

1 PubChemLite plus Collision Cross Section (CCS) values for 2 enhanced interpretation of non-target environmental data

3 Anjana Elapavalore^a, Dylan H. Ross^{b,c}, Valentin Grouès^a, Dagny Aurich^a,
4 Allison M. Krinsky^b, Sunghwan Kim^d, Paul A. Thiessen^d, Jian Zhang^d, James N. Dodds^e,
5 Erin S. Baker^e, Evan E. Bolton^{d*}, Libin Xu^{b*}, Emma L. Schymanski^{a*}

6 ^a Luxembourg Centre for Systems Biomedicine (LCSB), University of Luxembourg, 6 Avenue du Swing,
7 4367, Belvaux, Luxembourg. ORCID^s AE: [0000-0002-0295-6618](https://orcid.org/0000-0002-0295-6618); VG: [0000-0001-6501-0806](https://orcid.org/0000-0001-6501-0806) DA: [0000-0001-8823-0596](https://orcid.org/0000-0001-8823-0596); ELS: [0000-0001-6868-8145](https://orcid.org/0000-0001-6868-8145).

9 ^b Department of Medicinal Chemistry, University of Washington, Seattle, Washington 98195, United
10 States. ORCID^s DHR: [0009-0005-2943-2282](https://orcid.org/0009-0005-2943-2282); AMK: NA; LX: [0000-0003-1021-5200](https://orcid.org/0000-0003-1021-5200).

11 ^c Current Address: Biological Sciences Division, Pacific Northwest National Laboratory, Richland, WA,
12 USA.

13 ^d National Center for Biotechnology Information (NCBI), National Library of Medicine (NLM), National
14 Institutes of Health (NIH), Bethesda, MD, 20894, USA. ORCID^s SK: [0000-0001-9828-2074](https://orcid.org/0000-0001-9828-2074); PAT: [0000-0002-1992-2086](https://orcid.org/0000-0002-1992-2086); JZ: [0000-0002-6192-4632](https://orcid.org/0000-0002-6192-4632); EEB: [0000-0002-5959-6190](https://orcid.org/0000-0002-5959-6190).

16 ^e Department of Chemistry, University of North Carolina, Chapel Hill, North Carolina 27599, USA. ORCID^s
17 JND: [0000-0002-9702-2294](https://orcid.org/0000-0002-9702-2294) ESB: [0000-0001-5246-2213](https://orcid.org/0000-0001-5246-2213)

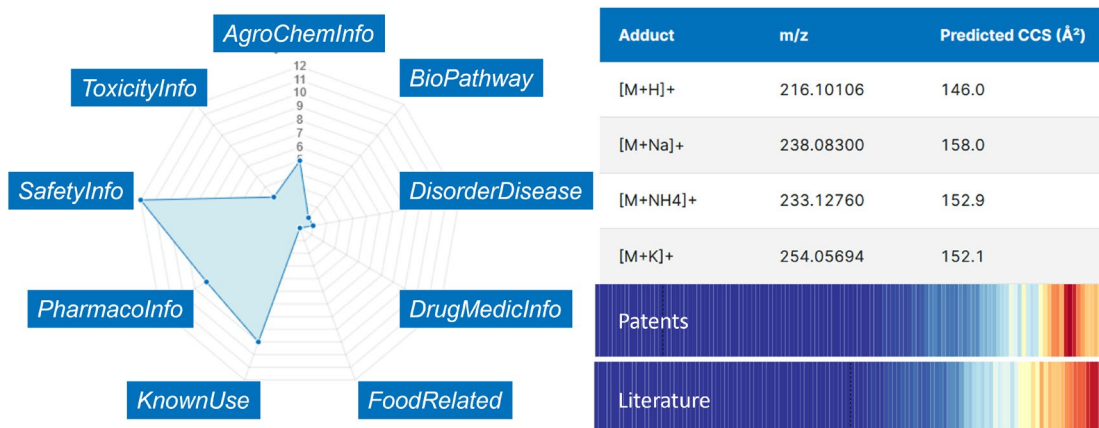
18 *Contact: EEB: bolton@ncbi.nlm.nih.gov, LX: libinxu@uw.edu and ELS: emma.schymanski@uni.lu

19 Abstract

20 Finding relevant chemicals in the vast (known) chemical space is a major challenge for environmental
21 and exposomics studies leveraging non-target high resolution mass spectrometry (NT-HRMS) methods.
22 Chemical databases now contain hundreds of millions of chemicals, yet many are not relevant. This
23 article details an extensive collaborative, open science effort to provide a dynamic collection of
24 chemicals for environmental, metabolomics and exposomics research, along with supporting
25 information about their relevance to assist researchers in the interpretation of candidate hits. The
26 PubChemLite for Exposomics collection is compiled from ten annotation categories within PubChem,
27 enhanced with patent, literature and annotation counts, predicted partition coefficient (logP) values, as
28 well as predicted collision cross section (CCS) values using CCSbase. Monthly versions are archived on
29 Zenodo under a CC-BY license, supporting reproducible research, and a new interface has been
30 developed, including the chemical stripes on patent and literature data, for researchers to browse the
31 collection. This article further describes how PubChemLite can support researchers in
32 environmental/exposomics studies, describes efforts to increase the availability of experimental CCS
33 values, and explores known limitations and potential for future developments. The data and code
34 behind these efforts are openly available. PubChemLite content can be explored at
35 <https://pubchemlite.lcsb.uni.lu>.

36 **Keywords:** non-target screening; identification; PubChemLite; exposomics; ion mobility; collision cross
37 section; PubChem

38 **Synopsis:** PubChemLite empowers environmental non-target screening data interpretation by
39 combining annotation content, patent, literature, logP and predicted collision cross section data.



40

41 *Table of Contents Graphic*

42 Introduction

43 Environmental and exposomics researchers are faced with the daunting task of determining which
44 chemicals, among other factors, may be either potentially detrimental or beneficial in the context of
45 human and/or environmental health. Non-target screening (NTS) methods leveraging high resolution
46 mass spectrometry (HRMS) approaches are now commonly used to explore complex samples due to
47 high sensitivity and selectivity plus improved availability of HRMS instruments^{1,2}. Ion mobility
48 spectrometry (IMS), which separates molecules based on their size and shape, is also increasingly
49 accessible. The calculated collision cross section (CCS) values serve as an additional parameter to
50 support identification in NTS^{3,4,5}. Nonetheless, the identification and - importantly - interpretation of
51 features detected during NTS is still challenging, while integration of IMS/CCS into workflows remains
52 poor, hindering the broader adoption of NTS¹. Identification in NTS primarily relies on mass spectral
53 libraries, suspect lists (lists of hundreds or thousands of chemicals that may occur in the samples) and
54 chemical databases, as recently reviewed elsewhere^{1,2}.

55 A diverse array of compound databases, often openly accessible, serve as primary sources of candidates
56 for identification in NTS. These range from hundreds of thousands to just over a million entries (*e.g.*,
57 HMDB⁶ with 220,945 metabolites and CompTox⁷ with 1,218,248 chemicals as of 20 Nov. 2024), through
58 to hundreds of millions of entries (*e.g.*, ChemSpider⁸, PubChem⁹ and the Chemical Abstract Services
59 (CAS) Registry¹⁰, with 129, 119 and 219 million chemicals as of 20 Nov. 2024, respectively). Since the CAS
60 Registry is licensed and ChemSpider introduced limitations to their application programming interface
61 (API) in 2018, PubChem has become the *de facto* standard large chemical database for open science-
62 based NTS methods. While PubChem, with >1000 sources, integrates the contents of many of the
63 smaller openly available databases, PubChem also includes tens of millions of entries that are neither
64 likely to be found in the environment, nor pertinent to the exposome. This hinders both the
65 performance and efficiency of NTS. Additionally, many of the chemicals in other potential sources for
66 NTS identification efforts, such as the Global Chemical Inventory (350,000 chemicals)¹¹ and various lists

67 from European regulators contributing to the NORMAN Suspect List Exchange (NORMAN-SLE)¹² include
68 large proportions of chemicals that have very little supporting evidence about their existence and
69 relevance, which makes interpretation of potential hits in NTS very challenging.

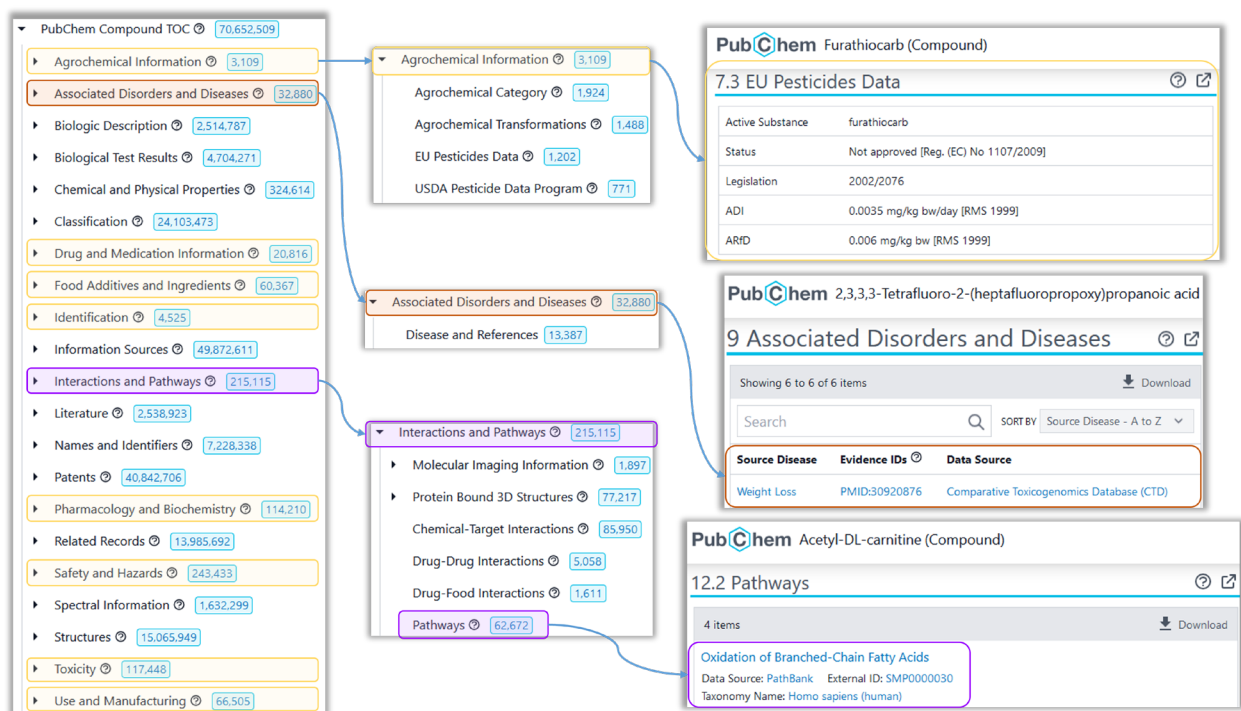
70 To mitigate these challenges, a subset of PubChem called PubChemLite was developed specifically to
71 streamline NTS identification and interpretation¹³. PubChemLite has been integrated into existing HR-
72 MS workflows, such as patRoom¹⁴ and MetFrag¹⁵. Although PubChemLite is familiar to many researchers
73 already, the original 2021 article¹³ was primarily technical. This article explains PubChemLite for an
74 environmental/exposomics audience and details extensions since the original publication, including the
75 development of an open experimental CCS pipeline in PubChem, integration of predicted CCS values in
76 PubChemLite to support IMS, and finally a new web interface (<https://pubchemlite.lcsb.uni.lu/>).

77 **Materials and Methods**

78 **Building PubChemLite**

79 The full technical details of PubChemLite are published elsewhere¹³. Briefly, PubChemLite is derived
80 from major categories relevant to environmental/exposomics applications appearing in the PubChem
81 Table of Contents (TOC) pages of the PubChem database¹⁶. The ten categories currently used to compile
82 PubChemLite (see Figure 1) are: Agrochemical Information (AgroChemInfo), Associated Disorders and
83 Diseases (DisorderDisease), Drug and Medication Information (DrugMedicInfo), Food Additives and
84 Ingredients (FoodRelated), Identification (Identification), Interactions and Pathways – Pathways subset
85 (BioPathway), Pharmacology and Biochemistry (PharmacolInfo), Safety and Hazards (SafetyInfo), Toxicity
86 (ToxicityInfo), Use and Manufacturing (KnownUse). These categories have remained consistent since the
87 original publication, except for the “Biomolecular Interactions and Pathways” category, which was
88 renamed by PubChem to “Interactions and Pathways” in 2022, then limited to the Pathways subset in
89 May 2023 (see Results and Discussion). The [input files](#) and code for the PubChemLite [build system](#) are
90 available on the Environmental Cheminformatics (ECI) [GitLab](#) repository (see Data Availability
91 Statement).

92 Any compound with one or more of the selected annotation categories is included. The matching
93 compounds (represented by PubChem Compound IDentifiers, CIDs) are aggregated by the first block of
94 the InChIKey into a primary entry and related CIDs, where the primary entry is the neutral or “parent”
95 form. Entries such as mixtures, disconnected substances, or those causing errors, such as some
96 transition metals, are excluded (see the [build system](#) code for details¹³). Chemical information (SMILES,
97 InChI, InChIKey, formula, mass), patent and literature (PubMed) counts plus predicted XlogP values are
98 retrieved in bulk using PubChem APIs. The chemical identifiers, mass and XlogP values correspond to the
99 “parent” (primary entry), while the annotation, patent and literature counts are aggregated across all
100 related CIDs. Importantly, the presence of an entry in PubChemLite means that at least one of these
101 annotation categories is available for each CID, with the resulting information publicly available on
102 PubChem to help interpret the relevance of the candidate, see Figure 1. Since the chemical content of
103 PubChem changes daily and annotation content weekly, PubChemLite is built and evaluated early Friday
104 mornings, following the weekly PubChem update cycle. New versions are currently released publicly
105 (typically last Friday of the month) on Zenodo (DOI: [10.5281/zenodo.5995885](https://doi.org/10.5281/zenodo.5995885) redirects to the latest
106 version).



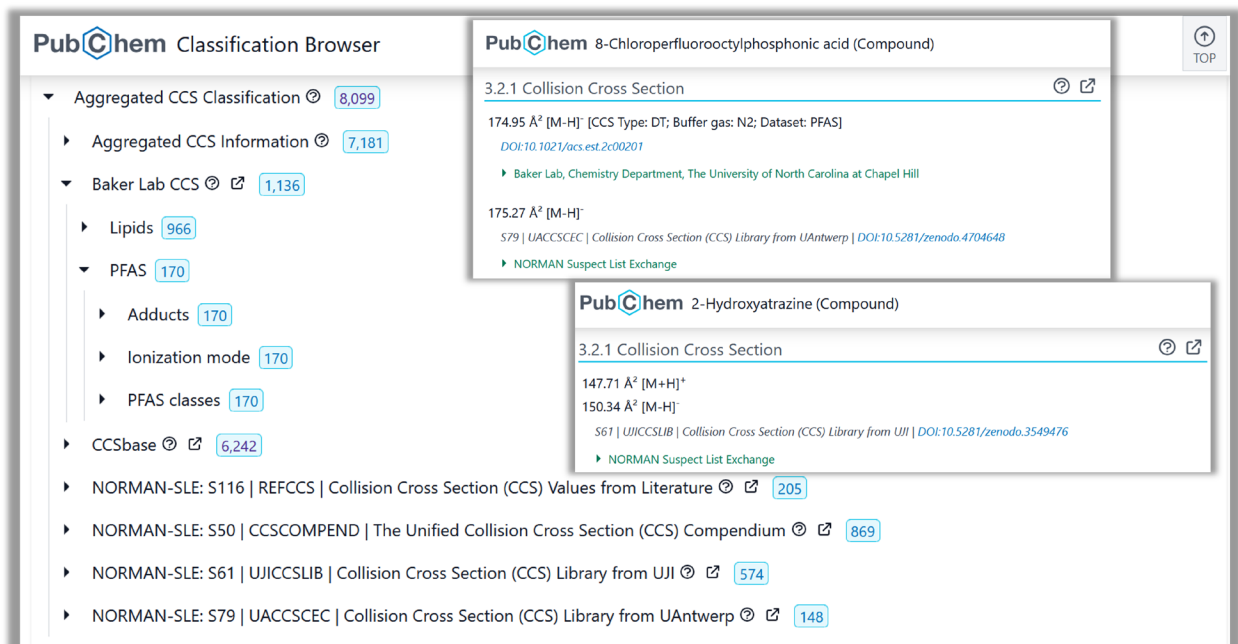
107

108 *Figure 1: PubChemLite categories in the PubChem Table of Contents (TOC), selected subcategories and*
 109 *associated annotation examples. Yellow shading denotes “environmental” categories (example CID*
 110 *47759), red the “exposomics” (example CID 114481) and purple the “metabolomics” sections (example*
 111 *CID 1).* For high resolution live images, please click the embedded hyperlinks.

112 Adding CCS Values

113 Although several methods to predict CCS values are available, few of these are suitable for a large
 114 database like PubChem, or even PubChemLite. The calculation time of quantum methods (e.g., MobCal¹⁷
 115 and ISICLE¹⁸) is prohibitive. Furthermore, since ISICLE only predicts values for C, H, N, O, P and S-
 116 containing compounds, ~40% of PubChemLite would remain uncovered. Machine learning (ML)
 117 prediction methods developed on experimental CCS datasets are much faster, with options including
 118 AllCCS¹⁹, CCSbase²⁰, SigmaCCS²¹, DeepCCS²², and CCS Predictor 2.0²³. Since initial calculations with
 119 CCSbase were promising, this became the method of choice to establish pipelines for PubChemLite.
 120 Calculations are performed using cs3db (the code base behind CCSbase²⁰), which has been re-trained
 121 with the updated experimental database in CCSbase using the same methodology²⁰, modified to run on
 122 PubChem internal systems. Following the monthly PubChemLite release, CCS value calculations (by
 123 PubChem) are triggered on the second day of the following month. The resulting file is transferred by
 124 FTP, merged into the PubChemLite files, quality-controlled, then released to Zenodo (DOI:
 125 [10.5281/zenodo.4081056](https://doi.org/10.5281/zenodo.4081056) redirects to the latest version).

126 To ensure that ML models have better coverage of environmentally relevant compounds to improve
 127 their predictions, part of this work involved establishing a pipeline to integrate experimental CCS values
 128 into PubChem. Currently PubChem contains CCS values from the Baker Lab^{24,25}, CCSbase²⁰ and four
 129 collections via the NORMAN-SLE¹²: S50 CCSCOMPEND^{26,27}, S61 UJICSLIB^{28,29}, S79 UACCSCEC^{3,30} and S116
 130 REFCCS³¹. These values are displayed on individual records in PubChem and navigable in the PubChem
 131 Classification Browser via the CCSbase, Baker Lab, NORMAN-SLE and the Aggregated CCS trees (see
 132 Figure 2).



133
 134 *Figure 2: Aggregated Collision Cross Section (CCS) Classification Tree in PubChem. Inset: Experimental*
 135 *CCS values in individual PubChem compound records for CI-PFOPA (CID 138395139) and the*
 136 *transformation product 2-hydroxyatrazine (CID 135398733). For high resolution live images, please click*
 137 *the embedded hyperlinks.*

138 This data can be retrieved from PubChem (code available on the [ECI GitLab](#)). The resulting compiled
 139 dataset is available on Zenodo (DOI: [10.5281/zenodo.6800138](https://doi.org/10.5281/zenodo.6800138) redirects to the latest version).

140 PubChemLite Web Interface

141 The PubChemLite web interface is developed as a plugin for the ELIXIR-Luxembourg [Data Catalog](#)^{32,33}. It
 142 is developed in Python, CSS, HTML and Javascript, using RDKit^{34,35} for structure depiction. For full details,
 143 see the [PubChemLite-web](#) code on GitLab.

144 The PubChemLite interface is based almost entirely on the information available in the archived
 145 PubChemLite-CCSbase CSV files (see Figure 3), except that additional synonyms (excluded from
 146 PubChemLite files for efficiency) are retrieved from the PubChem FTP site for better searchability.
 147 Additionally, records are enhanced with visualizations of the annotation categories and tables of the CCS
 148 and associated adduct mass values. Finally, the chemical stripes^{36,37,38} are included where available for
 149 both literature and patents (see Figure 4). The original chemical stripes R version was rewritten in
 150 Python for integration in [PubChemLite-web](#), with code available in both repositories^{38,39}.

PubChemLite EXPOSOMICS

Search

try [C10H14N2](#) [DUOANANYKXIQY-UHFFFAOYSA-N](#) [atrazine](#)
or [Explore](#)

Informative subset of PubChem relevant for various environmental, metabolomics, exposomics and mass spec. applications

Structural Information **2D Structure**

Molecular Formula
C₈H₁₄ClN₅

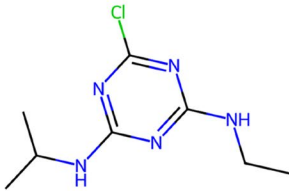
SMILES
CCNC1=NC(=NC(=N1)C)NC(C)C

InChI
InChI=1S/C8H14ClN5/c1-4-10-7-12-6(9)13-8(14-7)11-5(2)3/h5H,4H2,1-3H3,(H2,10,11,12,13,14)

InChIKey
MXWJVTQOROXGIU-UHFFFAOYSA-N

Compound name
6-chloro-4-N-ethyl-2-N-propan-2-yl-1,3,5-triazine-2,4-diamine

Related CIDs
[CID:2256](#) [CID:12306645](#) [CID:16213378](#)



8 Annotations Hits | 3616 References | 51422 Patents | 215.09378 Da Monoisotopic Mass | 2.6 XlogP (predicted)

151

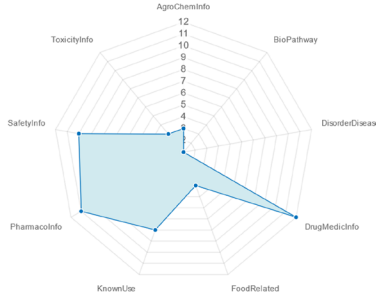
152 *Figure 3: PubChemLite web interface (composite image), compound view of Atrazine. For high resolution*
153 *live images, please click the embedded hyperlinks.*

PubChemLite EXPOSOMICS

Search

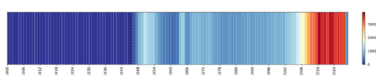
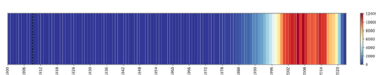
try [C10H14N2](#) [DUOANANYKXIQY-UHFFFAOYSA-N](#) [atrazine](#)
or [Explore](#)

Informative subset of PubChem relevant for various environmental, metabolomics, exposomics and mass spec. applications



Adduct	m/z	Predicted CCS (Å ²)
[M+H] ⁺	582.27298	228.3
[M+Na] ⁺	604.25492	229.2
[M+NH4] ⁺	599.29952	230.1
[M+K] ⁺	620.22886	228.0
[M-H] ⁻	580.25842	222.5
[M+Na-2H] ⁻	602.24037	243.9
[M] ⁺	581.26515	228.2
[M] ⁻	581.26625	228.2

m/z: mass to charge ratio of the adduct.
Predicted Collision Cross Section (CCS) values (Å²) per adduct calculated using [CCSbase](#).

Literature stripe  Patent stripe 

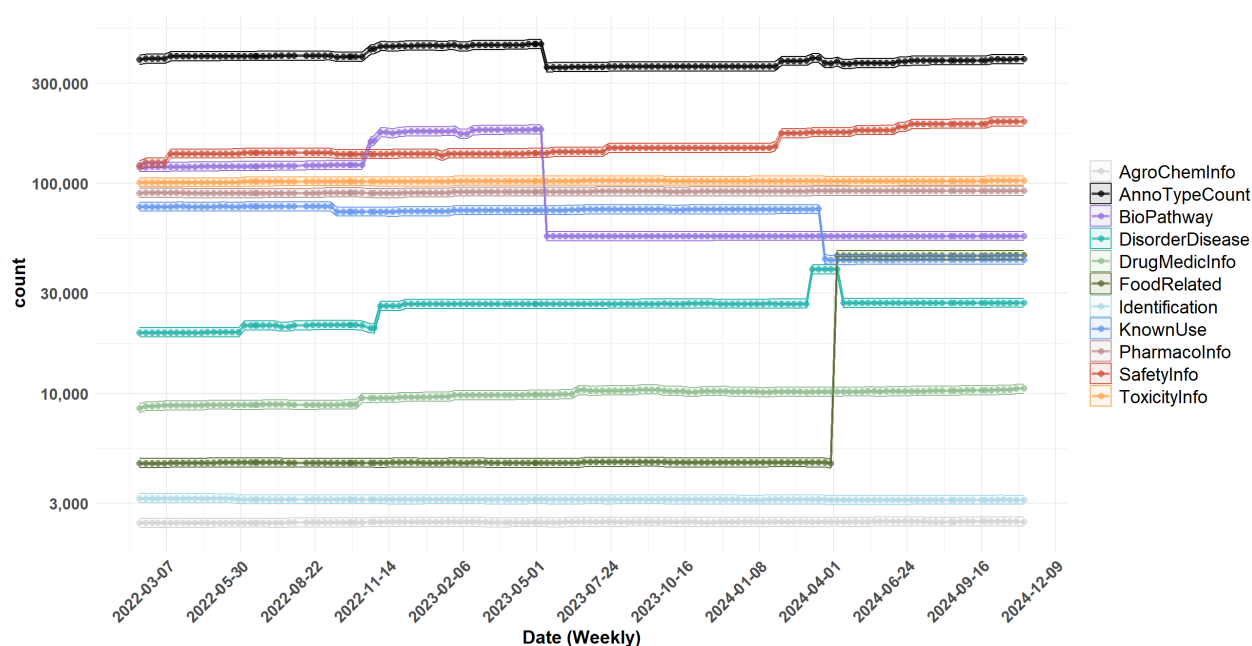
154

155 *Figure 4: PubChemLite web interface (composite image), view of additional data including annotations,*
156 *CCS values, patent and literature stripes for Streptomycin. For high resolution live images, please click*
157 *the embedded hyperlinks.*

158 Results and Discussion

159 PubChemLite Over Time

160 The performance of PubChemLite is monitored weekly with every build using the evaluation dataset of
161 977 compounds established in the original publication¹³. The ranking performance has been quite stable
162 over the three-year period, with median rank=1 of 794 (81.3%), 1-2 of 917 (93.9%), 1-5 of 960 (98.3%)
163 and 12 (1.2%) failures (compounds absent from PubChemLite due to lack of corresponding annotation).
164 The respective ranges [min;max] are rank=1 [788;800], 1-2 [909;922], 1-5 [955;963] and [10;15] failures.
165 The current performance is close to the median values: rank=1 of 797, 1-2 of 916, 1-5 of 960 and 11
166 failures and slightly better than the original publication (794, 912, 954, respectively, and 15 failures;
167 October 2020 version). The distribution of annotation content included in PubChemLite, including the
168 total number of entries between Feb. 2022 and Nov. 2024, is shown in Figure 5.



169
170 *Figure 5: PubChemLite annotation content (total and by category) between 4 Feb. 2022 and 3 Nov. 2024.*

171 Overall, despite PubChem increasing in content dramatically over that time, PubChemLite has remained
172 generally stable at ~380,000 entries, growing slowly over time. The systematic increase of the
173 BioPathway category (purple, Figure 5) starting in October 2022 introduced a number of irrelevant
174 candidates in preliminary results of non-target studies⁴⁰ and caused continual expansion; this content
175 was switched to the “Pathways” subcategory rather than the entire “Interactions and Pathways”
176 heading in May 2023 (see Figure 1), improving performance and interpretability of candidate hits. The
177 dramatic increase in FoodRelated information was due to the integration of FooDB⁴¹ into PubChem
178 annotation content; despite the large increase in that category, the overall candidate numbers remained
179 stable, indicating that many of these candidates already had other annotation content in PubChemLite.
180 The increase and then decrease of content in the DiseaseDisorder category in March-April 2024 was due
181 to an update of one data source that suddenly introduced many low-quality candidates - upon
182 contacting the contributor they checked their data and identified several issues; the fixes were
183 processed by 12 April (see Figure 5). Overall, the continuous monitoring and use of PubChemLite in
184 various NTS studies helps ensure relevance and usefulness for the community.

185 CCS Values

186 The experimental CCS data in PubChem currently (5 Nov. 2024) includes a total of 22,192 experimental
187 CCS values corresponding to 8099 unique compounds (CIDs). The contributions include 1554 CCS values
188 for 1136 CIDs from the Baker Lab^{24,25}; 17,187 CCS values for 6242 CIDs from CCSbase²⁰; and 3451 CCS
189 values corresponding to 869, 574, 148 and 205 CIDs from the NORMAN-SLE¹² collections S50
190 CCSCOMPEND^{26,27}, S61 UJICCSLIB^{28,29}, S79 UACCSECC^{3,30} and S116 REFCCS³¹, respectively. Information is
191 available for 98 adducts, where the most common adducts are [M+H]⁺ (7545 CCS values, 4278 CIDs), [M-
192 H]⁻ (4279 CCS values, 2064 CIDs), [M+Na]⁺ (4064 CCS values, 2831 CIDs), [M+K]⁺ (1179 CCS values, 1140
193 CIDs) and [M+H-H₂O]⁺ (1154 CCS values, 1113 CIDs).

194 Predicted CCS values calculated with CCSbase are available for all but 12 CIDs in PubChemLite, since
195 these could not be parsed with RDKit (the toolkit used in [cs3db](#)), despite being compatible with the
196 toolkits used in PubChem (OpenEye) and MetFrag (CDK). Toolkit compatibility has been explored
197 elsewhere recently⁴². Nevertheless, 12 failures out of 389,779 CIDs corresponds to only ~0.003% of the
198 dataset. In contrast, 11,373 (~2.9%) XlogP values are missing from the same file (values cannot be
199 calculated with the XlogP model⁴³ used by PubChem). The CCSbase model currently in use has been
200 trained on a slightly different set to the public CCSbase website and the numbers integrated in
201 PubChemLite may differ slightly. Since a major motivation to improve the availability of open
202 experimental CCS values was to increase the amount of relevant data for predictive systems such as
203 CCSbase, the pipelines presented in this article have been designed to allow upgrades of the CCSbase
204 model once they are ready. The predicted CCS values in PubChemLite have already been applied in NTS
205 studies⁴⁴.

206 Future Perspectives

207 PubChemLite has been used in several NTS applications; user feedback helps ensure the relevance for
208 environmental screening. Collaborative research activities have already identified less relevant content
209 (as described above), but also poor coverage (possible lack of annotation content) for compounds in
210 sediments. The chemical stripe integration helps researchers visualize and interpret the chemical history
211 over time³⁶, while the annotation summary pages display basic information available simply, providing
212 direct access to PubChem for the full content. The predicted CCS values will help improve NTS workflows
213 for IMS, while the growth of the experimental CCS dataset will help predictive models such as CCSbase
214 improve accuracy and relevance over time. PubChemLite is openly available
215 (<https://pubchemlite.lcsb.uni.lu>) - feedback is welcome (see contact page).

216 Declarations

217 Data Availability Statement

218 PubChemLite is compiled weekly from openly available files on the PubChem [FTP](#) site and is archived
219 monthly on Zenodo (DOI: [10.5281/zenodo.5995885](https://doi.org/10.5281/zenodo.5995885)). CCS values are added using open [cs3db](#) code and
220 the PubChemLite-CCS files are archived on Zenodo at DOI: [10.5281/zenodo.4081056](https://doi.org/10.5281/zenodo.4081056). The Zenodo links
221 redirect to the latest version. The code for the PubChemLite [build system](#), [inputs](#), [chemical stripes](#) and
222 [interface](#) are available on the Environmental Cheminformatics (ECI) [GitLab](#). All are available under open
223 licenses, see individual resources for details.

224 Acknowledgements

225 The authors acknowledge the earlier efforts of Todor Kondic (now at LDNS) to parts of this work, and
226 Christine Gallampo (Umea University) for her insights on the sediment NTS, as well as the
227 Environmental Cheminformatics, Bioinformatics Core, Xu lab, BakerLab and PubChem team members
228 and other colleagues and collaborators who contributed to this work indirectly via other collaborative
229 and scientific activities and discussions.

230 Funding

231 AE, DA and ELS acknowledge funding support from the Luxembourg National Research Fund (FNR) for
232 project A18/BM/12341006 (AE, DA, ELS), the University of Luxembourg Institute for Advanced Studies
233 (IAS) for the Audacity project “LuxTIME” (DA, ELS) and the European Union Research and Innovation
234 program Horizon Europe for PARC, Grant No. 101057014 (AE). The work of SK, PAT, JZ and EEB was
235 supported by the National Center for Biotechnology Information of the National Library of Medicine
236 (NLM), National Institutes of Health. JND and ESB would like to acknowledge funding support from the
237 National Institute of Environmental Health Sciences (P42 ES027704) and National Institute of General
238 Medical Sciences (R01 GM141277 and RM1 GM145416). LX acknowledges financial support from the
239 National Institute of Environmental Health Sciences, National Institutes of Health (R01 ES031927).

240 Author Contributions

241 AE: Data curation, Methodology, Software (PubChemLite build, evaluation), Validation, Writing original
242 draft preparation (joint), Writing review and editing. DHR: Methodology, Software (CCSbase, cs3db),
243 Validation, Writing review and editing. VG: Methodology, Software (PubChemLite-web), Visualization,
244 Writing review and editing. DA: Methodology, Software (chemical stripes), Visualization, Writing review
245 and editing. AMK: Methodology, Software (CCSbase). SK: Methodology, Software (PubChem-CCS
246 interface), Validation, Writing review and editing. PAT: Methodology, Software (PubChemLite,
247 PubChem-CCS interface, experimental CCS), Validation, Writing review and editing. JZ: Data curation,
248 Methodology, Software (PubChemLite, PubChem-CCS interface, experimental CCS), Validation, Writing
249 review and editing. JND: Data curation, Supervision, Writing review and editing. ESB: Data curation,
250 Project administration, Resources, Supervision, Writing review and editing. EEB: Conceptualization, Data
251 curation, Methodology, Project administration, Resources, Software (PubChemLite), Supervision,
252 Validation, Writing review and editing. LX: Conceptualization, Funding acquisition, Methodology, Project
253 administration, Resources, Software (CCSbase), Supervision, Writing review and editing. ELS:
254 Conceptualization, Data curation, Funding acquisition, Methodology, Project administration, Resources,
255 Software (PubChemLite, evaluation, experimental CCS), Supervision, Validation, Visualization, Writing
256 original draft preparation (joint), Writing review and editing.

257 Conflicts of Interest

258 The authors have no competing financial interests to declare.

259 Supporting Information

260 See Data Availability Statement.

261 References

- 262 (1) Hollender, J.; Schymanski, E. L.; Ahrens, L.; Alygizakis, N.; Béen, F.; Bijlsma, L.; Brunner, A. M.;
263 Celma, A.; Fildier, A.; Fu, Q.; Gago-Ferrero, P.; Gil-Solsona, R.; Haglund, P.; Hansen, M.; Kaserzon, S.;
264 Kruve, A.; Lamoree, M.; Margoum, C.; Meijer, J.; Merel, S.; Rauert, C.; Rostkowski, P.; Samanipour, S.;
265 Schulze, B.; Schulze, T.; Singh, R. R.; Slobodnik, J.; Steininger-Mairinger, T.; Thomaidis, N. S.; Togola, A.;
266 Vorkamp, K.; Vulliet, E.; Zhu, L.; Krauss, M. NORMAN Guidance on Suspect and Non-Target Screening in
267 Environmental Monitoring. *Environmental Sciences Europe* **2023**, *35* (1), 75.
268 <https://doi.org/10.1186/s12302-023-00779-4>.
- 269 (2) Lai, Y.; Koelmel, J. P.; Walker, D. I.; Price, E. J.; Papazian, S.; Manz, K. E.; Castilla-Fernández, D.;
270 Bowden, J. A.; Nikiforov, V.; David, A.; Bessonneau, V.; Amer, B.; Seethapathy, S.; Hu, X.; Lin, E. Z.; Jbebli,
271 A.; McNeil, B. R.; Barupal, D.; Cerasa, M.; Xie, H.; Kalia, V.; Nandakumar, R.; Singh, R.; Tian, Z.; Gao, P.;
272 Zhao, Y.; Froment, J.; Rostkowski, P.; Dubey, S.; Coufalíková, K.; Seličová, H.; Hecht, H.; Liu, S.; Udhani, H.
273 H.; Restituto, S.; Tchou-Wong, K.-M.; Lu, K.; Martin, J. W.; Warth, B.; Godri Pollitt, K. J.; Klánová, J.;
274 Fiehn, O.; Metz, T. O.; Pennell, K. D.; Jones, D. P.; Miller, G. W. High-Resolution Mass Spectrometry for
275 Human Exposomics: Expanding Chemical Space Coverage. *Environmental Science & Technology* **2024**, *58*
276 (29), 12784–12822. <https://doi.org/10.1021/acs.est.4c01156>.
- 277 (3) Belova, L.; Caballero-Casero, N.; Nuijs, A. L. N. van; Covaci, A. Ion Mobility-High-Resolution Mass
278 Spectrometry (IM-HRMS) for the Analysis of Contaminants of Emerging Concern (CECs): Database
279 Compilation and Application to Urine Samples. *Analytical Chemistry* **2021**, *93* (16), 6428–6436.
280 <https://doi.org/10.1021/acs.analchem.1c00142>.
- 281 (4) Celma, A.; Bade, R.; Sancho, J. V.; Hernandez, F.; Humphries, M.; Bijlsma, L. Prediction of
282 Retention Time and Collision Cross Section (CCSH+, CCSH-, and CCSNa+) of Emerging Contaminants
283 Using Multiple Adaptive Regression Splines. *Journal of Chemical Information and Modeling* **2022**, *62*
284 (22), 5425–5434. <https://doi.org/10.1021/acs.jcim.2c00847>.
- 285 (5) Song, X.-C.; Dreolin, N.; Canellas, E.; Goshawk, J.; Nerin, C. Prediction of Collision Cross-Section
286 Values for Extractables and Leachables from Plastic Products. *Environmental Science & Technology* **2022**,
287 *56* (13), 9463–9473. <https://doi.org/10.1021/acs.est.2c02853>.
- 288 (6) Wishart, D. S.; Guo, A.; Oler, E.; Wang, F.; Anjum, A.; Peters, H.; Dizon, R.; Sayeeda, Z.; Tian, S.;
289 Lee, B. L.; Berjanskii, M.; Mah, R.; Yamamoto, M.; Jovel, J.; Torres-Calzada, C.; Hiebert-Giesbrecht, M.;
290 Lui, V. W.; Varshavi, D.; Varshavi, D.; Allen, D.; Arndt, D.; Khetarpal, N.; Sivakumaran, A.; Harford, K.;
291 Sanford, S.; Yee, K.; Cao, X.; Budinski, Z.; Liigand, J.; Zhang, L.; Zheng, J.; Mandal, R.; Karu, N.; Dambrova,
292 M.; Schiöth, H. B.; Greiner, R.; Gautam, V. HMDB 5.0: The Human Metabolome Database for 2022.
293 *Nucleic Acids Research* **2022**, *50* (D1), D622–D631. <https://doi.org/10.1093/nar/gkab1062>.
- 294 (7) Williams, A. J.; Grulke, C. M.; Edwards, J.; McEachran, A. D.; Mansouri, K.; Baker, N. C.; Patlewicz,
295 G.; Shah, I.; Wambaugh, J. F.; Judson, R. S.; Richard, A. M. The CompTox Chemistry Dashboard: A
296 Community Data Resource for Environmental Chemistry. *Journal of Cheminformatics* **2017**, *9* (1), 61.
297 <https://doi.org/10.1186/s13321-017-0247-6>.
- 298 (8) Pence, H. E.; Williams, A. ChemSpider: An Online Chemical Information Resource. *Journal of*
299 *Chemical Education* **2010**, *87* (11), 1123–1124. <https://doi.org/10.1021/ed100697w>.

- 300 (9) Kim, S.; Chen, J.; Cheng, T.; Gindulyte, A.; He, J.; He, S.; Li, Q.; Shoemaker, B. A.; Thiessen, P. A.;
301 Yu, B.; Zaslavsky, L.; Zhang, J.; Bolton, E. E. PubChem 2023 Update. *Nucleic Acids Research* **2023**, *51*,
302 gkac956. <https://doi.org/10.1093/nar/gkac956>.
- 303 (10) American Chemical Society. CAS REGISTRY - The CAS Substance Collection, 2024.
304 <https://www.cas.org/cas-data/cas-registry> (accessed 2024-08-03).
- 305 (11) Wang, Z.; Walker, G. W.; Muir, D. C. G.; Nagatani-Yoshida, K. Toward a Global Understanding of
306 Chemical Pollution: A First Comprehensive Analysis of National and Regional Chemical Inventories.
307 *Environmental Science & Technology* **2020**, *54* (5), 2575–2584. <https://doi.org/10.1021/acs.est.9b06379>.
- 308 (12) Mohammed Taha, H.; Aalizadeh, R.; Alygizakis, N.; Antignac, J.-P.; Arp, H. P. H.; Bade, R.; Baker,
309 N.; Belova, L.; Bijlsma, L.; Bolton, E. E.; Brack, W.; Celma, A.; Chen, W.-L.; Cheng, T.; Chirsir, P.; Čirka, L.;
310 D’Agostino, L. A.; Djoumbou Feunang, Y.; Dulio, V.; Fischer, S.; Gago-Ferrero, P.; Galani, A.; Geueke, B.;
311 Głowacka, N.; Glüge, J.; Groh, K.; Grosse, S.; Haglund, P.; Hakkinen, P. J.; Hale, S. E.; Hernandez, F.;
312 Janssen, E. M.-L.; Jonkers, T.; Kiefer, K.; Kirchner, M.; Koschorreck, J.; Krauss, M.; Krier, J.; Lamoree, M.
313 H.; Letzel, M.; Letzel, T.; Li, Q.; Little, J.; Liu, Y.; Lunderberg, D. M.; Martin, J. W.; McEachran, A. D.;
314 McLean, J. A.; Meier, C.; Meijer, J.; Menger, F.; Merino, C.; Muncke, J.; Muschket, M.; Neumann, M.;
315 Neveu, V.; Ng, K.; Oberacher, H.; O’Brien, J.; Oswald, P.; Oswaldova, M.; Picache, J. A.; Postigo, C.;
316 Ramirez, N.; Reemtsma, T.; Renaud, J.; Rostkowski, P.; Rüdell, H.; Salek, R. M.; Samanipour, S.;
317 Scheringer, M.; Schliebner, I.; Schulz, W.; Schulze, T.; Sengl, M.; Shoemaker, B. A.; Sims, K.; Singer, H.;
318 Singh, R. R.; Sumarah, M.; Thiessen, P. A.; Thomas, K. V.; Torres, S.; Trier, X.; Wezel, A. P. van;
319 Vermeulen, R. C. H.; Vlaanderen, J. J.; Ohe, P. C. von der; Wang, Z.; Williams, A. J.; Willighagen, E. L.;
320 Wishart, D. S.; Zhang, J.; Thomaidis, N. S.; Hollender, J.; Slobodnik, J.; Schymanski, E. L. The NORMAN
321 Suspect List Exchange (NORMAN-SLE): Facilitating European and Worldwide Collaboration on Suspect
322 Screening in High Resolution Mass Spectrometry. *Environmental Sciences Europe* **2022**, *34* (1), 104.
323 <https://doi.org/10.1186/s12302-022-00680-6>.
- 324 (13) Schymanski, E. L.; Kondić, T.; Neumann, S.; Thiessen, P. A.; Zhang, J.; Bolton, E. E. Empowering
325 Large Chemical Knowledge Bases for Exposomics: PubChemLite Meets MetFrag. *Journal of*
326 *Cheminformatics* **2021**, *13* (1), 19. <https://doi.org/10.1186/s13321-021-00489-0>.
- 327 (14) Helmus, R.; Laak, T. L. ter; Wezel, A. P. van; Voogt, P. de; Schymanski, E. L. patRoom: Open
328 Source Software Platform for Environmental Mass Spectrometry Based Non-Target Screening. *Journal of*
329 *Cheminformatics* **2021**, *13* (1), 1. <https://doi.org/10.1186/s13321-020-00477-w>.
- 330 (15) Ruttkies, C.; Schymanski, E. L.; Wolf, S.; Hollender, J.; Neumann, S. MetFrag Relunched:
331 Incorporating Strategies Beyond in Silico Fragmentation. *Journal of Cheminformatics* **2016**, *8* (1), 3.
332 <https://doi.org/10.1186/s13321-016-0115-9>.
- 333 (16) NCBI/NLM/NIH. PubChem Table of Contents Classification Browser, 2024.
334 <https://pubchem.ncbi.nlm.nih.gov/classification/#hid=72> (accessed 2024-07-08).
- 335 (17) Ieritano, C.; Hopkins, W. S. Assessing Collision Cross Section Calculations Using MobCal-MPI with
336 a Variety of Commonly Used Computational Methods. *Materials Today Communications* **2021**, *27*,
337 102226. <https://doi.org/10.1016/j.mtcomm.2021.102226>.
- 338 (18) Colby, S. M.; Thomas, D. G.; Nuñez, J. R.; Baxter, D. J.; Glaesemann, K. R.; Brown, J. M.; Pirrung,
339 M. A.; Govind, N.; Teeguarden, J. G.; Metz, T. O.; Renslow, R. S. ISICLE: A Quantum Chemistry Pipeline for

- 340 Establishing in Silico Collision Cross Section Libraries. *Analytical Chemistry* **2019**, *91* (7), 4346–4356.
341 <https://doi.org/10.1021/acs.analchem.8b04567>.
- 342 (19) Zhou, Z.; Luo, M.; Chen, X.; Yin, Y.; Xiong, X.; Wang, R.; Zhu, Z.-J. Ion Mobility Collision Cross-
343 Section Atlas for Known and Unknown Metabolite Annotation in Untargeted Metabolomics. *Nature*
344 *Communications* **2020**, *11* (1), 4334. <https://doi.org/10.1038/s41467-020-18171-8>.
- 345 (20) Ross, D. H.; Cho, J. H.; Xu, L. Breaking Down Structural Diversity for Comprehensive Prediction of
346 Ion-Neutral Collision Cross Sections. *Analytical Chemistry* **2020**, *92* (6), 4548–4557.
347 <https://doi.org/10.1021/acs.analchem.9b05772>.
- 348 (21) Guo, R.; Zhang, Y.; Liao, Y.; Yang, Q.; Xie, T.; Fan, X.; Lin, Z.; Chen, Y.; Lu, H.; Zhang, Z. Highly
349 Accurate and Large-Scale Collision Cross Sections Prediction with Graph Neural Networks.
350 *Communications Chemistry* **2023**, *6* (1), 1–10. <https://doi.org/10.1038/s42004-023-00939-w>.
- 351 (22) Plante, P.-L.; Francovic-Fontaine, É.; May, J. C.; McLean, J. A.; Baker, E. S.; Laviolette, F.;
352 Marchand, M.; Corbeil, J. Predicting Ion Mobility Collision Cross-Sections Using a Deep Neural Network:
353 DeepCCS. *Analytical Chemistry* **2019**, *91* (8), 5191–5199.
354 <https://doi.org/10.1021/acs.analchem.8b05821>.
- 355 (23) Rainey, M. A.; Watson, C. A.; Asef, C. K.; Foster, M. R.; Baker, E. S.; Fernández, F. M. CCS
356 Predictor 2.0: An Open-Source Jupyter Notebook Tool for Filtering Out False Positives in Metabolomics.
357 *Analytical Chemistry* **2022**, *94* (50), 17456–17466. <https://doi.org/10.1021/acs.analchem.2c03491>.
- 358 (24) Kirkwood, K. I.; Christopher, M. W.; Burgess, J. L.; Littau, S. R.; Foster, K.; Richey, K.; Pratt, B. S.;
359 Shulman, N.; Tamura, K.; MacCoss, M. J.; MacLean, B. X.; Baker, E. S. Development and Application of
360 Multidimensional Lipid Libraries to Investigate Lipidomic Dysregulation Related to Smoke Inhalation
361 Injury Severity. *Journal of Proteome Research* **2022**, *21* (1), 232–242.
362 <https://doi.org/10.1021/acs.jproteome.1c00820>.
- 363 (25) Foster, M.; Rainey, M.; Watson, C.; Dodds, J. N.; Kirkwood, K. I.; Fernández, F. M.; Baker, E. S.
364 Uncovering PFAS and Other Xenobiotics in the Dark Metabolome Using Ion Mobility Spectrometry, Mass
365 Defect Analysis, and Machine Learning. *Environmental Science & Technology* **2022**, *56* (12), 9133–9143.
366 <https://doi.org/10.1021/acs.est.2c00201>.
- 367 (26) Picache, J. A.; Rose, B. S.; Balinski, A.; Leaptrot, K. L.; Sherrod, S. D.; May, J. C.; McLean, J. A.
368 Collision Cross Section Compendium to Annotate and Predict Multi-Omic Compound Identities. *Chemical*
369 *Science* **2019**, *10* (4), 983–993. <https://doi.org/10.1039/C8SC04396E>.
- 370 (27) Picache, J.; McLean, J. S50 CCSCOMPEND The Unified Collision Cross Section (CCS)
371 Compendium, 2019. <https://doi.org/10.5281/zenodo.2658162>.
- 372 (28) Celma, A.; Sancho, J. V.; Schymanski, E. L.; Fabregat-Safont, D.; Ibáñez, M.; Goshawk, J.;
373 Barknowitz, G.; Hernández, F.; Bijlsma, L. Improving Target and Suspect Screening High-Resolution Mass
374 Spectrometry Workflows in Environmental Analysis by Ion Mobility Separation. *Environmental Science &*
375 *Technology* **2020**, *54* (23), 15120–15131. <https://doi.org/10.1021/acs.est.0c05713>.
- 376 (29) Celma, A.; Fabregat-Safont, D.; Ibáñez, M.; Bijlsma, L.; Hernandez, F.; Sancho, J. V. S61 UJICCSLIB
377 Collision Cross Section (CCS) Library from UJI, 2019. <https://doi.org/10.5281/ZENODO.3549476>.

- 378 (30) Belova, L.; Caballero-Casero, N.; Nuijs, A. L. N. van; Covaci, A. S79 UACCSECC Collision Cross
379 Section (CCS) Library from UAntwerp, 2021. <https://doi.org/10.5281/ZENODO.4704648>.
- 380 (31) Muller, H.; Palm, E.; Schymanski, E. S116 REFCCS Collision Cross Section (CCS) Values from
381 Literature, 2024. <https://doi.org/10.5281/ZENODO.10932895>.
- 382 (32) Grouès, V.; Rocca-Serra, P.; Ded, V. Elixir-Luxembourg/Data-Catalog, 2023.
383 <https://github.com/elixir-luxembourg/data-catalog> (accessed 2024-08-04).
- 384 (33) Welter, D.; Rocca-Serra, P.; Grouès, V.; Sallam, N.; Ancien, F.; Shabani, A.; Asariardakani, S.;
385 Alper, P.; Ghosh, S.; Burdett, T.; Sansone, S.-A.; Gu, W.; Satagopam, V. The Translational Data Catalog -
386 Discoverable Biomedical Datasets. *Scientific Data* **2023**, *10* (1), 470. [https://doi.org/10.1038/s41597-](https://doi.org/10.1038/s41597-023-02258-0)
387 [023-02258-0](https://doi.org/10.1038/s41597-023-02258-0).
- 388 (34) Greg Landrum. RDKit: Open-Source Cheminformatics Software, 2024. <https://www.rdkit.org/>
389 (accessed 2024-08-04).
- 390 (35) Greg Landrum; Paolo Tosco; Brian Kelley; Ricardo Rodriguez; David Cosgrove; Riccardo Vianello;
391 sriniker; gedeck; Gareth Jones; NadineSchneider; Eisuke Kawashima; Dan Nealschneider; Andrew Dalke;
392 Matt Swain; Brian Cole; Samo Turk; Aleksandr Savelev; Alain Vaucher; Maciej Wójcikowski; Ichiru Take;
393 Vincent F. Scalfani; Daniel Probst; Kazuya Ujihara; Rachel Walker; Axel Pahl; guillaume godin; tadhurst-
394 cdd; Juuso Lehtivarjo; François Bérenger; Jonathan Bisson. Rdkit/Rdkit: 2024_03_5 (Q1 2024) Release,
395 2024. <https://doi.org/10.5281/ZENODO.591637>.
- 396 (36) Aurich, D.; Schymanski, E. L.; De Jesus Matias, F.; Thiessen, P. A.; Pang, J. Revealing Chemical
397 Trends: Insights from Data-Driven Visualization and Patent Analysis in Exposomics Research.
398 *Environmental Science & Technology Letters* **2024**, *11* (10), 1046–1052.
399 <https://doi.org/10.1021/acs.estlett.4c00560>.
- 400 (37) Arp, H. P. H.; Aurich, D.; Schymanski, E. L.; Sims, K.; Hale, S. E. Avoiding the Next Silent Spring:
401 Our Chemical Past, Present, and Future. *Environmental Science & Technology* **2023**, *57* (16), 6355–6359.
402 <https://doi.org/10.1021/acs.est.3c01735>.
- 403 (38) Aurich, D. Uniluxembourg / LCSB / Environmental Cheminformatics / Chemicalstripes · GitLab.
404 *GitLab*, 2024. <https://gitlab.com/uniluxembourg/lcsb/eci/chemicalstripes> (accessed 2024-08-04).
- 405 (39) Grouès, V. Uniluxembourg / LCSB / Environmental Cheminformatics / PubChemLite-Web ·
406 *GitLab*. *GitLab*, 2024. <https://gitlab.com/uniluxembourg/lcsb/eci/pubchemlite-web> (accessed 2024-08-
407 04).
- 408 (40) Talavera Andújar, B.; Mary, A.; Venegas, C.; Cheng, T.; Zaslavsky, L.; Bolton, E. E.; Heneka, M. T.;
409 Schymanski, E. L. Can Small Molecules Provide Clues on Disease Progression in Cerebrospinal Fluid from
410 Mild Cognitive Impairment and Alzheimer’s Disease Patients? *Environmental Science & Technology* **2024**,
411 *58*, 4181–4192. <https://doi.org/10.1021/acs.est.3c10490>.
- 412 (41) WishartLab. FooDB, 2024. <https://foodb.ca/> (accessed 2024-11-06).
- 413 (42) Barnabas, S. J.; Böhme, T.; Boyer, S. K.; Irmer, M.; Ruttkies, C.; Wetherbee, I.; Kondić, T.;
414 Schymanski, E. L.; Weber, L. Extraction of Chemical Structures from Literature and Patent Documents
415 Using Open Access Chemistry Toolkits: A Case Study with PFAS. *Digital Discovery* **2022**, *1* (4), 490–501.
416 <https://doi.org/10.1039/D2DD00019A>.

- 417 (43) Cheng, T.; Zhao, Y.; Li, X.; Lin, F.; Xu, Y.; Zhang, X.; Li, Y.; Wang, R.; Lai, L. Computation of
418 Octanol–Water Partition Coefficients by Guiding an Additive Model with Knowledge. *Journal of Chemical*
419 *Information and Modeling* **2007**, 47 (6), 2140–2148. <https://doi.org/10.1021/ci700257y>.
- 420 (44) Menger, F.; Celma, A.; Schymanski, E. L.; Lai, F. Y.; Bijlsma, L.; Wiberg, K.; Hernández, F.; Sancho,
421 J. V.; Ahrens, L. Enhancing Spectral Quality in Complex Environmental Matrices: Supporting Suspect and
422 Non-Target Screening in Zebra Mussels with Ion Mobility. *Environment International* **2022**, 170, 107585.
423 <https://doi.org/10.1016/j.envint.2022.107585>.