

1      **Development of Chemical Categories for Per- and Polyfluoroalkyl  
2      Substances (PFAS) and the Proof-of-Concept Approach to the  
3      Identification of Potential Candidates for Tiered Toxicological  
4      Testing and Human Health Assessment**

5      Grace Patlewicz<sup>a,\*</sup>, Richard Judson<sup>a</sup>, Antony Williams<sup>a</sup>, Tristan Butler<sup>b</sup>, Stan Barone Jr.<sup>b</sup>, Kelly  
6      Carstens<sup>a</sup>, John Cowden<sup>a</sup>, Jeffrey L. Dawson<sup>b</sup>, Sigmund Degitz<sup>a</sup>, Kellie Fay<sup>b</sup>, Tala R. Henry<sup>b,1</sup>, Anna  
7      Lowit<sup>b</sup>, Stephanie Padilla<sup>a</sup>, Katie Paul Friedman<sup>a</sup>, Martin B. Phillips<sup>b</sup>, David Turk<sup>b</sup>, John Wambaugh<sup>a</sup>,  
8      Barbara Wetmore<sup>a</sup>, Russell S. Thomas<sup>a</sup>

<sup>a</sup>Center for Computational Toxicology and Exposure (CCTE), US Environmental Protection Agency, Durham, USA,

<sup>b</sup>Office of Chemical Safety and Pollution Prevention (OSCPP), US Environmental Protection Agency, DC, USA,

---

9      **Abstract**

Per- and Polyfluoroalkyl substances (PFAS) are a class of manufactured chemicals that are in widespread use and many present concerns for persistence, bioaccumulation and toxicity. Whilst a handful of PFAS have been characterized for their hazard profiles, the vast majority have not been extensively studied. Herein, a chemical category approach was developed and applied to PFAS that could be readily characterized by a chemical structure. The PFAS definition as described in the Toxic Substances Control Act (TSCA) section 8(a)(7) rule was applied to the Distributed Structure-Searchable Toxicity (DSSTox) database to retrieve an initial list of 13,054 PFAS. Plausible degradation products from the 563 PFAS on the non-confidential TSCA Inventory were simulated using the Catalogic expert system, and the unique predicted PFAS degradants (2484) that conformed to the same PFAS definition were added to the list resulting in a set of 15,538 PFAS. Each PFAS was then assigned into a primary category using Organisation for Economic Co-operation and Development (OECD) structure-based classifications. The primary categories were subdivided into secondary categories based on a chain length threshold ( $\geq 7$  vs  $< 7$ ). Secondary categories were subcategorized using chemical fingerprints to achieve a balance between total number of structural categories vs. level of structural similarity within a category based on the Jaccard index. A set of 128 terminal structural categories were derived from which a subset of representative candidates could be proposed for potential data collection, considering the sparsity of relevant toxicity data within each category, presence on environmental monitoring lists, and the ability to identify plausible manufacturers/importers. Refinements to the approach taking into consideration ways in which the categories could be updated by mechanistic data and physicochemical property information are also described. This categorization approach may be used to form the basis of identifying candidates for data collection with related applications in QSAR development, read-across and hazard assessment.

10     **Keywords:** Per- and polyfluoroalkyl substances (PFAS), chemical categories, read-across, New  
11     Approach Methods (NAMs), data collection, Toxic Substances Control Act (TSCA)

---

\*Corresponding author  
Email address: patlewicz.grace@epa.gov (Grace Patlewicz)

<sup>1</sup>Retired

12    **1. Introduction**

13    **1.1. Background**

14    Per- and Polyfluoroalkyl substances (PFAS) are a large class of human-made chemicals that have been  
15    manufactured and used in a variety of industries since the 1940s<sup>1,2,3</sup>. PFAS have been or are currently  
16    being synthesized for a myriad of different uses, including adhesives, stain resistant coatings for  
17    clothes or furniture and fire retardants. In addition to consumer and industrial applications, PFAS  
18    are being released into the environment during manufacturing and use<sup>4</sup>. PFAS and products containing  
19    them are regularly disposed of in landfills or incinerated, which can also lead to further release into  
20    soil, groundwater, and air<sup>5,6</sup>. They are also found in biosolids from wastewater treatment facilities  
21    which have been spread onto agricultural fields<sup>7</sup>.

22    Characterizing the scope and scale of the 'PFAS class' has been challenging in the absence of a har-  
23    monized PFAS definition. Some publications have cited thousands of PFAS being in the environment  
24    (estimates range from 4700<sup>8</sup> to greater than 10,000<sup>9</sup>), but there is likely to be an increasing number  
25    given that analytical methods are continually being evolved to detect them. An Organisation for Eco-  
26    nomic Co-operation and Development (OECD) working group defined PFAS as 'fluorinated substances  
27    that contain at least one fully fluorinated methyl or methylene carbon atom (without any H/Cl/Br/I  
28    atom attached to it); that is, any chemical with at least a perfluorinated methyl group (-CF<sub>3</sub>) or a  
29    perfluorinated methylene group (-CF<sub>2</sub>-)<sup>8,10</sup>. This broad OECD definition would make estimates of a  
30    few thousand PFAS too low; however, the OECD working group also acknowledges that a chemistry  
31    definition of PFAS does not necessarily equate to how PFAS should be assessed in terms of their  
32    hazard profile or to what extent subcategorizations of PFAS are appropriate depending on differ-  
33    ent legislative frameworks. Indeed, if the OECD definition were applied to a large inventory such as  
34    the US EPA's Distributed Structure-Searchable (DSSTox) Database<sup>11</sup> estimates of the number of  
35    PFAS would be in the order of 30,000. For contrast, the [PubChem Classification Browser](#) has tagged  
36    over 7 million substances as meeting the OECD PFAS definition<sup>12</sup>. This would imply that any sub-  
37    stance containing a CF<sub>3</sub> would be classified as a "PFAS" even though it might fall within the purview  
38    of different regulatory frameworks. The US EPA's Office of Pollution Prevention and Toxics (OPPT)  
39    recently finalized a structural definition of PFAS applicable to several Toxic Substances Control Act  
40    (TSCA) activities: the Significant New Use Rule (SNUR) on PFAS designated as inactive on the TSCA  
41    inventory<sup>13</sup>, the TSCA Report and Recordkeeping Requirements for Perfluoroalkyl and Polyfluoroalkyl  
42    Substances rule (referred to herein as the TSCA section 8(a)(7) rule)<sup>14</sup> and EPA's recently released  
43    framework for assessing PFAS under TSCA's New Chemicals activities<sup>15</sup>. For these TSCA actions, a  
44    PFAS is defined as 'including at least one of three substructures: 1) R-(CF<sub>2</sub>)-CF(R')R", where both the  
45    CF<sub>2</sub> and CF moieties are saturated carbons; 2) R-CF<sub>2</sub>OCF<sub>2</sub>-R', where R and R' can either be F, O, or

46 saturated carbons; or 3)  $\text{CF}_3\text{C}(\text{CF}_3)\text{R}'\text{R}''$ , where R' and R'' can either be F or saturated carbons. This  
47 definition is narrower in scope than the OECD chemistry definition yet EPA estimates that it would  
48 still identify several thousand PFAS.

49 Of the many thousands of PFAS, few have been studied extensively in terms of their toxicity pro-  
50 file. Beyond a handful of closely-related perfluoroalkyl acids, perfluoroalkyl sulfonates, and perflu-  
51 roalkyl ethers (e.g., perfluoroocanoic acid (PFOA), perfluorooctane sulfonic acid (PFOS), and hexaflu-  
52 oropropylene oxide dimer acid, HFPO-DA), the vast majority of PFAS lack data to facilitate a robust  
53 characterization of their potential toxicity<sup>16</sup>. In an effort to address these data gaps, Congress di-  
54 rected EPA (15 USC 8962) to develop a process for prioritizing which PFAS or 'class' of PFAS should  
55 be subject to additional research efforts based on potential for human exposure, potential toxicity,  
56 and other available information. In response, the EPA published the [EPA National PFAS Testing Strat-](#)  
57 [egy](#) in October 2021 which describes EPA's approach to developing categories of PFAS and identifying  
58 substances for further data collection efforts.

59 The notion of a 'class' underpins grouping approaches which includes the concept of developing cat-  
60 egories to perform associated read-across. Rather than assessing each PFAS individually, closely  
61 related PFAS could be, in principle, grouped together into categories. Thus, in a category approach,  
62 not every PFAS needs to be tested for every single endpoint. Instead, the overall data for that cat-  
63 egory could potentially prove applicable to support a hazard assessment for other members of the  
64 category.

65 Grouping approaches have been in use in regulatory programmes for many years dating back to 1998  
66 when guidance was developed by the EPA in support of the [US High Production Volume \(HPV\) Chal-](#)  
67 [lenge Program](#)<sup>17</sup>. The concepts of grouping, categories and read-across are defined and extensively  
68 described in OECD's grouping guidance document, last revised in 2014<sup>18</sup> and presently undergoing  
69 revision. Moreover, the state of the art in read-across has also been described extensively in the  
70 literature; from workflows which outline the steps undertaken to develop category and analogue ap-  
71 proaches through to the evaluation, justification and documentation of any read-across predictions  
72 made<sup>19,20,21,22</sup>. More recently the notion of enhancing structure-based groupings with new approach  
73 methods (NAMs) has also been an evolving topic. For example, of keen interest is the extent to  
74 which structural categories can be further justified by NAM data by providing a mechanistic under-  
75 pinning<sup>20,21,23,22</sup>. NAMs are defined as any technology, methodology, approach, or combination that  
76 can provide information on chemical hazard and risk assessment without the use of animals, including  
77 in silico, in chemico, in vitro, and ex vivo approaches<sup>24</sup>. Of note, EPA has been leading a research  
78 programme to test a targeted set of ~150 PFAS through an array of different NAM approaches as  
79 part of a category approach<sup>25,26,27,28,23,29,30,31,32</sup>.

80 This study describes the approach taken to further refine a relevant PFAS landscape to EPA from  
81 which an initial set of structural categories were derived. The work here is a continuation of the  
82 initial categorization efforts described in the [EPA National PFAS Testing Strategy](#). For the cate-  
83 gories developed, data gaps were assessed to help identify which categories were particularly data  
84 poor (e.g., lacking relevant repeated dose toxicity data) and/or associated with known exposures and  
85 therefore would benefit from data collection or new test data generation (using both NAMs or tradi-  
86 tional approaches) to better characterize the category as a whole. The aims of this manuscript are  
87 as follows:

- 88 1. Summarize the process of constructing a PFAS landscape;
- 89 2. Profile the PFAS landscape to assign substances into broad structural categories in combination  
90 with chain length;
- 91 3. Evaluate the degree of structural similarity within each category and determine which cate-  
92 gories needed to be further subset to maximize their structural similarity whilst maintaining a  
93 pragmatic total number of categories;
- 94 4. Facilitate the identification of potential candidate PFAS for data collection by capturing addi-  
95 tional considerations such as availability of a known manufacturer/importer (who would be re-  
96 sponsible for conducting testing via TSCA); EPA Agency and/or State priorities, environmental  
97 monitoring information and structural diversity within the category;
- 98 5. Evaluate the categories based on their predicted physical state and physicochemical properties  
99 (a context of evaluating the similarity within the category and informing on potential technical  
100 limitations for testing);
- 101 6. Consider the utility of the structural categories developed in performing read-across, as well  
102 as refinements such as incorporating mechanistic and toxicokinetic data derived from NAMs.  
103 The mechanistic insights derived from EPA's parallel research effort on selected PFAS offer  
104 potential opportunities to refine the structurally-based categories developed.
- 105 7. Evaluate the feasibility of operationalizing the structural categories so that new PFAS can be  
106 profiled and assigned into one of the terminal categories developed.

107 **2. Methods**

108 **2.1. Defining the PFAS landscape**

109 To define the PFAS landscape for the purpose of this study, the DSSTox database<sup>11,33</sup> was searched  
110 using a series of structure-based queries that reflected the PFAS structural definition described  
111 earlier (see Section 1.1). DSSTox forms the basis of the EPA CompTox Chemicals Dashboard (re-  
112 ferred to herein as the Dashboard)<sup>11,33</sup> and comprises 1,218,248 substances (at the time of writ-

113 ing, May 2024, <https://comptox.epa.gov/dashboard/>). As a result of the search, 13,054 substances  
114 were identified as forming the initial PFAS landscape for this study. This landscape is available as a  
115 list published on the Dashboard at [PFAS8a7v3](#). This set was then cross referenced with the TSCA  
116 inventory (see Section 2.4.1) to identify matches. The TSCA inventory is the list of chemical sub-  
117 stances in commerce (manufactured, processed or imported) in the US since January 1975 that do  
118 not qualify for an exemption or exclusion under TSCA (TSCA inventory; Section 2.10.1). Note that  
119 EPA maintains two TSCA inventories - one that is publicly available and another that contains confi-  
120 dential business information (CBI) and is not publicly available. Those substances reported to EPA  
121 as in commerce since June 2006 are designated as "active" on the TSCA inventory. For each of  
122 the PFAS (active and inactive) listed on the publicly available TSCA inventory, degradation products  
123 were simulated using the biodegradation model, Catalogic 301C v13.18 within the commercial soft-  
124 ware tool, OASIS Catalogic v5.16.1.10 (University As Zlatarov, Laboratory of Mathematical Chem-  
125 istry, Bourgas, Bulgaria; <http://oasis-lmc.org/>). The intent was to enrich the landscape for PFAS  
126 likely to be found in the environment that originated from substances in commerce. The set of  
127 PFAS degradation products (2484) for the parent TSCA substances were added to the initial land-  
128 scape such that the final PFAS landscape used in this study comprised 15,538 substances. Note only  
129 degradation products meeting the PFAS definition were considered. Chemicals were represented  
130 by unique DSSTox Substance Identifiers (DTXSID)<sup>11</sup>, Simplified Molecular-Input-Line-Entry Sys-  
131 tem (SMILES) (<https://www.daylight.com/dayhtml/doc/theory/theory.smiles.html>), chemical names  
132 and CAS Registry Numbers (CASRN). International Chemical Identifier keys (InChIKeys), (hashed  
133 InChI)<sup>34</sup> were used as identifiers for the degradation products. Chemical substances in the DSSTox  
134 database have been curated and standardized to ensure correctness in chemical structure as well  
135 as their associations to chemical names and other identifiers such as CASRN. Examples of this cu-  
136 ration include checking for errors and mismatches in chemical structure formats and mapping to  
137 identifiers, as well as structure validation and/or standardization issues such as hyper-valency, tau-  
138 tomerism, etc<sup>11</sup>.

## 139 2.2. Biodegradation potential

140 Biodegradation predictions were made for PFAS in the landscape that were on the TSCA inven-  
141 tory using the Catalogic 301C v13.18 model within the commercial software tool, OASIS Catalogic  
142 v5.16.1.10. The biodegradation Catalogic 301C model simulates aerobic biodegradation under Ministry  
143 of International Trade and Industry, Japan (MITI) I (OECD 301C) test conditions. The modelled  
144 endpoint is the percentage of theoretical biological oxygen demand (BOD) on day 28. The underlying  
145 training set for the model comprises BOD data for 2618 substances - 745 of these were collected  
146 from the MITI I database and 804 were provided by National Institute of Technology and Evaluation

147 (NITE), Japan. A further 1069 substances that were proprietary were provided by NITE, Japan. The  
148 training set includes 797 readily biodegradable and 1821 not readily biodegradable substances. In ad-  
149 dition to BOD data, a second database underpinning the model comprised pathways for 845 organic  
150 substances, documented pathways for 649 chemicals were collected from the primary and secondary  
151 literature whereas pathways for 196 proprietary substances were provided by NITE, Japan. In brief,  
152 the Catalogic model comprises a metabolic simulator and an endpoint model. The microbial metabolism  
153 is simulated by a rule-based approach based on a set of hierarchically ordered transformations and  
154 a system of rules controlling the application of these transformations. Recursive application of the  
155 transformations allows for the simulation of metabolism and generation of biodegradation pathways.  
156 Calculation of the modelled endpoint is based on the simulated metabolic tree and the material balance  
157 of transformations used to build the tree. Predictions were made for all PFAS in the landscape that  
158 were on the non-confidential TSCA inventory (see Section 2.4.1 for more details). Prediction results  
159 containing the list of simulated metabolites (as SMILES) along with their parent DTXSID identifiers  
160 were exported as a text file. Prediction results were then processed in the following manner:

- 161 1. DTXSID identifiers were extracted for each parent substance and mapped to each metabolite.  
162 This ensured for a given parent, all metabolites could be readily associated with its corresponding  
163 parent substance.
- 164 2. A new identifier was then created for the metabolites based on the parent DTXSID identi-  
165 fier. That is to say, the first listed metabolite simulated for parent DTXSID9065256 would be  
166 tagged as DTXSID9065256\_m\_1 and so on.
- 167 3. InChIKeys were then generated for all SMILES, parents and simulated metabolites. Use of  
168 InChIKeys provided an unambiguous means of structurally representing the substance (rather  
169 than using SMILES that are potentially non-unique) and enabled subsequent associations to be  
170 derived between substances. The processed results were saved for subsequent analysis.

171 Many degradation products were found to be common across parent substances. Grouping by  
172 InChIKeys created a set of unique degradation products. These were filtered to remove non PFAS  
173 degradation products or those not meeting the PFAS definition. A final step involved cross matching  
174 the degradates against the starting landscape to remove any duplicate entries i.e. if a degradate was  
175 a substance already captured. These steps resulted in a set of 2484 degradation products that were  
176 then added to the starting landscape of 13,054 substances.

177 To explore the coverage and relevance of the MITI training set (the non-proprietary portion) within  
178 the Catalogic 301C model relative to the PFAS on the TSCA inventory substances, a comparison was  
179 performed to assess the overlap in structural space as characterized by Morgan chemical finger-

prints<sup>35</sup> (see Section 2.3.4 for details on chemical fingerprint generation). In the latter case, this structural space was projected onto a 2-dimensional (2D) scatterplot (see Figure A1) using a Uniform Manifold Approximation and Projection (UMAP) to facilitate visualization<sup>36</sup>. This is a dimensionality reduction technique that assumes available data samples are evenly distributed across a topological space (manifold) which can be approximated from these finite data samples and mapped to a lower dimensional space. In essence, UMAP learns the manifold in the high dimensional space, in this case, these are the 1024 chemical fingerprints and aims to find a 2D representation of the same manifold. The Catalogic model also provided an indication of whether any of the PFAS profiled were part of the training set as well as whether they were within the structural domain of applicability<sup>37</sup>.

### 189 2.3. Profiling PFAS into structural categories

This study aimed to develop a hierarchy of PFAS categories starting with a handful of large, diverse categories that could be subcategorized into more structurally similar categories based on other considerations (e.g., chain length and chemical fingerprints). The conceptual workflow for creating the PFAS structural categories is summarized in Figure 1 and the details of each step are described in turn.

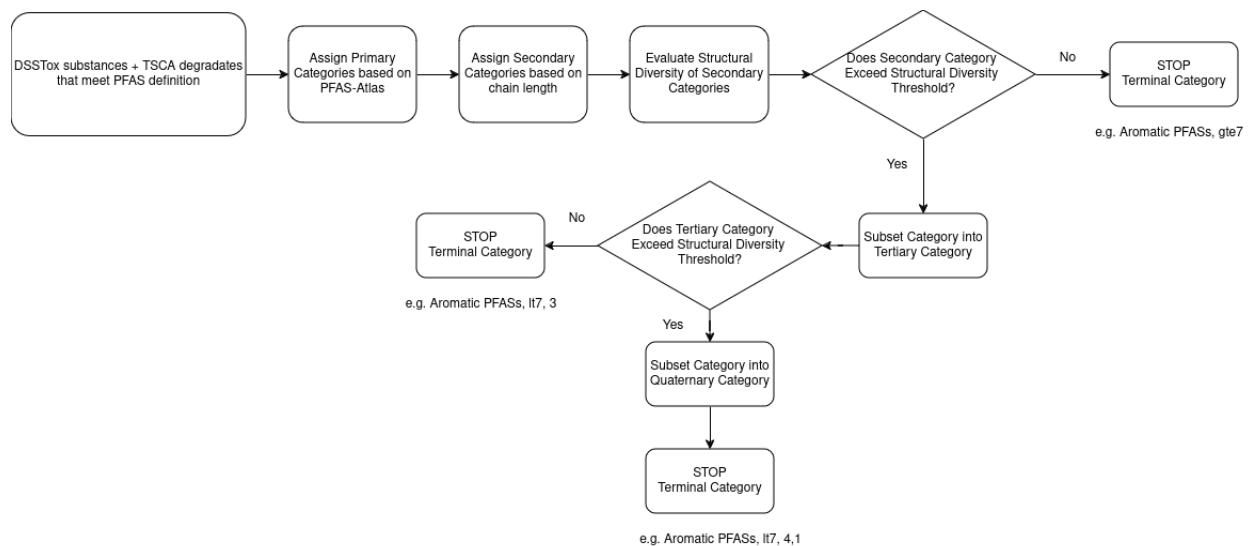


Figure 1: Conceptual workflow for generating PFAS structural categories

#### 195 2.3.1. Primary structural categories

Primary categories were derived by profiling the PFAS landscape of 15,538 substances through the PFAS subgroup classification tool developed by Su et al.<sup>38</sup> called **PFAS-Atlas**. As described in Su et al.<sup>38</sup>, PFAS can be classified into one of at least four broad classes:

- Perfluoroalkyl acids (PFAA)

- 200        • PFAA precursors  
201        • Polyfluoroalkyl acids  
202        • Other PFAS

203        In practice, when substances are batch processed by PFAS-Atlas, a first class and second class are  
204        assigned. The PFAS-Atlas first class corresponds to the 4 primary categories described above but  
205        additionally delineates linear substances from cyclic substances to create 8 primary categories. The  
206        tool also designates any substance that is not a PFAS based by the tool as 'Not PFAS'. The PFAS-  
207        Atlas second class subdivides each of the first classes further, for example in the PFAA precursors  
208        categories, hydrofluoroethers would be classified separately from semi-fluorinated alkenes or perflu-  
209        oroalkane sulfonyl fluorides (PASFs). For this study, a hybrid approach was used taking into account  
210        membership size. If the number of substances in a PFAS-Atlas first class designation exceeded 300  
211        substances, its second class designation was used. This was performed to limit the number of primary  
212        structural categories, and avoiding large membership sizes such as too many substances falling into  
213        the 'Other PFAS' category. The 300 membership threshold was chosen after manual inspection of  
214        the membership counts following substance assignment into the PFAS-Atlas first and second class  
215        designations. The PFAS-Atlas first class designation was used as the initial primary category assign-  
216        ment for PFAAs, whereas a handful of second class assignments were used for PFAA precursors and  
217        Polyfluoroalkyl acids. The PFAS-Atlas second class designation was used in lieu of "Other PFAS" ex-  
218        cept when the number of substances were very low. The PFAS-Atlas classification tree is an update of  
219        the original PFAS-Map database framework developed by some of the same authors<sup>39</sup> (and which had  
220        been used in the original NTS) but is more closely aligned with the OECD Terminology 2021 guidance<sup>8</sup>  
221        with some modifications.

222        The PFAS landscape was also processed through the OPEn structure-activity/property Relationship  
223        App (OPERA) v2.9 tool<sup>40</sup> (<https://github.com/kmansouri/OPERA>) to derive QSAR-READY SMILES  
224        and selected physicochemical property predictions (as discussed in Section 4.3). QSAR-READY  
225        SMILES are standardized SMILES where salts and stereochemistry are removed. QSAR-READY  
226        SMILES were used to facilitate the processing of substances through PFAS-Atlas. Substances  
227        without QSAR-READY SMILES or which could not be computationally resolved were assigned  
228        as "unclassified". Substances assigned as "Not PFAS" by PFAS-Atlas were also re-assigned as  
229        "unclassified".

230        The category assignments used in this study are captured in Table 1.

Table 1: List of PFAS-Atlas class assignments and the corresponding primary categories used in this study. Si PFAS refer to Silicon PFAS, HFCs refer to hydrofluorocarbons, PASF-based substances refer to perfluoroalkane sulfonyl fluorides, PFAA refer to perfluoroalkyl acids and PolyFACs are polyfluoroalkyl alcohols. Substances that could not be positively categorized by PFAS-Atlas were denoted as 'other'.

PFAS-Atlas first class	Primary category
PFAAs	PFAAs
PFAAs, cyclic	PFAAs, cyclic
PFAA precursors	PFAA precursors
PFAA precursors	PASF-based substances
PFAA precursors	n:2 fluorotelomer-based substances
PFAA precursors	HFCs
PFAA precursors, cyclic	PFAA precursors, cyclic
Polyfluoroalkyl acids	Polyfluoroalkyl acids
Polyfluoroalkyl acids	PolyFCA derivatives
Polyfluoroalkyl acids,cyclic	Polyfluoroalkyl acids,cyclic
Other PFAS	Other PFAS
Other PFAS	Aromatic PFASs
	Si PFASs
	Polyfluoroalkanes
	others
Other PFAS,cyclic	Other PFAS,cyclic
Other PFAS,cyclic	others, cyclic
Not PFAS	Unclassified

<sup>231</sup> 2.3.2. Secondary structural categories

<sup>232</sup> It is hypothesized that the length of the contiguous fluorinated carbon chain influences differences  
<sup>233</sup> in toxicity as well as the length of time the chemical spends in the body and environment. This sup-  
<sup>234</sup> position draws from experiences with PFAAs<sup>41,42</sup>. Due to the potential importance of chain length in  
<sup>235</sup> the toxicity, persistence and bioaccumulation of PFAS, secondary structural categories were defined  
<sup>236</sup> using a carbon chain length threshold.

<sup>237</sup> Chain length determination

<sup>238</sup> The maximum number of contiguous  $CF_2$  groups in a chain was determined for all 15,538 sub-  
<sup>239</sup> stances. This was achieved by iterating through a range of  $CF_2$  chain lengths (from 1-30) for each  
<sup>240</sup> substance in turn and determining its longest chain length. For instance, Perfluorosebacamide

241 [DTXSID40380015] contains 8 contiguous  $\text{CF}_2$  units; hence, its chain length was denoted as 8.  
242 PFOA [DTXSID8031865] had a maximum chain length of 7 whereas PFOS [DTXSID3031864] had a  
243 maximum chain length of 8. For PFOA, although there are 8 carbons in its backbone, the 8th is part  
244 of the carboxyl group whereas in PFOS, there are 8  $\text{CF}_2$  groups plus the sulfonate group.

245 For the current analysis, the chain length threshold was set at 7 ( $\geq 7$  vs  $< 7$ ) as representative of  
246 a "long chain" PFAS. The chain length threshold is broadly consistent with the EPA's 2009 [PFAS ac-](#)  
247 [tion plan](#). A PFAS with a maximum number of contiguous  $\text{CF}_2$  number greater than or equal to 7 was  
248 denoted "gte7". Using this threshold, both PFOS and PFOA would be assigned to the "gte7" sec-  
249 ondary category. A PFAS with a maximum number of contiguous  $\text{CF}_2$  groups less than 7 was denoted  
250 "lt7". Defining chain lengths for PFAS with non-contiguous chains or branching is less straightfor-  
251 ward but has been evaluated in more detail by Richard et al <sup>43,44</sup> through the development of new  
252 PFAS specific chemical fingerprints, so-named PFAS ToxPrints, as an extension of the logic used  
253 to develop the original ToxPrints that had been defined for a broader chemistry <sup>45</sup>. A secondary  
254 category was thus denoted by its PFAS-Atlas assignment, akin to the primary OECD structural clas-  
255 sification and a carbon chain length threshold e.g., 2,2,3,3,4,4,5,5,6,6,7,7,8,8,9,9,9-Heptadecafluoro-  
256 N,N-diphenylnonanamide [DTXSID90896196] would thus be described as belonging to the "Aromatic  
257 PFASs, gte7" secondary category (see Figure 1).

### 258 2.3.3. Derivation of terminal structural categories

259 The underlying motivation for the study was to identify categories that would balance maximiz-  
260 ing structural similarity that could permit read-across within those categories versus pragmatism in  
261 terms of total number of categories. Too many categories with very few substances renders the  
262 approach less generalizable, too few categories could result in extrapolating between substances  
263 that were not sufficiently similar. To that end, an objective threshold was needed to determine how  
264 granular categories needed to be to manage this trade-off and ensure that the categorization was ac-  
265 tionable. An objective threshold was developed, described in Section 2.3.6, that compared structural  
266 similarity within a category relative to the structural similarity between different categories.

### 267 2.3.4. Chemical fingerprints

268 Morgan chemical fingerprints <sup>35</sup> were calculated for all substances within each secondary category  
269 using the open-source Python library RDkit <sup>46</sup>, with a radius of 3 and a bit-length of 1024. The fin-  
270 gerprint data was represented as bit vectors where presence of structural features were denoted  
271 by 1 and absence by 0. These fingerprint (FP) files were stored for subsequent processing. Morgan  
272 fingerprints also known as Extended Connectivity Fingerprints (ECFPs) are widely used in machine  
273 learning applications for cheminformatics <sup>47</sup> especially when ranking diverse structures by similarity.

274 These circular fingerprints map the molecular environment of every atom.

275 2.3.5. Chemical similarity

276 Pairwise distance matrices were calculated for each secondary category. These were generated by  
277 using the chemical fingerprint files as inputs and computing the Jaccard distance for each pair of  
278 substances. The Jaccard distance captures the proportion of FP bits between 2 substances that  
279 differ<sup>48</sup>. The Jaccard distance ranges from 0 to 1 where 0 would indicate zero distance (or high  
280 similarity) and 1 would indicate high distance (or low similarity). Distance matrices were computed  
281 for all secondary categories and stored for subsequent processing. These are referred to as 'within  
282 category' distance matrices in Section 2.3.6.

283 2.3.6. Objective distance threshold

284 The rationale underpinning the objective distance threshold was based on the expectation that the  
285 variance in the distribution of the pairwise distances for each secondary category representing the  
286 'within category' similarity would be lower than distributions of the pairwise distances between dif-  
287 ferent secondary categories ('between category'). The 'within category' distances had already been  
288 computed as described in Section 2.3.5.

289 'Between category' combinations aimed to identify categories that did not share the same primary  
290 category root. A list of all possible binary combinations of secondary categories was created using the  
291 names of the secondary categories, "Aromatic PFASs, lt7" and "PFAA precursors, gte7" is an example  
292 of such a binary combination. These were then filtered to remove secondary categories that shared  
293 the same primary category root (i.e., a combination such as "PFAAs, lt7" and "PFAAs, gte7" would  
294 be excluded from consideration as a 'between category'). Chemical fingerprint datasets for each  
295 binary combination were created by combining the secondary category chemical fingerprint datasets.  
296 Pairwise distance matrices were then derived for the combined category set. These matrices were  
297 filtered to retain only the pairwise distances between the starting secondary categories.

298 The empirical cumulative distribution functions (ECDFs) of the pairwise distances were calculated  
299 for each secondary category (see Figure A2 for a plot of the ECDFs). ECDFs were also derived for  
300 the 'between category' combinations (see Figure A3 for a plot of the first 10 ECDFs). The ECDFs  
301 permitted a visual inspection of the range of the pairwise distances across all secondary categories  
302 as well as across all 'between category' combinations. Based on visual inspection of the ECDFs, the  
303 median value for each distribution was selected as the summary metric.

304 Probability density functions of median values from all within and between secondary categories  
305 were plotted to explore their overlap. The 15<sup>th</sup> percentile of the 'between categories' distribution

306 was selected, by reference to the density plot, as the threshold to determine whether a secondary  
307 category merited further subcategorization. A secondary category was only subcategorized if the  
308 median of its 'within category' pairwise distance distribution exceeded this threshold.

309 2.3.7. Deriving terminal categories

310 Secondary categories that exceeded the threshold were subcategorized using agglomerative hierar-  
311 chical clustering. The condensed form of the pairwise distance matrix computed for each secondary  
312 category that exceeded the threshold was used as an input into a hierarchical clustering using Ward's  
313 method<sup>49</sup>. Ward's method is a criterion that minimizes the total within-cluster variance. For each  
314 secondary category, the dendrogram was plotted and the number of first-generation clusters was set  
315 as the maximum cluster number. Clusters were labelled as 1,2,3 etc. Each of the clustering results  
316 were combined into one table which was then merged with the starting table of primary and secondary  
317 categories.

318 The next generation of categories, quaternary categories, would then be processed in the same  
319 manner to determine whether any exceeded the objective threshold and needed to be subcategorized  
320 further as already described. In practice, a maximum of two generations of subcategorizations were  
321 performed, with the expectation that this would balance the structural similarity within the category  
322 relative to total number of terminal categories.

323 Secondary categories or tertiary categories which did not exceed the threshold were ultimately  
324 denoted as the terminal category. Thus, a terminal category could be tagged as "Aromatic PFASs,  
325 gte7", effectively a secondary category, or could be tagged as "Aromatic PFASs, lt7, 3", a tertiary  
326 category or following two iterations of subcategorization would be tagged as "Aromatic PFASs, lt7, 4,  
327 1" (see Figure 1). Note: in the data files and figures, terminal categories without one or two iterations  
328 of subcategorization are denoted as "Aromatic PFASs, gte7, nan, nan" or "Aromatic PFASs, lt7, 3, nan"  
329 where "nan" represents a null value.

330 2.3.8. Identification of centroid substances

331 For each terminal category, a single substance was identified that was nominally representative  
332 of the category. This substance was the computed centroid calculated from the Jaccard pairwise  
333 distance matrices (see Section 2.3.5). The sum of the pairwise distances across all substances for  
334 a given structural category was computed and the substance with the minimum value was denoted as  
335 the centroid (i.e., this substance would have the lowest distance from all other category members).  
336 Technically, this calculation gives rise to the medoid of a cluster. However, for the purposes of this  
337 analysis and for consistency with the NTS, the term centroid is used to denote it as the 'central'

338 substance within the category. Distances of all category members relative to the centroid substance  
339 were also computed.

340 **2.3.9. Identification of additional representative substances**

341 Since a number of the terminal categories were large in size (e.g. greater than 300 members), a single  
342 substance would be potentially insufficient to both characterize the category and its potential hazard  
343 profile. In an effort to address this limitation, the MaxMinPicker approach, as implemented within  
344 the RDKit Python library, was applied to identify additional substances which would in turn capture  
345 the breadth and diversity of each terminal category<sup>50</sup>. The MaxMin approach is a well-established  
346 algorithm for dissimilarity-based compound selection that has been applied in drug discovery for many  
347 years. The reader is referred to Snarey et al<sup>51</sup> for a comparison of the different algorithms. The  
348 MaxMinPicker approach proceeds as follows:

349 1. Molecular descriptors are generated for all substances. In this case, the Morgan fingerprints  
350 calculated for all substances within a terminal category represented the candidate pool whereas the  
351 pre-computed centroid equated to the initial seed.

352 2. From the substances in the terminal category, the substance that had the maximum value for its  
353 minimum distance to the picked set (initially this would be just the centroid) would then be identified.  
354 This substance would be the most distant one to those already picked so it would be transferred to  
355 the 'picked set' (now centroid + 1).

356 3. An iteration back to step 2 would then be performed until the desired number of substances were  
357 picked.

358 The MaxMinPicker was applied to all terminal categories containing more than 5 members to identify  
359 the next 3 most diverse substances within a category (centroid + up to 3 additional substances). The  
360 intention of identifying additional diverse substances was to help bound the domain of the structural  
361 category. The identification of 3 most diverse substances was chosen out of convenience to provide  
362 an actionable number of additional substances.

363 A systematic evaluation of the relationship between the number of diverse substances that could  
364 be identified relative to the structural diversity within each terminal category was also undertaken.  
365 This was performed as follows, first the ranked order by diversity of all members within a terminal  
366 category was computed. Then the pairwise distance matrices derived in Section 2.3.5 were filtered  
367 by the diverse substances, starting from the centroid, centroid plus first diverse chemical through to  
368 the complete set of category members. At each step the mean minimum distance was recorded. This  
369 enabled the construction of a matrix to capture the mean of the minimum pairwise distances relative

370 to the number of diverse chemicals selected. The normalized cumulative sum of all the mean mini-  
371 mum distances was then computed. This provided a means of evaluating the proportion of structural  
372 diversity that was captured as a function of the number of MaxMin substances identified.

373 This calculation provides for two objective assessments namely:

- 374 1. the amount of structural diversity captured by the 3 diverse picks originally identified; and  
375 2. the number of diverse substances that would need to be identified (if practical resources were  
376 not a limiting factor) to capture a specified level of structural diversity. For example, how many  
377 substances would need to be identified if capturing a specific percentage of the structural diversity  
378 within a terminal category was desired, 80% is presented here merely for illustrative purposes.

379 **2.4. Facilitating the identification of potential candidates for data collection**

380 To facilitate the identification of potential candidate PFAS for data collection, availability of a  
381 known manufacturer/importer, EPA Agency and/or State priorities, environmental monitoring infor-  
382 mation were evaluated as additional considerations. These are described in turn.

383 **2.4.1. Qualitative exposure and release designations**

384 Several qualitative designations were added to the landscape to identify substances for which expo-  
385 sure could be plausible, including their TSCA inventory status, production volumes per TSCA's Chemi-  
386 cal Data Reporting (CDR) rule, State/EPA Region priorities, as well as physical state and physicochem-  
387 ical properties.

388 The non-confidential (Non-CBI) TSCA Inventory active and inactive lists were downloaded from the  
389 Dashboard (see [TSCA\\_ACTIVE\\_NCTI\\_0224](#), [TSCA\\_INACTIVE\\_NCTI\\_0224](#)) and combined into one  
390 large set. Substances within this inventory included both Chemical Abstract Service (CAS) Registry  
391 Number, Chemical Abstracts (CA) Index Name, and DSSTox substance identifier (DTXSID). These  
392 were matched by the DTXSID identifiers already captured in the PFAS landscape. Substances were  
393 tagged as 'inactive', 'active' or 'unclassified'. Note the predicted degradation products of substances  
394 tagged as either "inactive" or "active" had been used to augment the PFAS landscape as already de-  
395 scribed in Section [2.1](#).

396 The 2020 CDR data was downloaded from the public EPA web address (<https://www.epa.gov/chemical->  
397 data-reporting/access-cdr-data). The CDR data comprises information for a set of 8660 substances.  
398 DTXSID identifiers were available for 8017 of these substances when using the Batch search  
399 functionality within the Dashboard. A tag was created for CDR2020 status if a PFAS had a 2020  
400 CDR record. National Aggregated Production Volume (National Agg PV) data was also extracted  
401 to highlight how this was distributed across primary categories. Since some of the production

402 volume (PV) data was numeric and some represented in numeric ranges, the PV data was summa-  
403 rized into one of 10 different ranges (<25,000 lbs, 25,000-<100,000 lbs, 100,000-<500,000 lbs,  
404 500,000-<1,000,000 lbs, <1,000,000 lbs, 1,000,000 -<10,000,000 lbs, 1,000,000-< 20,000,000 lbs,  
405 20,000,000-< 100,000,000 lbs, 50,000,000-<100,000,000 lbs, 100,000,000-< 1,000,000,000 lbs).

406

407 Various EPA Regions or States have identified PFAS of interest based on validated analytical meth-  
408 ods or for environmental monitoring purposes. The data sources captured as part of the EPA's PFAS  
409 Analytic Tools website (<https://echo.epa.gov/trends/pfas-tools#data>) were used to construct lists  
410 of such PFAS. The specific data sources were Discharge Monitoring Data, Drinking Water (State)  
411 Data, Drinking Water (Unregulated Contaminant Monitoring Rule (UCMR)) Data, Environmental Media  
412 Data, Production Data, Toxics Release Inventory (TRI) Data - Waste Managed, TRI Data - On-Site,  
413 TRI Data - Off-Site and Production Data (all accessed 7th April 2024). Discharge Monitoring data  
414 is collected by virtue of the National Pollutant Elimination System permit. Drinking Water Data com-  
415 prises UCMR and State level monitoring data. Environmental Media data comprises ambient sampling  
416 data reported by federal, state, tribal and local governments, academic and non-governmental orga-  
417 nizations, and individuals that are submitted to the Water Quality Portal (WQP). Production data  
418 entails information reported under the Chemical Data Reporting (CDR) Rule under TSCA. TRI tracks  
419 the management of certain toxic chemicals that may pose a threat to human health or the environ-  
420 ment by more than 21,000 facilities throughout the US and its territories. The National Defense  
421 Authorization Act of Fiscal year 2020 (NDAA) added certain PFAS to the TRI list and provided a  
422 framework for the ongoing listing of additional PFAS.

423 Identifiers were extracted from these source files and searched against the Dashboard to map  
424 to DTXSID records. The set of identifiers (Names and CASRN) within the entire PFAS landscape  
425 were also queried against PubMed, the National Library of Medicine's citation index for biomedical  
426 literature, to determine whether studies for a substance might have been reported in the literature.  
427 The article counts were obtained using the same queries as used within the Abstract Sifter v7.5<sup>52</sup>.

428 Expected routes of exposure and presence in environmental media are dependent on the phys-  
429 ical state and physicochemical properties. Physicochemical properties were predicted using the  
430 open-source OPERA v2.9 tool<sup>40</sup> for all substances with QSAR-READY SMILES (as discussed in  
431 Section 2.3.1). The properties predicted were melting point, boiling point, Henry's Law constant,  
432 water solubility and vapour pressure. Physical state was predicted at 25 deg C using the predicted  
433 values of melting point and boiling point. Gases had boiling points less than 25 deg C, solids had melting  
434 points greater than or equal to 25 deg C and liquids had melting points less than 25 deg C and boiling  
435 points greater than or equal to 25 deg C. These are the guiding principles underpinning the EPA's

436 Sustainable Futures Framework guidance (see [Interpretative Guidance Document](#)). Whilst physical  
437 properties are continuously distributed, and cutoff values are necessarily arbitrary, there is utility  
438 in grouping substances into broad categories as a way to acknowledge the practicalities of testing  
439 and human exposure under "typical" (i.e. room temperature and atmospheric pressure) conditions. A  
440 water solubility threshold of 0.5 mg/L was used to denote whether a substance was soluble/insoluble  
441 whereas a vapour pressure threshold of 75 mmHg determined volatility and a HLC threshold of  
442 0.1 atm m<sup>3</sup>/mol highly volatile. Based on these properties, each substance was assigned into 1 of  
443 4 "physical state and physicochemical designations" (from A-D). Designation A covered substances  
444 that were insoluble solids, designation B captured both soluble solids and soluble non-volatile liquids,  
445 whereas C tagged soluble volatile liquids/insoluble liquids and soluble gases. Designation D assigned  
446 substances as insoluble gases or highly volatile gases. Substances that could not be assigned into  
447 one of these 4 designations were tagged as 'not determined'. For each of the terminal structural  
448 categories, Morgan fingerprint representations were projected into two dimensions using UMAP to  
449 facilitate visualization<sup>36</sup>. The projections were plotted as 2D kernel density distributions overlaid  
450 with physical state and physicochemical designation information to help explore the extent to which  
451 members were assigned to the same designation and therefore had a consistent profile across a  
452 given terminal category.

453 Each of these respective qualitative designations were then matched to the PFAS landscape to  
454 provide another attribute for consideration when identifying potential candidates for data collection.

#### 455 2.4.2. Constrained PFAS landscape

456 One of the limitations of the identification of centroids and additional diverse substances was that  
457 they might yet not yield feasible candidates for data collection due to the lack of assignable manufac-  
458 turer/importer. This was articulated as a potential challenge in the National PFAS Testing Strategy.  
459 To address this practical constraint, the same process of computing centroids, identifying additional  
460 diverse substances and evaluating their structural diversity coverage was also performed using the  
461 terminal categories as a basis as described in Section [2.3.7](#) but constraining the landscape to only  
462 those substances on the public TSCA inventory and specifically those substances that were actives  
463 on the public TSCA inventory. Constraining the landscape would allow identification of substances for  
464 data collection that were already in commerce and/or could be more readily procured.

#### 465 2.5. Evaluation of variance of in vivo toxicity within terminal categories

466 Ultimately, read-across of data within categories could be performed such that the hazard profile  
467 of the category is adequate without needing to test a significant number of category members. To  
468 evaluate the feasibility of performing read-across within the terminal categories derived, an explo-

469 ration of the distribution of in vivo points of departure (PODs) within and across terminal categories  
470 was performed for the oral route of exposure.

471 2.5.1. Variance of in vivo PODs across and within terminal categories

472 From ToxValDB version 9.5, the [Toxicity Values Database](#), all studies where 'oral' was the route of  
473 exposure were extracted. Only records where a point of departure (POD) was reported as a NOEL,  
474 NOAEL, NEL, NOAEC, LOAEL, LOEL, LOAEC, LEL and where the dose units were expressed as mg/kg-  
475 bw/day or mg/kg were retrieved. Study types were also restricted to the following: 'short-term', 'sub-  
476 chronic', 'chronic', 'developmental', 'reproduction', 'reproduction developmental', '28-day' as captured  
477 in the 'study type' field within the database. Species were standardized into one of 'rat', 'mouse',  
478 'rabbit', 'dog', 'hamster' or 'guinea pig'. Effect levels were harmonized consistent with the approach  
479 taken by Aurisano et al.<sup>53</sup> where non-cancer effects vs. reproductive/developmental effects were  
480 processed separately. Records with sub-acute or sub-chronic as the study type were extrapolated  
481 to chronic using a subchronic-to-chronic factor of 2 and a subacute-to-chronic factor of 5<sup>2</sup>. LOAEL  
482 effect level types were extrapolated to NOAELs by dividing by an extrapolation factor of 3. Effect  
483 levels for all records were extrapolated to humans by dividing reported effect values by conversion  
484 factors based on the average body of weight of humans relative to the average body weight of the  
485 test species. NOAELs were extrapolated to human equivalent Benchmark Dose (BMDh) values based  
486 on assigned conceptual model depending on the critical effect reported. The calculated BMDh was  
487 based on the mean of 2 assigned conceptual models. In Aurisano et al<sup>53</sup>, the 25<sup>th</sup> percentile of the  
488 fitted log-normal distributed (using the mean BMD and standard deviation (sd) BMD) was calculated  
489 to derive a POD per substance. The sd was set to the median sd of all records in cases where the  
490 number of study records was less than 5.

491 The derived BMDh values were then merged with the PFAS substances from the landscape. The  
492 summary values provided an estimate of the POD for each substance and the expected level of variation  
493 across and within categories. Box and whisker plots were created to reflect the distribution  
494 of the PODs across the terminal categories for general non-cancer and repro/developmental effects  
495 for the oral route of exposure. Strip plots were overlaid to show the variation of chain length across  
496 a given terminal category for general non-cancer effects.

497 2.6. Qualitative mechanistic and toxicokinetic designations

498 A summary of the NAM testing being undertaken for ~150 PFAS was described in Patlewicz et  
499 al.<sup>23</sup>. See Houck et al.<sup>26</sup> for results from various nuclear receptor and oxidative stress targeted

---

<sup>2</sup>Note: the approach taken in Aurisano et al. is not necessarily consistent with uncertainty factor selection information provided in EPA's 2002 Review of the Reference Dose and Reference Concentration Processes and EPA's 2022 ORD Staff Handbook for Developing IRIS Assessments

500 assays, Houck et al.<sup>27</sup> for 12 human primary cell-based assay models of pathophysiology including  
501 immunosuppression, Carstens et al.<sup>25</sup> for the developmental neurotoxicity assays, Degitz et al.<sup>32</sup> for  
502 the thyroid pathway assays and, for toxicokinetic information, Smeltz et al.<sup>29</sup> and Kreutz et al.<sup>28</sup>.  
503 The manuscript for the remaining data stream (zebrafish developmental toxicity) is in under internal  
504 review (Britton et al., in prep).

505 In addition to the NAM testing, a quality control (QC) evaluation of the chemical stock solutions  
506 was undertaken to confirm PFAS analyte presence and stability<sup>30</sup>. This evaluation was warranted  
507 given recent reports of certain PFAS degrading in the aprotic solvent dimethyl sulfoxide (DMSO),  
508 readily used as the solvent of choice in HTS<sup>54,55</sup>. Two hundred and five PFAS selected based on  
509 criteria described in Patlewicz et al.<sup>23</sup> were evaluated using low resolution tandem mass spectro-  
510 metric detection strategies to confirm presence of intended analyte, evaluate analyte stability and  
511 presence of isomers, and verify stock concentrations for a subset for which commercially available  
512 verified standards were available. Ultimately 57 PFAS failed QC evaluation, with three exhibiting  
513 degradation in DMSO and the remainder not detected as present, likely due to volatilization. The  
514 pass/fail scores and informational flags as described in Smeltz et al.<sup>30</sup>, and can be downloaded from  
515 the following figshare url, [https://epa.figshare.com/articles/dataset/Chemistry\\_Dashboard\\_Data\\_Analytical\\_QC\\_for\\_PFAS/22118099](https://epa.figshare.com/articles/dataset/Chemistry_Dashboard_Data_Analytical_QC_for_PFAS/22118099).

517 For each of the NAM data streams, substances were tagged with a qualitative flag to indicate  
518 the class of mechanistic information that could be derived from the associated assay outcome (e.g.,  
519 estrogen receptor activity from a nuclear receptor assay) and an expert-derived qualitative level of  
520 confidence associated with the outcome (high confidence of activity, medium confidence of activity  
521 or low concern). Only NAM results from substances that passed QC were carried forward. These  
522 flags were considered as an additional line of evidence to determine whether a terminal category  
523 might merit being split based on its mechanistic or toxicokinetic information or to inform what types  
524 of higher order testing might be most impactful for a given substance drawn from said terminal  
525 category. The derivation of the flags are described in more detail in Judson et al. in prep. Confidence  
526 scores across the NAM flags were standardized as appropriate to facilitate visualizations across data  
527 streams. Each flag could take on one of three values, low concern, medium or high confidence, color  
528 coded as blue, yellow and red. The immune flag was the exception. It was binary in nature and only  
529 gave rise to a low concern and medium confidence value. The flag categories are summarized below  
530 in Table 2. One final parameter computed was a TK half-life bin score using the machine learning  
531 model developed by Dawson et al.<sup>56</sup>. Predictions were scored from 1-4 where 1 signified a half-life  $\leq$   
532 12 h, 2 = 12 h-1 week, 3 = 1 week-2 months, and 4  $\geq$  2 months. Predictions were generated for humans  
533 assuming a oral dosing regimen and aggregated by maximum half-life score so that there was a single

<sup>534</sup> prediction per substance.

Table 2: Summary of NAM Flag Rationales

Endpoint	Low Concern (Blue)	Medium Confidence (Yellow)	High Confidence (Red)
Nuclear Receptors	No nuclear receptor activity	Activity against at least one of the receptors ER, PPARA, PPARG, PPARD, NFE2L2, PXR, RARG, RXRB at the level of one or more samples in one assay.	Activity in the yellow medium concern that is confirmed in at least one sample in 2 orthogonal assays
DNT	No activity or activity was only observed at the highest concentration related to cytotoxicity	Low number of hits which demonstrated selective bioactivity	Moderate to high bioactivity (as measured by hitcall) and demonstrated selective bioactivity (activity below cytotoxicity AC50 as measured by AUC) and median AC50 < 10 $\mu$ M
Zebrafish	Development was normal in all larvae	Test results were equivocal or if less than 50% of the larvae were affected	Positive activity (i.e., elicited death, non-hatching, or malformations in at least 50% of the animals)
Thyroid	No activity greater than 50% of the model inhibitors/binders	Activity greater than 50% of the model inhibitors/binder, but the concentration necessary to result in this activity was 2 orders of magnitude higher than the model inhibitors/binders	AC50s that were within 2 orders of magnitude of the model inhibitors/binders
Immune	Selectivity scores less than 0.25 log <sub>10</sub> $\mu$ M	Selectivity scores of greater than 0.25 log <sub>10</sub> $\mu$ M	

Endpoint	Low Concern (Blue)	Medium Confidence (Yellow)	High Confidence (Red)
TK Plasma Binding (TK_PlasBind)	TK_PlasBInd_High: Plasma protein binding higher than 50% of non-PFAS chemicals (f <sub>up</sub> < 0.11) (this corresponds to 25 <sup>th</sup> percentile of PFAS (f <sub>up</sub> <0.10)	TK_PlasBInd_Higher: Plasma protein binding higher than 50% of PFAS chemicals (f <sub>up</sub> < 0.0109)	TK_PlasBInd_Highest: Plasma protein binding higher than 75% of PFAS (f <sub>up</sub> < 0.0039)
TK Intrinsic Clear- ance (TK_Metab)	TK_Metab_Moderate: Clint in upper 75 <sup>th</sup> percentile of exp PFAS data (Clint>5.97 ul/min/million cells). Max Clint = 49.86	TK_Metab_Slow: Clint<5.97 ul/min/million heps. (lower 75 <sup>th</sup> percentile)	TK_Metab_Stable: Stable in in vitro hepatocyte incubation (Clint = 0 or Clint pvalue > 0.05)
TK half-life predic- tions	TK_Struc_Endo	Non-fluorinated structure is similar to endogenous chemicals. More likely to be a transporter substrate.	
	category 1 or 2 to denote half-life $\leq$ 12 h, or 12 h-1 week	category 3 = 1 week-2 months	category 4 $\geq$ 2 months

535 Qualitative observations of the consistency of the various flags across all the tested substances  
 536 and within terminal categories were made. The Fisher's exact test was used to compute an odds ratio  
 537 and associated p value for each PFAS ToxPrint<sup>44</sup> relative to a NAM flag that had been converted  
 538 into a binary scale. This enrichment analysis was comparable with the methodology discussed in Wang  
 539 et al.<sup>57</sup>. A PFAS ToxPrint was considered enriched if it had an odds ratio greater than or equal  
 540 to 3, an one-sided Fishers exact p-value less than 0.05 (probability value of the odds ratio being  
 541 greater than 1) and the number of true positives equal or greater than 3. The set of 'enriched' PFAS  
 542 ToxPrints were then used to profile the entire PFAS landscape to assign potential predicted NAM  
 543 flags. Comparisons were made of the actual NAM flags and their predictions to evaluate performance  
 544 metrics (sensitivity and specificity). TK half-life predictions were generated for all substances and  
 545 the outcomes binarized where substances with a bin category of 4 were assigned a 1, and any other  
 546 bin category lower was assigned a 0.

547 **2.7. Operationalizing the terminal categories for re-use**

548 In the absence of a model to predict the terminal category, the categorization would need to be  
549 re-run for each new set of PFAS. To operationalize the terminal categories for practical use, a ma-  
550 chine learning approach was used to develop a model that could be used to profile a new PFAS and  
551 assign it to its most likely terminal category. A random forest classifier (RFC) as implemented in  
552 the python library scikit-learn<sup>58</sup> was used to predict assignment of substances into one of the final  
553 terminal categories developed. Morgan fingerprints generated earlier (as discussed in Section 2.3.4)  
554 in conjunction with primary category and chain length were combined with the final terminal category  
555 names for all substances. Terminal categories with less than 10 members were aggregated together  
556 into one miscellaneous category. The dataset was then split into a training 80% and test 20% split  
557 using a random stratification approach based on the terminal category labels. A dummy classifier was  
558 applied first to establish a baseline. Then an initial RFC with default settings was assessed within a  
559 5-fold stratified cross validation (CV) procedure to evaluate initial performance using balanced accu-  
560 racy as a metric. This was performed using a pipeline within scikit-learn where the primary category  
561 names were treated as category features and passed into an OrdinalEncoder whereas chain length and  
562 Morgan fingerprints were first standardized before being passed to the RFC. A randomized search  
563 was then undertaken as part of a nested 5-fold CV to identify the best parameters and evaluate test  
564 CV performance. The resulting model was then applied to the test set that had been held out to  
565 evaluate performance. Finally, the model was refitted to the entire dataset.

566 **3. Data analysis software and code**

567 Data processing was conducted using the Anaconda distribution of Python 3.9 and associated li-  
568 braries. Jupyter Notebooks, scripts and datasets will be made available on github and Figshare.

569 **4. Results and discussion**

570 **4.1. Primary and secondary structural categories**

571 The PFAS landscape following application of the TSCA section 8(a)(7) rule to DSSTox resulted in a  
572 dataset comprising 13,054 substances plus 2484 degradation products for a total of 15,538.

573 Minimal structural overlap was found between the 1549 accessible training set substances (of which  
574 1429 substances could be resolved into structures) from the Catalogic model and the TSCA inventory  
575 substances. Figure A1 depicts a UMAP plot for the MITI training set substances relative to the TSCA  
576 substances using Morgan chemical fingerprints as inputs. There were 12 TSCA substances of the  
577 dataset that were part of the training set. Only 29% of the PFAS TSCA substances were tagged as

578 being within the structural domain of the model. In view of this, the degradation products simulated  
579 (comprising 16% of the landscape) should be interpreted with caution until additional experimental  
580 data are collected and new models developed.

581 A chain length could not be computed for 14 substances due to issues with resolving structures within  
582 RDKit. Only one of the 14 substances, DTXSID20153820, could be resolved since the chain length  
583 failed on account of the chain length exceeding the range used to calculate the maximum values. The  
584 chain length of this substance was manually annotated. The remaining 13 substances were dropped  
585 from further consideration. There were 31 substances that were tagged as "Not PFAS" based on the  
586 OECD structure definitions used within PFAS-Atlas. All were reassigned to the "unclassified" primary  
587 category. Figure 2 is a bar chart showing the number of PFAS within each secondary category. The  
588 final PFAS landscape used for the remainder of the analysis comprised 15,525 substances.

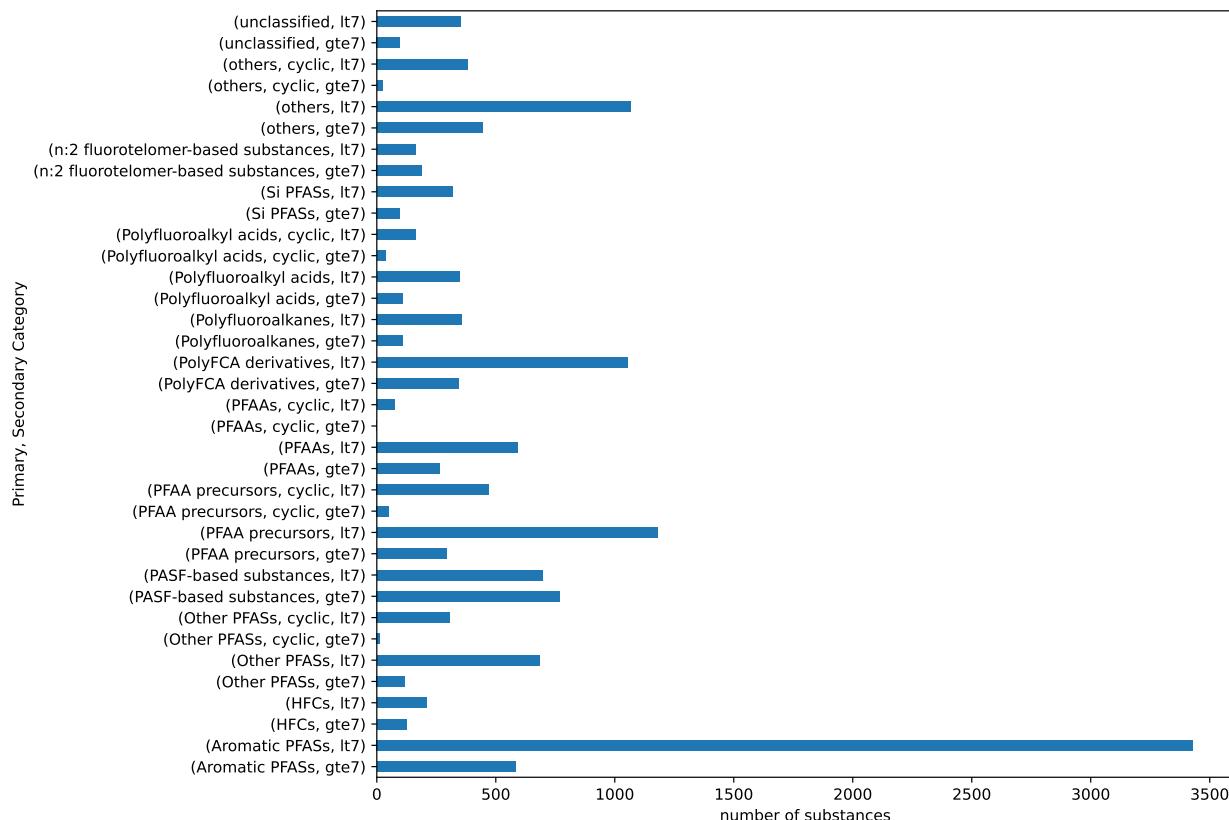


Figure 2: Bar chart showing the number of substances within each secondary category, ordered by primary category root. Methods to define primary and secondary categories are outlined in Sections 2.31 and 2.32. Lt7, chain length less than 7; gte7, chain length greater than or equal to 7.

589 Across the more than 15,000 PFAS substances evaluated, twenty-two percent (3430) of the sub-  
590 stances fell into the "Aromatics PFAS, lt7" category. In addition, 1066 substances fell into the "Oth-

591 ers, It7" and 1180 in the "PFAA precursors, It7" secondary categories. This represents a potential  
592 limitation of using broad definitions represented by the OECD primary categories themselves. The  
593 smallest secondary category was the "Others PFAS, cyclic, gte7" with 14 members whereas "PFAAs,  
594 cyclic, gte7" was a singleton.

595 A chemotype ToxPrint enrichment was explored following the approach outlined in Wang et al.<sup>57</sup> but  
596 using the PFAS specific ToxPrints developed in Richard et al<sup>44</sup> (see Section [Appendix A.1](#) for method-  
597 ological details). This was an effort to identify whether there were specific structural features that  
598 might be helpful in subcategorizing those primary categories with the largest memberships namely  
599 (i.e., "Aromatic PFASs", "PFAA precursors" or "unclassified"). The most enriched features for the  
600 "unclassified" category included fluorotelomer chains and sulfonic acid functional groups whereas al-  
601 cohols and carbonyls featured as functional groups for the "PFAA precursors". No specific features  
602 were enriched for the "Aromatic PFAAs" categories. However, where there were enriched features,  
603 these were not determined to be sufficiently distinctive to justify creation of additional primary  
604 categories.

605 Structural similarity was evaluated within and between secondary categories to determine which sec-  
606 ondary categories required further subcategorization (as discussed in Section [2.3.6](#) of the Methods).  
607 Figure 3 shows the two distributions of the median pairwise distance distributions in the between  
608 and within secondary category combinations. The objective distance threshold derived by taking the  
609 15<sup>th</sup> percentile of the median pairwise distances from the between categories combinations resulted  
610 in a value of 0.8.

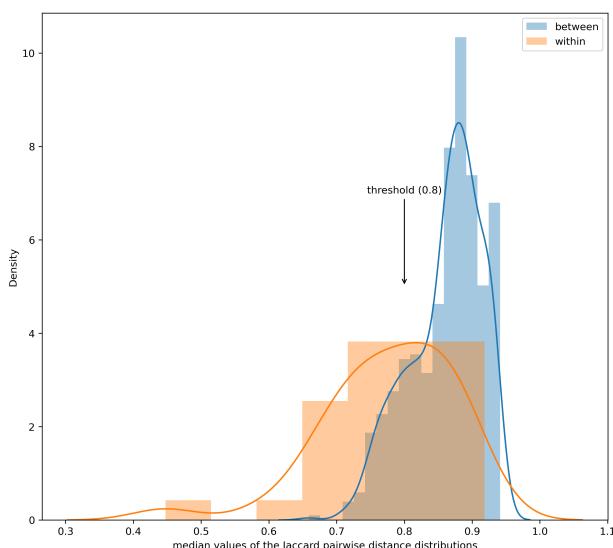


Figure 3: Probability density functions of the median Jaccard pairwise distance distributions for within (orange) and between (blue) secondary categories. Orange and blue graphed lines represent the fits to the probability density distributions.

611 Based on the threshold, 16 secondary categories (Table 3) were found to exceed the value that  
 612 would render them subject to further subcategorization. The sixteen secondary categories included  
 613 the "Aromatic PFASs", "PFAA precursors", "Others" and "PolyFCA derivatives". These categories are  
 614 of little surprise given their membership sizes were the largest out of all the secondary combina-  
 615 tions; hence, these categories were expected to be the most diverse in terms of their structural  
 616 makeup. Since the "PFAAs, cyclic" category only comprised 1 substance, this was excluded from any  
 617 subcategorization.

Table 3: List of secondary categories exceeding the threshold and their corresponding median pairwise distances (rounded to 2 decimal places)

Primary-Secondary Categories	Median Within-Category Pairwise distance
others, cyclic, lt7	0.92
PFAA precursors, cyclic, lt7	0.90
Other PFASs, cyclic, lt7	0.89
Aromatic PFASs, lt7	0.88
Other PFASs, lt7	0.88
unclassified, lt7	0.87
others, lt7	0.87
Polyfluoroalkyl acids, cyclic, lt7	0.86
PFAA precursors, lt7	0.86
Other PFASs, cyclic, gte7	0.85
Polyfluoroalkanes, lt7	0.85
PolyFCA derivatives, lt7	0.83
PFAAs, lt7	0.82
PFAAs, cyclic, lt7	0.82
Polyfluoroalkyl acids, lt7	0.81
HFCs, lt7	0.81

618 Figure 4 shows the membership following the first generation of clusters being created for the 16  
 619 secondary categories that exceeded this objective threshold.

620 Following creation of the next generation categories, there were 23 tertiary categories that met  
 621 or exceeded the threshold and were subcategorized further. The root primary categories were pre-  
 622 dominantly from the "Aromatic PFASs", "Other PFASs", and "PFAA precursors" categories. Figure 5  
 623 reflects the quaternary categories for the 23 that were subset further.

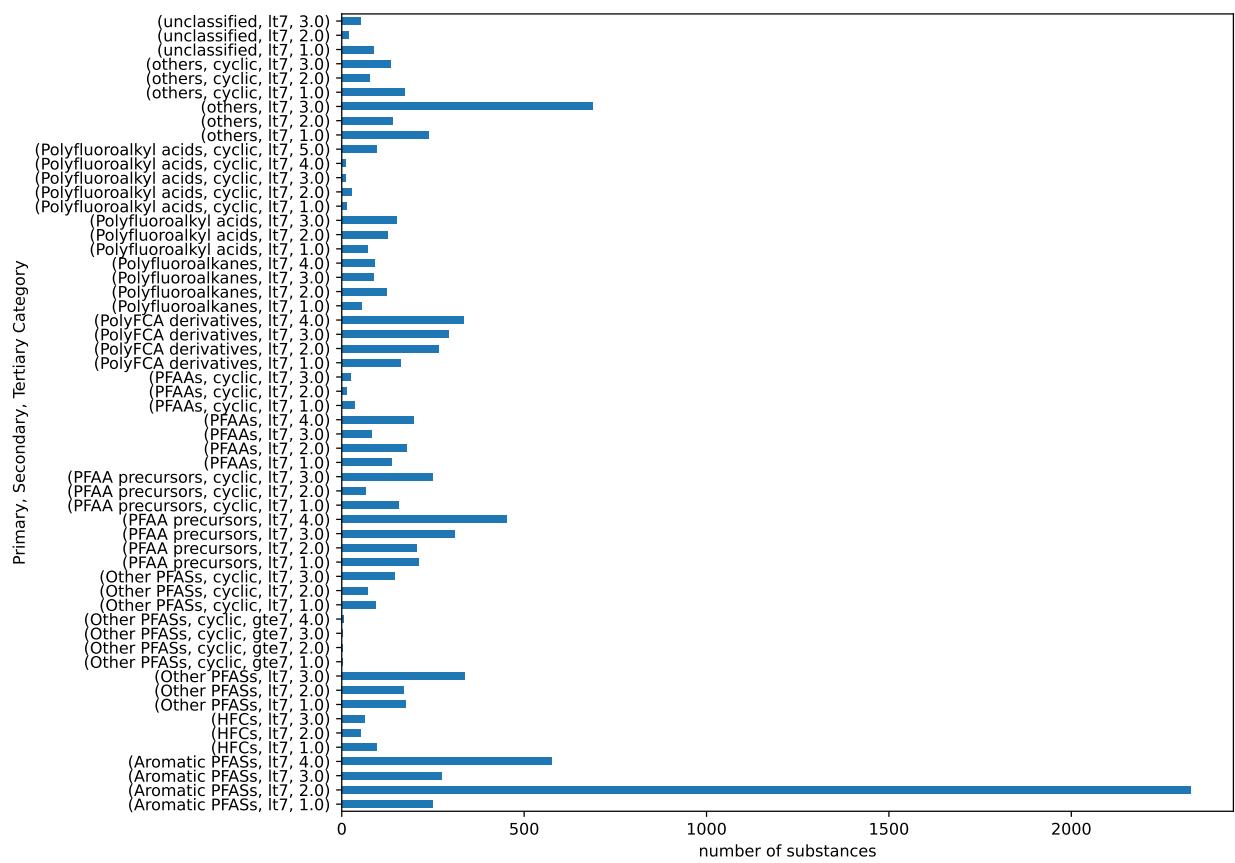


Figure 4: Bar chart showing the number of substances within each tertiary category, ordered by primary and secondary category roots. Methods to define tertiary categories are outlined in Section 2.7. Lt7, chain length less than 7; gte7, chain length greater than or equal to 7.

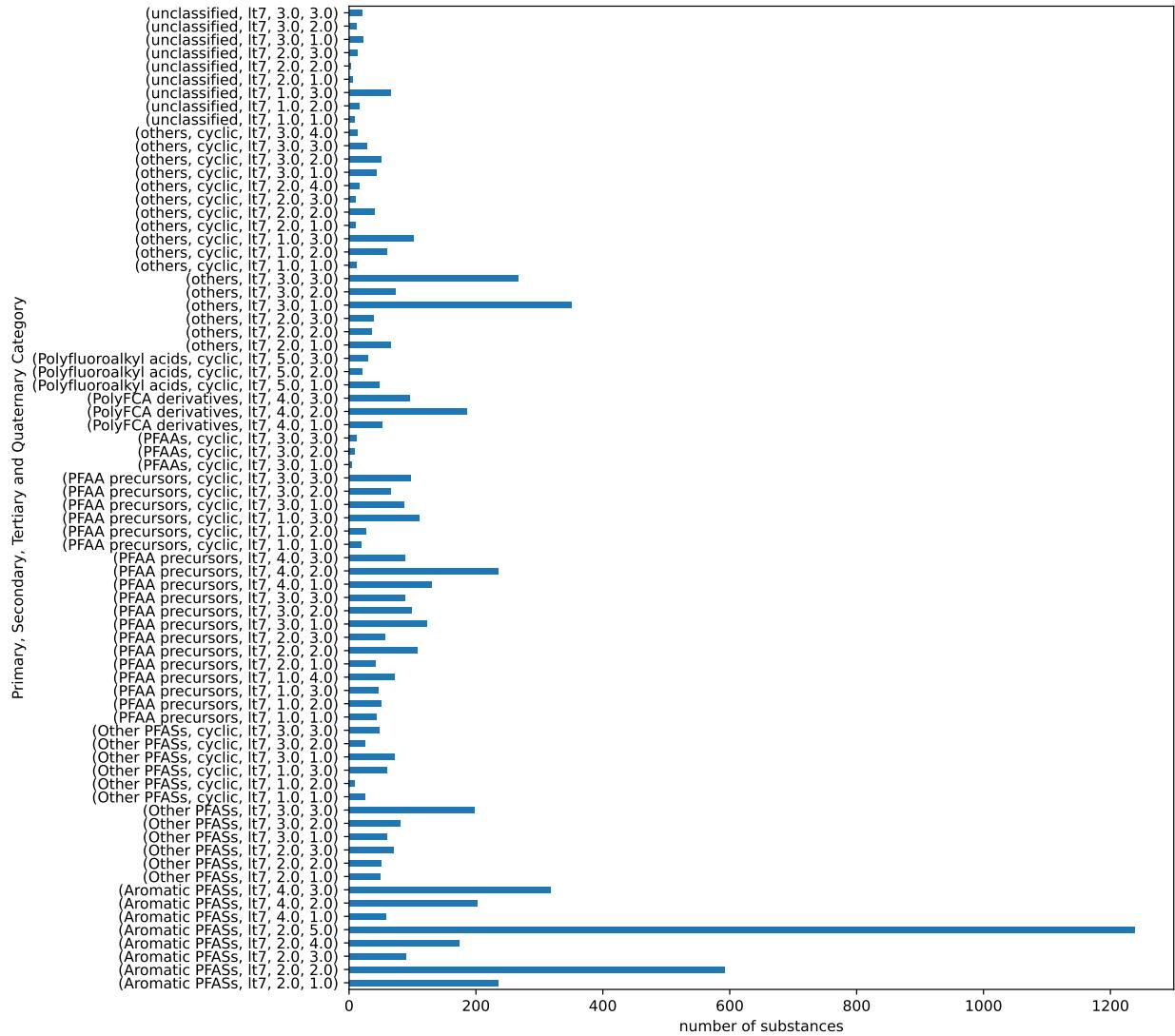


Figure 5: Bar chart showing the number of substances within each quaternary category, ordered by primary, secondary, and tertiary category roots. Methods to define quaternary categories are outlined in Section 2.7. Lt7, chain length less than 7; gte7, chain length greater than or equal to 7.

624 Terminal categories were defined as either secondary or tertiary categories that did not exceed  
625 the threshold, as well as all quaternary categories. A total of 128 terminal categories (127 categories  
626 + 1 singleton (PFAAs, cyclic, gte7)) were ultimately derived. This represented a trade-off in terms  
627 of the final number of terminal categories that was a practical number to characterize the landscape  
628 of PFAS balanced with maximizing structural similarity within the categories themselves. The full  
629 list of 15,525 substances together with their terminal category assignments are provided as supple-  
630 mentary information. Structural similarity within categories did increase following subcategorization,  
631 Figure A4 shows the ECDFs of several terminal categories which are left shifted relative to the  
632 original ECDFs for the secondary categories (Figure A2), i.e. the pairwise distance range decreases.

633 4.2. Selection of representative substances

634 Whilst centroids were selected as the most representative substance from each terminal category,  
635 there was a recognition that a single chemical was unlikely to capture the breadth of diversity within  
636 a category. Additional substances to capture the breadth and structural diversity relied on the  
637 MaxMinPicker method<sup>50</sup>. This method was used to select up to 3 further substances in addition to  
638 the centroid. A total of 484 substances were selected using this approach for 121 of the terminal  
639 categories. Terminal categories with 5 or fewer members did not result in any additional substances  
640 being selected (beyond the centroid) by the approach. Table 4 lists the 7 terminal categories which  
641 had insufficient membership to apply the MaxMinPicker approach.

Table 4: Terminal categories for which the MaxMin approach was not undertaken.

Terminal category	Membership
Other PFASs, cyclic, gte7, 1	2
Other PFASs, cyclic, gte7, 2	4
Other PFASs, cyclic, gte7, 3	2
PFAAs, cyclic, gte7	1
PFAAs, cyclic, lt7, 3, 1	4
unclassified, lt7, 2.0, 1.0	5
unclassified', lt7, 2.0, 2.0	2

642 To evaluate the proportion of structural diversity captured by the selected representative sub-  
643 stances, the normalized cumulative minimum distance was calculated as a function of the number of  
644 substances selected using the MaxMinPicker method as discussed in Section 2.3.9. There were 15  
645 terminal categories, out of the 121 terminal categories for which diverse substances were selected,  
646 where picking 3 substances captured at least 50% of the structural diversity (shown in Table 5).

Table 5: Terminal categories for which 3 representative substance selections capture more than 50% of the structural diversity.

Terminal category	Number of chemicals for 80% structural diversity	Cumulative % of Structural Diversity	Terminal Category size
Other PFASs, cyclic, gte7, 4	3	84.03	6
Polyfluoroalkyl acids, cyclic, lt7, 3	3	83.67	11
unclassified, lt7, 1.0, 1.0	4	67.35	8
Other PFASs, cyclic, lt7, 1.0, 2.0	4	67.14	9
PFAAs, cyclic, lt7, 3.0, 2.0	5	64.87	9
others, cyclic, lt7, 2.0, 3.0	5	63.47	10
others, cyclic, lt7, 3.0, 4.0	5	62.24	13
PFAAs, cyclic, lt7, 2.0	6	59.02	14
others, cyclic, lt7, 2.0, 1.0	5	58.79	10
PFAA precursors, cyclic, lt7, 1.0, 1.0	6	58.03	19
PFAAs, cyclic, lt7, 3.0, 3.0	6	57.47	12
Polyfluoroalkyl acids, cyclic, lt7, 4.0	6	55.93	12
unclassified, lt7, 3.0, 2.0	6	55.75	11
others, cyclic, lt7, 1.0, 1.0	6	55.17	12
Polyfluoroalkyl acids, cyclic, lt7, 1.0	6	54.9	15

647 Notes: Column 1 represents the number of substances that would be required to capture 80% of  
648 the structural diversity in the category, Cumulative % of Structural Diversity represents the normal-  
649 ized cumulative minimum distance for up to 3 selected diverse substances. Note the 80% used as a  
650 threshold is purely for illustrative purposes only.

651 For the largest terminal category, "Aromatic PFASs, It7, 2.0, 5.0", selecting up to 3 diverse sub-  
652 stances only captured 0.8% of the structural diversity. In order to capture 80% of the structural  
653 diversity for this terminal category, 528 substances would need to be selected for data collection.  
654 The number of substances that needed to be selected from each terminal category to capture 80%  
655 of the structural diversity varied from 3 (as shown above in Table 5) to 528 with the median number  
656 being 30.

657 Figure 6 shows the curves of the number of diverse selections as a function of the percentage  
658 normalized cumulative minimum distances for 10 representative terminal categories. These vary in  
659 steepness showing how quickly or not the structural diversity coverage converges with number of  
660 diverse selections depending on the terminal category of interest. Figure 7 highlights the difference  
661 in the number of diverse chemicals that would be needed to capture a minimum structural diversity  
662 across each terminal category.

663 Figure 8 attempts to summarize the tradeoff of the number of diverse chemicals (thus the centroids  
664 and MaxMin) as a function of % structural diversity captured across the terminal categories.

665 The diverse selections identified earlier for the terminal categories reflects a pragmatism in terms  
666 of identifying a potential candidate list of substances. As discussed later in Section 4.6, the struc-  
667 tural diversity captured forms one of the considerations in selecting candidates for additional data  
668 collection relative to those terminal categories that are data poor or contain substances that are on  
669 the TSCA inventory.

#### 670 4.3. Evaluation of physical state and physicochemical consistency within terminal categories

671 In order to determine the nature of further data collection activities and understanding the poten-  
672 tial presence in different environmental media, physical state and physicochemical information was  
673 determined for the PFAS landscape as far as possible. For the 15,525 substances in the PFAS land-  
674 scape, 431 substances (2.8%) could not be assigned into any specific physical state and physicochem-  
675 ical designation owing to a lack of predicted physicochemical property information. The designations  
676 of the remaining substances are shown in Table 6.

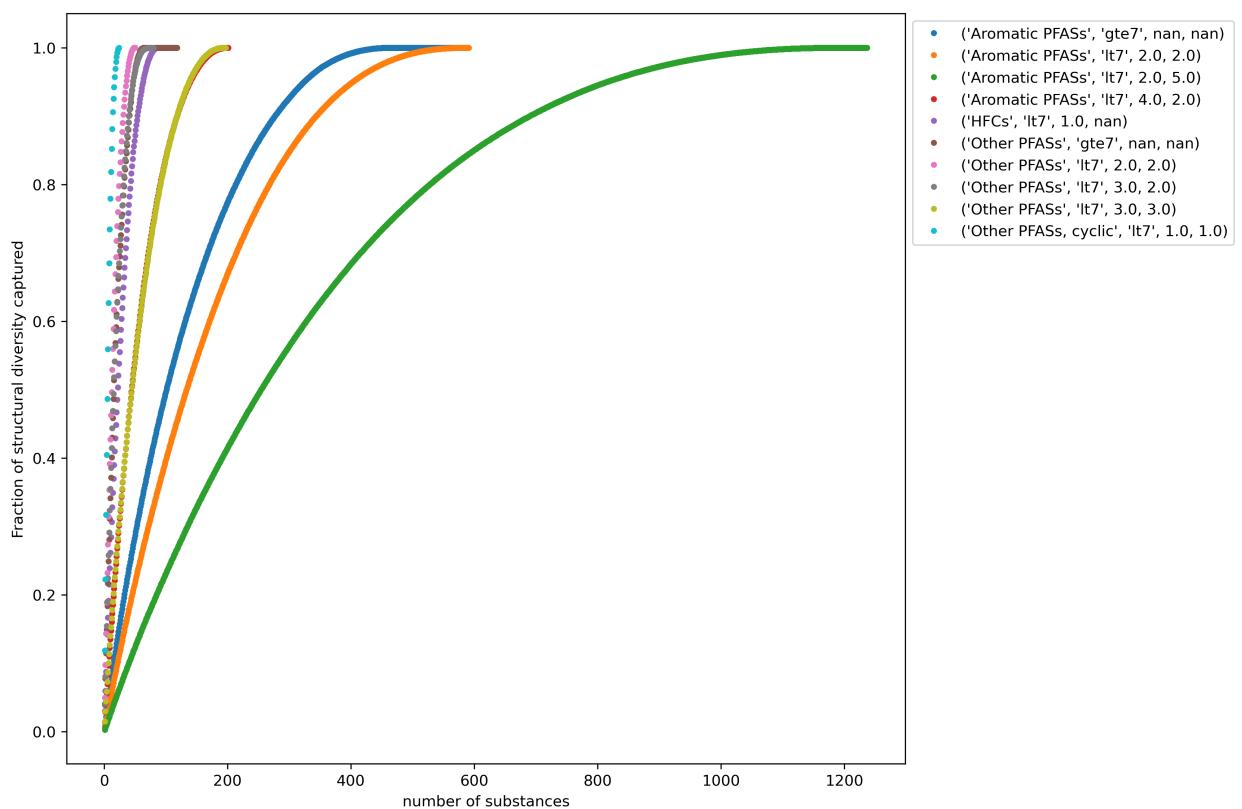


Figure 6: For a selection of terminal categories, the extent to which the fraction of structural diversity is captured relative to number of diverse chemicals selected varies.

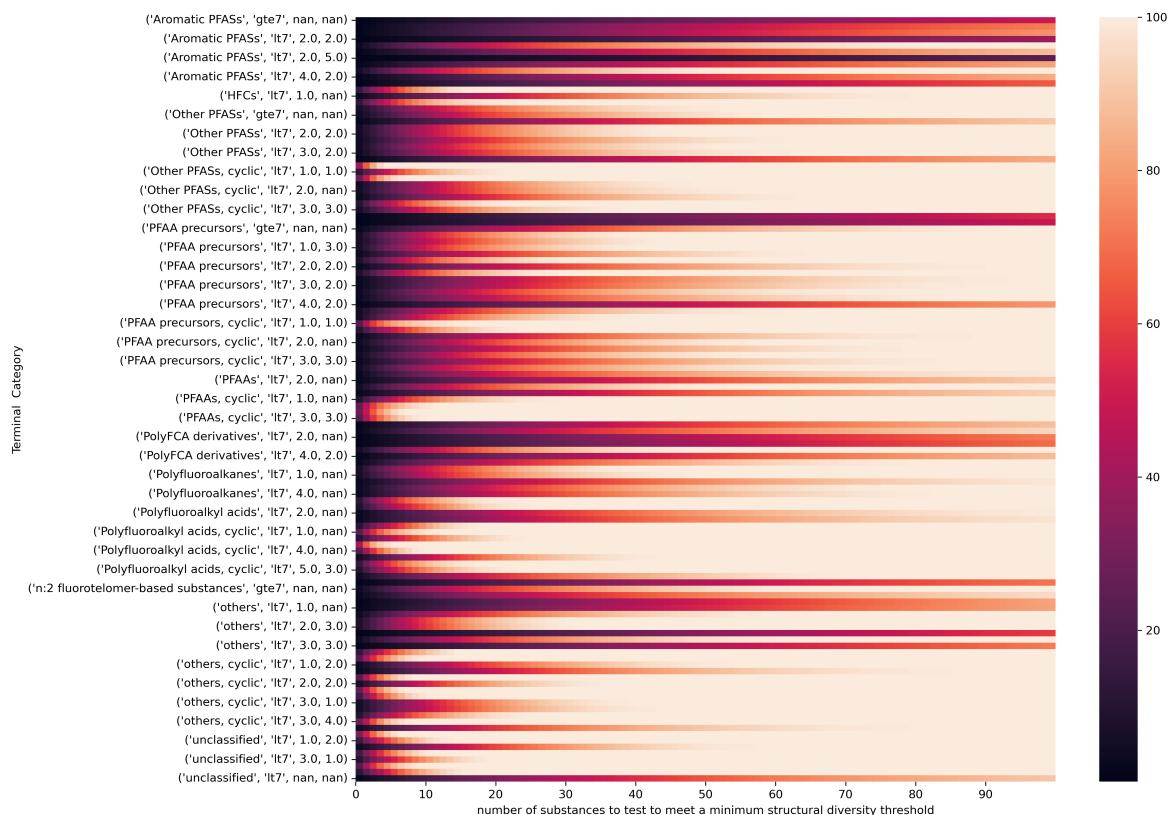


Figure 7: Heatmap showing the number of diverse substances that would need to be selected to achieve a specific minimum % structural diversity coverage across each terminal category. A selection of terminal categories are plotted to highlight how the number of diverse substances needed varies. The legend color corresponds to the number of substances needed as presented on the x axis.

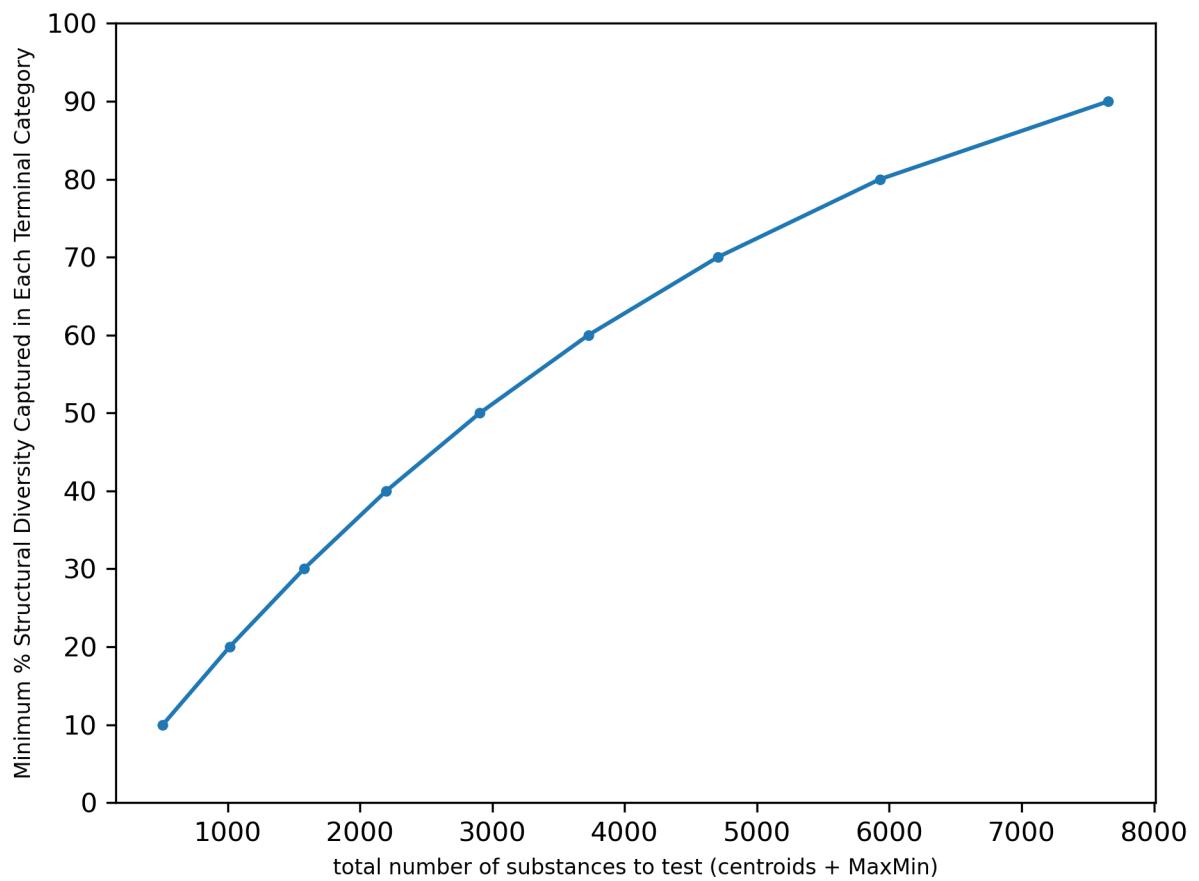


Figure 8: Lineplot showing the number of diverse substances that would need to be selected to achieve a specific minimum % structural diversity coverage across all the terminal categories. To achieve 80% structural diversity, a total of 5929 substances would need to be selected across the various terminal categories.

Table 6: Number (percentage) of substances assigned to each physical state and physicochemical designation.

Physical state and physicochemical designation	Full landscape	TSCA active constrained landscape
A (insoluble solids)	2060 (13.2%)	25 (12.6%)
B (soluble solids and soluble non-volatile liquids)	9824 (63.3%)	71 (35.7%)
C (soluble volatile liquids/insoluble liquids and soluble gases)	3115 (20%)	85 (42.7%)
D (insoluble gases or highly volatile gases)	95 (0.6%)	10 (5%)
No designation	431 (2.8%)	8 (4%)

677 The high percentage of substances in designation C additionally raises questions about the compatibility  
 678 with most NAM-based systems. Across the terminal categories, there was a general trend of  
 679 number of different designations increasing with size in category membership (see Figure A5). Figure  
 680 9 shows an example of one of the most diverse and largest terminal categories "Aromatic PFASs,  
 681 It7, 2.0, 5.0" which comprises 1238 members and spans 3 of the 4 designations. Although substances  
 682 predominantly lie within designation B, there is no discernible separation between the designations  
 683 across the structural category as characterized by Morgan fingerprints. In contrast all 96 substances  
 684 belonging to terminal category "PolyFCA derivatives, It7, 4.0, 3.0" fell into designation B (figure not  
 685 shown) whereas the 58 substances in "Aromatic PFASs, It7, 4.0, 1.0" fell into designations A and B.  
 686 There was a positive association between how structurally similar a terminal category was and the  
 687 consistency in physical state and physicochemical profile observed (as reflected by the designations).  
 688 However, the Morgan fingerprints could not resolve all the differences. For the selection of potential  
 689 candidates for data collection, the physical state and physicochemical profile remains an important  
 690 consideration in concert with the structural diversity described in Section 4.2.

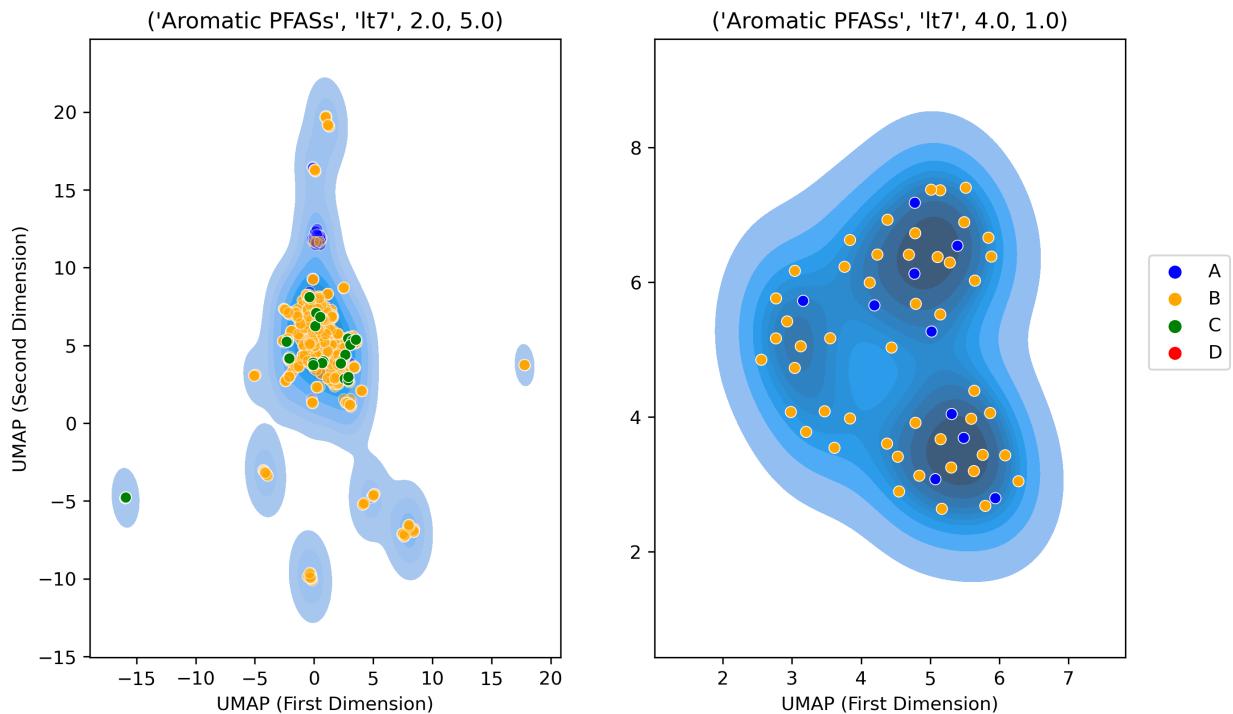


Figure 9: UMAP projections for terminal category a) "Aromatic PFASs, It7, 2.0, 5.0" and b) "Aromatic PFASs, It7, 4.0, 1.0" using Morgan chemical fingerprints with physical state and physicochemical designations A-D overlaid.

#### 691 4.4. Variation of POD values across and within terminal categories

692 Ultimately, the terminal categories are intended to facilitate a read-across for human health as-  
 693 sessment. To explore the feasibility of this further, the 25<sup>th</sup> percentile values of oral BMDhs were  
 694 calculated using available non-cancer data for 55 substances and repro/developmental toxicity data  
 695 for 35 substances. The distributions were plotted in a series of box plots. In vivo toxicity data were  
 696 available for at least one chemical in 28 of the 128 terminal categories across the two study types (28  
 697 for non-cancer, 19 for repro/developmental). The available data allowed preliminary trends for termi-  
 698 nal categories to be observed where the primary root was Aromatic PFASs, PASF-based substances,  
 699 PFAAs, Polyfluoroalkanes and Polyfluoroalkyl acids categories (see Figure A6 in the supplementary  
 700 information for the boxplots for both study types).

701 Figure 10 shows boxplot and strip plots for the oral non-cancer studies only. It appears that sub-  
 702 stances at each end of the spectrum of chain length within a category tended to exhibit lower toxicity,  
 703 i.e., their aggregate POD is higher. The spread of POD values within a category with greater diversity  
 704 in chain length tend to span ~ 1-2 orders of magnitude.

705 Although the available toxicity data are limited, there does appear to be some separation in the  
 706 potency distributions between terminal categories based on a common primary root. Inspection of

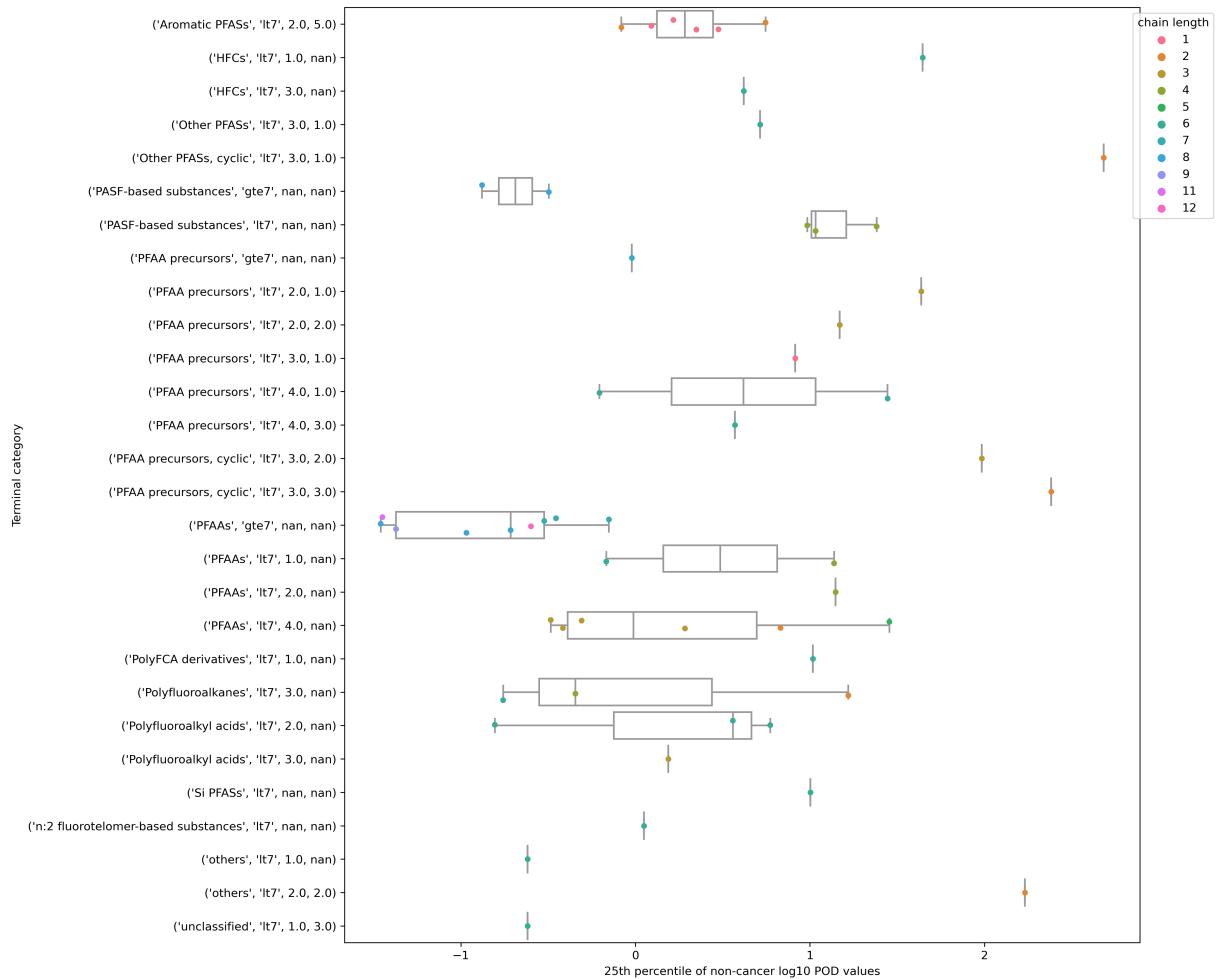


Figure 10: Boxplots showing the spread of the 25<sup>th</sup> percentile of the oral non-cancer log<sub>10</sub> POD values across and within terminal categories bounded by the carbon chain number. The box in the boxplot reflects the quartiles of the dataset, whilst the whiskers extend to + 1.5 \* inter-quartile range (IQR). Outliers are shown as points if they exceed 1.5 \* IQR.

707 Figure 10 does show a shift in potency values between the PFAA categories with a left shift for  
 708 those substances in the gte7 category vs the majority of the PFAA lt7 categories. A similar shift  
 709 was observed for the PASF-based categories. However, the relatively large spread for some of the  
 710 terminal categories suggests that additional refinement beyond structural similarity and chain length  
 711 (such as factoring in toxicokinetic information) will likely be needed for some terminal categories  
 712 prior to broader application in a read-across context.

#### 713 4.5. Qualitative mechanistic and toxicokinetic designations

714 There were six data streams with qualitative flags assigned for the ~150 PFAS tested as part  
 715 of the research project described in Patlewicz et al<sup>23</sup> namely: 1) nuclear receptor assays (NR); 2)  
 716 developmental toxicity (zebrafish testing); 3) DNT (developmental neurotoxicity); 4) thyroid toxicity;  
 717 5) immunosuppression (BioMAP assays); and 6) toxicokinetics (TK). Figure 11 profiles all the NAM flags  
 718 across the different technologies together with a stock QC flag<sup>30</sup> (Pass (red)) and a qc\_httk flag (Pass  
 719 (red)).

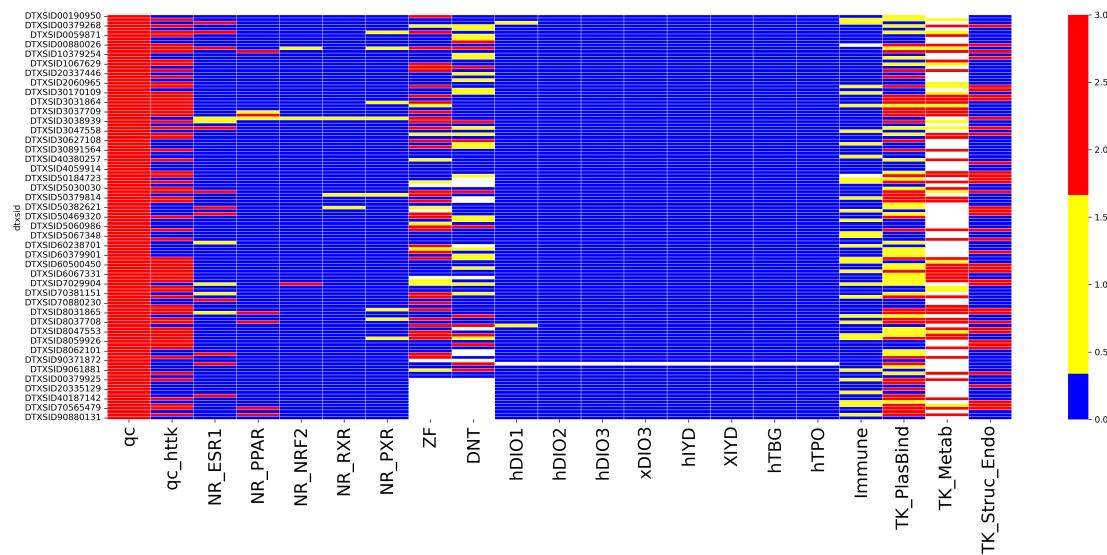


Figure 11: Heatmap of NAMs flags for the ~150 PFAS substances (~120 of which passed analytical QC (qc)) tested as part of the research programme described in Patlewicz et al<sup>23</sup>. Y axis tick labels do not capture all substances, only every 3<sup>rd</sup> substance by DTXSID is shown. qc = analytic QC and qc-httk analytical QC for TK; ESR1 = Estrogen Receptor 1; PPAR = peroxisome proliferator- activated receptor; NRF2 = nuclear factor erythroid 2-related factor 2; RXR= retinoid X receptor; PXR=pregnane X receptor; ZF = zebrafish; DNT = developmental neurotoxicity; DIO1, DIO2, DIO3 = Type 1,2,3 deiodinase; IYD = iodotyrosine deiodinase; TBG = thyroxine binding globulin; TPO = thyroid peroxidase.No data were represented as null values (white colored), data available but no flag identified as denoted a 0 (blue colored), 1 denoted a medium confidence flag (yellow colored) and 2 was associated with a high confidence flag (colored in red) consistent with the descriptions described in Table 2.

720 From Figure 11, the first two columns represent the quality control (QC) information. The next  
 721 5 columns represent the NR data. The next 2 columns represent the developmental toxicity (ZF)

assay and the DNT assay. The next 8 columns represent the thyroid assay outcomes followed by the integrated immunotoxicity flag from the BioMap assays. The last 3 columns represent the TK flags.

A semi quantitative analysis was performed by computing which PFAS ToxPrints were enriched for each NAM flag (excluding the TK\_Metab and TK\_Struc\_Endo flags since this information was captured using the predictions from the Dawson et al.<sup>56</sup> model). The full set of enriched ToxPrints are provided in the supplementary information. Those enriched PFAS ToxPrints were then used to profile the entire landscape to provide a predicted NAM flag profile. Performance metrics were derived to compare the actual and predicted NAM flag (see Table 7). Performance appeared weakest for the DNT and Immune flags which had the fewest number of ToxPrints that were enriched. It is worth noting that only ~150 chemicals (of which 124 substances overlapped with the PFAS landscape) were tested and the performance of these ToxPrint signatures may change as more substances are tested.

Table 7: Performance metrics for enriched PFAS ToxPrints

comparison	ROC_AUC_score	Sensitivity	Specificity
[NR_ESR1, pred_NR_ESR1]	0.67	0.47	0.87
[NR_PPAR, pred_NR_PPAR]	0.78	0.88	0.68
[NR_NRF2, pred_NR_NRF2]	0.82	1.00	0.64
[NR_PXR, pred_NR_PXR]	0.76	1.00	0.53
[ZF, pred_ZF]	0.66	0.59	0.74
[DNT, pred_DNT]	0.56	0.14	0.98
[Immune, pred_Immune]	0.55	0.10	0.99
[TK_PlasBind, pred_TK_PlasBind]	0.71	0.78	0.63

Predictions of half-life in humans were made for all the 15,525 substances. 69% of the substances were predicted to have the slowest half-life of greater than 2 months (bin 4), 14 % in the next slowest bin (bin 3, 1 week–2 months) and the remaining 17% in the second fastest bin (bin 2, 12 h–1 week).

To demonstrate the integration of the mechanistic and TK-related data with the structural cate-

737 gories, two terminal categories were clustered based on the predicted mechanistic and TK enrichment  
 738 flags as shown in Figure 12. Terminal categories “Aromatic PFASs, gte7” and “PFAAs, lt7, 4.0” were  
 739 first profiled against the enriched PFAS ToxPrints in conjunction with the predicted half-lives and  
 740 then clustered to show the potential NAM profile across the entire category and how consistent it was  
 741 across the terminal category members. Terminal category “Aromatic PFASs, gte7” had a far more  
 742 consistent profile whereas “PFAAs, lt7, 4.0” showed some differences in the half-life predictions  
 743 which might warrant further additional substances to be identified for data collection to reflect the  
 744 TK diversity.

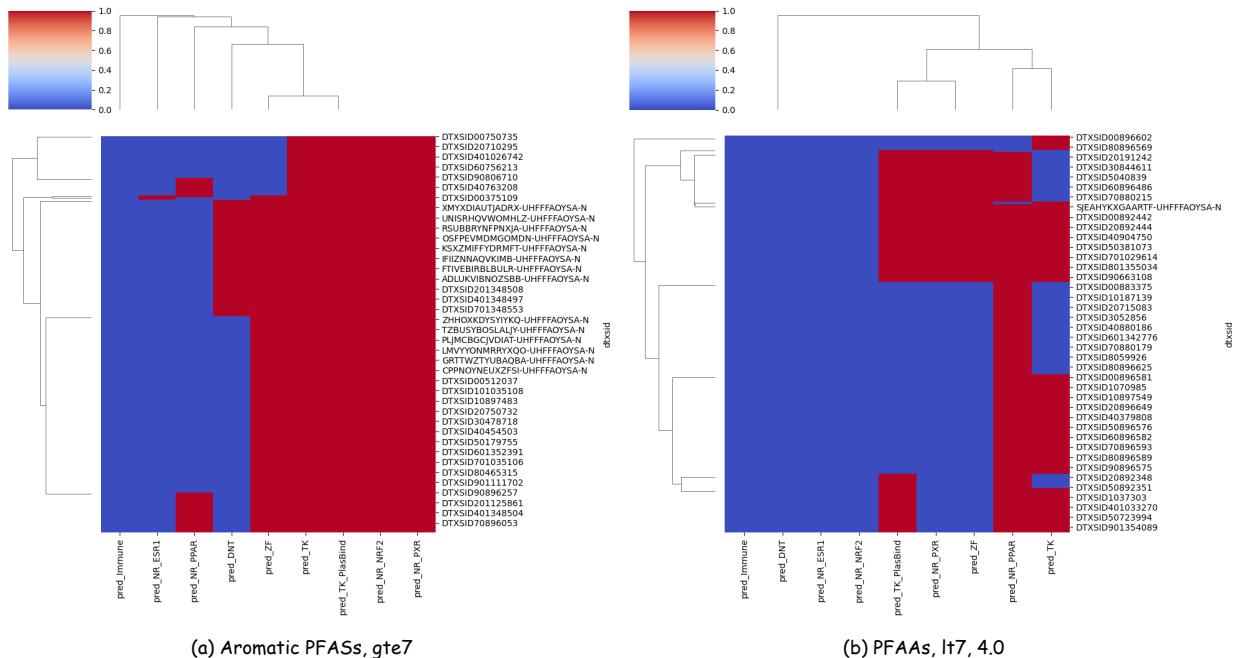


Figure 12: Clustermaps for terminal category “Aromatic PFASs, gte7” and “PFAAs, lt7, 4.0” to illustrate their concordance across predicted NAM profiles.

#### 745 4.6. Potential application to support the National PFAS Testing Strategy (NTS)

746 There are several considerations that come into play when identifying potential candidates for data  
 747 collection in concert with the landscape defined. To make the NTS actionable, one consideration  
 748 was to limit the landscape to one that was constrained by the TSCA active inventory to increase the  
 749 feasibility of being able to identify a manufacturer/importer of the substance. A second enables the  
 750 tradeoff between the number of diverse substances to select vs capturing the structural diversity  
 751 to be more practically addressed. Herein, the scope of terminal categories represented by the full  
 752 TSCA and TSCA active inventory and the impact this had in terms of capturing structural diversity  
 753 was evaluated. Finally, a proposal was outlined that considers how the terminal categories could be

754 triaged to initially focus on terminal categories which were either data poor or contained members  
755 that represented large exposure sources.

756 4.6.1. Constraining the landscape to the TSCA active inventory

757 4.6.1.1. TSCA inventory.

758 Of the substances in the PFAS landscape, only 563 substances were identified to be on the TSCA  
759 inventory, of which 237 were 'active' and the remaining 326 'inactive'. Active and inactive refers to  
760 the EPA's designation of whether a substance is active in US commerce based on the rule requiring  
761 industry to report chemicals manufactured or imported or processed in the US over a 10 year period  
762 ending 21st June 2016. There were 384 substances in the full landscape that matched a degradant  
763 of a TSCA substance. Figure 13 shows a bar chart of the membership of the terminal categories and  
764 how that differs when considering TSCA inventory status (overall or by active TSCA only).

765 The largest category memberships when constrained by presence on the TSCA inventory reflected  
766 the "PASF-based substances, lt7", "PASF-based substances, lt7" and "PFAA precursors, gte7" cate-  
767 gories. Across the terminal categories, 63% of the categories (80 out of the 128 categories) contain  
768 members on the TSCA inventory. If only categories containing substances that are on the TSCA  
769 active inventory are considered, then the number of terminal categories decreases to 60, i.e., 47%  
770 coverage. Some of the categories where there were no examples on the TSCA inventory were fairly  
771 large in size, examples include several of the Aromatic PFASs categories with 174-592 members as  
772 well as the PolyFCA derivatives and Polyfluoroalkyl acids with 186 and 151 members respectively.

773 4.6.1.2. Selection of representative substances in the constrained TSCA active inventory.

774 Centroids were computed for the 60 terminal categories containing substances that were on the  
775 active TSCA inventory. For 14 of these terminal categories, membership exceeded 5, which permitted  
776 the MaxMinPicker approach to be applied to identify further analogues. An additional 56 analogues  
777 were selected from this constrained landscape. Figure 14 shows the overlap in substances (centroids  
778 and diverse) across the unconstrained and the TSCA active constrained landscapes. The minimal  
779 overlap between the sets highlights the limitations of using a constrained landscape, i.e., one which  
780 does not represent the breadth of the PFAS chemistry. However, the substances on the TSCA  
781 active inventory represent those substances that are currently in commerce in the US and potentially  
782 represent the largest exposure source. It is worth noting that the overlap in the venn diagram in  
783 terms of exact substance may not reflect the overlap in structurally similar substances.

784 An evaluation of the structural diversity captured using the centroids and additional MaxMin sub-  
785 stances relative to the number of substances that would need to be selected to attain 80% structural

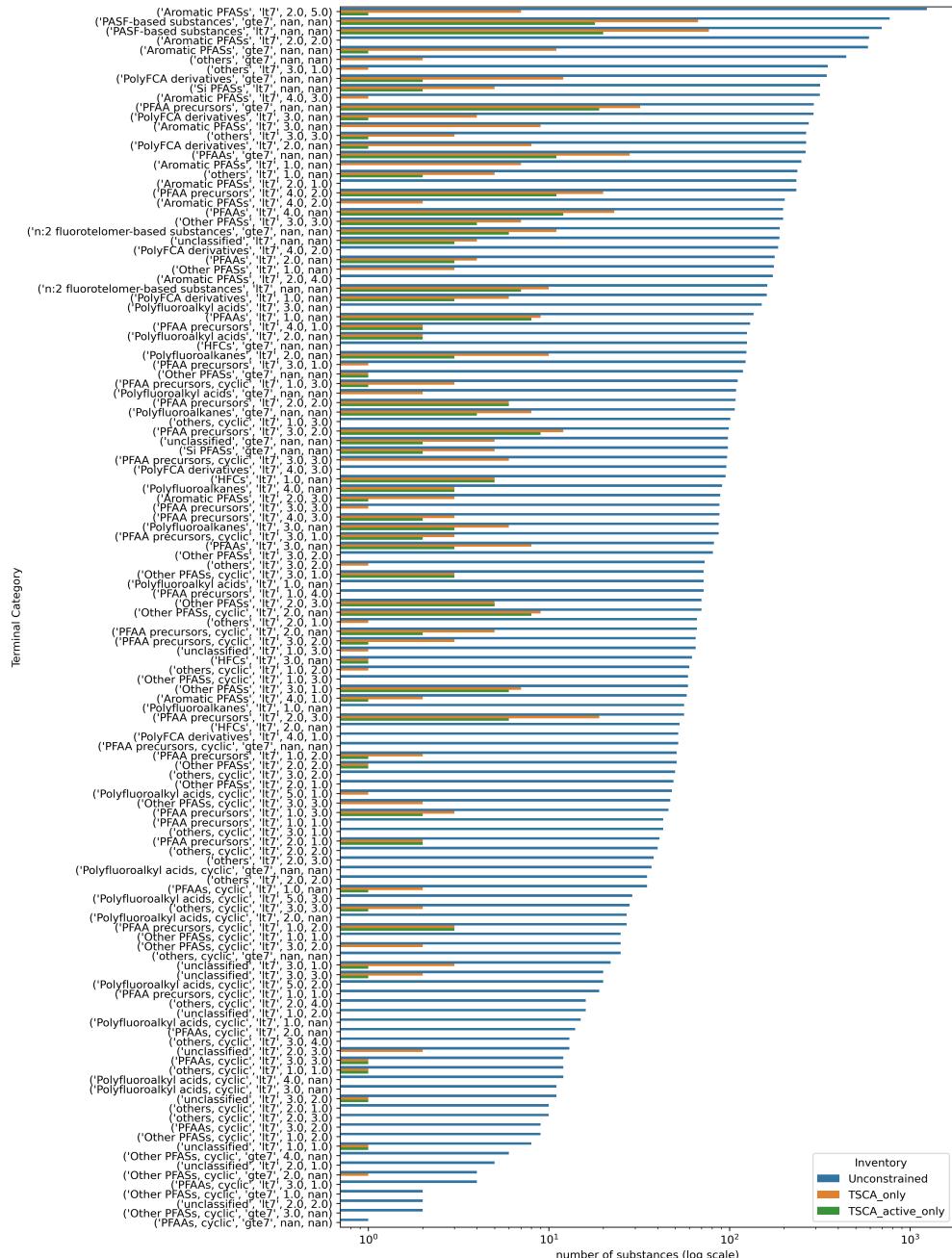


Figure 13: Bar chart showing membership of terminal categories and how that differs when constrained by TSCA inventory or TSCA active inventory.

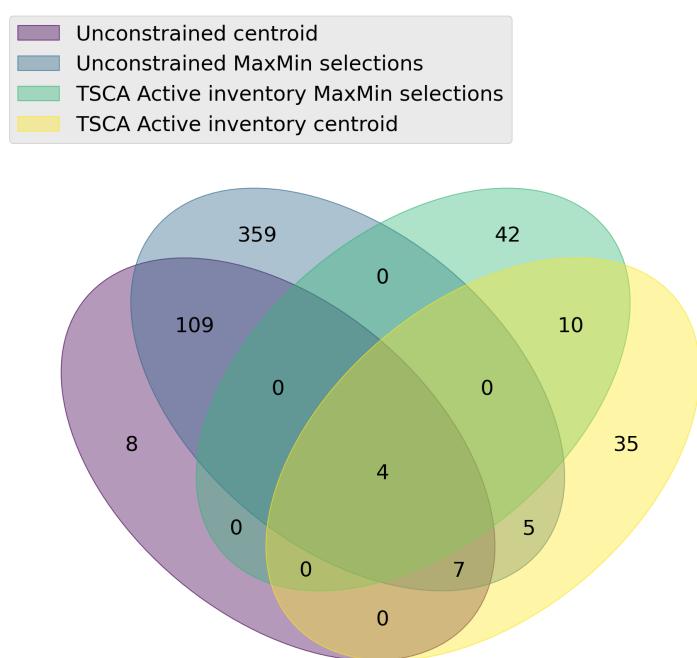


Figure 14: Venn diagram showing the overlap in substances based on whether they were identified as additional diverse picks or centroids in the unconstrained PFAS landscape and that constrained by the TSCA active inventory.

786 diversity coverage was also undertaken in the same manner as had been performed for the full land-  
 787 scape. For the 13 of the 14 categories where the MaxMin approach had been applied, the diverse picks  
 788 originally selected captured more than 50% of the structural diversity as shown in Table 8. This is  
 789 not so surprising given the TSCA active set substantially limited the terminal category size and in  
 790 turn their diversity.

Table 8: Terminal categories from the constrained TSCA active landscape where the MaxMin approach had been applied. Terminal categories for which 3 representative substance selections capture more than 50% of the structural diversity

Terminal category	Number of chemicals for 80% structural diversity	Cumulative % of structural diversity	Terminal category size
n:2	4	79.39	7
fluorotelomer-based substances, lt7			
Other PFASs, lt7, 3.0, 1.0	3	87.82	6
Other PFASs, cyclic, lt7, 2.0	4	75.73	8
PASF-based substances, gte7	6	50.43	18
PFAA precursors, gte7	4	66.79	19
PFAA precursors, lt7, 2.0, 2.0	3	83.96	6
PFAA precursors, lt7, 2.0, 3.0	2	97.06	6
PFAA precursors, lt7, 3.0, 2.0	5	66.84	9
PFAA precursors, lt7, 4.0, 2.0	5	63.06	11
PFAAs, gte7	3	90	11
PFAAs, lt7, 1.0	3	81.61	8
PFAAs, lt7, 4.0	5	59.97	12
n:2	1	100	6
fluorotelomer-based substances, gte7			

791 Notes: 'Number of chemicals for 80% structural diversity' represents the number of diverse sub-  
792 stances that would need to be selected to capture 80% of the structural diversity, 'Cumulative % of  
793 Structural Diversity' reflects the structural diversity captured by the up to 3 diverse substances  
794 already made and 'Terminal Category size' reflects the size of the terminal category if constrained  
795 by the availability of TSCA active substances.

796 For the largest terminal category, "PASF-based substances, lt7", 3 diverse substance selections  
797 captured 47.72% of the structural diversity. In order to capture 80% of the structural diversity  
798 for this terminal category, 7 substances would need to be selected for additional data collection.  
799 The number of substances to select from each terminal category to capture 80% of the structural  
800 diversity varied from 1 to 7 with the median number being 4. Across the entire TSCA active space,  
801 considering the 14 categories where a MaxMin approach could be applied - 55 substances would need  
802 to be selected to capture a 80% structural diversity. In order to capture a minimum of 80% structural  
803 diversity across all the TSCA active categories, at least 101 substances (the centroids for categories  
804 where no MaxMin had been applied + MaxMin) would be ideally selected for data collection. Figure 15  
805 summarizes the structural diversity attained across all TSCA active terminal categories.

#### 806 4.6.2. Proof of Concept Workflow: Identifying potential candidates for data collection

807 The availability of toxicity data across different study types and the presence of substances within  
808 different monitoring lists was arrayed across the terminal categories. All oral and inhalation studies  
809 from ToxValDB 9.5 were first retrieved. There were 76 substances with data for one or more of the  
810 study types which were then matched on the basis of DTXSID with substances in the PFAS landscape.  
811 The resulting table was then transformed to produce a table where columns represented different  
812 study types, rows were substances and cells were labelled 1 if data for a specific study type existed  
813 for a specific substance and 0 if no data existed. The qualitative lists reflecting various priorities,  
814 environmental detection/discharged, and availability of analytical methods etc. as described in Sec-  
815 tion 2.4.1 were compiled together and transformed into a table where rows represented substances,  
816 columns represented the different list sources and cells were populated with a 1 or 0 to denote pres-  
817 ence or absence on a specific list. There were 448 substances identified across these lists but only  
818 198 unique substances which were then matched with substances in the PFAS landscape. CDR status  
819 tags and Pubmed count tags were then added to the PFAS landscape. Columns representing the tox-  
820 icty study types and various lists were grouped by terminal category to produce a new table which  
821 reflected presence or absence of information (denoted by 1 or 0). Study quality was not considered  
822 - only the availability of publicly available toxicity data. The set of terminal categories were filtered  
823 to retain only those terminal categories which contained members on the TSCA active inventory (60  
824 terminal categories). Figure 16 provides a perspective of this information, namely the toxicity data

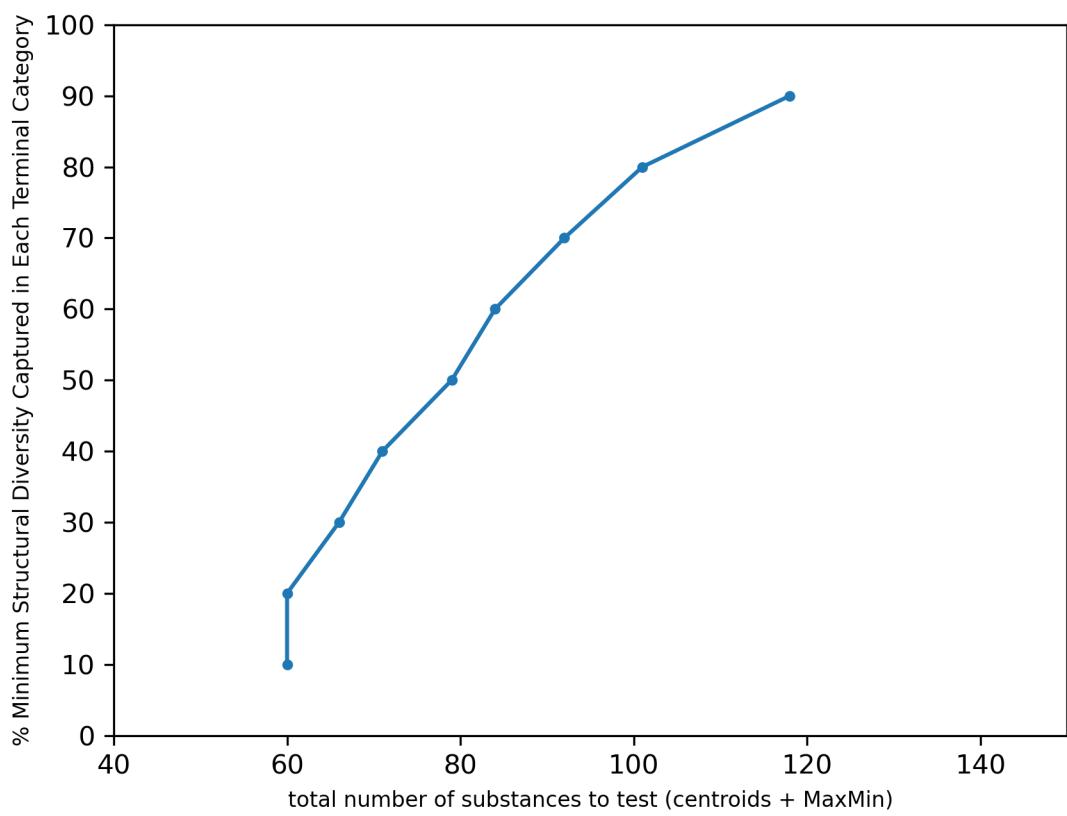


Figure 15: Lineplot of the TSCA active constrained terminal categories as a function of number of diverse substances selected and the minimum % structural diversity captured across all terminal categories.

825 sparsity across the categories that fall within the scope of the TSCA active inventory as well as dif-  
826 ferent environmental monitoring efforts or discussed in the literature. The PFAAs categories and  
827 their subcategorizations show up with data entries which is largely unsurprising, given the extent  
828 to which PFOA and PFOS have been studied. Note: The figure provides a landscape perspective of  
829 the data coverage across broad structural categories which may not entirely align with toxicological  
830 classes.

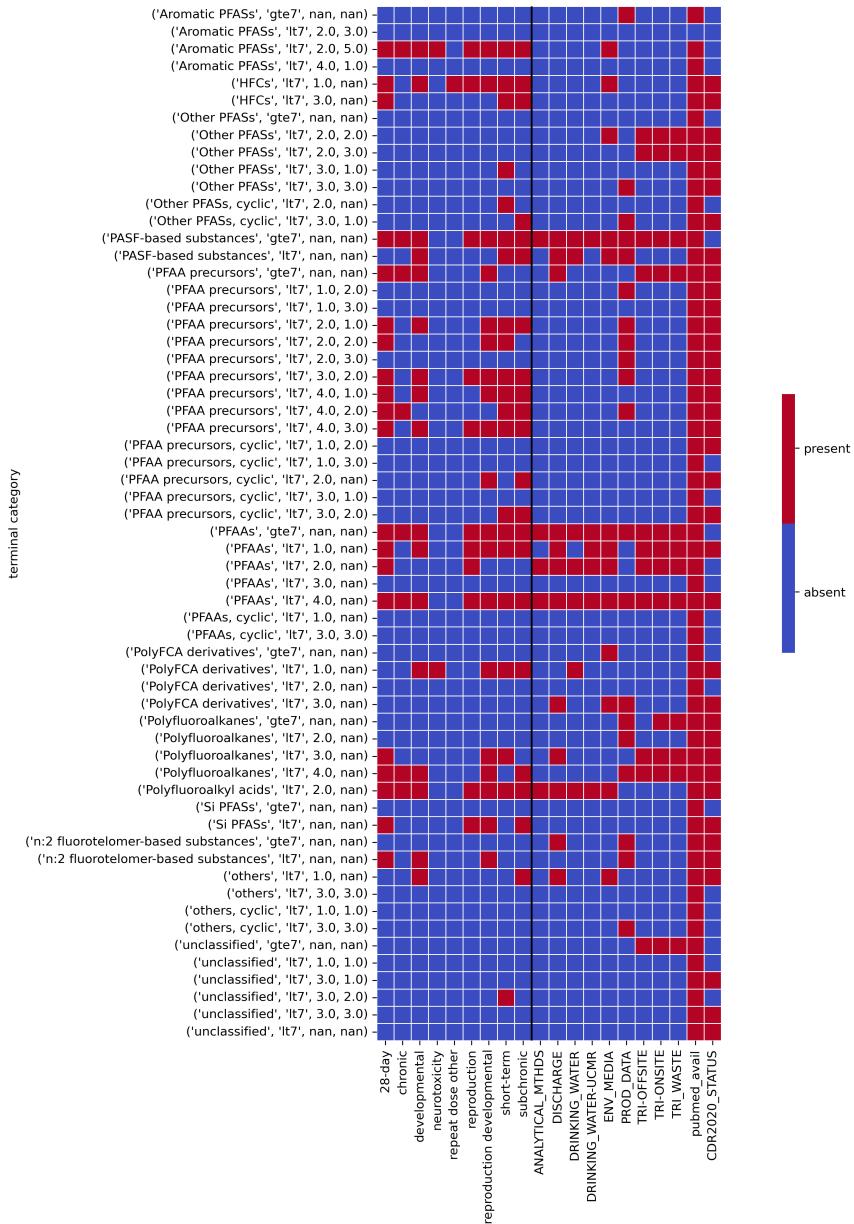


Figure 16: Heatmap of toxicity data availability and qualitative exposure and release designations. Notes: pubmed\_avail is a tag to denote presence or absence of articles indexed in Pubmed. PROD-Data = Production data, DISCHARGE = Discharge Monitoring data, DRINKING\_WATER = Drinking Water (State) Data, DRINKING\_WATER-UCMR = Drinking Water data comprising Unregulated Contaminant Monitoring Rule data and State level monitoring data, ENV\_MEDIA = Environmental Media data, TRI\_Waste = Toxics Release Inventory (TRI) Data Waste Managed, TRI\_On-Site = On Site TRI Data, TRI\_Off-Site = Off Site TRI Data, Analytical\_Mthds = PFAS with Validated Analytical Methods 533 and 537.

831      Each of the earlier sections in of themselves highlight different lines of evidence that can inform the  
 832      identification of potential candidates for data collection. Note: The qualitative release designations  
 833      are not intended to be exhaustive but could be refined to factor other relevant information data

834 streams such as existing epidemiological studies etc. Here, an attempt was made to demonstrate how  
 835 these steps can be integrated together to triage terminal categories and their potential candidates  
 836 for subsequent tiered data collection efforts (Figure 17). Step 1 is to consider a given terminal  
 837 category and determine whether it meets the condition of being a 'data poor category'. Data-poor  
 838 in this context was to consider whether this was a category that did not contain any members for  
 839 which repeated dose toxicity data existed (by the oral or inhalation route and with a reported NOAEL,  
 840 LOAEL, LOEL, NOEL, NEL or LEL value). Note in this study, only repeated dose toxicity data within  
 841 the publicly available ToxVal was considered. There were 94 terminal categories out of the 128 total  
 842 number of categories that met this condition.

843 The next step was to focus on terminal categories that overlapped with those which contained sub-  
 844 stances that were on the TSCA inventory.

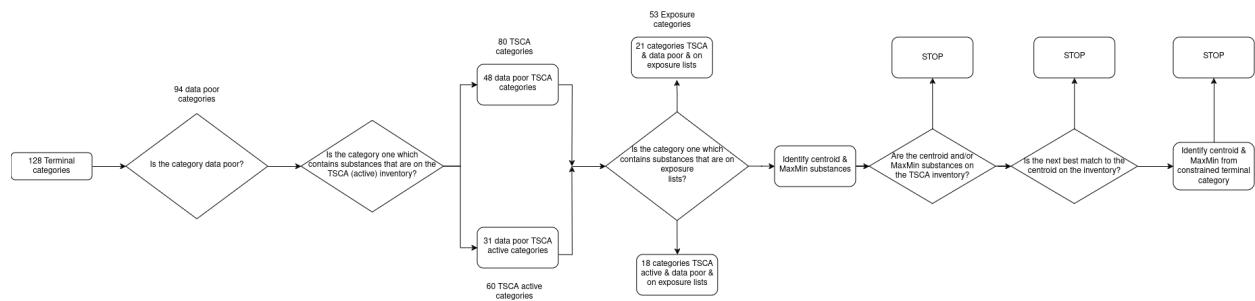


Figure 17: Workflow to highlight the main steps involved in prioritizing potential candidate selection for data collection for a given terminal category.

845 There were 80 terminal categories that contained substances that were on the TSCA inventory of  
 846 which 60 terminal categories contained substances that were on the TSCA active inventory. Of the  
 847 TSCA categories, 48 also satisfied the condition of being a 'data poor' category. In contrast, 31 of  
 848 the TSCA active categories were 'data poor'. The following step was to consider terminal categories  
 849 that contained substances that were on different environmental monitoring (EM) lists. There were  
 850 53 terminal categories that contained substances that were on one or more monitoring lists or 117  
 851 if Pubmed article availability was taken into account. Of these EM terminal categories, 21 were also  
 852 overlapping with data poor TSCA categories or 18 of the data poor TSCA active categories. The 18  
 853 terminal categories included "Aromatic PFASs, gte7", "Other PFASs, lt7, 2.0, 2.0", "PFAA precursors,  
 854 lt7, 1.0, 2.0", "PolyFCA derivatives, gte7", "Polyfluoroalkanes, gte7" and "unclassified, lt7, 2.0, 1.0".  
 855 Note: Only 2 of these categories met the condition to apply the MaxMin approach - taking into  
 856 account those categories for which 3 substances would need to be selected to achieve 80% structural  
 857 diversity, the remaining 16 categories would be limited to selecting the centroid.

858 For a category that satisfied all these conditions, the next step would be identify the representa-

859 tive substances characterizing the category (namely the centroid and MaxMin substances and check  
860 whether any were on the TSCA inventory). If none of these were on the inventory, then the next  
861 step would be to check whether the next closest match to the centroid was on the inventory. If not,  
862 the next steps would be to identify the centroid and MaxMin substances from either the TSCA con-  
863 strained landscape or the TSCA active constrained landscape for that terminal category. Figure 17  
864 summarizes these steps in a conceptual workflow.

865 For illustrative purposes, terminal category "PFAA precursors, lt7, 2.0, 3.0" was identified that met  
866 the conditions of being a data poor category, containing members on the TSCA active inventory and  
867 containing members on various environmental monitoring lists, discharge and TRI lists. This terminal  
868 category comprises 56 members. If the category were constrained by TSCA active substances only,  
869 the category size would be reduced to 6 members of which 2 substances would capture 80% of its  
870 structural diversity. The centroid, DTXSID70884511 was on the TSCA inventory but the TSCA  
871 active centroid DTXSID60880406 could be selected. Figure 18 shows an UMAP projection<sup>36</sup> with  
872 the centroid, MaxMin and TSCA centroid substances shown for illustrative purposes to highlight  
873 their relative positions in the structural space captured within the terminal category.

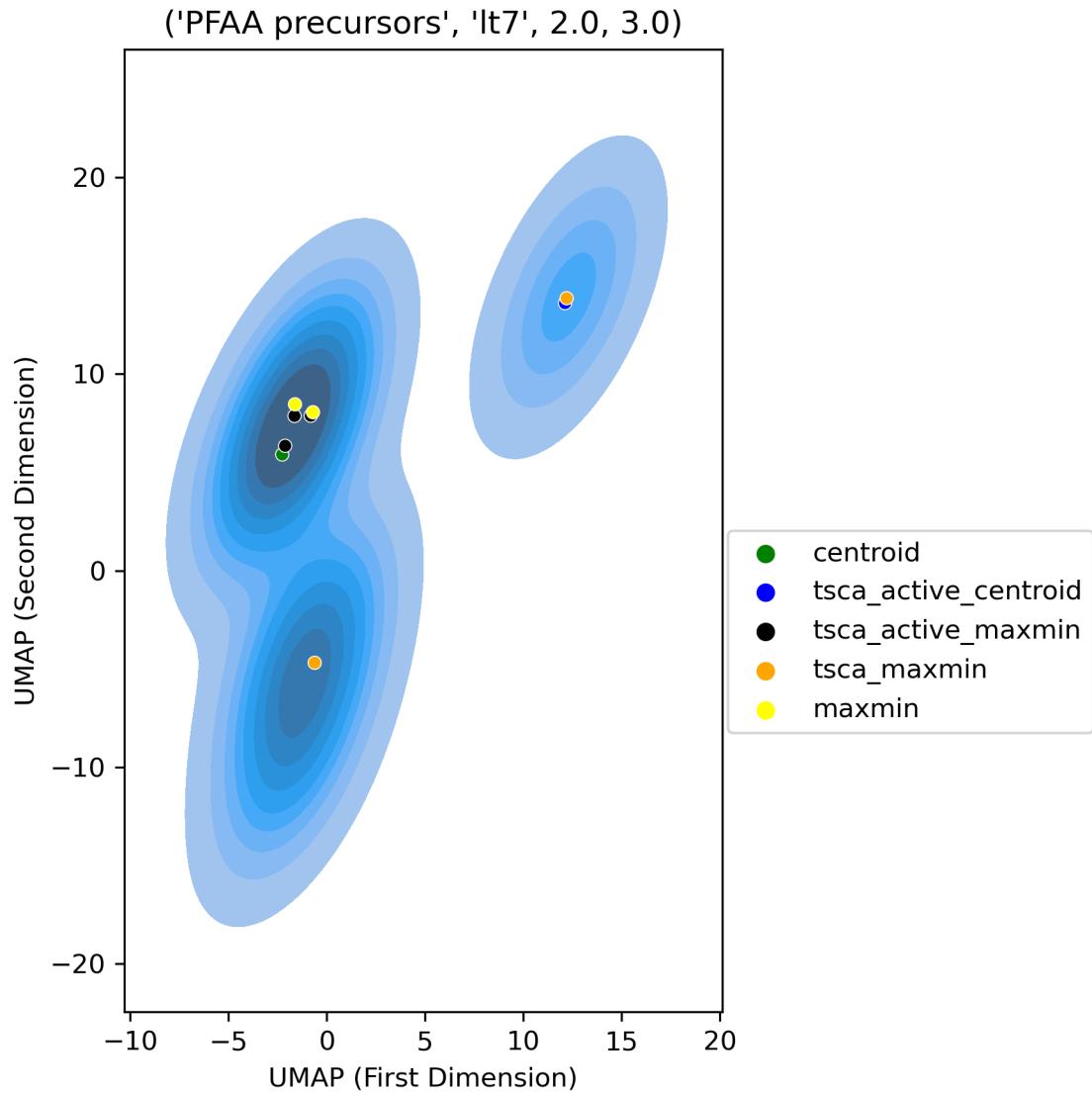


Figure 18: UMAP projection of terminal category "PFAA precursors, It7, 2.0, 3.0" with its (TSCA active) centroid and MaxMin substances shown.

<sup>874</sup> 4.7. Operationalizing the terminal categories for re-use

<sup>875</sup> There were 11 terminal categories that had memberships of less than 10 which were aggregated  
<sup>876</sup> into one miscellaneous group. A random forest classifier was trained to assignment membership of  
<sup>877</sup> a substance into one of 118 categories (117 terminal categories + a miscellaneous category). The  
<sup>878</sup> features used were Morgan fingerprints, chain length and primary category assignments whereas the  
<sup>879</sup> labels were the terminal category names. A random forest classifier with default settings applied as  
<sup>880</sup> part of a 5-fold stratified CV procedure gave rise to a mean balanced accuracy of 0.808 (std 0.014).

881 A randomized search CV procedure using a range of different hyperparameters found the highest  
882 balanced accuracy (BA) to be with 400 trees, a minimum split size of 2, a minimum number of data  
883 points allowed in a leaf node to be 2, a maximum features to be the square root of the number of  
884 features, and a balanced subsample class weight. The CV mean balanced accuracy was determined to  
885 be 0.845 (std 0.013). This model was then applied to the test set that had been held out. The balanced  
886 accuracy of the test set was determined to be 0.857. The BA varied across the terminal categories  
887 with a median value of 0.973 and a minimum value of 0.50. The worse performing categories were the  
888 unclassified categories namely "unclassified, lt7, 3.0, 1.0", "unclassified, lt7, 3.0, 2.0", "unclassified,  
889 lt7, 3.0, 3.0" and "unclassified, lt7, 1.0, 3.0" all of which had BA less than or equal to 0.5. There were 8  
890 categories with a BA less than 0.75 (of these 4 had a BA less than or equal to 0.5). The top 5 highest  
891 performing categories were "Aromatic PFASs, lt7, 2.0, 3.0", "others, cyclic, lt7, 3.0, 3.0", "others,  
892 cyclic, lt7, 1.0, 1.0", "PFAA precursors, cyclic, lt7, 1.0, 1.0" and "PFAAs, cyclic, lt7, 2.0". The full test  
893 set performance scores and the feature importances are provided in the supplementary information.  
894 Figure A8 shows the BA and recall scores across the terminal categories.

## 895 5. Conclusions

896 EPA was directed by Congress to develop a process for prioritizing which PFAS or classes of PFAS  
897 should be subject to additional research efforts based on potential for human exposure, toxicity, and  
898 other available information. Herein, we describe an approach that can be used to create a relevant  
899 PFAS landscape using the TSCA section 8(a)(7) rule definition to continue the efforts initiated in the  
900 National PFAS Testing Strategy.

901 A landscape of 13,054 PFAS substances was first created based on the structural definitions out-  
902 lined in the TSCA rule. The landscape was then augmented with 2484 simulated degradation products  
903 of PFAS substances on the TSCA inventory using the Catalogic expert system. Adding simulated  
904 degradates was intended to enrich the landscape by substances that might be expected to be found  
905 in the environment from existing substances in commerce. The simulated degradation products were  
906 derived from an expert system which includes training set substances that are PFAS; however a full  
907 characterization of the model relative to the PFAS landscape was not feasible as some of the training  
908 set was proprietary in nature. For the portion of training set substances that could be evaluated -  
909 there was a minimal overlap in datasets as shown in Figure A1. The robustness of the simulated degra-  
910 dation products is a limitation in the approach and requires additional work but a pragmatic one given  
911 the absence of data to refine and improve the model further.

912 The 15,525 substances in the PFAS landscape were then grouped into 128 terminal structural cat-  
913 egories based on a stepwise process that combined OECD functional categories, chain length and

914 structural similarity. The use of a chain length threshold of 7 was a pragmatic choice to help iden-  
915 tify persistent long chain substances, though the subcategorization using this threshold may be best  
916 suited for straight chain linear PFAS. Some of the terminal categories were very large and structurally  
917 diverse whilst others showed much greater structural similarity highlighting both the complexity of  
918 the PFAS structural landscape and local areas of homogeneity.

919 A method was developed to select the most representative substance (centroid) and other sub-  
920 stances (MaxMin) to capture the structural diversity of each terminal category and guide data col-  
921 lection efforts. A substantial number of representative and diverse substances (~6000) would be  
922 required to capture 80% percent of structural diversity in the terminal categories for the uncon-  
923 strained landscape. Significantly fewer representative and diverse substances (101) would be required  
924 to capture 80% percent of structural diversity in the terminal categories for the TSCA constrained  
925 landscape (though if ToxVal data availability was factored in, this would reduce to 76 substances).  
926 The difference in utilizing the unconstrained and TSCA constrained landscape highlights the chal-  
927 lenges in data collection to address future and theoretical data gaps versus those data gaps that  
928 exist amongst substances known to be in commerce.

929 Publicly available in vivo data across the terminal categories were used to evaluate whether read-  
930 across could be potentially viable based on the variation of the in vivo data itself. A 25<sup>th</sup> percentile  
931 of the derived human equivalent BMDs served as surrogate POD value for a given substance. Not  
932 all terminal categories were associated with toxicity data but for those categories, the following in-  
933 sights were noted; substances with very short or long carbon chain length within a category tended to  
934 exhibit lower toxicities (i.e., higher PODs), but the spread of PODs within a category could be large  
935 particularly for diverse categories based on carbon chain length, spanning 1-2 orders of magnitude or  
936 more. In addition to the shift in potency between terminal categories containing longer vs shorter  
937 chain lengths, there was also a shift between terminal categories with different functional groups  
938 e.g. Aromatic PFASs tended to be less potent vs. PFAAs. The variability of in vivo PODs from tradi-  
939 tional toxicity tests within some of the terminal categories suggests that structural considerations  
940 may not be sufficient for performing read-across without additional data collection.

941 Information from NAMs were layered on the terminal categories to help identify potentially distinct  
942 mechanistic and TK-related subgroups. The NAM information was intended to help refine terminal  
943 categories, guide candidate selection, and inform data collection efforts. Currently, the terminal  
944 categories and candidates for data collection are primarily identified based on chemical structural  
945 considerations; however, other factors such as toxicokinetics, hazard, and modes-of-action are also  
946 important when considering a category and read-across approach<sup>22</sup>. To incorporate these other fac-  
947 tors, the current process for selecting representative and diverse substances (i.e., centroid, MaxMin)

948 could also be applied to the mechanistic and TK-related subgroups to select candidates for data collec-  
949 tion as well as inform the types of tests that may be useful in characterizing the hazards associated  
950 with a specific terminal category. For example, subgroups predicted to have endocrine-related activ-  
951 ity may benefit from developmental and reproductive tests, whilst those predicted to have significant  
952 cross-species TK differences may have limited benefits from rodent-based TK studies. While the cur-  
953 rent NAM data and enrichment flags are limited, the performance of these models will improve over  
954 time as additional substances are tested. Similar approaches have been proposed or incorporated  
955 into case studies, but not operationalized for large groups of chemicals<sup>20,59,60</sup>. Information from  
956 environmental measurements and release, traditional toxicity data, and chemical properties (physical  
957 state and physicochemical properties) were also layered on the terminal categories to help identify  
958 priorities for further data collection efforts.

959 The landscape of PFAS substances is substantially large and diverse with limited human health data.  
960 A category approach enables strategic data collection with the longer-term goal of enabling read-  
961 across within a particular category. In an effort to address this goal, a stepwise systematic process  
962 was developed to group the substances using a combination of OECD primary categories, sequential flu-  
963 orinated carbon chain length, and structural similarity. The process attempted to balance maximizing  
964 structural similarity within each category relative to a manageable number of categories. Within each  
965 category, representative substances to capture the structural diversity were identified to guide data  
966 collection efforts. A substantial number of substances were required to capture a large percentage  
967 of structural diversity in the terminal categories for the full PFAS landscape, whilst a significantly  
968 fewer number were needed to capture the structural diversity for the TSCA active constrained land-  
969 scape. The difference in utilizing the full and TSCA active constrained PFAS landscape highlights the  
970 challenges in data collection to address future and theoretical data gaps versus those data gaps that  
971 exist amongst substances known to be in commerce. To assist in prioritizing the categories for data  
972 collection, information from environmental measurements and release, traditional toxicity data, and  
973 exposure considerations based on chemical properties were incorporated to focus on those categories  
974 of greatest need. The variability in POD values from existing traditional in vivo toxicity tests within  
975 some of the terminal categories suggests that the structural considerations may not be sufficient  
976 for performing read-across without additional data collection. TK information is another important  
977 factor that can help resolve the variability observed. It should also be noted that POD values are  
978 one of many considerations - a reference dose effect will likely vary across similar chemicals within  
979 a category even if they exhibit similar hazard profiles. Information from NAMs were used to help  
980 refine the terminal categories based on potentially distinct mechanistic and TK-related subgroups  
981 and inform the types of data collection activities that may be required. Finally, a machine learning

982 modelling approach was applied in an attempt to build a predictive model to operationalize the terminal  
983 categories developed, such that new PFAS not already captured in the landscape could be profiled  
984 and assigned to their most probable terminal category. The methods developed for categorizing the  
985 PFAS landscape, selecting representative substances, refining categories based on mechanistic and  
986 TK information, and prioritizing categories for data collection provide a robust foundation to aid  
987 EPA in addressing the significant challenges associated with evaluating the environmental and human  
988 health impacts of this class of chemicals. The methods and associated categories are flexible in  
989 accommodating additional data as it is generated and may evolve as the scientific knowledge grows.

## 990 **6. Acknowledgements**

991 The authors thank Drs Meghan Tierney, Louis (Gino) Scarano, Charles Lowe and Nathaniel Charest  
992 for their thoughtful reviews of this manuscript. We also thank the anonymous reviewers for their  
993 comments.

## 994 **Funding**

995 The work presented in this manuscript was solely supported by appropriated funds of the US Envi-  
996 ronmental Protection Agency (US EPA).

## 997 **Disclaimer**

998 The approach presented and the scientific views discussed in this manuscript are those of the au-  
999 thors and do not necessarily reflect final policies of the US Environmental Protection Agency (US  
1000 EPA). Any mention of trade names, manufacturers or products does not imply an endorsement by the  
1001 U.S. Government or the EPA. The EPA and its employees do not endorse any commercial products,  
1002 services, or enterprises.

## 1003 **References**

- 1004 [1] Z. Wang, J. C. DeWitt, C. P. Higgins, I. T. Cousins, [A Never-Ending Story of Per- and Polyfluoroalkyl Substances \(PFASs\)?](#),  
1005 Environmental Science & Technology 51 (5) (2017) 2508-2518, number: 5 Publisher: American Chemical Society. [doi:  
10.1021/acs.est.6b04806](https://doi.org/10.1021/acs.est.6b04806).  
1006 URL <https://doi.org/10.1021/acs.est.6b04806>
- 1007 [2] J. Glüge, M. Scheringer, I. T. Cousins, J. C. DeWitt, G. Goldenman, D. Herzke, R. Lohmann, C. A. Ng, X. Trier, Z. Wang, [An  
1008 overview of the uses of per- and polyfluoroalkyl substances \(PFAS\)](#), Environmental Science: Processes & Impacts 22 (12)  
1009 (2020) 2345-2373, publisher: The Royal Society of Chemistry. [doi:10.1039/D0EM00291G](https://doi.org/10.1039/D0EM00291G).  
1010 URL <https://pubs.rsc.org/en/content/articlelanding/2020/em/d0em00291g>

- 1012 [3] L. G. T. Gaines, Historical and current usage of per- and polyfluoroalkyl substances (PFAS):  
1013 A literature review, American Journal of Industrial Medicine 66 (5) (2023) 353-378, \_eprint:  
1014 <https://onlinelibrary.wiley.com/doi/pdf/10.1002/ajim.23362>. doi:10.1002/ajim.23362.  
1015 URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/ajim.23362>
- 1016 [4] D. J. Wallis, K. E. Barton, D. R. U. Knappe, N. Kotlarz, C. A. McDonough, C. P. Higgins, J. A. Hoppin, J. L. Adgate, Source  
1017 apportionment of serum PFASs in two highly exposed communities, The Science of the Total Environment 855 (2023)  
1018 158842. doi:10.1016/j.scitotenv.2022.158842.
- 1019 [5] Y. Chen, H. Zhang, Y. Liu, J. A. Bowden, T. M. Tolaymat, T. G. Townsend, H. M. Solo-Gabriele, Evaluation of per- and  
1020 polyfluoroalkyl substances (PFAS) in leachate, gas condensate, stormwater and groundwater at landfills, Chemosphere  
1021 318 (2023) 137903. doi:10.1016/j.chemosphere.2023.137903.  
1022 URL <https://www.sciencedirect.com/science/article/pii/S0045653523001704>
- 1023 [6] J. Li, B. Xi, G. Zhu, Y. Yuan, W. Liu, Y. Gong, W. Tan, A critical review of the occurrence, fate and treatment of per-  
1024 and polyfluoroalkyl substances (PFASs) in landfills, Environmental Research 218 (2023) 114980. doi:10.1016/j.envres.  
1025 2022.114980.  
1026 URL <https://www.sciencedirect.com/science/article/pii/S0013935122023076>
- 1027 [7] N. Bolan, B. Sarkar, M. Vithanage, G. Singh, D. C. W. Tsang, R. Mukhopadhyay, K. Ramadass, A. Vinu, Y. Sun, S. Ramanayaka,  
1028 S. A. Hoang, Y. Yan, Y. Li, J. Rinklebe, H. Li, M. B. Kirkham, Distribution, behaviour, bioavailability and remediation of poly-  
1029 and per-fluoroalkyl substances (PFAS) in solid biowastes and biowaste-treated soil, Environment International 155 (2021)  
1030 106600. doi:10.1016/j.envint.2021.106600.  
1031 URL <https://www.sciencedirect.com/science/article/pii/S0160412021002257>
- 1032 [8] OECD, Reconciling Terminology of the Universe of Per- and Polyfluoroalkyl Substances: Recommendations and Practical  
1033 Guidance, Series on Risk Management No. 61., Tech. rep. (2021).
- 1034 [9] N. Gaber, L. Bero, T. J. Woodruff, The Devil they Knew: Chemical Documents Analysis of Industry Influence on PFAS  
1035 Science, Annals of Global Health 89 (1) 37. doi:10.5334/aogh.4013.  
1036 URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10237242/>
- 1037 [10] Z. Wang, A. M. Buser, I. T. Cousins, S. Demattio, W. Drost, O. Johansson, K. Ohno, G. Patlewicz, A. M. Richard, G. W.  
1038 Walker, G. S. White, E. Leinala, A New OECD Definition for Per- and Polyfluoroalkyl Substances, Environmental Science  
1039 & Technology 55 (23) (2021) 15575-15578, number: 23 Publisher: American Chemical Society. doi:10.1021/acs.est.  
1040 1c06896.  
1041 URL <https://doi.org/10.1021/acs.est.1c06896>
- 1042 [11] C. M. Grulke, A. J. Williams, I. Thillanadarajah, A. M. Richard, EPA's DSSTox database: History of development of a cu-  
1043 rated chemistry resource supporting computational toxicology research, Computational Toxicology (Amsterdam, Nether-  
1044 lands) 12 (Nov. 2019). doi:10.1016/j.comtox.2019.100096.
- 1045 [12] E. L. Schymanski, J. Zhang, P. A. Thiessen, P. Chirsir, T. Kondic, E. E. Bolton, Per- and Polyfluoroalkyl Substances (PFAS)  
1046 in PubChem: 7 Million and Growing, Environmental Science & Technology (Oct. 2023). doi:10.1021/acs.est.3c04855.
- 1047 [13] U. EPA, Per- and Poly-Fluoroalkyl Chemical Substances Designated as Inactive on the TSCA Inventory: Significant New  
1048 Use Rule (2023).  
1049 URL <https://www.federalregister.gov/documents/2023/01/26/2023-01156/per--and-poly-fluoroalkyl-chemical->  
1050 substances-designated-as-inactive-on-the-tscas-inventory
- 1051 [14] U. EPA (2023). [link].  
1052 URL <https://www.federalregister.gov/documents/2023/10/11/2023-22094/toxic-substances-control-act-reporting->  
1053 and-recordkeeping-requirements-for-perfluoroalkyl-and
- 1054 [15] U. EPA (2024). [link].

- 1055 URL <https://www.federalregister.gov/documents/2024/01/11/2024-00412/per--and-poly-fluoroalkyl-chemical-substances-designated-as-inactive-on-the-tsca-inventory>
- 1056 [16] L. M. Carlson, M. Angrish, A. V. Shirke, E. G. Radke, B. Schulz, A. Kraft, R. Judson, G. Patlewicz, R. Blain, C. Lin, N. Vetter, C. Lemeris, P. Hartman, H. Hubbard, X. Arzuaga, A. Davis, L. V. Dishaw, I. L. Druwe, H. Hollinger, R. Jones, J. P. Kaiser, L. Lizarraga, P. D. Noyes, M. Taylor, A. J. Shapiro, A. J. Williams, K. A. Thayer, Systematic Evidence Map for Over One Hundred and Fifty Per- and Polyfluoroalkyl Substances (PFAS), *Environmental Health Perspectives* 130 (5) (2022) 056001. doi:10.1289/EHP10343.
- 1061 [17] U. EPA, *Status and future directions of the high production volume challenge program* (2004).
- 1062 URL <https://nepis.epa.gov/Exe/ZyNET.exe/P1004QXK.TXT?ZyActionD=ZyDocument&Client=EPA&Index=2000+Thru+2005&Docs=&Query=&Time=&EndTime=&SearchMethod=1&TocRestrict=n&Toc=&TocEntry=&QField=&QFieldYear=&QFieldMonth=&QFieldDay=&IntQFieldOp=0&ExtQFieldOp=0&XmlQuery=&File=D%3A%5Czyfiles%5CIndex%20Data%5C00thru05%5CTxt%5C00000021%5CP1004QXK.txt&User=ANONYMOUS&Password=anonymous&SortMethod=h%7C-&MaximumDocuments=1&FuzzyDegree=0&ImageQuality=r75g8/r75g8/x150y150g16/i425&Display=hpfr&DefSeekPage=x&SearchBack=ZyActionL&Back=ZyActionS&BackDesc=Results%20page&MaximumPages=1&ZyEntry=1&SeekPage=x&ZyPURL#>
- 1063 [18] OECD, *Guidance on Grouping of Chemicals, Second Edition OECD Series on Testing and Assessment, No 194*, OECD Publishing (2017). doi:<https://doi.org/10.1787/9789264274679-en>.
- 1064 URL <https://www.oecd.org/publications/guidance-on-grouping-of-chemicals-second-edition-9789264274679-en.htm>
- 1065 [19] M. T. D. Cronin, *Chapter 1:An Introduction to Chemical Grouping, Categories and Read-Across to Predict Toxicity*, in: Chemical Toxicity Prediction, 2013, pp. 1-29. doi:10.1039/9781849734400-00001.
- 1066 URL <https://pubs.rsc.org/en/content/chapter/bk9781849733847-00001/978-1-84973-384-7>
- 1067 [20] S. E. Escher, H. Kamp, S. H. Bennekou, A. Bitsch, C. Fisher, R. Graepel, J. G. Hengstler, M. Herzler, D. Knight, M. Leist, U. Norinder, G. Ouédraogo, M. Pastor, S. Stuard, A. White, B. Zdrail, B. van de Water, D. Kroese, *Towards grouping concepts based on new approach methodologies in chemical hazard assessment: the read-across approach of the EU-ToxRisk project*, *Archives of Toxicology* 93 (12) (2019) 3643-3667, number: 12. doi:10.1007/s00204-019-02591-7.
- 1068 URL <http://link.springer.com/10.1007/s00204-019-02591-7>
- 1069 [21] G. Patlewicz, M. T. Cronin, G. Helman, J. C. Lambert, L. E. Lizarraga, I. Shah, *Navigating through the minefield of read-across frameworks: A commentary perspective*, *Computational Toxicology* 6 (2018) 39-54. doi:10.1016/j.comtox.2018.04.002.
- 1070 URL <https://linkinghub.elsevier.com/retrieve/pii/S2468111318300331>
- 1071 [22] G. Patlewicz, I. Shah, *Towards systematic read-across using Generalised Read-Across (GenRA)*, *Computational Toxicology* 25 (2023) 100258. doi:10.1016/j.comtox.2022.100258.
- 1072 URL <https://www.sciencedirect.com/science/article/pii/S2468111322000469>
- 1073 [23] G. Patlewicz, A. M. Richard, A. J. Williams, R. S. Judson, R. S. Thomas, *Towards reproducible structure-based chemical categories for PFAS to inform and evaluate toxicity and toxicokinetic testing*, *Computational Toxicology* 24 (2022) 100250. doi:10.1016/j.comtox.2022.100250.
- 1074 URL <https://www.sciencedirect.com/science/article/pii/S246811132200038X>
- 1075 [24] A. O. Stucki, T. S. Barton-Maclaren, Y. Bhuller, J. E. Henriquez, T. R. Henry, C. Hirn, J. Miller-Holt, E. G. Nagy, M. M. Perron, D. E. Ratzlaff, T. J. Stedeford, A. J. Clippinger, *Use of new approach methodologies (nams) to meet regulatory requirements for the assessment of industrial chemicals and pesticides for effects on human health*, *Frontiers in Toxicology* 4 (2022). doi:10.3389/ftox.2022.964553.
- 1076 URL <https://www.frontiersin.org/articles/10.3389/ftox.2022.964553>
- 1077 [25] K. E. Carstens, T. Freudenrich, K. Wallace, S. Choo, A. Carpenter, M. Smeltz, M. S. Clifton, W. M. Henderson, A. M.

- 1098 Richard, G. Patlewicz, B. A. Wetmore, K. Paul Friedman, T. Shafer, [Evaluation of Per- and Polyfluoroalkyl Substances](#)  
1099 ([PFAS](#)) In Vitro Toxicity Testing for Developmental Neurotoxicity, *Chemical Research in Toxicology* 36 (3) (2023) 402-  
1100 419, publisher: American Chemical Society. [doi:10.1021/acs.chemrestox.2c00344](#).  
1101 URL <https://doi.org/10.1021/acs.chemrestox.2c00344>
- 1102 [26] K. A. Houck, G. Patlewicz, A. M. Richard, A. J. Williams, M. A. Shobair, M. Smeltz, M. S. Clifton, B. Wetmore, A. Medvedev,  
1103 S. Makarov, Bioactivity profiling of per- and polyfluoroalkyl substances (PFAS) identifies potential toxicity pathways  
1104 related to molecular structure, *Toxicology* 457 (2021) 152789. [doi:10.1016/j.tox.2021.152789](#).
- 1105 [27] K. A. Houck, K. P. Friedman, M. Feshuk, G. Patlewicz, M. Smeltz, M. S. Clifton, B. A. Wetmore, S. Velichko, A. Berenyi,  
1106 E. L. Berg, Evaluation of 147 perfluoroalkyl substances for immunotoxic and other (patho)physiological activities through  
1107 phenotypic screening of human primary cells, *ALTEX* 40 (2) (2023) 248-270. [doi:10.14573/altex.2203041](#).
- 1108 [28] A. Kreutz, M. S. Clifton, W. M. Henderson, M. G. Smeltz, M. Phillips, J. F. Wambaugh, B. A. Wetmore, [Category-Based](#)  
1109 [Toxicokinetic Evaluations of Data-Poor Per- and Polyfluoroalkyl Substances \(PFAS\) using Gas Chromatography Coupled](#)  
1110 [with Mass Spectrometry](#), *Toxics* 11 (5) (2023) 463. [doi:10.3390/toxics11050463](#).  
1111 URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10223284/>
- 1112 [29] M. Smeltz, J. F. Wambaugh, B. A. Wetmore, [Plasma Protein Binding Evaluations of Per- and Polyfluoroalkyl Substances for](#)  
1113 [Category-Based Toxicokinetic Assessment](#), *Chemical Research in Toxicology* (May 2023). [doi:10.1021/acs.chemrestox.](#)  
1114 [3c00003](#).
- 1115 [30] M. G. Smeltz, M. S. Clifton, W. M. Henderson, L. McMillan, B. A. Wetmore, Targeted Per- and Polyfluoroalkyl substances  
1116 ([PFAS](#)) assessments for high throughput screening: Analytical and testing considerations to inform a PFAS stock quality  
1117 evaluation framework, *Toxicology and Applied Pharmacology* 459 (2023) 116355. [doi:10.1016/j.taap.2022.116355](#).
- 1118 [31] T. E. Stoker, J. Wang, A. S. Murr, J. R. Bailey, A. R. Buckalew, High-Throughput Screening of ToxCast PFAS Chemical  
1119 Library for Potential Inhibitors of the Human Sodium Iodide Symporter, *Chemical Research in Toxicology* 36 (3) (2023)  
1120 380-389. [doi:10.1021/acs.chemrestox.2c00339](#).
- 1121 [32] S. J. Degitz, J. H. Olker, J. S. Denny, P. P. Degoeij, P. C. Hartig, M. C. Cardon, S. A. Eytcheson, J. T. Haselman, S. A.  
1122 Mayasich, M. W. Hornung, [In vitro screening of per- and polyfluorinated substances \(pfas\) for interference with seven](#)  
1123 [thyroid hormone system targets across nine assays](#), *Toxicology in Vitro* 95 (2024) 105762. [doi:https://doi.org/10.](#)  
1124 [1016/j.tiv.2023.105762](#).  
1125 URL <https://www.sciencedirect.com/science/article/pii/S0887233323002114>
- 1126 [33] A. J. Williams, C. M. Grulke, J. Edwards, A. D. McEachran, K. Mansouri, N. C. Baker, G. Patlewicz, I. Shah, J. F. Wambaugh,  
1127 R. S. Judson, A. M. Richard, [The CompTox Chemistry Dashboard: a community data resource for environmental chemistry](#),  
1128 *Journal of Cheminformatics* 9 (1) (2017) 61. [doi:10.1186/s13321-017-0247-6](#).
- 1129 [34] S. R. Heller, A. McNaught, I. Pletnev, S. Stein, D. Tchekhovskoi, [InChI, the IUPAC International Chemical Identifier](#),  
1130 *Journal of Cheminformatics* 7 (1) (2015) 23. [doi:10.1186/s13321-015-0068-4](#).  
1131 URL <https://doi.org/10.1186/s13321-015-0068-4>
- 1132 [35] D. Rogers, M. Hahn, [Extended-Connectivity Fingerprints](#), *Journal of Chemical Information and Modeling* 50 (5) (2010)  
1133 742-754, publisher: American Chemical Society. [doi:10.1021/ci100050t](#).  
1134 URL <https://doi.org/10.1021/ci100050t>
- 1135 [36] L. McInnes, J. Healy, J. Melville, Umap: Uniform manifold approximation and projection for dimension reduction (2020).  
1136 [arXiv:1802.03426](https://arxiv.org/abs/1802.03426).
- 1137 [37] S. Dimitrov, G. Dimitrova, T. Pavlov, N. Dimitrova, G. Patlewicz, J. Niemela, O. Mekyan, A stepwise approach for defining  
1138 the applicability domain of sar and qsar models., *J Chem Inf Model* 45 (2005) 839-849. [doi:10.1021/ci0500381](#).
- 1139 [38] A. Su, Y. Cheng, C. Zhang, Y.-F. Yang, Y.-B. She, K. Rajan, An artificial intelligence platform for automated pfas subgroup  
1140 classification: A discovery tool for pfas screening, *Science of the Total Environment* 921 (2024) 171229. [doi:10.1016/](#)

- 1141       j.scitotenv.2024.171229.
- 1142 [39] A. Su, K. Rajan, A database framework for rapid screening of structure-function relationships in PFAS chemistry, *Scientific Data* 8 (1) (2021) 14, number: 1. doi:10.1038/s41597-021-00798-x.
- 1143 [40] K. Mansouri, C. M. Grulke, R. S. Judson, A. J. Williams, OPERA models for predicting physicochemical properties and environmental fate endpoints, *Journal of Cheminformatics* 10 (1) (2018) 10. doi:10.1186/s13321-018-0263-1.
- 1144 [41] W. S. Chambers, J. G. Hopkins, S. M. Richards, *A Review of Per- and Polyfluorinated Alkyl Substance Impairment of Reproduction*, *Frontiers in Toxicology* 3 (2021).
- 1145       URL <https://www.frontiersin.org/articles/10.3389/ftox.2021.732436>
- 1146 [42] K. Sznajder-Katarzyńska, M. Surma, I. Cieślik, *A Review of Perfluoroalkyl Acids (PFAAs) in terms of Sources, Applications, Human Exposure, Dietary Intake, Toxicity, Legal Regulation, and Methods of Determination*, *Journal of Chemistry* 2019 (2019) e2717528, publisher: Hindawi. doi:10.1155/2019/2717528.
- 1147       URL <https://www.hindawi.com/journals/jchem/2019/2717528/>
- 1148 [43] A. M. Richard, H. Hidle, G. Patlewicz, A. J. Williams, *Identification of Branched and Linear Forms of PFOA and Potential Precursors: A User-Friendly SMILES Structure-based Approach*, *Frontiers in Environmental Science* 10 (2022).
- 1149       URL <https://www.frontiersin.org/article/10.3389/fenvs.2022.865488>
- 1150 [44] A. M. Richard, R. Lougee, M. Adams, H. Hidle, C. Yang, J. Rathman, T. Magdziarz, B. Bienfait, A. J. Williams, G. Patlewicz, A New CSRML Structure-Based Fingerprint Method for Profiling and Categorizing Per- and Polyfluoroalkyl Substances (PFAS), *Chemical Research in Toxicology* 36 (3) (2023) 508-534. doi:10.1021/acs.chemrestox.2c00403.
- 1151 [45] C. Yang, A. Tarkhov, J. Marusczyk, B. Bienfait, J. Gasteiger, T. Kleinoeder, T. Magdziarz, O. Sacher, C. H. Schwab, J. Schwoebel, L. Terfloth, K. Arvidson, A. Richard, A. Worth, J. Rathman, New publicly available chemical query language, CSRML, to support chemotype representations for application to data mining and modeling, *Journal of Chemical Information and Modeling* 55 (3) (2015) 510-528. doi:10.1021/ci500667v.
- 1152 [46] G. L. Landrum, *RDKit: Open-source cheminformatics*;
- 1153       URL <http://www.rdkit.org>
- 1154 [47] N. M. O'Boyle, R. A. Sayle, *Comparing structural fingerprints using a literature-based similarity benchmark*, *Journal of Cheminformatics* 8 (1) (2016) 36. doi:10.1186/s13321-016-0148-0.
- 1155       URL <https://doi.org/10.1186/s13321-016-0148-0>
- 1156 [48] J. W. Raymond, C. J. Blankley, P. Willett, *Comparison of chemical clustering methods using graph- and fingerprint-based similarity measures*, *Journal of Molecular Graphics and Modelling* 21 (5) (2003) 421-433, type: Journal Article. doi: [https://doi.org/10.1016/S1093-3263\(02\)00188-2](https://doi.org/10.1016/S1093-3263(02)00188-2).
- 1157       URL <https://www.sciencedirect.com/science/article/pii/S1093326302001882>
- 1158 [49] J. H. Ward, *Hierarchical Grouping to Optimize an Objective Function*, *Journal of the American Statistical Association* 58 (301) (1963) 236-244, publisher: Taylor & Francis \_eprint: <https://www.tandfonline.com/doi/pdf/10.1080/01621459.1963.10500845>. doi:10.1080/01621459.1963.10500845.
- 1159       URL <https://www.tandfonline.com/doi/abs/10.1080/01621459.1963.10500845>
- 1160 [50] M. Ashton, J. Barnard, F. Casset, M. Charlton, G. Downs, D. Gorse, J. Holliday, R. Lahana, P. Willett, *Identification of Diverse Database Subsets using Property-Based and Fragment-Based Molecular Descriptions*, *Quantitative Structure-Activity Relationships* 21 (6) (2002) 598-604, \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/qsar.200290002>. doi:10.1002/qsar.200290002.
- 1161       URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/qsar.200290002>
- 1162 [51] M. Snarey, N. K. Terrett, P. Willett, D. J. Wilton, *Comparison of algorithms for dissimilarity-based compound selection*, *Journal of Molecular Graphics & Modelling* 15 (6) (1997) 372-385. doi:10.1016/s1093-3263(98)00008-4.
- 1163 [52] N. Baker, T. Knudsen, A. Williams, *Abstract Sifter: a comprehensive front-end system to PubMed*, *F1000Research* 6

- 1184 (2017) *Chem Inf Sci*-2164. doi:10.12688/f1000research.12865.1.
- 1185 URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5801564/>
- 1186 [53] N. Aurisano, O. Jollet, W. A. Chiu, R. Judson, S. Jang, A. Unnikrishnan, M. B. Kosnik, P. Fantke, **Probabilistic points of**  
1187 **departure and reference doses for characterizing human noncancer and developmental/reproductive effects for 10,145**  
1188 **chemicals**, *Environmental Health Perspectives* 131 (3) (2023) 037016. arXiv:<https://ehp.niehs.nih.gov/doi/pdf/10.1289/EHP11524>.  
1189 URL <https://ehp.niehs.nih.gov/doi/abs/10.1289/EHP11524>
- 1190 [54] H. K. Liberatore, S. R. Jackson, M. J. Strynar, J. P. McCord, **Solvent Suitability for HFPO-DA ("GenX" Parent Acid) in**  
1191 **Toxicological Studies**, *Environmental Science & Technology Letters* 7 (7) (2020) 477-481, publisher: American Chemical  
1192 Society. doi:10.1021/acs.estlett.0c00323.  
1193 URL <https://doi.org/10.1021/acs.estlett.0c00323>
- 1194 [55] C. Zhang, A. C. McElroy, H. K. Liberatore, N. L. M. Alexander, D. R. U. Knappe, **Stability of Per- and Polyfluoroalkyl**  
1195 **Substances in Solvents Relevant to Environmental and Toxicological Analysis**, *Environmental Science & Technology* 56 (10)  
1196 (2022) 6103-6112, publisher: American Chemical Society. doi:10.1021/acs.est.1c03979.  
1197 URL <https://doi.org/10.1021/acs.est.1c03979>
- 1198 [56] D. E. Dawson, C. Lau, P. Pradeep, R. R. Sayre, R. S. Judson, R. Tornero-Velez, J. F. Wambaugh, **A machine learning model to**  
1199 **estimate toxicokinetic half-lives of per- and polyfluoro-alkyl substances (pfas) in multiple species**, *Toxics* 11 (2) (2023).  
1200 doi:10.3390/toxics11020098.  
1201 URL <https://www.mdpi.com/2305-6304/11/2/98>
- 1202 [57] J. Wang, D. R. Hallinger, A. S. Murr, A. R. Buckalew, R. R. Lougee, A. M. Richard, S. C. Laws, T. E. Stoker, High-throughput  
1203 screening and chemotype-enrichment analysis of ToxCast phase II chemicals evaluated for human sodium-iodide sym-  
1204 porter (NIS) inhibition, *Environment International* 126 (2019) 377-386. doi:10.1016/j.envint.2019.02.024.
- 1205 [58] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg,  
1206 J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, G. Duchesnay, Scikit-learn: Machine learning in python,  
1207 *The Journal of Machine Learning Research* 12 (null) (2011) 2825-2830, number: null.
- 1208 [59] F. Webster, M. Gagné, G. Patlewicz, P. Pradeep, N. Trefiak, R. Judson, T. Barton-Maclaren, Predicting estrogen receptor  
1209 activation by a group of substituted phenols: An integrated approach to testing and assessment case study., *Regul Toxicol*  
1210 *Pharmacol.* 106 (2019) 278-291. doi:10.1016/j.yrtph.2019.05.017.
- 1211 [60] OECD, Case study on the use of integrated approaches for testing and assessment IATA for estrogenicity of the substi-  
1212 tuted phenols., OECD Publishing (2018).
- 1213

1214 **Appendix A. Supplementary information**

1215 **Appendix A.1. Evaluating the feasibility of subdividing the primary categories**

1216 Corina Symphony on the command line (licensed from Molecular Networks GmbH and Altamira LLC)  
1217 was used to compute the 129 PFAS ToxPrints<sup>44</sup>. The Fisher's exact test was used to compute an  
1218 odds ratio and associated p value for each PFAS ToxPrint relative to the OECD primary category  
1219 designation. This was comparable with the methodology discussed in Wang et al.<sup>57</sup>. A PFAS Tox-  
1220 Print was considered enriched if it had an odds ratio greater than or equal to 3, a one-sided Fish-  
1221 ers exact p-value less than 0.05 (probability value of the odds ratio being greater than 1) and the  
1222 number of True Positives (TP) was determined to be greater than or equal to 3. For the "unclas-  
1223 sified" primary category, the top enriched ToxPrints were PFAS bond and chain features including  
1224 pfas\_bond:S(=O)O\_sulfonicAcid\_acyclic\_(chain)\_SCF, pfas\_chain:FT\_n1\_OP, pfas\_chain:FT\_n2\_OP  
1225 the latter represent fluorotelomer chains with either 1 or 2 CH<sub>2</sub> units and an organophosphorus  
1226 terminus. On the otherhand, the "PFAA precursors" had alcohols and carbonyls as enriched func-  
1227 tional groups (pfas\_bond:COH\_alcohol\_pri-alkyl\_CF, pfas\_bond:CC(=O)C\_ketone\_generic\_CF). The  
1228 intention was to explore whether certain types of features were specifically enriched in these broad  
1229 primary categories to consider subcategorizing them to reduce the starting membership. The full set  
1230 of enrichments for all primary categories are provided as a separate data file.

1231 Supplementary Figures

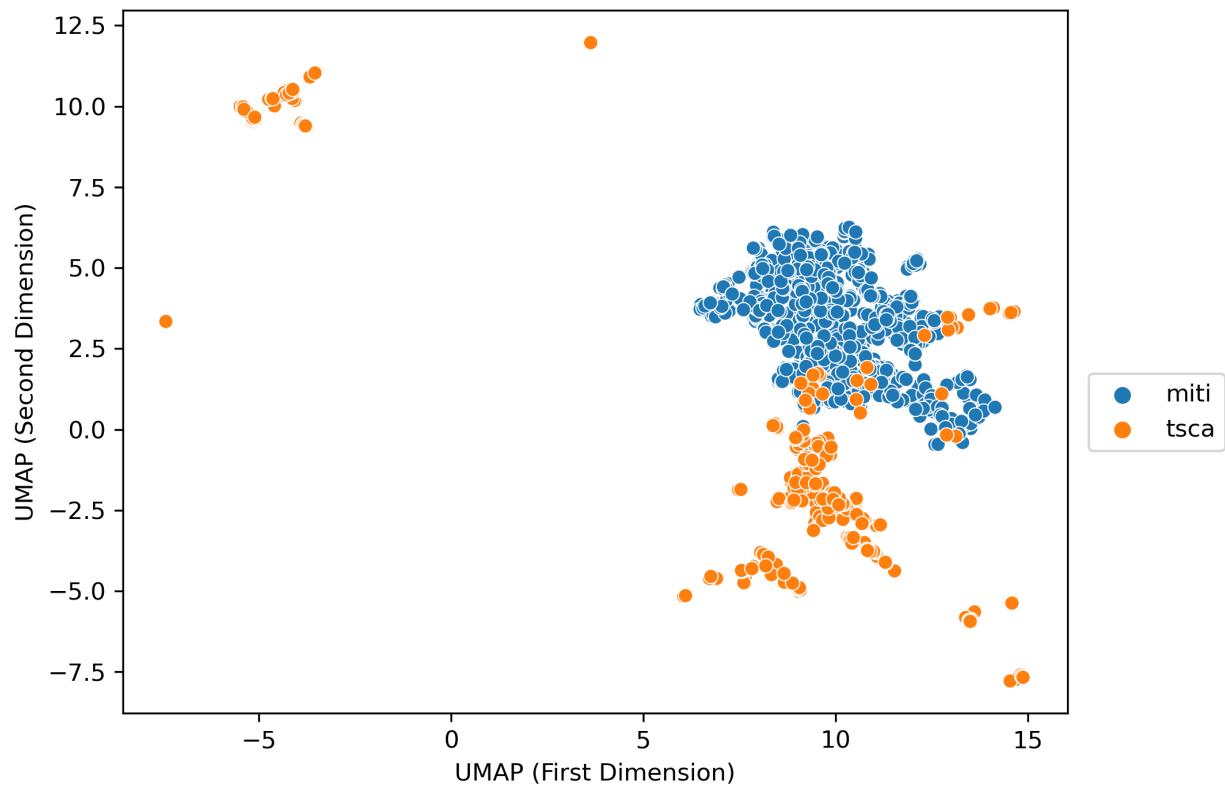


Figure A1: Overlap of MITI training data substances with TSCA substances using Morgan chemical fingers and represented in a UMAP plot

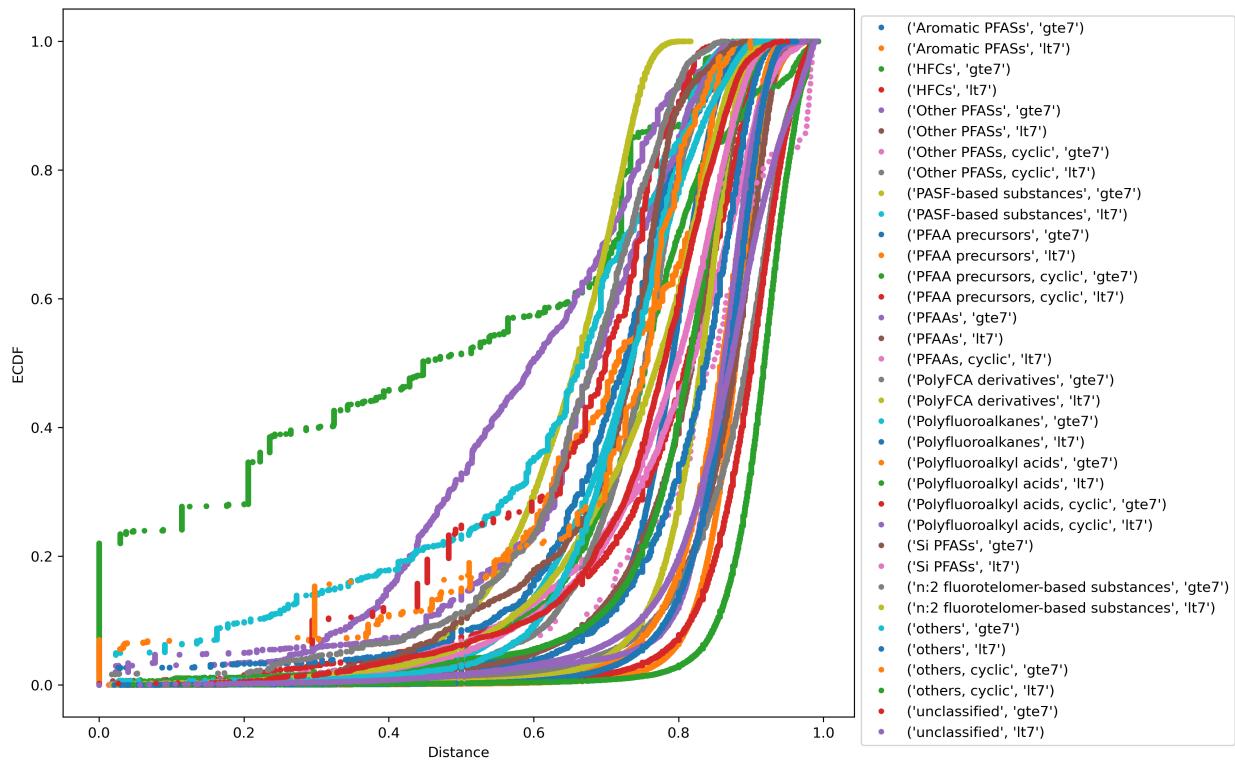


Figure A2: ECDFs of the within categories based on the chain length threshold of 7

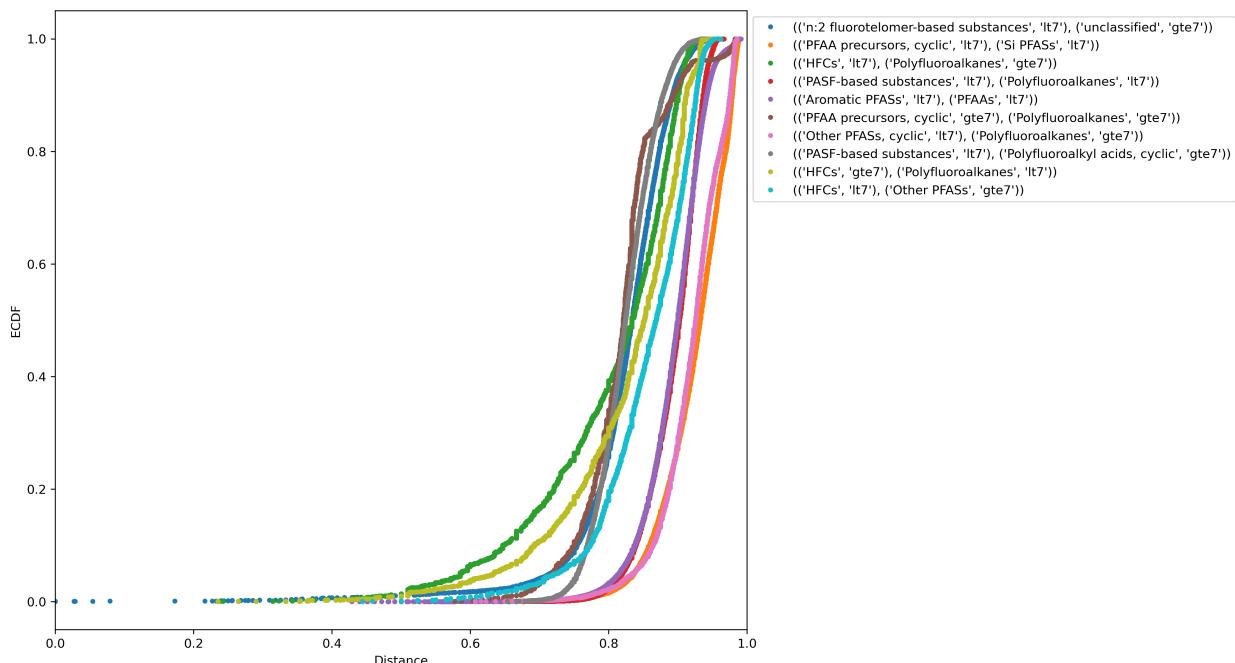


Figure A3: ECDFs for selected between category combinations for carbon chain length categories

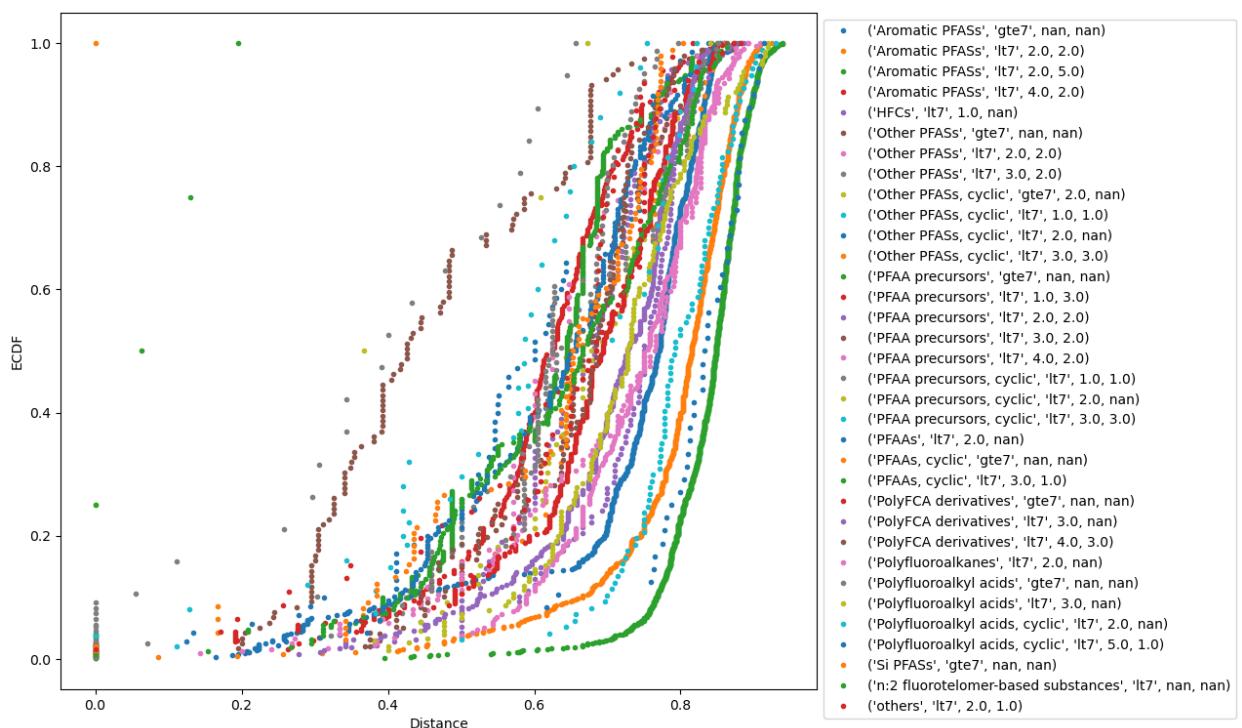


Figure A4: EDCFs for selected terminal categories to demonstrate left shift in pairwise distance

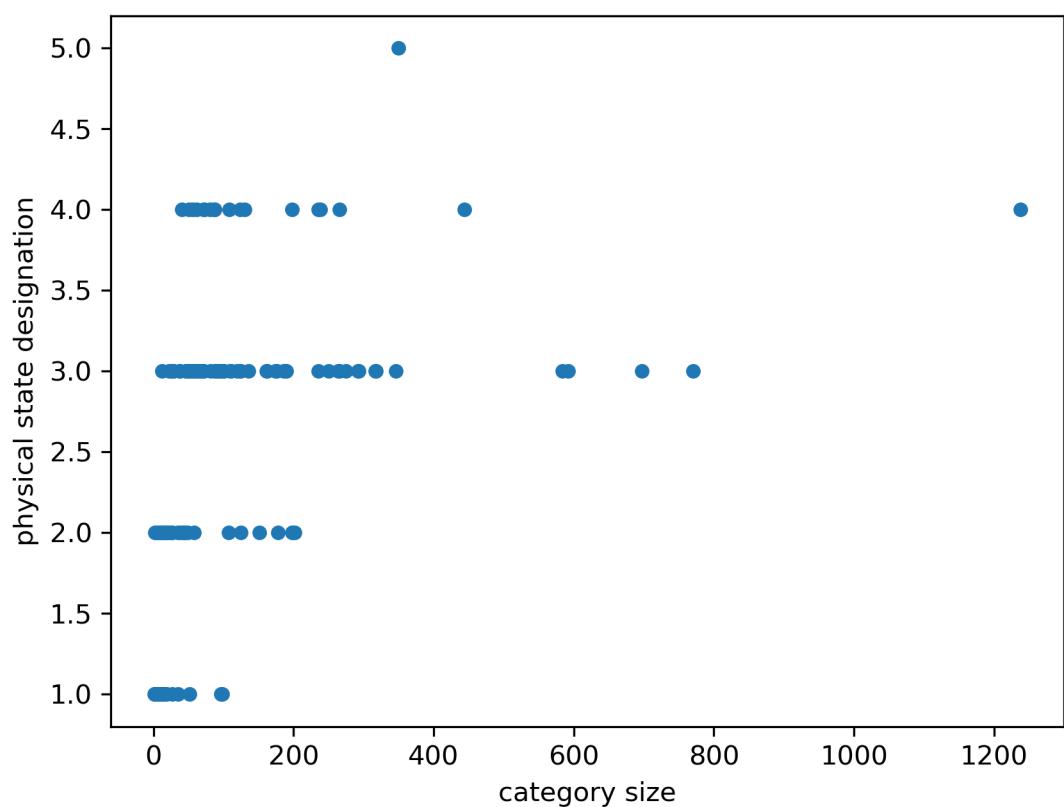
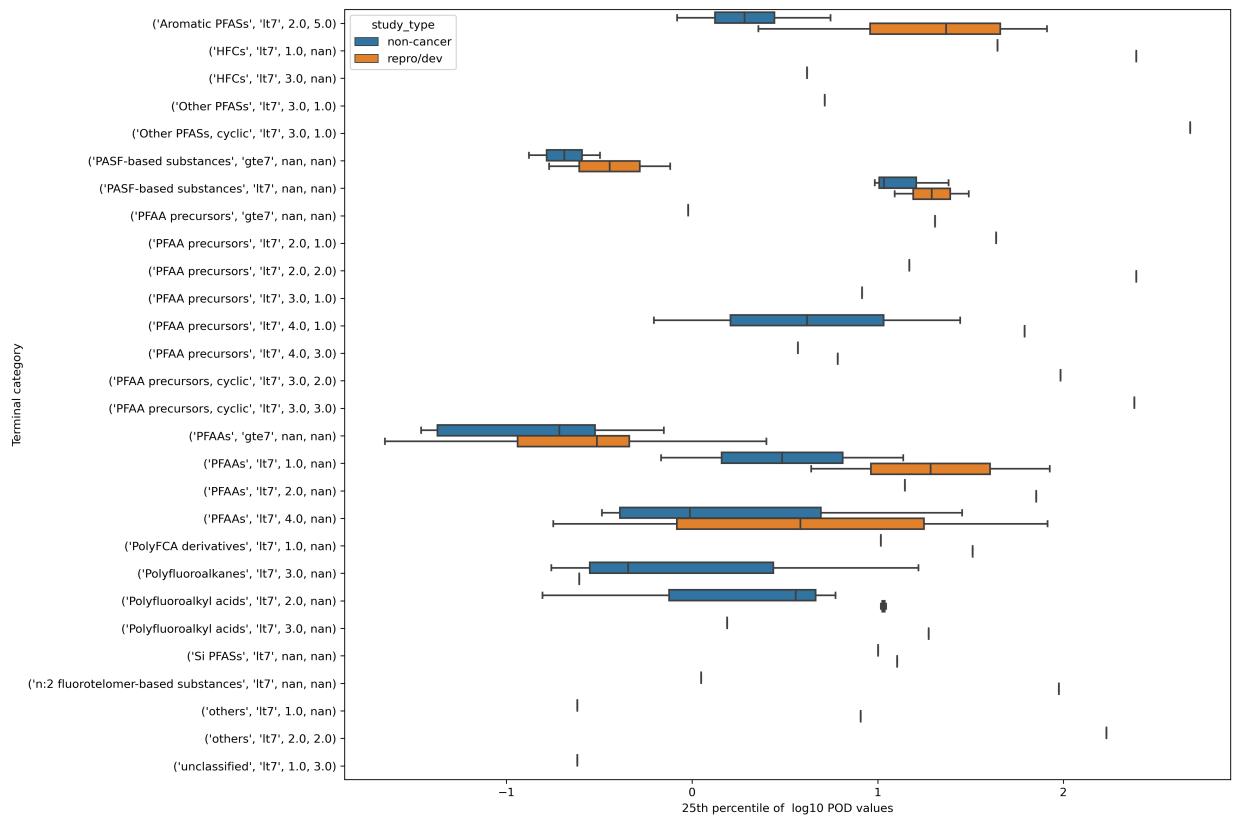


Figure A5: Correlation between terminal categories with large membership size and the number of designations represented amongst their memberships



**Figure A6:** Boxplots of the variation of 25<sup>th</sup> percentiles of point of departure values (PODs) from non-cancer and repro/developmental studies. The box in the boxplot reflects the quartiles of the dataset, whilst the whiskers extend to + 1.5 \* inter-quartile range (IQR). Outliers are shown as points if they exceed 1.5 \* IQR. The repro/developmental boxplot is shown below the non-cancer boxplot for a given terminal category.

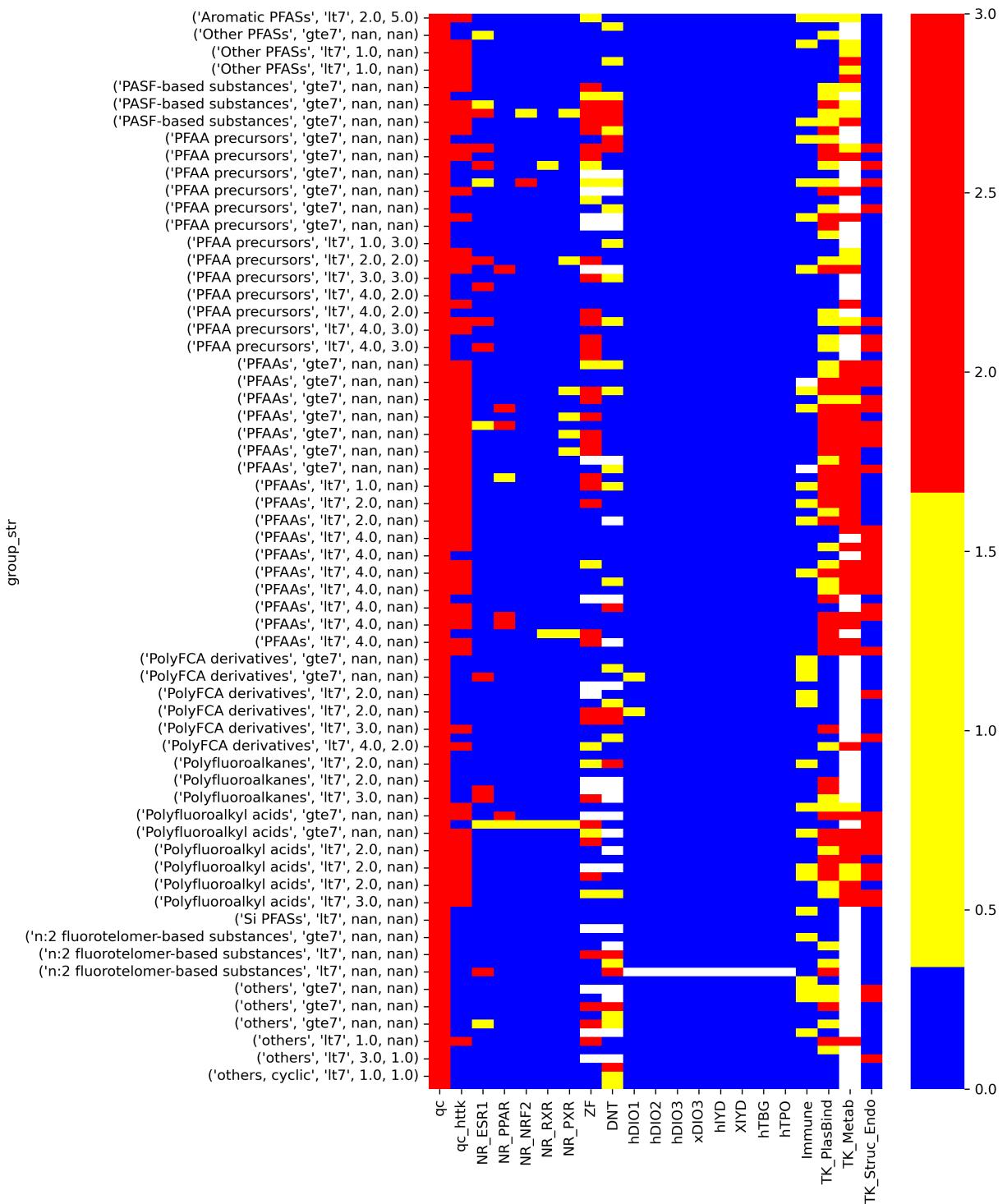


Figure A7: Heatmap of NAM flags for substances tested that overlap with the PFAS inventory

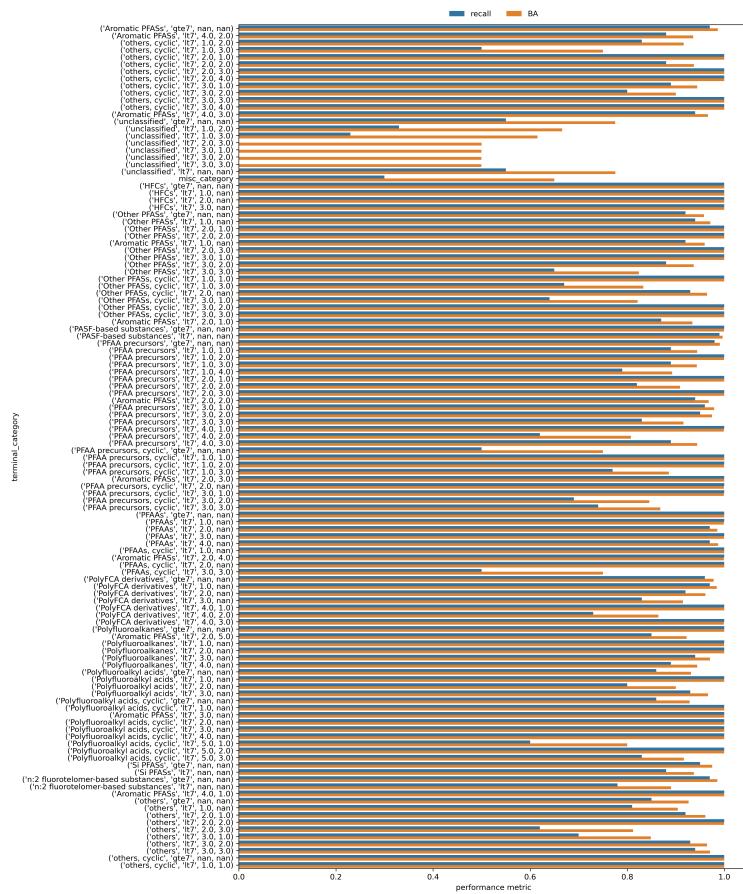


Figure A8: Barplot of the performance scores for the random forest classification model as applied to the hold out set of terminal categories