

1 **SUPPORTING INFORMATION**

2
3 **Prioritizing Chemical Candidates from Non-Targeted Analysis Using Metadata, Spectral**
4 **Similarity, and Hazard Scoring within INTERPRET NTA**

5
6 Alex Chao^{1*}, Jeffrey M. Minucci^{2*}, Troy M. Ferland^{1,3}, E. Tyler Carr^{1,4}, Greg Janesch^{1,4}, Safia
7 Rizwan^{1,4}, Heather D. Whitehead¹, Tommy Cathey⁵, Shirley Pu^{1,3}, Laura D. Brunelle^{3,6}, Angela L.
8 Batt⁶, Jon R. Sobus^{1*}, and Antony J. Williams^{1*}

9
10 ¹United States Environmental Protection Agency, Office of Research and Development, Center
11 for Computational Toxicology and Exposure, 109 TW Alexander Dr., Research Triangle Park, NC
12 27711, United States

13
14 ²United States Environmental Protection Agency, Office of Research and Development, Center
15 for Public Health and Environmental Assessment, 109 TW Alexander Dr., Research Triangle Park,
16 NC 27711, United States

17
18 ³Oak Ridge Institute for Science and Education (ORISE) Participant, Oak Ridge, TN 37831, USA

19
20 ⁴Oak Ridge Affiliated Universities (ORAU) Student Services Contractor, 109 T.W Alexander
21 Drive, Research Triangle Park, NC 27711, United States

22
23 ⁵General Dynamics Information Technology, 109 TW Alexander Dr., Research Triangle Park, NC
24 27711, United States

25
26 ⁶United States Environmental Protection Agency, Office of Research and Development, Center
27 for Environmental Solutions and Emergency Response, 26 W Martin Luther King Dr., Cincinnati,
28 OH 45268, United States

29
30 † Authors contributed equally

31
32 *Authors to whom correspondence should be addressed:

33 Alex Chao: chao.alex@epa.gov

34 Jeffrey Minucci: minucci.jeffrey@epa.gov

35 Jon Sobus: sobus.jon@epa.gov

36 Antony Williams: williams.antony@epa.gov

47	Supporting Information Table of Contents
48	
49	Section 1.0 – Experimental Details
50	
51	Section 1.1 – INTERPRET NTA and Associated Database Details
52	
53	Section 1.2 – Data Analysis and Results Scoring
54	
55	Section 2.0 – Input Requirements for MS ² Workflow
56	
57	Section 3.0 – Chemical metadata scoring process
58	
59	Section 4.0 – Supplemental Figures
60	
61	Section 5.0 – Interactive Chemical Results Visualizations
62	
63	Section 5.1 – Interactive Scatterplot of Chemical Results
64	
65	Section 5.2 – Interactive Bar Chart and Grid Table of Chemical Results
66	
67	

68 Section 1.0 – Experimental Details

69

70 The following sections provide further experimental details associated with this study.

71

72

73 Section 1.1 – INTERPRET NTA and Associated Database Details

74

75 *INTERPRET NTA Calculation Engine*

76 INTERPRET NTA (hereafter referred to as I-NTA) is a module of US EPA’s Analytical Methods
77 and Open Spectra (AMOS) web application. The backend components of I-NTA (also known as
78 the Calculation Engine) were developed in Python using Django for application programming
79 interfaces (APIs) and the Dask library¹ for distributed and asynchronous processing. When a user
80 submits their data and desired settings to an I-NTA workflow on AMOS, a request is sent to the
81 backend API that triggers job submission to the Dask scheduler. Once the job is distributed to a
82 Dask worker, the worker completes the required calculations and data retrieval, and the output is
83 stored in a central MongoDB database for later retrieval by AMOS to serve to the user (**Figure**
84 **S1**). The various components of I-NTA’s backend are orchestrated and deployed via Kubernetes
85 on an US EPA Enterprise Amazon Web Services (AWS) environment. The source code for I-NTA
86 is publicly available and can be found at https://github.com/quanted/nta_app.

87

88 *DSSTox MS-Ready Structures*

89 The I-NTA module, and more generally the AMOS application, uses DSSTox as the underlying
90 chemical substance database, facilitating chemical lookup based on measured accurate mass or
91 predicted molecular formula.² DSSTox has been accessible to the public for almost a decade via
92 the web-based application titled the “CompTox Chemicals Dashboard” (hereafter, the
93 “Dashboard”), available at <https://comptox.epa.gov/dashboard>.³ The Dashboard has been used to
94 great effect for NTA projects, with multiple demonstrations showing benefits of metadata-based
95 candidate prioritization.^{4,5} The majority of users have taken advantage of capabilities for structure-
96 based queries of “MS-Ready” chemical structures via batch interface searches.^{6,7} An MS-Ready
97 structure refers to a chemical structure that has been specifically formatted and processed to
98 accurately represent how a molecule would be observed by mass spectrometry, meaning it is de-
99 salted, de-solvated, and separated into individual components (when in mixtures). The MS-Ready
100 structures associated with the DSSTox structural content are used as input for the generation of *in*
101 *silico* spectra (see *CFM-ID in silico Spectra* section), which have been shown to aid compound
102 identification in our previous NTA works.^{8,9} In recent years, NTA support provided to the
103 community by EPA’s publicly accessible tools has been limited to the Dashboard (via advanced
104 searching capabilities) and open datasets provided as downloads (i.e., *in silico* predicted spectral
105 data, DSSTox MS-Ready structures, etc.). I-NTA development has paralleled the efforts described
106 above to finally integrate all tested and proven approaches for chemical identification and
107 candidate prioritization into an application that also performs necessary NTA data quality review.

108

109 *AMOS Chemical Metadata*

110 Chemical metadata has been shown to assist the identification of “known unknowns” via ranking
111 of candidate structures.^{4,5,10} In previous work, we have used the following metadata: 1) number

112 of associated “data sources”; 2) number of publications; and 3) number of patents. Data sources
113 are equivalent to the “chemical lists” provided in the Dashboard
114 (<https://comptox.epa.gov/dashboard/chemical-lists>). At the time of writing there are 503 lists, each
115 associated with a specific collection of chemicals registered into the DSSTox database. Lists can
116 represent a specific regulation (e.g., BIOSOLIDS, [https://comptox.epa.gov/dashboard/chemical-](https://comptox.epa.gov/dashboard/chemical-lists/CWA311HS)
117 [lists/CWA311HS](https://comptox.epa.gov/dashboard/chemical-lists/CWA311HS)), a collection of chemicals based on specific formula rules (e.g., PFAS
118 chemicals, <https://comptox.epa.gov/dashboard/chemical-lists/PFASSTRUCTV5>), or data
119 collections associated with a publication (e.g., suspect screening of chemicals in consumer
120 products, <https://comptox.epa.gov/dashboard/chemical-lists/EPACONS>).¹¹⁻¹³ Regarding
121 publications and patents, PubChem and PubMed counts have historically been harvested and
122 stored within the Dashboard.⁴ However, metadata refresh rates have lagged behind ongoing new
123 chemical registration (in support of NTA projects and other agency activities). Relevant metadata
124 have therefore been added to the AMOS database and are updated on a more frequent basis using
125 the PubChem API.¹⁴ The inclusion of chemical metadata into AMOS, along with analytical
126 method documents, chemical fact sheets (harvested from domestic and international collections),
127 and experimental spectra, has provided expanded capabilities for candidate ranking via I-NTA
128 interactive visualizations (*vide infra*).

129

130 *CFM-ID in silico Spectra*

131 When available, chemical fragmentation (here MS²) spectra allow for higher confidence
132 identifications in certain NTA studies.¹⁵ Limitations in the number of chemicals with available
133 experimental MS² spectra have led to the generation of *in silico* MS² spectra that provide wider
134 chemical coverage in NTA experiments. Our previous work documented the generation of *in silico*
135 MS² spectra for the entirety of chemicals within the DSSTox database using the CFM-ID
136 algorithms.^{8, 16} It further reported on the performance of this “CFM-ID database” for accurately
137 identifying chemicals in an NTA study.^{5, 9, 17} At the time of writing, the CFM-ID database contains
138 predictions generated via CFM-ID for 1.04 million MS-Ready chemical structures; these
139 predictions cover ESI+, ESI-, and EI spectra. While the work reported here demonstrates the utility
140 of the predicted ESI spectra for chemical identification, future work is intended to similarly
141 highlight value of the predicted EI spectra.

142

143 *CHM Hazard and Exposure Values*

144 In our earliest NTA work, we used *in vitro* bioactivity data and exposure predictions to prioritize
145 chemicals tentatively identified in household dust and drinking water.^{18, 19} At that time, our
146 cheminformatics applications and underlying databases were in their infancy. Over the past decade
147 they have advanced significantly to include additional exposure and toxicity information (available
148 databases include Chemical and Products Database [CPDat], Multimedia Monitoring Database
149 [MMDB], Toxicity Forecasting Database [ToxCast], Toxicity Value Database [ToxVal], and
150 Toxicity Reference Database [ToxRef]).²⁰⁻²⁴ While the majority of these data are available via the
151 Dashboard in some form, I-NTA primarily utilizes data assembled into the Cheminformatics
152 Hazard Module (CHM), which is part of EPA’s publicly available Cheminformatics Modules
153 (<https://hcd.rtpnc.epa.gov/#/hazard>). CHM data are most desirable for NTA projects, as refreshes
154 are more frequent, with different data streams updated at least annually, allowing more current

155 chemical information to be utilized for candidate prioritization. CHM is a web-based application
156 (expanding on the work of Vegosen and Martin)²⁵ that delivers access to human health,
157 ecotoxicity, environmental fate, and exposure information (**Table S4**). It has previously been
158 applied to perform priority-based ranking of chemical candidates identified by NTA. CHM
159 integrates both experimental and predicted (based on molecular structure) toxicity data, along with
160 predicted fate and exposure values, to assess potential risks associated with a given chemical
161 substance. As discussed below, the collective information available from the CHM (hereafter
162 termed “hazard values” for brevity) can be used to prioritize NTA suspects for further
163 investigation.

164

165 **Section 1.2 – Data Analysis and Results Scoring**

166

167 *NTA Dataset for Proof-of-Concept Demonstration*

168 Study details and initial findings were first described in Brunelle *et al.* (2024).²⁶ The data used
169 here are based on the subsequent I-NTA QA/QC demonstration work of Sobus *et al.* (2025).²⁷
170 Briefly, passive organic chemical integrative samplers (POCIS) were deployed for 27-31 days
171 during one of three sampling campaigns – October 2014, April 2015, or August 2015. POCIS
172 samples were collected at one of four locations, including upstream of a wastewater treatment
173 plant, within the wastewater effluent mixing zone, at the intake tap of a drinking water treatment
174 plant, and at the drinking water tap of the drinking water treatment plant. Samples were extracted
175 using methanol and analyzed using ultra-performance liquid chromatography (UPLC) coupled
176 with HRMS in ESI+ and ESI- modes.

177

178 *MS¹ data processing*

179 The generation of preprocessed MS¹ feature matrices (one each for ESI+ and ESI- data) using
180 MZmine 3 was previously described.²⁷ The outputs from MZmine 3 were used as inputs for the I-
181 NTA MS¹ workflow and are provided in Sobus *et al.* (2025).²⁷ For the current work, MS¹
182 workflow parameters exactly follow that of Sobus *et al.* (2025), but with additional parameters
183 selected for chemical database searching. Specifically, MS¹ features that met QA/QC criteria were
184 searched by mass (7 ppm tolerance) against the MS-Ready masses of DSSTox structures,
185 retrieving a list of chemical candidates with their associated PubChem (sources, patents, and
186 articles) and PubMed (articles) counts. These candidates were searched by DTXSID against CHM
187 in order to retrieve available hazard values. Full analysis parameters for I-NTA MS¹ data
188 processing are given in **Table S5**.

189

190 *Additional database searches*

191 Chemical results from the I-NTA MS¹ workflow output were batch searched, via DTXSID, within
192 AMOS to retrieve counts of fact sheets, analytical methods, and experimental spectra associated
193 with each chemical candidate. DTXSIDs were also searched against the Dashboard to identify
194 presence on 50 chemical lists associated with water as a matrix (**Table S6**).

195

196 *MS² data processing*

197 MS² spectral data were preprocessed and aligned via Progenesis QI (Nonlinear Dynamics, v2.3)
198 and exported in .msp format.²⁶ The export files were inputs for the I-NTA MS² workflow, which
199 accepts .msp and .mgf files exported from either preprocessed aligned data or individual sample
200 runs (as described in **SI Section 2.0**). Within the MS² workflow, input files were parsed to extract
201 a list of individual MS² features, where each MS² feature is an MS² spectrum with an associated
202 precursor mass, retention time, and adduct assignment. When present, adduct assignments were
203 used to calculate a neutral mass from the precursor mass (see **Table S7**); otherwise, neutral masses
204 were calculated with the assumption of [M+H]⁺ and [M-H]⁻ ions for positive and negative mode
205 data, respectively. MS² features were deduplicated using neutral mass and retention time windows
206 of 10 ppm and 0.2 min, respectively; for any duplicate MS² features, the instance with the highest
207 total sum intensity of fragment peaks was retained. Neutral masses of deduplicated MS² features
208 were searched against the CFM-ID database using a 10-ppm mass search window. Chemical
209 candidates were retrieved with associated *in silico* spectra at collision energy (CE) levels of 10,
210 20, and 40 eV. Full analysis parameters for I-NTA MS² data processing are given in **Table S5**.

211

212 *Scoring INTERPRET NTA Outputs*

213

214 *Chemical metadata scoring*

215 Metadata categories used for candidate scoring include (1) PubChem source counts (i.e., the
216 number of depositors in PubChem associated with a chemical)²⁸, (2) PubChem patent counts, (3)
217 PubChem article counts, (4) PubMed article counts, (5) AMOS experimental spectra counts, (6)
218 AMOS analytical method counts, (7) AMOS fact sheet counts, and (8) Dashboard water list counts
219 (see **Table S6** for considered Dashboard lists). A total metadata score was calculated and assigned
220 to each MS-Ready candidate returned from the I-NTA mass query. Scores were calculated via a
221 three-step process (see **SI section 3.0** for an illustrative example). First, for each MS-Ready
222 candidate with multiple associated substances, a sum value across substances was separately
223 calculated for each metadata category. Second, each metadata category score was normalized to a
224 value between 0 and 1 at the feature level (by dividing each candidate score by the max score for
225 that feature). Third, normalized metadata values for each candidate assigned to a given feature
226 were summed into a total metadata score, with a potential range of 0 to 8 (a value of 8 would
227 indicate a feature candidate had the highest counts in all eight metadata categories).

228

229 *MS² in silico spectra scoring*

230 Scoring between MS² features and corresponding candidates considered all three CE levels of
231 predicted *in silico* spectra, following the method of Chao *et al.* (2020).⁹ Specifically, predicted
232 spectra were scored against experimental spectra using a cosine dot-product algorithm²⁹, with
233 fragment peaks matched on a 0.02 Da window, and the three similarity scores for each candidate
234 (from 10, 20, and 40 eV predictions) summed into a total score. For each MS² feature, total scores
235 for all candidates were normalized into quotient scores by dividing each candidate score by the
236 highest candidate score for that feature.

237

238 *Hazard value scoring*

239 Hazard value scoring was based on CHM hazard levels (i.e., very high, high, medium, low,
240 inconclusive, no data) and authority source designations (i.e., authoritative, screening, QSAR
241 model) across 20 endpoints (**Table S4**). Following the procedures of Newmeyer *et al.* (2024),
242 hazard levels and authority source designations were converted to numeric values and used to
243 calculate a quality-adjusted hazard (QAH) score for each candidate substance.³⁰ An accompanying
244 completeness score was also calculated at the substance level. Briefly, average hazard and quality
245 scores were calculated for each candidate substance, which were then multiplied to create the QAH
246 score. Completeness scores were calculated by dividing the number of endpoints with data for
247 each candidate substance by the total number of endpoints (n=20). Scores at the chemical
248 substance level were collapsed to the MS-Ready chemical structure level by taking only the highest
249 QAH score, and associated completeness score, amongst all associated chemical substances.

250

251 *Ratio calculations*

252 Using data for all *known chemical features* (i.e., those associated with the 77 *known chemicals*),
253 distributions of ratios were calculated using metadata field values, the overall metadata scores, the
254 MS² scores, and the QAH scores. For each *known chemical feature*, exactly one correct candidate
255 exists, and a series of incorrect candidates exists. Thus, for each pairing of correct candidate and
256 incorrect candidate, ratio values were calculated as the correct candidate score divided by the
257 incorrect candidate score. Zero values for any score category did not allow for meaningful ratio
258 calculations. Thus, imputations were performed in instances where scores were zero. Imputed
259 values were set as the lowest nonzero value for a given score category divided by the square root
260 of two. No imputation was performed if the correct candidate for a *known chemical feature* was
261 the only candidate with a nonzero score.

262

263

264

265 **Section 2.0 – Input Requirements for MS² Workflow**

266
267 The I-NTA MS² workflow currently accepts .mgf and .msp formats as input files. Shown below
268 are examples of file formats for both, as well as a description of how they are parsed by I-NTA
269 when read in.

270
271 *Example MGF format file*

```
272  
BEGIN IONS  
PEPMASS=105.0423685  
CHARGE=1+  
TITLE=EntactEnv_Pos_MS1_Dust1IDA_01.d, MS/MS of 105.0423685 1+ at 0.04745 mins  
RTINSECONDS=2.847  
86.5515288395944      10.08654  
86.5781883146149      11.06618  
87.4602539495444      12.05921  
103.348974721758      10.08654  
103.996592740677      13.04167  
110.133284589395      13.11111  
126.802093185387      12.20833  
136.071037424052      11.08654  
144.893286462799      13.125  
154.376888744377      10.075  
157.374623516955      11.16667  
158.92422332612       10  
185.559080046904      14.1  
194.428548644229      38.80357  
214.063863547247      10.47115  
216.511031353531      11.125  
218.149310353634      12.06618  
221.991928276189      11.00961  
263.123839285519      11.68056  
318.116692565632      11.03571  
332.775479092216      13.00735  
337.009514276059      16.00833  
347.614828778035      23.07143  
386.061720707835      11.00833  
466.52480913612       13.11111  
474.667515677988      13  
END IONS
```

273
274
275 - “BEGIN IONS” – Indicates the beginning of a new MS² spectrum. This will trigger I-NTA
276 to create a new MS² feature for which to populate with the relevant spectrum/feature
277 information.

278
279 - “PEPMASS” – The numeric value in this line will be saved as the precursor mass for the
280 MS² feature.
281 o Example value: 105.0423685

282
283 - “RTINSECONDS” – The numeric value in this line will be saved as the retention time (in
284 seconds).
285 o Example value: 2.847

286

- 287 - Any numeric value at the beginning of the line – Indicates a line that has a fragment m/z,
 288 fragment intensity pair. The numeric values in the line are split, where the first is added to
 289 a list of fragment m/z values and the second is added to a list of fragment intensity values.
 290 ○ Example value: [86.5515288395944], [10.08654]
 291
 292 - “END IONS” – This indicates the end of an MS² spectrum in the input file. The precursor
 293 mass, retention time, and list of fragment m/z and intensity pairs are saved as an MS²
 294 feature in a list.
 295

296 *Example MSP format file*
 297

```
Name: Unknown (0.72_101.9637n)
Precursor_type: [M+H-H2O]+
PrecursorMZ: 84.9604
Comment: 0.72_101.9637n
Num Peaks: 25
53.0199 80.6951
61.0080 2561.4187
61.5094 102.5312
62.0240 190.7465
62.5263 46.6838
63.0048 47.1720
66.0197 485.2085
68.0167 46.0241
68.9532 51.1917
68.9831 824.8460
70.0131 2110.7041
70.5147 92.6136
70.9805 170.6099
71.0303 130.9467
72.0108 49.6613
73.5328 1239.6155
74.0334 178.1040
77.0206 159.1474
79.0187 499.7516
81.5214 2441.2715
82.0148 2040.3814
82.5377 1487.0464
83.0386 287.4457
83.5182 47.6695
84.9604 16908.8770
```

- 298
 299
 300 - “Name” – Indicates the beginning of a new MS² spectrum. This will trigger I-NTA to create
 301 a new MS² feature for which to populate with the relevant spectrum/feature information.
 302
 303 - “Precursor_type” – Where present, the text value in this line will be saved as the adduct
 304 type of the MS² feature. If this line is not present, the adduct type is assumed to be [M+H]⁺
 305 or [M-H]⁻ for positive or negative mode, respectively.
 306 ○ Example value: “[M+H-H2O]”

- 307 - “PrecursorMZ:” – The numeric value in this line will be saved as the precursor mass for
308 the MS² feature. If there is an adduct type present for the MS² feature, the precursor mass
309 is adjusted (See **Table S7** for specific mass adjustments performed by I-NTA)
310 ○ Example value: 84.9604 + 18.010565 = 102.970965
311
- 312 - “Comment:” – The first numeric value in this line prior to the underscore will be saved as
313 the retention time (in minutes).
314 ○ Example value: 0.72
315
- 316 - “Num Peaks:” – The numeric value in this line represents if there are fragment peaks
317 associated with the spectrum. A “has_fragments” Boolean is created checking if this
318 number is greater than zero.
319
- 320 - Any numeric value at the beginning of the line – Indicates a line that has a fragment m/z,
321 fragment intensity pair. The numeric values in the line are split, where the first is added to
322 a list of fragment m/z values and the second is added to a list of fragment intensity values.
323 ○ Example value: [53.0199], [80.6951]
324
- 325 - Blank space line: This indicates the end of an MS² spectrum in the input file. If the MS²
326 spectrum has any fragments (has_fragments == True), then the precursor mass, retention
327 time, and list of fragment m/z and intensity pairs are saved as an MS² feature in a list.
328
329

330 **Section 3.0 – Chemical metadata scoring process**

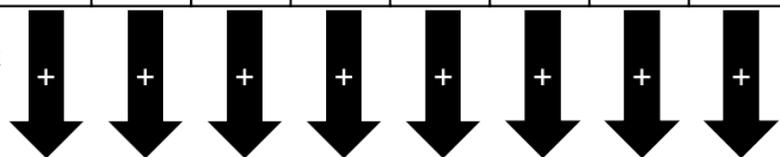
331
 332 The following text describes how overall metadata scores were calculated for chemical structures
 333 from their associated substance-related metadata retrieved from DSSTox.

334
 335 Step 1: Collapse chemical substances (identified by a DTXSID) into chemical structures
 336 (identified by a DTXCID) by summing metadata fields for each substance.

All Chemical Substances Mapped from Chemical Structure DTXCID0027983

Chemical Substance Name	DTXSID	PubChem Patents	PubChem Articles	PubMed Articles	PubChem Sources	AMOS spectra	AMOS methods	AMOS fact sheets	Dash-board Water lists
Acesulfame potassium	DTXSID1030606	15842	813	286	21	12	15	1	1
Acesulfame	DTXSID0048006	27072	752	346	31	75	8	0	3
Aspartame acesulfame	DTXSID30909991	111	0	0	10	0	0	0	0
Calcium bis(6-methyl-2,2-dioxo-2H-1,2lambda~6~,3-oxathiazin-4-olate)	DTXSID50968888	25	283	283	5	0	0	0	0

Sum within each metadata field



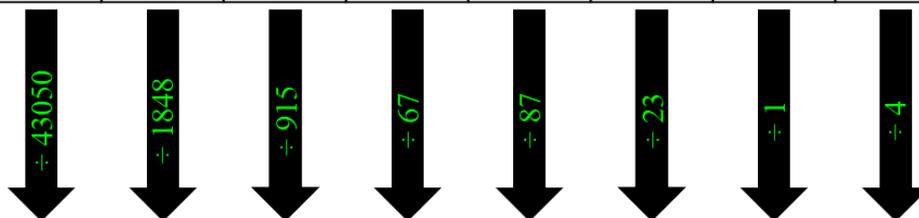
Chemical Structure (DTXCID)	PubChem Patents	PubChem Articles	PubMed Articles	PubChem Sources	AMOS spectra	AMOS methods	AMOS fact sheets	Dash-board Water lists
DTXCID0027983	43050	1848	915	67	87	23	1	4

338
 339
 340
 341
 342
 343
 344
 345
 346
 347

348 Step 2: Within each feature, normalize all the metadata field counts for all associated chemical
 349 structures by dividing by the maximum metadata field count.
 350
 351

DTXCID	PubChem Patents	PubChem Articles	PubMed Articles	PubChem Sources	AMOS spectra	AMOS methods	AMOS fact sheets	Dash-board Water lists
DTXCID0027983	43050	1848	915	67	87	23	1	4
DTXCID10618574	11	0	0	2	0	0	0	0
DTXCID20585208	14	0	0	2	0	0	0	0
DTXCID60300210	1	8	0	6	3	0	0	0
DTXCID70618575	0	0	0	2	0	0	0	0
DTXCID80401596	20	0	0	3	0	0	0	0

Divide values by
 maximum value
 within each
 metadata field



DTXCID	PubChem Patents	PubChem Articles	PubMed Articles	PubChem Sources	AMOS spectra	AMOS methods	AMOS fact sheets	Dash-board Water lists
DTXCID0027983	1	1	1	1	1	1	1	1
DTXCID10618574	0.00026	0	0	0.030	0	0	0	0
DTXCID20585208	0.00033	0	0	0.030	0	0	0	0
DTXCID60300210	0.00002	0.0043	0	0.090	0.034	0	0	0
DTXCID70618575	0.00000	0	0	0.030	0	0	0	0
DTXCID80401596	0.00046	0	0	0.045	0	0	0	0

352
 353
 354
 355
 356
 357
 358
 359
 360
 361
 362
 363
 364
 365
 366
 367
 368
 369
 370
 371
 372

373 Step 3: Sum across the eight metadata field normalized values for each chemical structure
 374 (DTXCID) to create an overall metadata score for each chemical structure.

375
 376

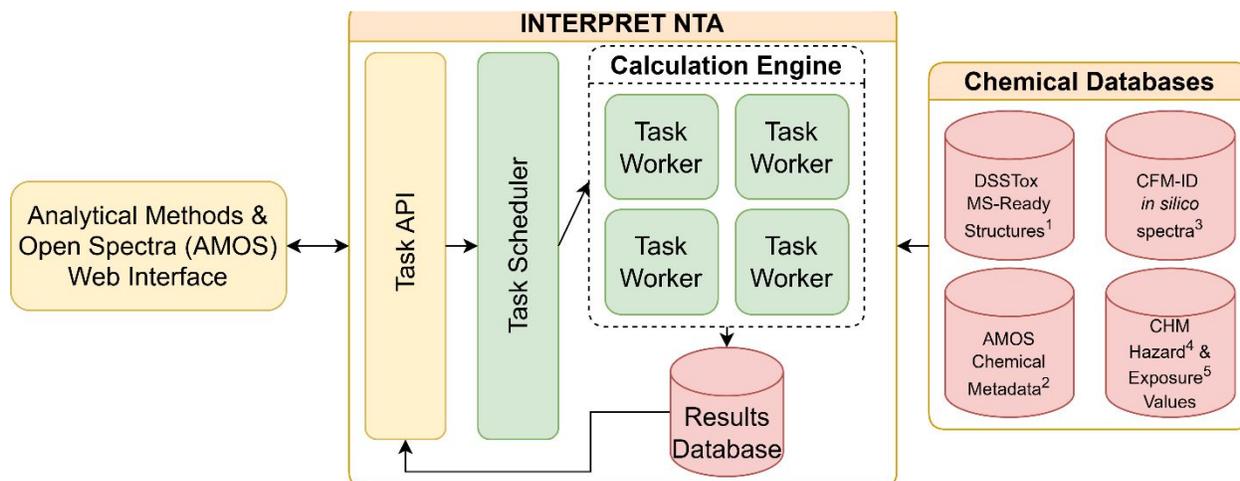
DTXCID	PubChem Patents	PubChem Articles	PubMed Articles	PubChem Sources	AMOS spectra	AMOS methods	AMOS fact sheets	Dash-board Water lists	<i>Sum values for all fields</i>	Overall Metadata Score
DTXCID0027983	1	1	1	1	1	1	1	1		8
DTXCID10618574	0.00026	0	0	0.030	0	0	0	0		0.030
DTXCID20585208	0.00033	0	0	0.030	0	0	0	0		0.030
DTXCID60300210	0.00002	0.0043	0	0.090	0.034	0	0	0		0.128
DTXCID70618575	0.00000	0	0	0.030	0	0	0	0		0.030
DTXCID80401596	0.00046	0	0	0.045	0	0	0	0		0.045

377
 378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404

405 Section 4.0 – Supplemental Figures

406

407



408

409

410 ¹ Includes MS-Ready chemical structures linked to parent substances and chemical identifiers^{3, 7}

411 ² Includes fact sheets, methods, patent counts, publication counts, and open spectra³¹

412 ³ Includes predicted electron ionization (EI) and electrospray ionization (ESI) spectra based on the CFM-ID algorithm^{8, 9, 17}

413 ⁴ Includes hazard (measured and predicted) scores and authority levels (Table S4)^{20-25, 30, 32-34}

414 ⁵ Includes exposure predictions based on the systematic empirical evaluation of models (SEEM3) framework³⁵

415

416

417

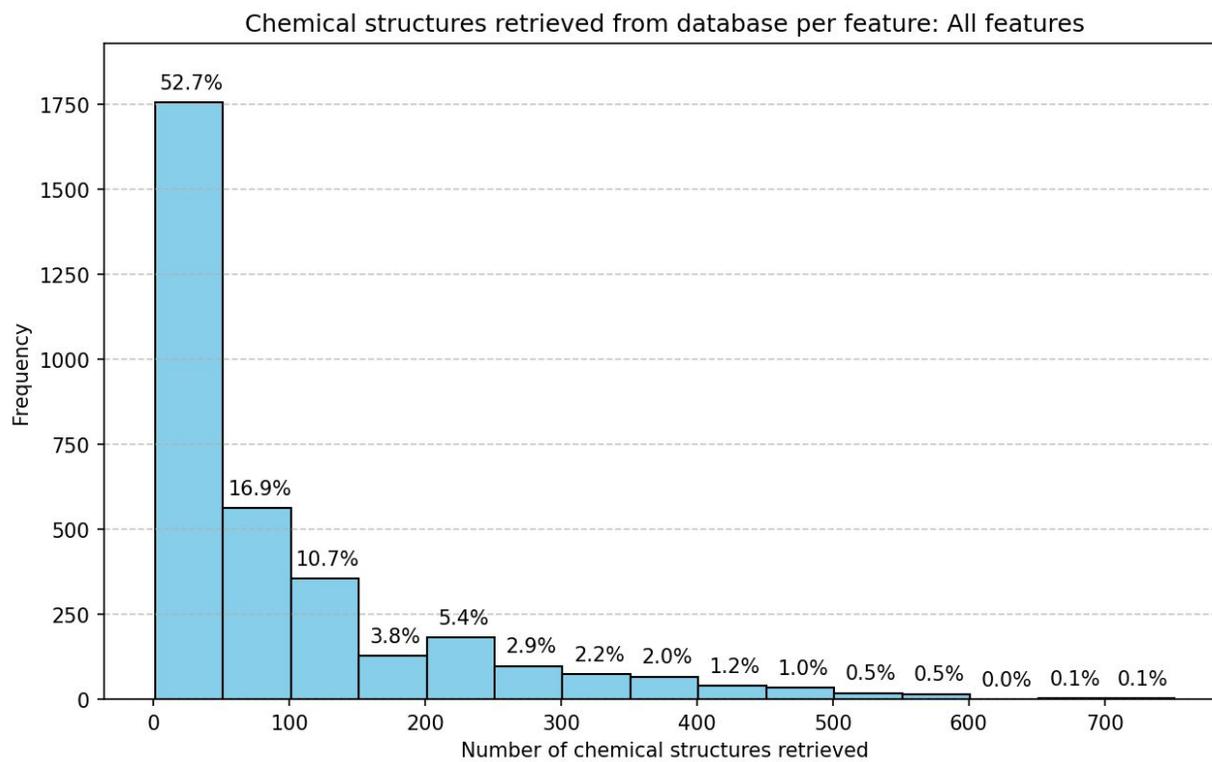
418 **Figure S1:** Application diagram illustrating the various components of the INTERPRET NTA

419 application and its database connections.

420

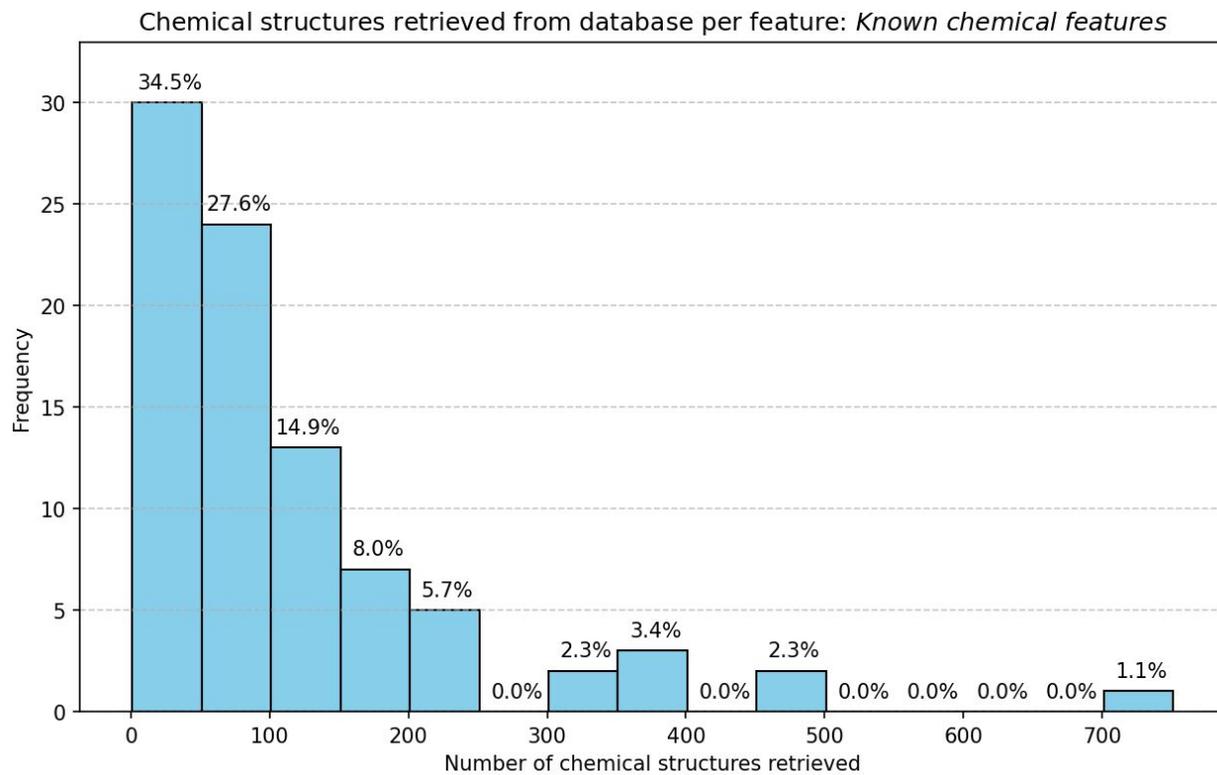
421

422



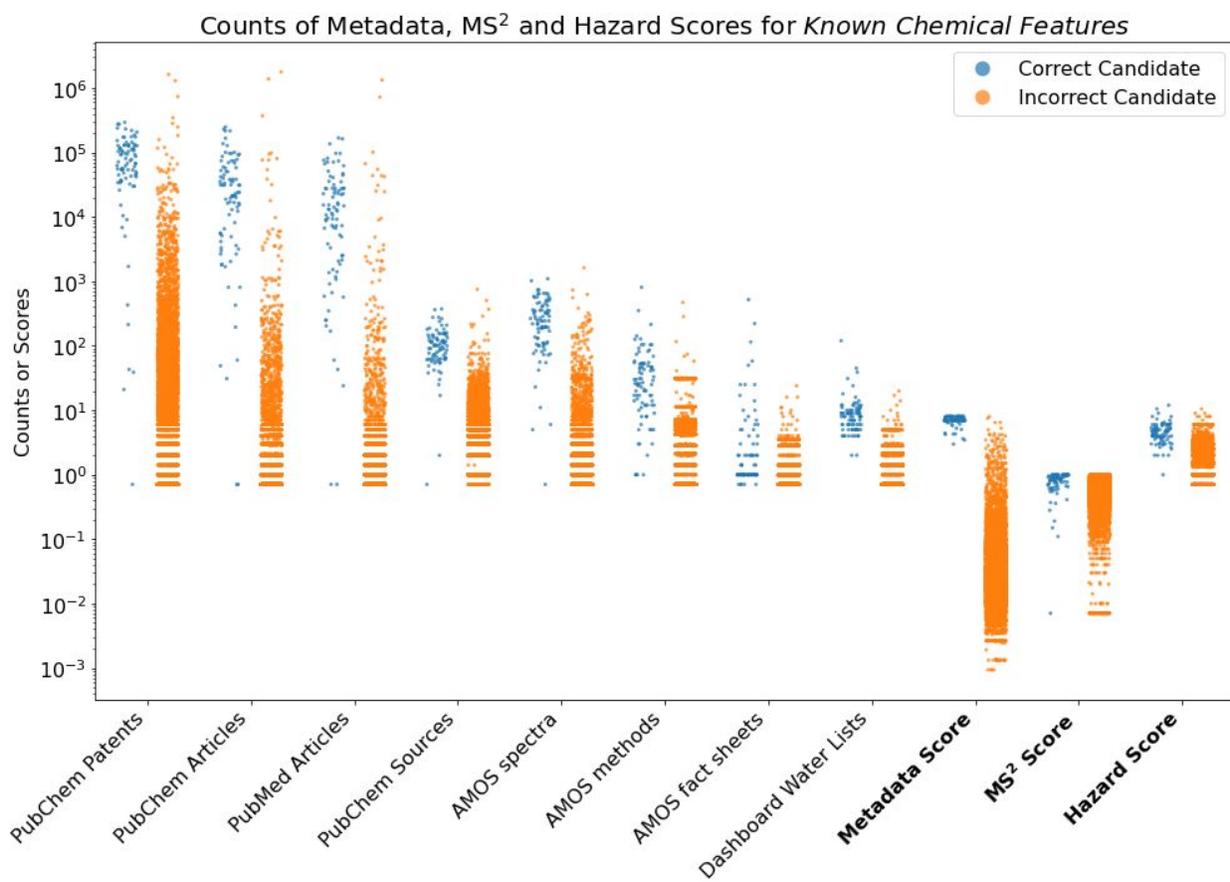
423
424
425
426
427
428

Figure S2a. The number of chemical structures retrieved for all features based on mass searches against the DSSTox chemical database.



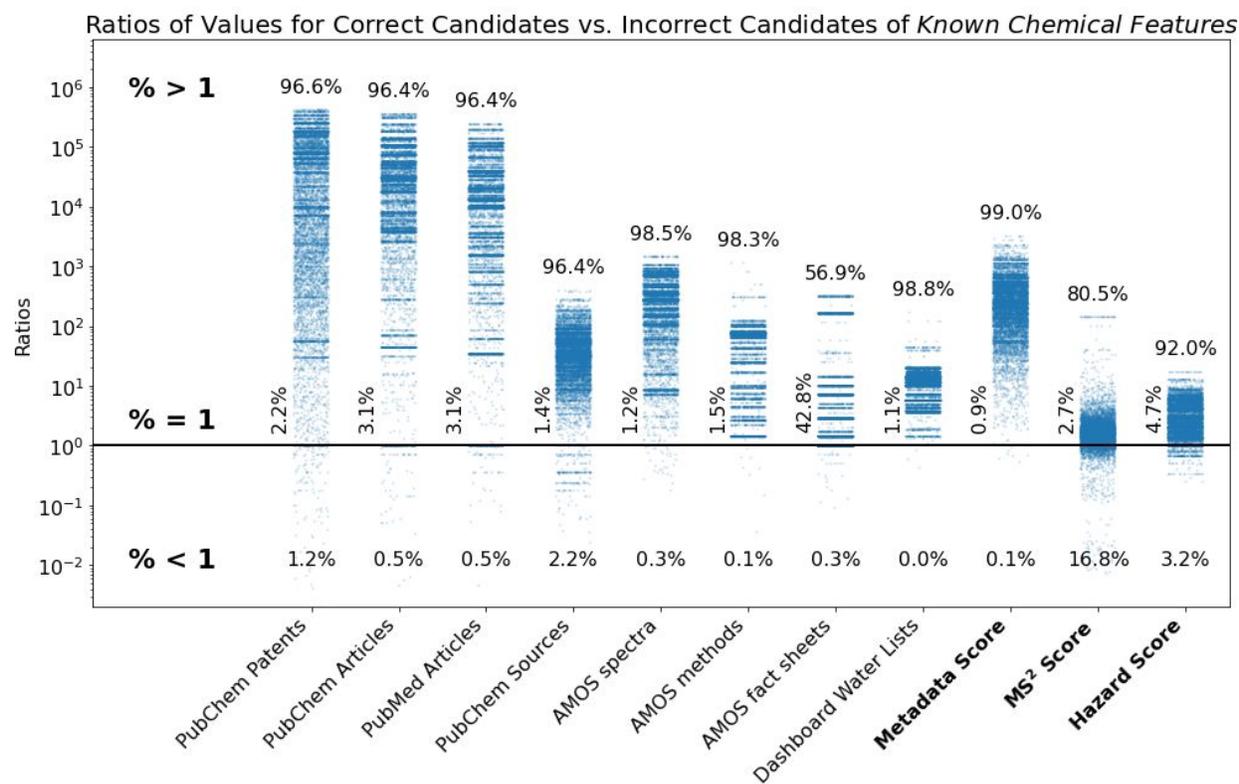
429
430
431
432
433
434
435

Figure S2b. The number of chemical structures retrieved for *known chemical features* based on mass searches against the DSSTox chemical database.



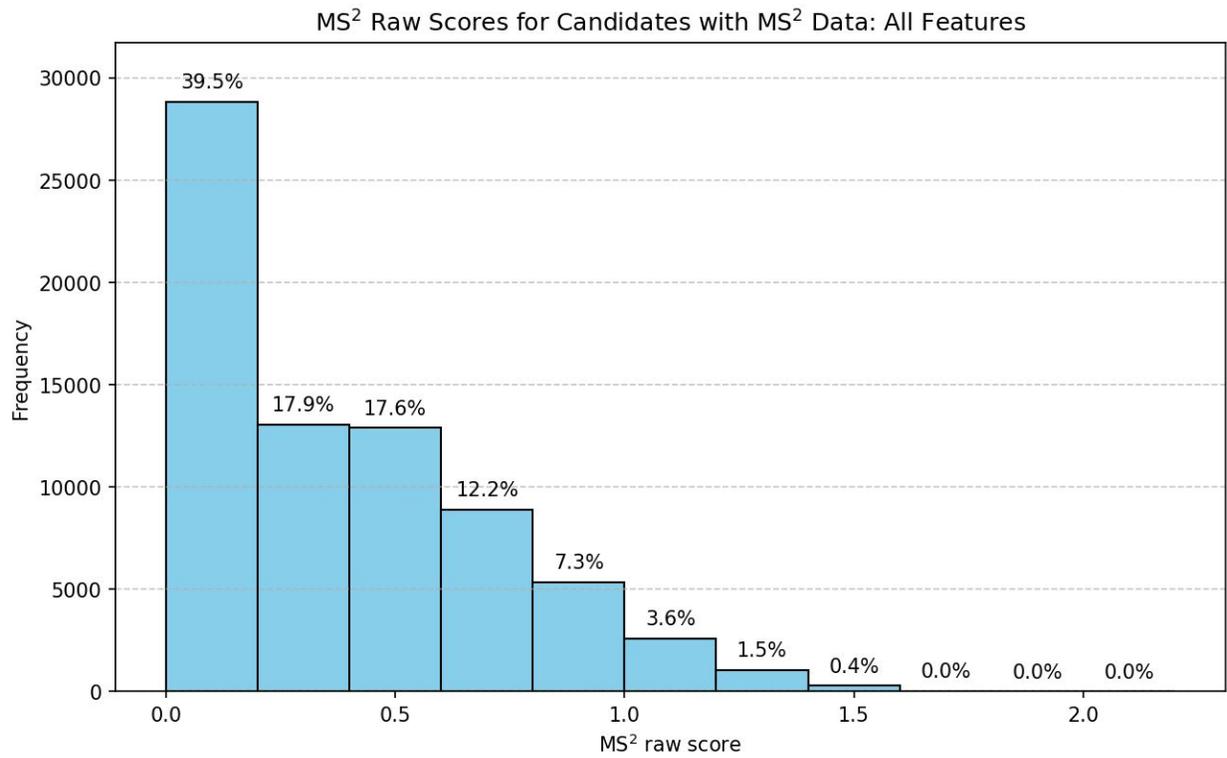
436
 437
 438
 439
 440
 441
 442

Figure S3. Absolute counts for chemical metadata fields, total metadata score, MS² score, and hazard score. Distributions are separated between correct and incorrect candidates for *known chemical features*.



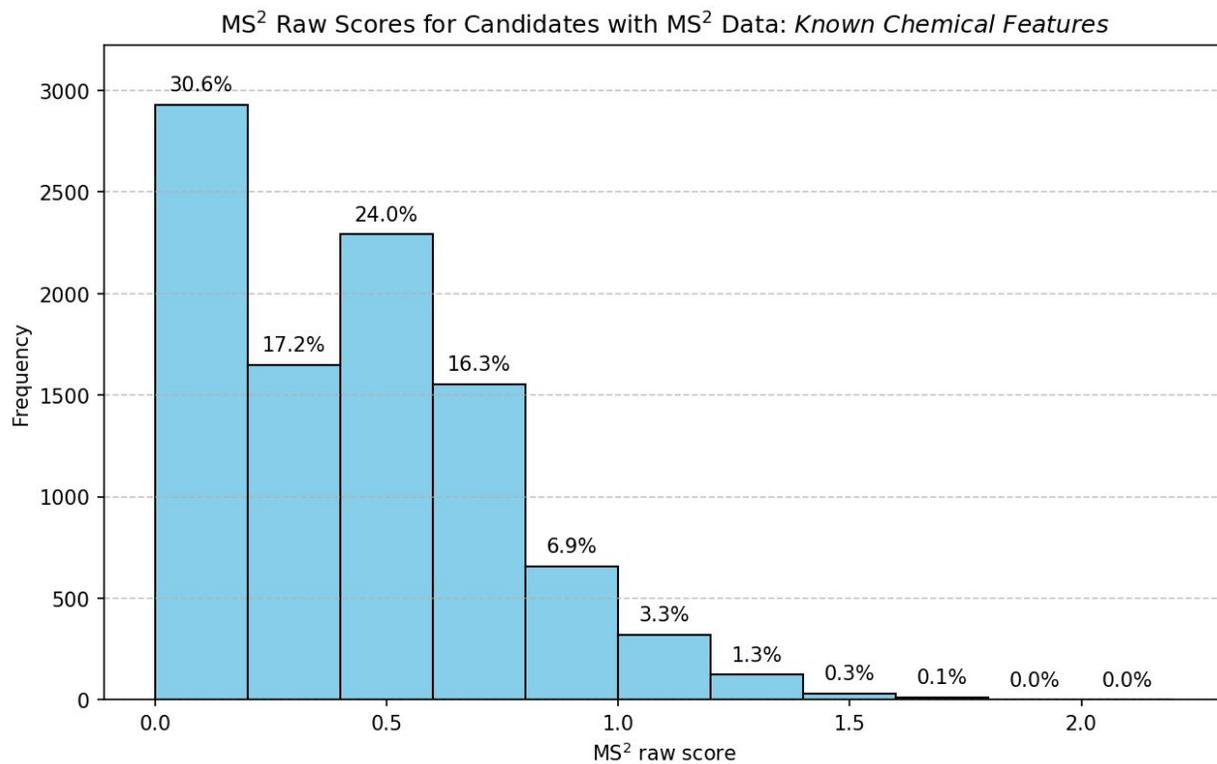
443
 444
 445 **Figure S4.** Score ratios between the correct candidate and all associated incorrect candidates for
 446 each *known chemical feature* (n=87). Ratios greater than one indicate higher scores for correct
 447 candidates. Percentage values indicate the proportion of ratios greater than, equal to, and less than
 448 one.

449
 450
 451



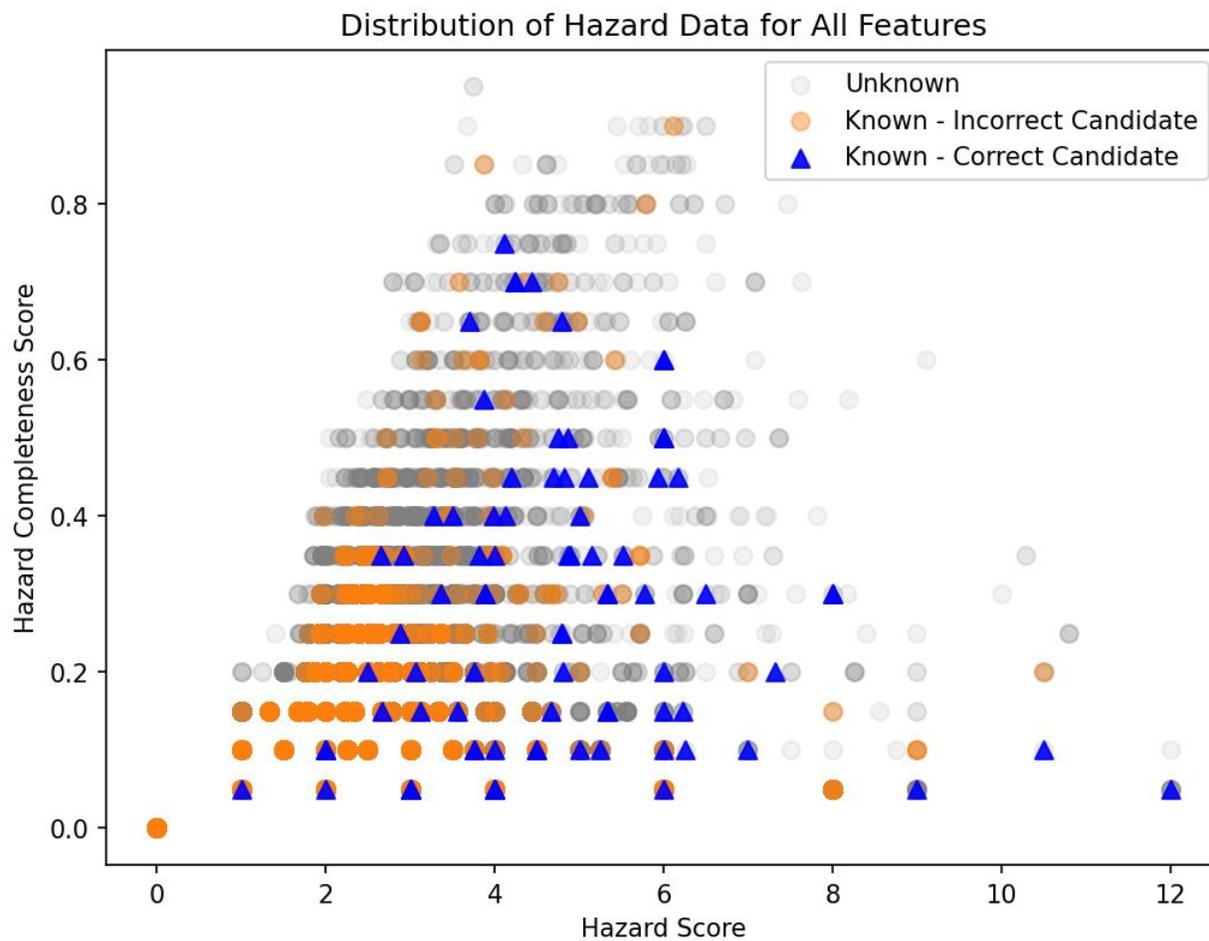
452
 453
 454
 455
 456
 457
 458
 459

Figure S5a. Distribution of MS² raw scores returned for all candidates of all features with acquired MS² data.



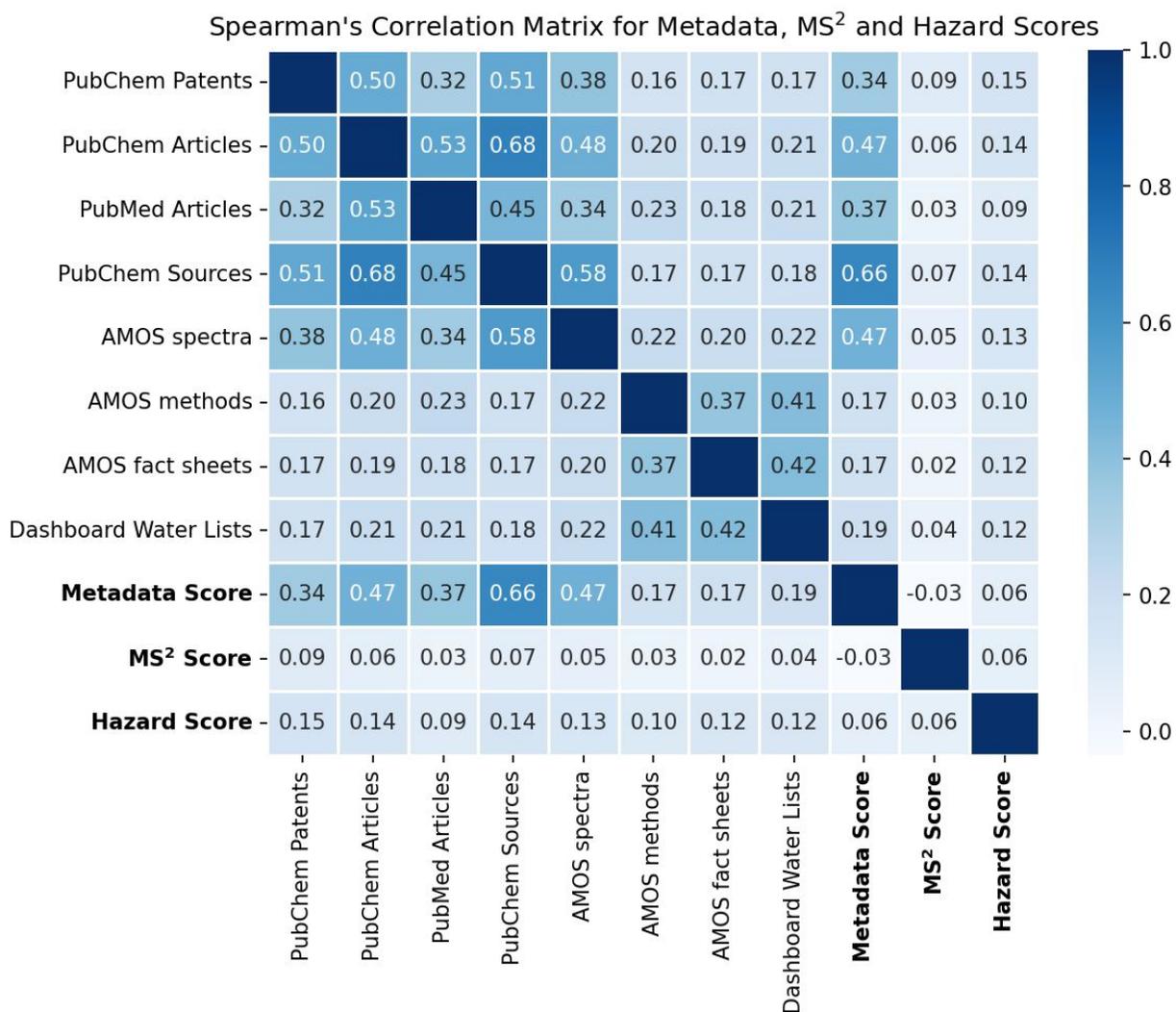
460
461
462
463
464
465
466
467

Figure S5b. Distribution of MS² raw scores returned for all candidates of *known chemical features* with acquired MS² data.



468
 469
 470
 471
 472
 473
 474
 475
 476

Figure S6. Distribution of calculated hazard scores for candidates of unknown chemical features (grey), and correct (blue) and incorrect (orange) candidates of *known chemical features*.



477
 478 **Figure S7.** Spearman correlation matrix results for absolute counts of metadata fields, and the total
 479 metadata, MS², and hazard scores, with Spearman correlation values for each field/score pair.
 480
 481
 482
 483
 484
 485
 486
 487
 488
 489
 490
 491
 492
 493
 494

495 Section 5.0 – Interactive Chemical Results Visualizations

496
497 The following sections describe the usage and functionality of interactive visualizations developed
498 within I-NTA to display, control, and interrogate chemical results generated by I-NTA. These
499 visualizations are generated using data from the chemical results sheet output of the I-NTA’s MS¹
500 workflow. Visualizations were developed in JavaScript and utilized the D3.js and AG Grid
501 JavaScript libraries.

502
503

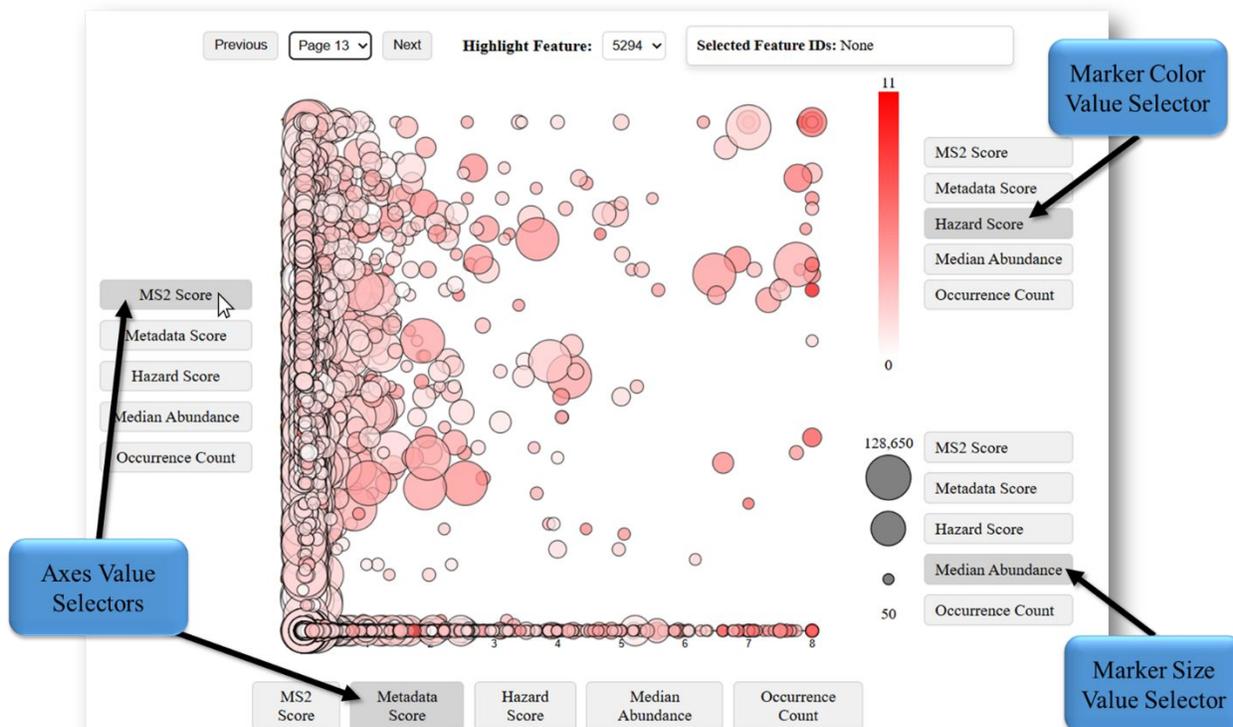
504 Section 5.1 – Interactive Scatterplot of Chemical Results

505

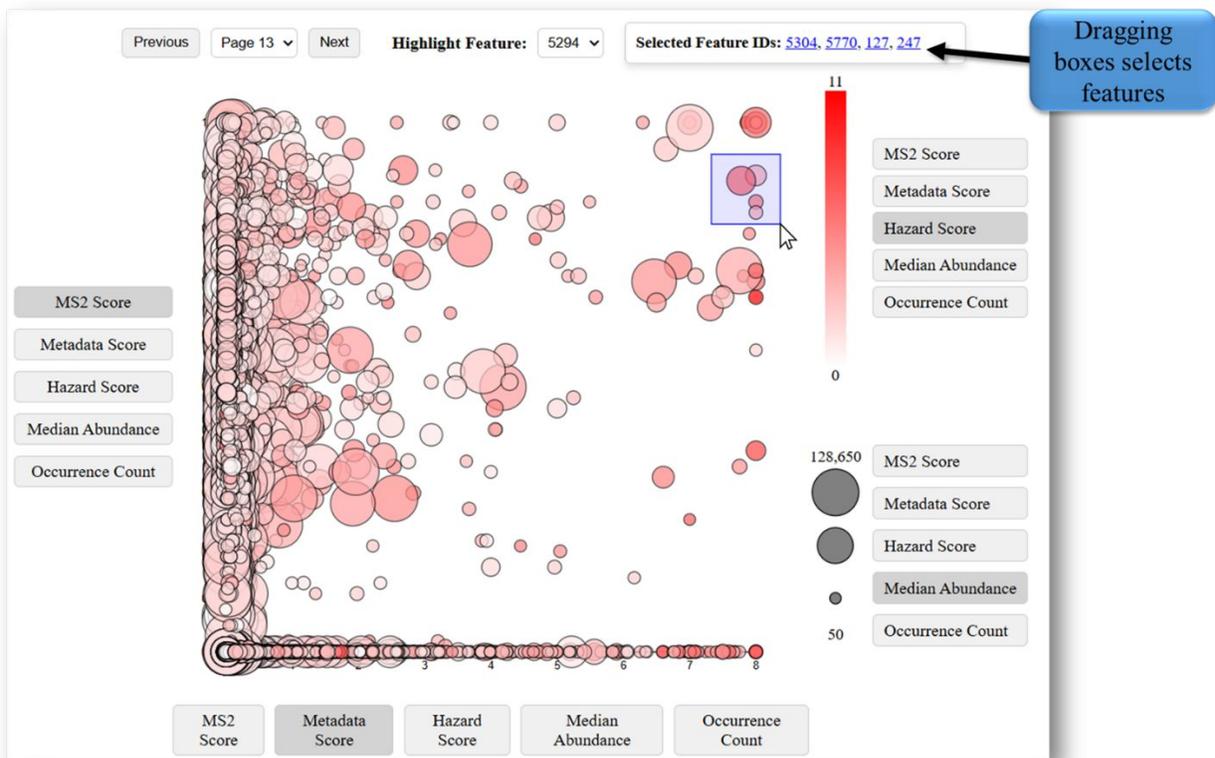


506
507 An interactive scatterplot displays the chemical results for chemical structure candidates retrieved
508 for NTA study features, where each circle marker on the plot represents a single chemical structure
509 candidate. A subset of all study features is displayed on a given page for purposes of clarity and
510 visualization speed. Users may toggle through pages of features and their respective candidates
511 using either “Previous” and “Next” buttons, or with a dropdown menu.

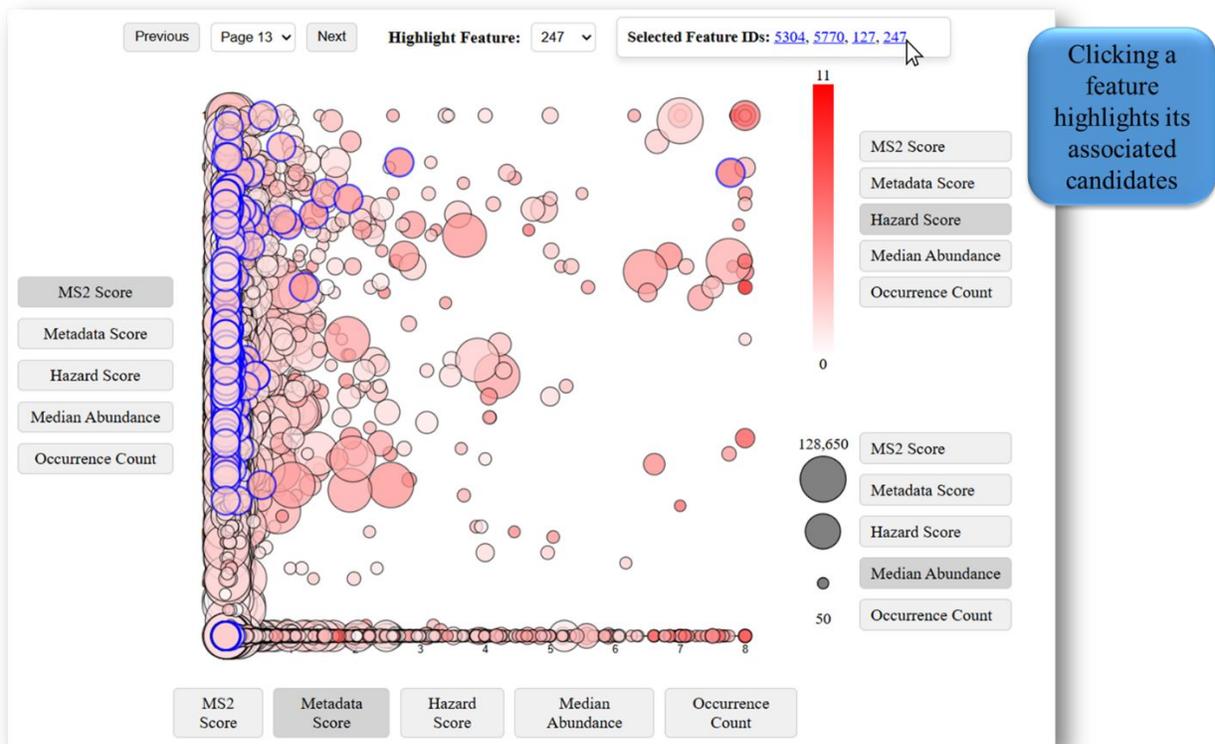
512
513



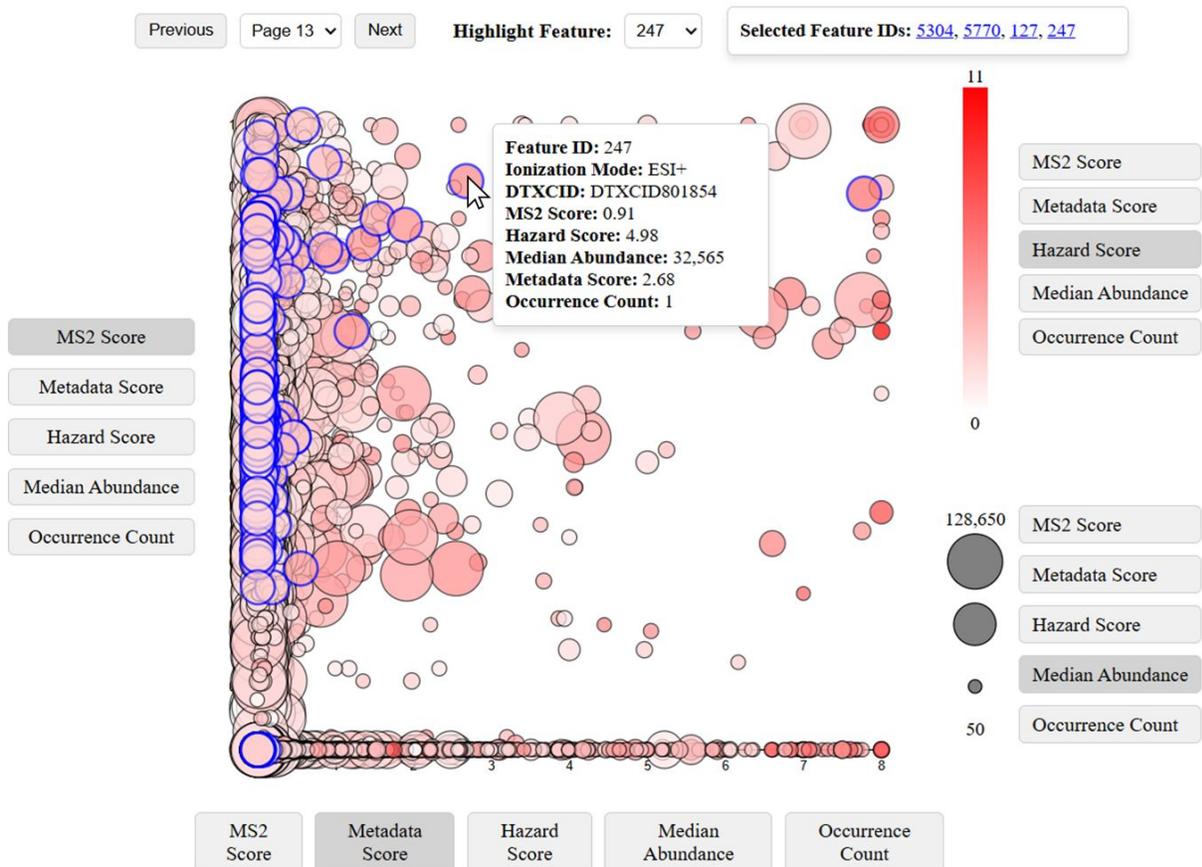
514
 515 Value selectors for the x-axis and y-axis allow users to select what chemical or feature data to plot
 516 on either plot axis, ultimately controlling how candidates are sorted and grouped in space. Value
 517 selectors for the marker color and marker size allow users to select what chemical or feature data
 518 is represented by the color and size of the marker, respectively, ultimately controlling what markers
 519 stand out visually wherever they are located.
 520
 521
 522



523
 524 After configuring parameters of the scatterplot, users may select candidates/markers and their
 525 respective features by dragging a box around data markers. This reveals features associated with
 526 the selected markers in the “Selected Feature IDs” box.
 527
 528
 529

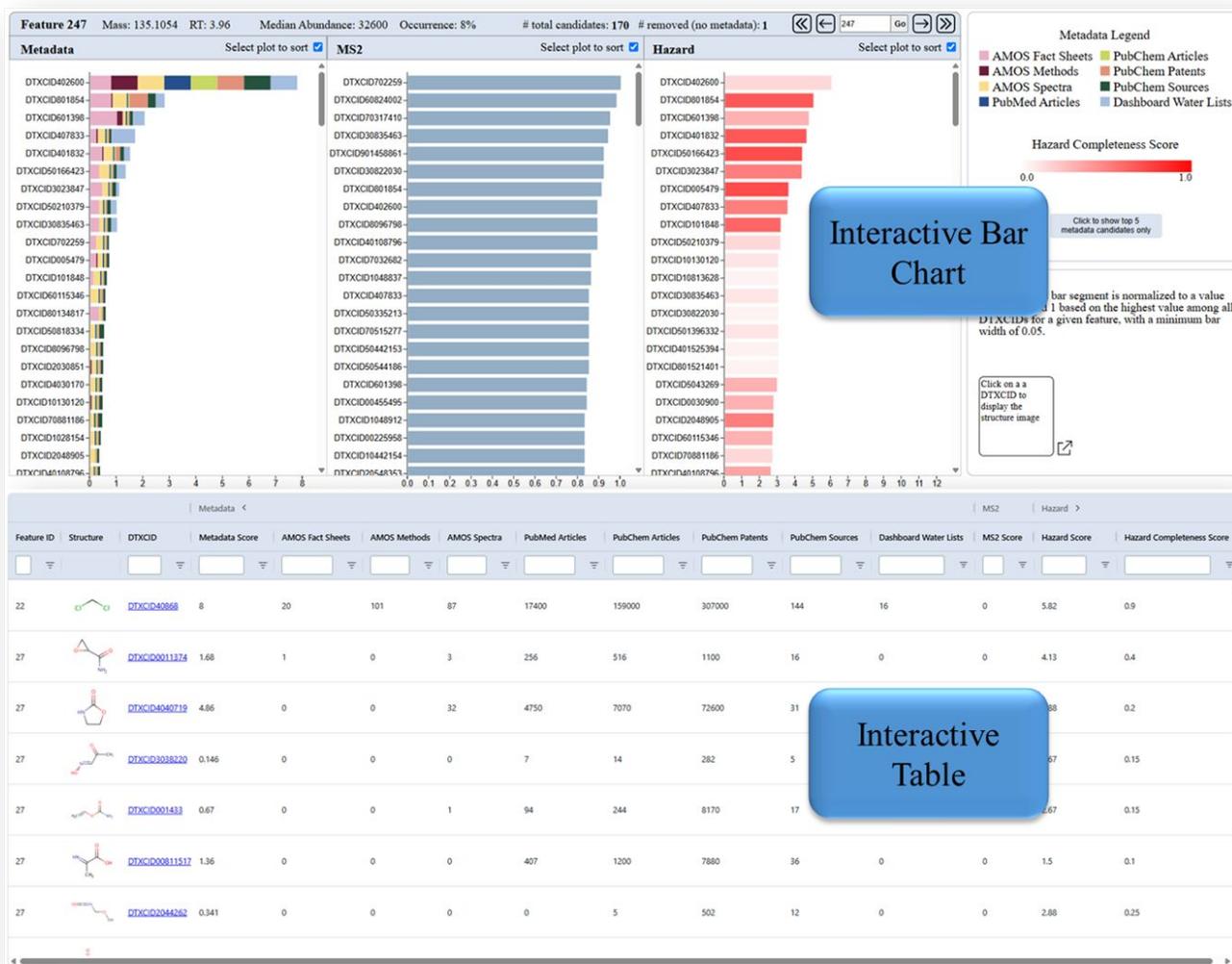


530
 531 A feature may be selected by either clicking on its ID in the “Selected Feature IDs” box or by
 532 selecting it within the “Highlight Features” dropdown menu. Selecting a feature highlights all the
 533 candidates associated with the feature.
 534
 535
 536



537
 538 Finally, mousing over any marker in the plot brings up a textbox containing the relevant feature
 539 and candidate information for that specific candidate. Altogether, the user can drive how the data
 540 are sorted spatially and highlighted via color and size to visually determine the specific data points
 541 that are of interest. To explore further, the user can mouseover markers to see underlying feature
 542 and candidate information.

Section 5.2 – Interactive Bar Chart and Grid Table of Chemical Results



A combined interactive bar chart and grid table visualization allows for further investigation of the underlying data once features and candidates of interest have been identified from the previous interactive scatterplot. The bar chart section displays chemical information graphically, whereas the grid table contains the data underlying the visuals.

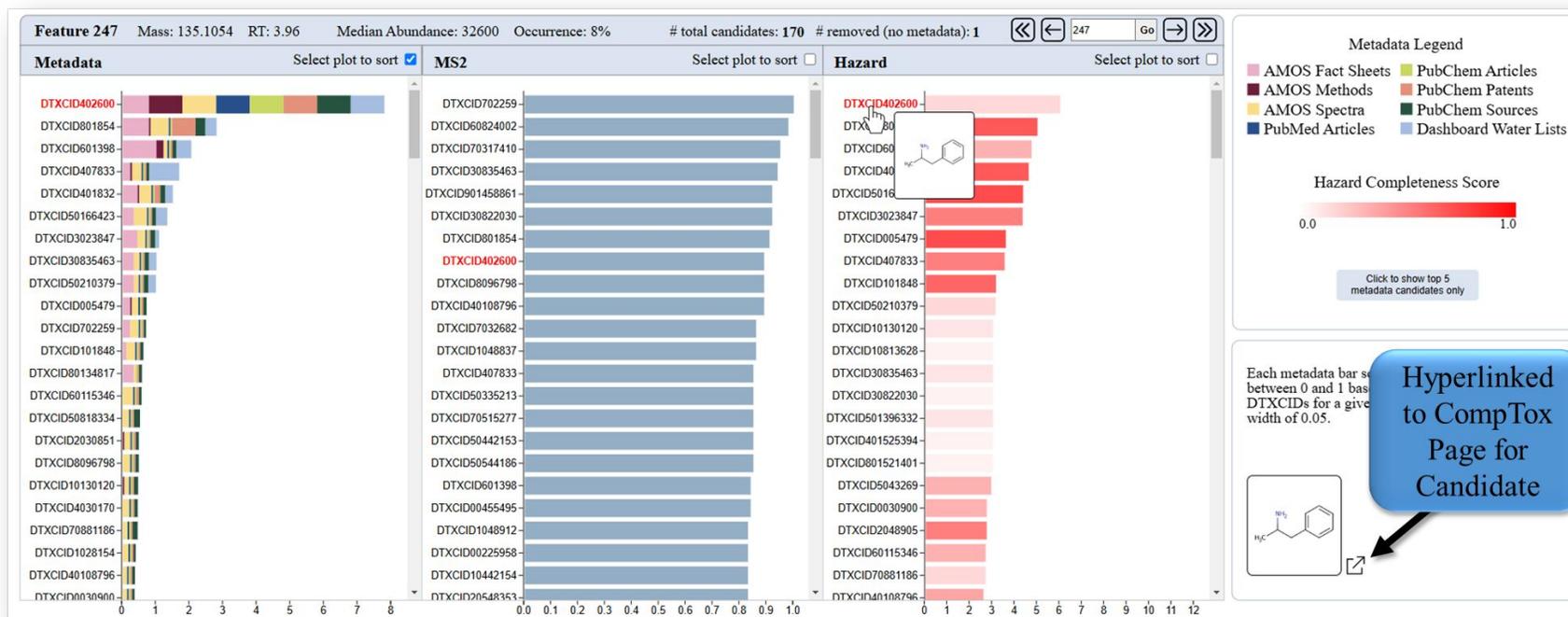


In the title block of the bar chart plot, following the feature ID, mass, retention time, median abundance and occurrence percentage, the total number of candidates for that feature and the number of those candidates with no associated metadata are displayed (NOTE: candidates with no metadata do not appear in the bar plot). There is also a feature selection tool to scroll through the features, as well as to select a specific feature ID. When a feature is selected, three bar charts are populated for the associated candidate structures of that feature; a metadata, MS², and hazard bar chart. Each bar chart's bar represents that respective category's score value. The metadata bar chart displays the normalized contributions from each metadata field broken out by color, and the hazard bar chart displays the completeness score of the hazard score value based on coloration as well. To ensure that metadata bar segments are wide enough to mouse-over, segments with a normalized value of less than 0.05 are adjusted to have a width of 0.05. A bar chart legend on the right explains the colors for these two bar charts.



The user can select the metadata fields to display in the metadata bar chart by interacting with the metadata legend. As a metadata field is selected or de-selected, the bars will re-sort in descending order of the sum of the selected metadata scores only.

Using the “Click to show top 5 metadata candidates” button beneath the legends, the user can choose to toggle between displaying all candidates for a feature or displaying only the five candidates with the highest total metadata scores. Clicking the button to show only the top 5 candidates will cause the grid to populate with only the top 5 candidates for each feature.



A mouseover of any candidate DTXCID on the y-axis will display a tooltip containing the candidate structure, and also highlight the candidate in all three bar charts. Clicking on a candidate DTXCID will bring up the structure of the candidate in a hyperlinked diagram on the lower right that will load the CompTox Chemicals Dashboard page for the specified candidate. Additionally, a mouseover of any bar segment in the plot displays the actual value for that metadata/MS²/hazard field.

The screenshot shows a data grid with columns: Feature ID, Structure, DTXCID, Metadata Sc..., AMOS Fact Sheets, AMOS Methods, AMOS Spectra, PubMed Articles, PubChem Articles, PubChem Patents, PubChem Sources, Dashboard Water Lists, MS2 Score, Hazard Score, and Hazard Completeness Score. A dropdown menu is open over the 'Metadata Sc...' column, showing options: Equals, Does not equal, Greater than, Greater than or equal to, Less than, Less than or equal to, Between, Blank, and Not blank. A blue box labeled 'Filtering options for each field' points to this menu. Another blue box labeled 'Sort fields on column header click' points to the 'PubChem Sources' column header.

Feature ID	Structure	DTXCID	Metadata Sc...	AMOS Fact Sheets	AMOS Methods	AMOS Spectra	PubMed Articles	PubChem Articles	PubChem Patents	PubChem Sources	Dashboard Water Lists	MS2 Score	Hazard Score	Hazard Completeness Score		
247		DTXCID402600	7.78				69	163	73900	101000	172000	264	9	0.89	6	0.15
247		DTXCID801854	2.68				1	2500	1200					0.91	4.98	0.65
247		DTXCID601398	1.89				15	91	707					0.84	4.72	0.3
247		DTXCID407833	1.54	2	4			58	1940					0.85	3.52	0.5
247		DTXCID401832	1.36	4	1	58	98	742	28500	40	2			0.81	4.59	0.65
247		DTXCID50166423	1.21	3	0	63	16	145	5880	32	3			0.67	4.34	0.7
247		DTXCID3023847	0.965	4	0	37	36	316	5310	39	1			0.79	4.32	0.5

The grid table portion of this visualization contains all the underlying data used to generate the bar charts. This includes all the counts for each metadata field, the MS² quotient scores, and each individual hazard endpoint's hazard and quality score. Numerical filters can be applied to any metadata column. The grid can be sorted in ascending or descending order of any of the data field values by clicking on the column headers. Sort and filter changes applied to the grid table will be reflected in the bar plot.

REFERENCES

- (1) *Dask: Library for dynamic task scheduling*; Dask Development Team, 2016. <http://dask.pydata.org>.
- (2) Grulke, C. M.; Williams, A. J.; Thillanadarajah, I.; Richard, A. M. EPA's DSSTox database: History of development of a curated chemistry resource supporting computational toxicology research. *Computational Toxicology* **2019**, *12*, 100096. DOI: <https://doi.org/10.1016/j.comtox.2019.100096>.
- (3) Williams, A. J.; Grulke, C. M.; Edwards, J.; McEachran, A. D.; Mansouri, K.; Baker, N. C.; Patlewicz, G.; Shah, I.; Wambaugh, J. F.; Judson, R. S.; 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.
- (4) McEachran, A. D.; Sobus, J. R.; Williams, A. J. Identifying known unknowns using the US EPA's CompTox Chemistry Dashboard. *Analytical and Bioanalytical Chemistry* **2017**, *409* (7), 1729-1735. DOI: 10.1007/s00216-016-0139-z.
- (5) McEachran, A. D.; Chao, A.; Al-Ghoul, H.; Lowe, C.; Grulke, C.; Sobus, J. R.; Williams, A. J. Revisiting Five Years of CASMI Contests with EPA Identification Tools. In *Metabolites*, 2020; Vol. 10.
- (6) Lowe, C. N.; Williams, A. J. Enabling High-Throughput Searches for Multiple Chemical Data Using the U.S.-EPA CompTox Chemicals Dashboard. *Journal of Chemical Information and Modeling* **2021**, *61* (2), 565-570. DOI: 10.1021/acs.jcim.0c01273.
- (7) McEachran, A. D.; Mansouri, K.; Grulke, C.; Schymanski, E. L.; Ruttkies, C.; Williams, A. J. "MS-Ready" structures for non-targeted high-resolution mass spectrometry screening studies. *Journal of Cheminformatics* **2018**, *10* (1), 45. DOI: 10.1186/s13321-018-0299-2.
- (8) McEachran, A. D.; Balabin, I.; Cathey, T.; Transue, T. R.; Al-Ghoul, H.; Grulke, C.; Sobus, J. R.; Williams, A. J. Linking *in silico* MS/MS spectra with chemistry data to improve identification of unknowns. *Scientific Data* **2019**, *6* (1), 141. DOI: 10.1038/s41597-019-0145-z.
- (9) Chao, A.; Al-Ghoul, H.; McEachran, A. D.; Balabin, I.; Transue, T.; Cathey, T.; Grossman, J. N.; Singh, R. R.; Ulrich, E. M.; Williams, A. J.; et al. *In silico* MS/MS spectra for identifying unknowns: a critical examination using CFM-ID algorithms and ENTACT mixture samples. *Analytical and Bioanalytical Chemistry* **2020**, *412* (6), 1303-1315. DOI: 10.1007/s00216-019-02351-7.
- (10) Little, J. L.; Williams, A. J.; Pshenichnov, A.; Tkachenko, V. Identification of "Known Unknowns" Utilizing Accurate Mass Data and ChemSpider. *Journal of the American Society for Mass Spectrometry* **2012**, *23* (1), 179-185. DOI: 10.1007/s13361-011-0265-y.
- (11) Richman, T.; Arnold, E.; Williams, A. J. Curation of a list of chemicals in biosolids from EPA National Sewage Sludge Surveys & Biennial Review Reports. *Scientific Data* **2022**, *9* (1), 180. DOI: 10.1038/s41597-022-01267-9.
- (12) Gaines, L. G. T.; Sinclair, G.; Williams, A. J. A proposed approach to defining per- and polyfluoroalkyl substances (PFAS) based on molecular structure and formula. *Integrated Environmental Assessment and Management* **2023**, *19* (5), 1333-1347. DOI: 10.1002/ieam.4735.
- (13) Phillips, K. A.; Yau, A.; Favela, K. A.; Isaacs, K. K.; McEachran, A.; Grulke, C.; Richard, A. M.; Williams, A. J.; Sobus, J. R.; Thomas, R. S.; et al. Suspect Screening Analysis of Chemicals in Consumer Products. *Environmental Science & Technology* **2018**, *52* (5), 3125-3135. DOI: 10.1021/acs.est.7b04781.
- (14) *Programmatic Access*. PubChem, 2025. <https://pubchem.ncbi.nlm.nih.gov/docs/programmatic-access> (accessed 2025 April 10th).

- (15) Schymanski, E. L.; Jeon, J.; Gulde, R.; Fenner, K.; Ruff, M.; Singer, H. P.; Hollender, J. Identifying Small Molecules via High Resolution Mass Spectrometry: Communicating Confidence. *Environmental Science & Technology* **2014**, *48* (4), 2097-2098. DOI: 10.1021/es5002105.
- (16) Allen, F.; Greiner, R.; Wishart, D. Competitive fragmentation modeling of ESI-MS/MS spectra for putative metabolite identification. *Metabolomics* **2015**, *11* (1), 98-110. DOI: 10.1007/s11306-014-0676-4.
- (17) Wang, F.; Liigand, J.; Tian, S.; Arndt, D.; Greiner, R.; Wishart, D. S. CFM-ID 4.0: More Accurate ESI-MS/MS Spectral Prediction and Compound Identification. *Analytical Chemistry* **2021**, *93* (34), 11692-11700. DOI: 10.1021/acs.analchem.1c01465.
- (18) Rager, J. E.; Strynar, M. J.; Liang, S.; McMahan, R. L.; Richard, A. M.; Grulke, C. M.; Wambaugh, J. F.; Isaacs, K. K.; Judson, R.; Williams, A. J.; et al. Linking high resolution mass spectrometry data with exposure and toxicity forecasts to advance high-throughput environmental monitoring. *Environment International* **2016**, *88*, 269-280. DOI: 10.1016/j.envint.2015.12.008.
- (19) Newton, S. R.; McMahan, R. L.; Sobus, J. R.; Mansouri, K.; Williams, A. J.; McEachran, A. D.; Strynar, M. J. Suspect screening and non-targeted analysis of drinking water using point-of-use filters. *Environmental Pollution* **2018**, *234*, 297-306. DOI: 10.1016/j.envpol.2017.11.033.
- (20) Dionisio, K. L.; Phillips, K.; Price, P. S.; Grulke, C. M.; Williams, A. J.; Biryol, D.; Hong, T.; Isaacs, K. K. The Chemical and Products Database, a resource for exposure-relevant data on chemicals in consumer products. *Scientific Data* **2018**, *5* (1), 180125. DOI: 10.1038/sdata.2018.125.
- (21) Isaacs, K. K.; Wall, J. T.; Williams, A. R.; Hobbie, K. A.; Sobus, J. R.; Ulrich, E.; Lyons, D.; Dionisio, K. L.; Williams, A. J.; Grulke, C.; et al. A harmonized chemical monitoring database for support of exposure assessments. *Scientific Data* **2022**, *9* (1), 314. DOI: 10.1038/s41597-022-01365-8.
- (22) Richard, A. M.; Judson, R. S.; Houck, K. A.; Grulke, C. M.; Volarath, P.; Thillainadarajah, I.; Yang, C.; Rathman, J.; Martin, M. T.; Wambaugh, J. F.; et al. ToxCast Chemical Landscape: Paving the Road to 21st Century Toxicology. *Chemical Research in Toxicology* **2016**, *29* (8), 1225-1251. DOI: 10.1021/acs.chemrestox.6b00135.
- (23) Feshuk, M.; Kolaczkowski, L.; Watford, S.; Paul Friedman, K. ToxRefDB v2.1: update to curated in vivo study data in the Toxicity Reference Database. *Frontiers in Toxicology* **2023**, *5*, Mini Review.
- (24) Wall, T.; Sayre, R.; Smith, D.; Winter, S.; Groover, M.; Hope, J.; Webb, A.; Friedman, K.; Feshuk, M.; Williams, A.; et al. Development of the Toxicity Values Database, ToxValDB: A Curated Resource for Experimental and Derived Human Health-Relevant Toxicity Data. *Submitted* **2025**.
- (25) Vegosen, L.; Martin, T. M. An automated framework for compiling and integrating chemical hazard data. *Clean Technologies and Environmental Policy* **2020**, *22* (2), 441-458. DOI: 10.1007/s10098-019-01795-w.
- (26) Brunelle, L. D.; Batt, A. L.; Chao, A.; Glassmeyer, S. T.; Quinete, N.; Alvarez, D. A.; Kolpin, D. W.; Furlong, E. T.; Mills, M. A.; Aga, D. S. De facto Water Reuse: Investigating the Fate and Transport of Chemicals of Emerging Concern from Wastewater Discharge through Drinking Water Treatment Using Non-targeted Analysis and Suspect Screening. *Environmental Science & Technology* **2024**, *58* (5), 2468-2478. DOI: 10.1021/acs.est.3c07514.
- (27) Sobus, J. R.; Sayre-Smith, N. A.; Chao, A.; Ferland, T. M.; Minucci, J. M.; Carr, E. T.; Brunelle, L. D.; Batt, A. L.; Whitehead, H. D.; Cathey, T.; et al. Automated QA/QC reporting for

non-targeted analysis: a demonstration of “INTERPRET NTA” with de facto water reuse data. *Analytical and Bioanalytical Chemistry* **2025**. DOI: 10.1007/s00216-025-05771-w.

(28) *Data Sources*. PubChem, <https://pubchem.ncbi.nlm.nih.gov/source> (accessed 2025 April 10th).

(29) Stein, S. E.; Scott, D. R. Optimization and testing of mass spectral library search algorithms for compound identification. *Journal of the American Society for Mass Spectrometry* **1994**, *5* (9), 859-866. DOI: 10.1016/1044-0305(94)87009-8.

(30) Newmeyer, M. N.; Lyu, Q.; Sobus, J. R.; Williams, A. J.; Nachman, K. E.; Prasse, C. Combining Nontargeted Analysis with Computer-Based Hazard Comparison Approaches to Support Prioritization of Unregulated Organic Contaminants in Biosolids. *Environmental Science & Technology* **2024**, *58* (27), 12135-12146. DOI: 10.1021/acs.est.4c02934.

(31) Janesch, G.; Carr, E. T.; Sivasupramaniam, S.; Charest, N.; Williams, A. J. Applying Cheminformatics to Develop a Structure Searchable Database of Analytical Methods. *Submitted*. **2025**.

(32) Sloop, J. T.; Chao, A.; Gundersen, J.; Phillips, A. L.; Sobus, J. R.; Ulrich, E. M.; Williams, A. J.; Newton, S. R. Demonstrating the Use of Non-targeted Analysis for Identification of Unknown Chemicals in Rapid Response Scenarios. *Environ Sci Technol* **2023**, *57* (8), 3075-3084. DOI: 10.1021/acs.est.2c06804 From NLM.

(33) Brueck, C. L.; Xin, X.; Lupolt, S. N.; Kim, B. F.; Santo, R. E.; Lyu, Q.; Williams, A. J.; Nachman, K. E.; Prasse, C. (Non)targeted Chemical Analysis and Risk Assessment of Organic Contaminants in Darkibor Kale Grown at Rural and Urban Farms. *Environmental Science & Technology* **2024**, *58* (8), 3690-3701. DOI: 10.1021/acs.est.3c09106.

(34) Ring, C. L.; Arnot, J. A.; Bennett, D. H.; Egeghy, P. P.; Fantke, P.; Huang, L.; Isaacs, K. K.; Jolliet, O.; Phillips, K. A.; Price, P. S.; et al. Consensus Modeling of Median Chemical Intake for the U.S. Population Based on Predictions of Exposure Pathways. *Environmental Science & Technology* **2019**, *53* (2), 719-732. DOI: 10.1021/acs.est.8b04056.

(35) Wambaugh, J. F.; Bare, J. C.; Carignan, C. C.; Dionisio, K. L.; Dodson, R. E.; Jolliet, O.; Liu, X.; Meyer, D. E.; Newton, S. R.; Phillips, K. A.; et al. New approach methodologies for exposure science. *Current Opinion in Toxicology* **2019**, *15*, 76-92. DOI: <https://doi.org/10.1016/j.cotox.2019.07.001>.