6 UGT1A1	P19224 P22309	PMID:23462933	Data from ChEMBL - IDs (predecessor successor enzyme): C...
	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...

Figure 4. An example of biological transformations constructed according to the proposed Transformations template⁴⁴ (using Name and SMILES), taking a subset of transformations from NORMAN-SLE dataset MetXBioDB³⁵ (from BioTransformer³⁶) and the ChEMBL⁴³ datasets on PubChem; both datasets have some degree of enzyme, gene, and/or protein information available.

in Figures 5 and 6, one can see the benefits of arranging data in FAIR templates. Figure 5 is an example of a resulting Transformation entry in PubChem, while Figure 6 can be created automatically in CDK Depict using simple code in R to create annotated reaction SMILES from the fields shown in Figure 3 only.

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, and Reusable). When authors provide this content as Supplementary Material, 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 FAIRness of Supplementary Material greatly

8.11 Transformations



7 items View More Rows & Details

Download

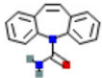
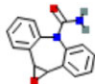
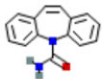
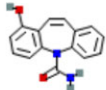
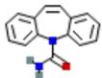
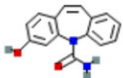
SORT BY ⌵ Please Choose One ⌵					
Predecessor	Predecessor Name	Successor	Successor Name	Transformation	Enzyme
	carbamazepine		Carbamazepine 10,11-epoxide	Epoxidation of 1,2-disubstituted alkene / Human Phase I	CYP3A4 CYP2C8
	carbamazepine		9-Hydroxycarbamazepine	Aromatic hydroxylation of fused benzene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

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

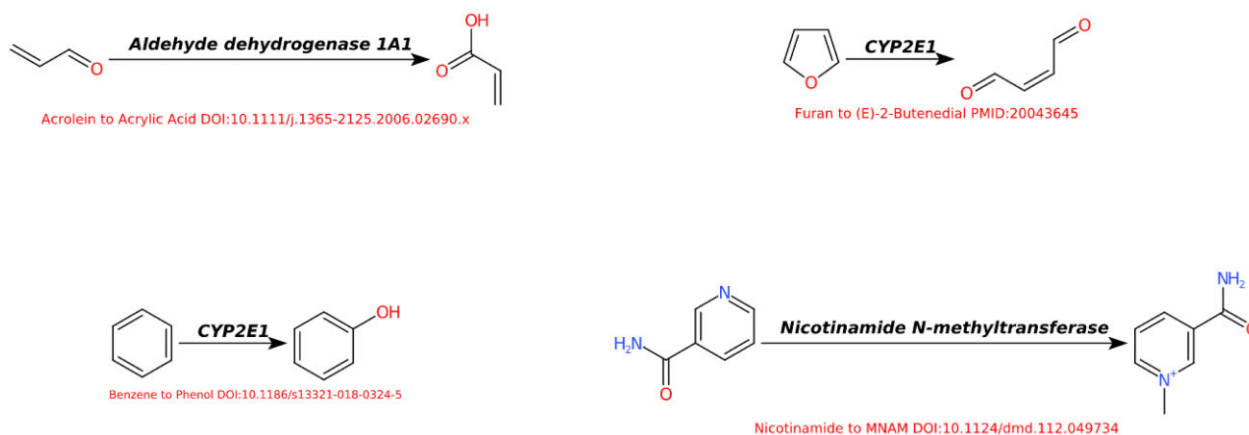


Figure 6. Example reactions corresponding with the last four rows of Figure 3, automatically created and depicted with CDK Depict⁵⁶ (<https://www.simolecule.com/cdkdepict/depict.html>) directly from template content shown in Figure 3 (SMILES, Name, Enzyme, and Reference_ID fields).

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.”

Acknowledgments

We gratefully acknowledge discussions with the entire PubChem team, especially Jian (Jeff) Zhang and Tiejun Cheng for their joint work on the transformations, as well as Ben Shoemaker, Paul Thiessen, Siqian He, and Asta Gindulyte. We also gratefully acknowledge discussions with Egon Willighagen and the editorial team at the Journal of Cheminformatics (surrounding the lead-up article to this article), and many collaborators who have worked on depositions within PubChem and the NORMAN-SLE. Special mentions go to Frank Menger (SLU, Sweden) and Lidia Belova (University of Antwerp, Belgium), for testing and depositing data using earlier versions of the transformations template

(SLUPESTPS and UACCSCEC, respectively). We are also grateful to Anca Baesu (McGill University, Canada) and Parviel Chirsir (University of Luxembourg), as well as Noelia Ramirez and colleagues (URV, Tarragona, Spain) for their testing and contributions using the existing templates (REFTPS and THSTPS, respectively).

Supplementary materials

Supplementary material is available at *Exposome* online. 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.

Funding

E.E.B. is funded by the Intramural Research Program of the National Library of Medicine, National Institutes of Health. E.L.S. acknowledges funding support from the Luxembourg National Research Fund (FNR) for project A18/BM/12341006.

Conflict of interest statement

The authors declare no competing interests.

References

- Wild CP. Complementing the genome with an “exposome”: The outstanding challenge of environmental exposure measurement in molecular epidemiology. *Cancer Epidemiol Biomarkers Prev* 2005;14:1847–1850. doi:10.1158/1055-9965.EPI-05-0456
- Miller GW. Exposome: A new field, a new journal. *Exposome* 2021;1. doi:10.1093/exposome/osab001
- 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>
- 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>
- Vermeulen R, Schymanski EL, Barabási AL, Miller GW. The exposome and health: Where chemistry meets biology. *Science*. 2020;367:392–396. doi:10.1126/science.aay3164
- Miller GW, Jones DP. The nature of nurture: Refining the definition of the exposome. *Toxicol Sci*. 2014;137:1–2. doi:10.1093/toxsci/kft251
- Kim S, Chen J, Cheng T, et al. PubChem 2019 update: Improved access to chemical data. *Nucleic Acids Res*. 2019;47:D1102–D1109. doi:10.1093/nar/gky1033
- Davis AP, Grondin CJ, Johnson RJ, et al. Comparative toxicogenomics database (CTD): Update 2021. *Nucleic Acids Res*. 2021;49:D1138–D1143. doi:10.1093/nar/gkaa891
- Williams AJ, Grulke CM, Edwards J, et al. The CompTox chemistry dashboard: A community data resource for environmental chemistry. *J Cheminform*. 2017;9:61. doi:10.1186/s13321-017-0247-6
- 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 Process Impacts*. 2019;21:1426–1445. doi:10.1039/C9EM00068B
- Baker NC, Schymanski EL, Williams AJ. Literature neurotoxins: Excel Macro File. FigShare doi:10.23645/epacomptox.7334603
- Schymanski EL, Bolton EE. FAIR chemical structures in the Journal of Cheminformatics. *J Cheminform*. 2021;13:50. doi:10.1186/s13321-021-00520-4
- Peter Suber. *Open Access Overview (Definition, Introduction)*. Accessed July 3, 2021. <http://legacy.earlham.edu/~peters/fos/overview.htm>
- GO FAIR. *FAIR Principles*. Published 2021. Accessed March 23, 2021. <https://www.go-fair.org/fair-principles/>
- Wilkinson MD, Dumontier M, Aalbersberg IJ, et al. Comment: The FAIR guiding principles for scientific data management and stewardship. *Sci Data*. 2016;3:1–9. doi:10.1038/sdata.2016.18
- Heller S, McNaught A, Stein S, Tchekhovskoi D, Pletnev I. InChI—the worldwide chemical structure identifier standard. *J Cheminform*. 2013;5:7. doi:10.1186/1758-2946-5-7
- NCBI/NLM/NIH. *PubChem Documentation: Substance SDF Submission*. Published 2021. Accessed March 23, 2021. https://pubchem.ncbi.nlm.nih.gov/upload/docs/examples/substance_submission.sdf
- NCBI/NLM/NIH. *Chemical Structure Data Template (CSV)*. Published 2021. Accessed May 9, 2021. https://ftp.ncbi.nlm.nih.gov/pubchem/Other/Submissions/Chemical_Structure_Data_Template.csv
- LCSB-ECI, Krier J, Schymanski E, et al. S68 | HSDBTPS | Transformation Products Extracted from HSDB Content in PubChem. Published online June 11, 2020. doi:10.5281/zenodo.3827487
- NCBI/NLM/NIH. *PubChem Identifier Exchange*. Published 2021. Accessed March 23, 2021. <https://pubchem.ncbi.nlm.nih.gov/idxchange/idxchange.cgi>
- United States Environmental Protection Agency. *CompTox Batch Search*. Published 2021. Accessed March 23, 2021. https://comptox.epa.gov/dashboard/dsstoxdb/batch_search
- Krier J, Singh RR, Kondić T, et al. Discovering pesticides and their TPs in Luxembourg waters using open cheminformatics approaches. *Environ Int*. 2022;158:106885. doi:10.1016/j.envint.2021.106885
- Montáns FJ, Chinesta F, Gómez-Bombarelli R, Kutz JN. Data-driven modeling and learning in science and engineering. *Comptes Rendus Mécanique*. 2019;347:845–855. doi:10.1016/j.crme.2019.11.009
- Clauset A, Larremore DB, Sinatra R. Data-driven predictions in the science of science. *Science*. 2017;355:477–480. doi:10.1126/science.aal4217
- Liu K, Lee C, Singer G, et al. Enzyme-based chemical identification for metabolomics. *FASEB J*. 2021;35:fasebj.2021.35.S1.04277. doi:10.1096/fasebj.2021.35.S1.04277
- Ross DH, Seguin RP, Krinsky AM, Xu L. High-throughput measurement and machine learning-based prediction of collision cross sections for drugs and drug metabolites. *Bioinformatics* 2021. doi:10.1101/2021.05.13.443945
- Schymanski EL, Kondić T, Neumann S, Thiessen PA, Zhang J, Bolton EE. Large chemical knowledge bases for exposomics: PubChemLite meets MetFrag. *J Cheminform*. 2021;13:19. doi:10.1186/s13321-021-00489-0

28. NORMAN Network. NORMAN Suspect List Exchange. Accessed June 9, 2019. <https://www.norman-network.com/nds/SLE/>
29. NORMAN Network. NORMAN Suspect List Exchange on Zenodo. NORMAN Suspect List Exchange: Zenodo Community. Accessed June 9, 2019. <https://zenodo.org/communities/norman-sle/>
30. NORMAN Network, NCBI/NLM/NIH. NORMAN SLE Classification Browser. Accessed May 7, 2020. <https://pubchem.ncbi.nlm.nih.gov/classification/#hid=101>
31. Kiefer K, Müller A, Singer H, Hollender J. S60 | SWISSEST19 | Swiss Pesticides and Metabolites from Kiefer 2019. Published online November 17, 2019. <http://doi.org/10.5281/zenodo.3544760>
32. Kiefer K, Müller A, Singer H, Hollender J. New relevant pesticide transformation products in groundwater detected using target and suspect screening for agricultural and urban micropollutants with LC-HRMS. *Water Res.* 2019;165:114972. doi:10.1016/j.watres.2019.114972
33. Schollee J, Schymanski E. S66 | EAWAGTPS | Parent-Transformation Product Pairs from Eawag. Published online April 23, 2020. doi:10.5281/zenodo.3754448
34. Schollée JE, Schymanski EL, Stravs MA, Gulde R, Thomaidis NS, Hollender J. Similarity of high-resolution tandem mass spectrometry spectra of structurally related micropollutants and transformation products. *J Am Soc Mass Spectrom.* 2017;28:2692–2704. doi:10.1007/s13361-017-1797-6
35. Djoumbou-Feunang Y, Schymanski E, Zhang J, Wishart DS. S73 | METXBIODB | Metabolite Reaction Database from BioTransformer. Published online November 5, 2020. doi:10.5281/zenodo.4056560
36. Djoumbou-Feunang Y, Fiamoncini J, Gil-de-la-Fuente A, Greiner R, Manach C, Wishart DS. BioTransformer: A comprehensive computational tool for small molecule metabolism prediction and metabolite identification. *J Cheminform.* 2019;11:2. doi:10.1186/s13321-018-0324-5
37. Schymanski E. S74 | REFTPS | Transformation products and reactions from literature. Published online December 2020. doi:10.5281/zenodo.4318838
38. Menger F, Boström G. S78 | SLUPESTTPS | Pesticides and TPs from SLU, Sweden. Published online May 10, 2021. doi:10.5281/zenodo.4687924
39. Menger F, Boström G, Jonsson O, et al. Identification of pesticide transformation products in surface water using suspect screening combined with national monitoring data. *Environ Sci Technol.* 2021;55:10343–10353. doi:10.1021/acs.est.1c00466
40. Belova L, Caballero-Casero N, van Nuijs ALN, Covaci A. Ion mobility-high-resolution mass spectrometry (IM-HRMS) for the analysis of contaminants of emerging concern (CECs): Database compilation and application to urine samples. *Anal Chem.* 2021;93:6428–6436. doi:10.1021/acs.analchem.1c00142
41. Belova L, Caballero-Casero N, van Nuijs Alexander LN, Covaci A. S79 | UACCSECC | Collision Cross Section (CCS) Library from UAntwerp. Published online May 10, 2021. doi:10.5281/zenodo.4704648
42. Merino C, Vinaixa M, Ramirez N. S81 | THSTPS | Thirdhand Smoke Specific Metabolites. Published online September 2, 2021. doi:10.5281/zenodo.5394629
43. Gaulton A, Hersey A, Nowotka M, et al. The ChEMBL database in 2017. *Nucleic Acids Res.* 2017;45:D945–D954. doi:10.1093/nar/gkw1074
44. NCBI/NLM/NIH. Transformations data template (CSV). Published 2021. Accessed May 25, 2021. https://ftp.ncbi.nlm.nih.gov/pubchem/Other/Submissions/Transformations_Template.csv
45. McDonald AG, Boyce S, Tipton KF. ExplorEnz: The primary source of the IUBMB enzyme list. *Nucleic Acids Res.* 2009;37(Database issue):D593–D597. doi:10.1093/nar/gkn582
46. Bairoch A. The ENZYME database in 2000. *Nucleic Acids Res.* 2000;28:304–305. doi:10.1093/nar/28.1.304
47. Chang A, Jeske L, Ulbrich S, et al. BRENDA, the ELIXIR core data resource in 2021: New developments and updates. *Nucleic Acids Res.* 2021;49:D498–D508. doi:10.1093/nar/gkaa1025
48. Brown GR, Hem V, Katz KS, et al. Gene: A gene-centered information resource at NCBI. *Nucleic Acids Res.* 2015;43(Database issue):D36–D42. doi:10.1093/nar/gku1055
49. Benson DA, Cavanaugh M, Clark K, et al. GenBank. *Nucleic Acids Res.* 2018;46(D1):D41–D47. doi:10.1093/nar/gkx1094
50. The UniProt Consortium. UniProt: The universal protein knowledgebase. *Nucleic Acids Res.* 2017;45(D1):D158–D169. doi:10.1093/nar/gkw1099
51. Sayers EW, Beck J, Bolton EE, et al. Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.* 2021;49(D1):D10–D17. doi:10.1093/nar/gkaa892
52. International DOI Foundation. *Frequently Asked Questions about the DOI® System.* Accessed September 7, 2021. <https://www.doi.org/faq.html>
53. Helmus R, ter Laak TL, van Wezel AP, de Voogt P, Schymanski EL. patRoon: Open source software platform for environmental mass spectrometry based non-target screening. *J Cheminform.* 2021;13(1). doi:10.1186/1020-00477-w
54. Griss J, Jones AR, Sachsenberg T, et al. The mzTab data exchange format: Communicating mass-spectrometry-based proteomics and metabolomics experimental results to a wider audience. *Mol Cell Proteomics.* 2014;13:2765–2775. doi:10.1074/mcp.O113.036681
55. Hoffmann N, Rein J, Sachsenberg T, et al. mzTab-M: A data standard for sharing quantitative results in mass spectrometry metabolomics. *Anal Chem.* 2019;91:3302–3310. doi:10.1021/acs.analchem.8b04310
56. Mayfield J. CDK Depict Web Interface. Accessed December 6, 2021. <https://simolecule.com/cdkdepict/depict.html>
57. NCBI/NLM/NIH. *PubChem Submissions Template Folder.* Published 2021. Accessed May 25, 2021. <https://ftp.ncbi.nlm.nih.gov/pubchem/Other/Submissions/>