



## Two-Stage Machine Learning-Based Approach to Predict Points of Departure for Human Noncancer and Developmental/Reproductive Effects

Kvasnicka, Jacob; Aurisano, Nicolò; von Borries, Kerstin; Lu, En-Hsuan; Fantke, Peter; Jolliet, Olivier; Wright, Fred A.; Chiu, Weihsueh A.

*Published in:*  
Environmental Science and Technology

*Link to article, DOI:*  
[10.1021/acs.est.4c00172](https://doi.org/10.1021/acs.est.4c00172)

*Publication date:*  
2024

*Document Version*  
Peer reviewed version

[Link back to DTU Orbit](#)

*Citation (APA):*  
Kvasnicka, J., Aurisano, N., von Borries, K., Lu, E-H., Fantke, P., Jolliet, O., Wright, F. A., & Chiu, W. A. (in press). Two-Stage Machine Learning-Based Approach to Predict Points of Departure for Human Noncancer and Developmental/Reproductive Effects. *Environmental Science and Technology*.  
<https://doi.org/10.1021/acs.est.4c00172>

---

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

1  
2  
3  
4  
5  
6 1 Two-Stage Machine Learning-Based Approach  
7  
8 2 to Predict Points of Departure for Human Non-  
9  
10 3 cancer and Developmental/Reproductive Effects  
11  
12  
13 4

15 5 Jacob Kvasnicka<sup>1</sup>, Nicolò Aurisano<sup>2</sup>, Kerstin von Borries<sup>2</sup>, En-Hsuan Lu<sup>1</sup>, Peter  
16 6 Fantke<sup>2</sup>, Olivier Jolliet<sup>2</sup>, Fred A. Wright<sup>3</sup>, Weihsueh A. Chiu\*<sup>1</sup>  
17  
18

19 7 <sup>1</sup>Department of Veterinary Physiology and Pharmacology, Interdisciplinary  
20 8 Faculty of Toxicology, Texas A&M University, College Station, Texas, United  
21 9 States  
22

23  
24 10 <sup>2</sup>Quantitative Sustainability Assessment, Department of Environmental and  
25 11 Resource Engineering, Technical University of Denmark, Lyngby, Denmark  
26

27 12 <sup>3</sup>Departments of Statistics and Biological Sciences and Bioinformatics Research  
28 13 Center, North Carolina State University, Raleigh, North Carolina, United States  
29

30 14 Corresponding author: Weihsueh A. Chiu, [wchiu@tamu.edu](mailto:wchiu@tamu.edu)  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**ABSTRACT**

Chemical points of departure (PODs) for critical health effects are crucial for evaluating and managing human health risks and impacts from exposure. However, PODs are unavailable for most chemicals in commerce due to a lack of *in vivo* toxicity data. We therefore developed a two-stage machine learning (ML) framework to predict human-equivalent PODs for oral exposure to organic chemicals based on chemical structure. Utilizing ML-based predictions for structural/physical/chemical/toxicological properties from OPERA 2.9 as features (Stage 1), ML models using random forest regression were trained with human-equivalent PODs derived from *in vivo* datasets for general noncancer effects (n = 1,791) and reproductive/developmental effects (n = 2,228), with robust cross-validation for feature selection and estimating generalization errors (Stage 2). These two-stage models accurately predicted PODs for both effect categories, with cross-validation-based root-mean-squared errors less than an order of magnitude. We then applied one or both models to 34,046 chemicals expected to be in the environment, revealing several thousand chemicals of *moderate* concern and several hundred chemicals of *high* concern for health effects at estimated median population exposure levels. Further application can expand by orders of magnitude the coverage of organic chemicals that can be evaluated for their human health risks and impacts.

**Keywords:** QSAR model, machine learning, toxicity prediction, chemical risk assessment, high-throughput screening, life cycle impact assessment (LCIA)

**Synopsis:** Most chemicals lack toxicity data related to human health. This study uses machine learning to fill this gap, greatly expanding the ability to characterize chemical risks and impacts.

## 41 INTRODUCTION

42 Determining a chemical's point of departure (POD) is crucial to evaluating and managing  
43 health risks and toxicity impacts associated with chemical exposure. The POD is the starting  
44 point along the dose-response curve for extrapolating health risks to relevant exposure levels that  
45 may be encountered in the general population.<sup>1</sup> A variety of impact and risk assessment  
46 frameworks, such as contaminated site remediation, life cycle impact assessment (LCIA),  
47 chemical alternatives assessment (CAA), and health-based risk screening, heavily rely on  
48 PODs.<sup>2,3</sup> These PODs are primarily developed in regulatory or other authoritative assessments by  
49 agencies, such as the United States Environmental Protection Agency (U.S. EPA), that  
50 synthesize available toxicity data from *in vivo* studies and identify the "critical" or "most-  
51 sensitive" endpoint for characterizing health effects. However, due to the resource-intensive  
52 nature of these assessments, such authoritative PODs are available for less than 1,000 chemicals,  
53 which is a tiny fraction of the more than 150,000 commercial chemicals to which humans may  
54 be exposed.<sup>4,5</sup> Consequently, most of these chemicals lack comprehensive human health  
55 assessments and are not included in impact and risk assessment tools, such as USEtox.<sup>6</sup>

56 To partially address the lack of availability of authoritative assessments, a number of  
57 open-source databases compiling publicly available experimental *in vivo* toxicity data required  
58 for POD derivation have emerged, such as the U.S. EPA's Toxicity Value Database  
59 (ToxValDB)<sup>7</sup> and the European Chemicals Agency's International Uniform Chemical  
60 Information Database (IUCLID; <https://iuclid6.echa.europa.eu/>). These databases have enabled  
61 researchers to derive "surrogate" PODs, through rigorous curation and statistical approaches, as a  
62 proxy for PODs that would be selected in an authoritative assessment.<sup>8</sup> However, even though  
63 use of these databases increases the availability of PODs by an order of magnitude to about ten

1  
2  
3 64 thousand chemicals, the remaining gap underscores the need for a high-throughput approach to  
4  
5 65 develop surrogate PODs in the absence of *in vivo* data.  
6

7  
8 66 “New approach methods” (NAMs), including *in vitro* and computational (*in silico*)  
9  
10 67 approaches, have emerged as promising, high-throughput alternatives to animal testing, while  
11  
12 68 also addressing ethical concerns regarding animal use. A prime example of *in silico* NAMs is  
13  
14 69 QSAR modeling (Quantitative Structure-Activity Relationship). QSAR models commonly use  
15  
16 70 machine learning (ML) to predict biological activity based on chemical structure information.  
17  
18  
19 71 Applications of QSAR modeling have substantially expanded the availability of toxicologically  
20  
21 72 relevant data. For example, Mansouri et al. developed a collection of open-source ML models  
22  
23 73 known as “OPERA” [Open (Quantitative) Structure-activity/property Relationship App].<sup>9,10</sup>  
24  
25 74 These models predict structural and physical-chemical properties, environmental fate metrics,  
26  
27 75 acute toxicity, and toxicokinetic endpoints for hundreds of thousands of chemicals. Many of  
28  
29 76 these predictions are available through open-source web platforms such as the CompTox  
30  
31 77 Chemistry Dashboard by U.S. EPA,<sup>11</sup> and the National Toxicology Program (NTP) Integrated  
32  
33 78 Chemical Environment (ICE).<sup>12</sup>  
34  
35  
36

37  
38 79 Previous studies have also developed QSAR models to predict PODs. For instance, the  
39  
40 80 models developed by Wignall et al. (2018) included those that predict PODs, such as benchmark  
41  
42 81 doses (BMDs) and No Observed Adverse Effect Levels (NOAELs), using training data from  
43  
44 82 several hundred chemicals with available authoritative human health assessments (n=137 for  
45  
46 83 BMDs and n=487 for NOAELs).<sup>4</sup> For these PODs, the Wignall et al. (2018) models explained  
47  
48 84 between 28% and 45% of the variance, with mean absolute errors of 0.93-1.13 log<sub>10</sub>-units.  
49  
50  
51 85 Pradeep et al. (2020) used a similar approach to predict effect levels for specific species-study  
52  
53 86 type combinations in ToxValDB, with training sets ranging in size from <100 to over 3600 and a  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 87 wide range of performance depending on the study type.<sup>13</sup> Combining all study types, they  
4  
5 88 achieved an  $R^2$  of 0.53 and RMSE of 0.71 in  $\log_{10}$ -units, but their approach does not provide  
6  
7  
8 89 surrogate PODs that reflect the “critical” or “most-sensitive” endpoints for characterizing health  
9  
10 90 effects. Thus, there remains a substantial gap in the availability of surrogate PODs for a wider  
11  
12 91 range of chemicals.

13  
14 92 Conventional ML-based QSAR models often rely on hundreds of molecular descriptors  
15  
16 93 as features.<sup>4,13</sup> While these descriptors can enable accurate predictions, and many have good  
17  
18 94 structural interpretability, it can be challenging to explain their toxicological importance to  
19  
20 95 practitioners and decision-makers. Recognizing this challenge, the Organisation for Economic  
21  
22 96 Co-operation and Development’s (OECD) *(Q)SAR Assessment Framework*<sup>14</sup> includes a key  
23  
24 97 “mechanistic interpretation” criterion for evaluating a QSAR model, defined as “how the  
25  
26 98 rationale behind a (Q)SAR model is consistent with or accounts for the knowledge related to the  
27  
28 99 predicted property.” This guidance highlights the importance of QSAR models that not only  
29  
30  
31 100 predict accurately but also provide insights into their underlying scientific basis to enhance their  
32  
33 101 utility and trustworthiness. Thus, in accordance with the OECD report suggesting preference for  
34  
35 102 a “physical-chemical interpretation (if possible) that is consistent with a known mechanism of  
36  
37 103 biological action,” we posit that the structural/physical/chemical/toxicological properties that are  
38  
39 104 available in OPERA, such as water solubility and bioconcentration factor, are more easily  
40  
41 105 understood by a typical practitioner than typical chemoinformatic descriptors, and offer a path  
42  
43 106 towards more “understandable” machine learning.

44  
45 107 Building on prior efforts, this study aimed to expand the coverage of chemicals with  
46  
47 108 toxicity values that can be used as a surrogate for human-equivalent noncancer PODs for oral  
48  
49 109 exposure in the absence of *in vivo* data. Our objectives were threefold:  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 110 1. Develop and evaluate a two-stage QSAR modeling framework that incorporates  
111 an intermediate layer of structural/physical/chemical/toxicological properties as  
112 features.
- 113 2. Generate an extended set of oral surrogate PODs, with quantified model  
114 prediction errors based on cross-validation, for a wide range of chemicals.
- 115 3. Apply this framework to a large dataset of chemicals observed in the  
116 environment, assessing potential health risks using the margin of exposure as a  
117 metric.

118 Following Aurisano et al. (2023),<sup>8</sup> we differentiated between reproductive/developmental and  
119 nonreproductive/developmental effects (“general noncancer effects”).<sup>3,15</sup> The surrogate PODs  
120 from this study can be integrated into various chemical management and exposure and impact  
121 assessment frameworks for health-based risk screening, LCIA, CAA for chemical substitution,  
122 and exposure and risk prioritization.<sup>3,16,17</sup>

## 123 METHODS

124 To address the stated objectives, we developed a two-stage ML framework. The first  
125 stage derives ML-based predictions for structural/physical/chemical/toxicological properties that  
126 are readily interpretable. The second stage leverages these properties as features in a separate  
127 ML model to predict surrogate PODs. **Figure 1A** illustrates the conceptual framework, while  
128 **Figure 1B** shows an overview of the model development, evaluation, and application. The  
129 conceptual framework comprises the following steps:

- 130 1. Select and identify chemicals for modeling.
- 131 2. Standardize chemical structures to make them “QSAR-ready.”
- 132 3. Run prior QSAR models for feature extraction (Stage 1).

1  
2  
3 133 4. Clean and parse the QSAR predictions to obtain raw features.  
4

5 134 5. Apply these features in a modeling pipeline to predict PODs (Stage 2).  
6

7  
8 135 All ML algorithms for predicting PODs were implemented using Python 3.9, leveraging open-

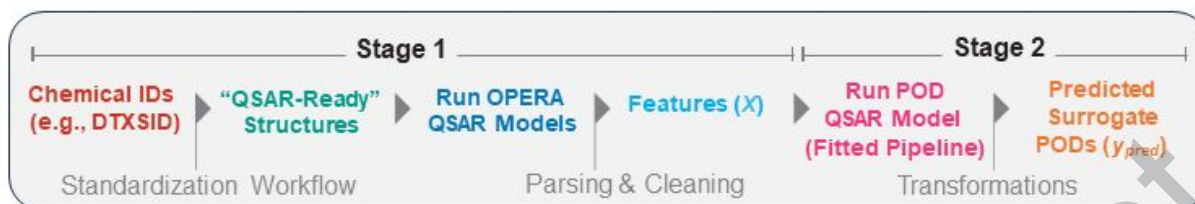
9  
10 136 source libraries such as scikit-learn 1.2.2.<sup>18</sup> The source code, results, and input files associated

11  
12 137 with this study are openly available in a GitHub repository at <https://github.com/jkvasnicka/Two->

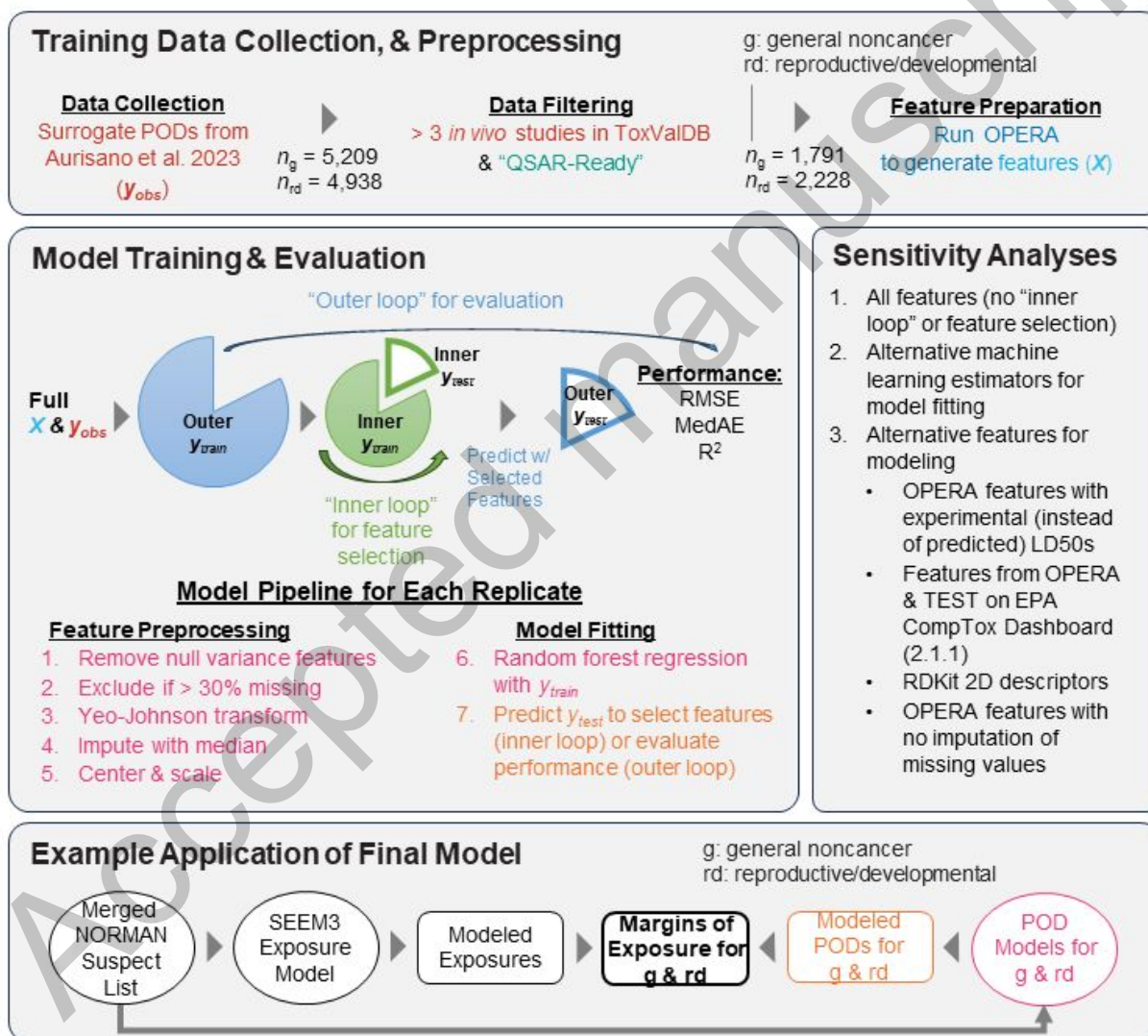
13  
14  
15 138 [Stage-ML-Oral-PODs](#).  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



## A Conceptual Framework: Two-Stage QSAR Model



## B Model Development, Evaluation, & Application



139  
140 **Figure 1.** Overview of the two-stage machine learning framework for predicting points of departure. (A) Conceptual  
141 framework; (B) Model development, evaluation, and application. The surrogate points of departure were obtained  
142 from Table S5 of Aurisano et al. (2023).<sup>8</sup> Features were extracted from predictions by OPERA 2.9.<sup>9,10</sup> **Figures S1-**  
143 **S2** provide an overview of the model training and evaluation. Exposure estimates were obtained from SEEM3 by

1  
2  
3 144 Ring et al. (2019).<sup>19</sup> Application chemicals were expected to occur in the environment and lacked *in vivo* points of  
4 145 departure.<sup>20,21</sup> Note: ML, machine learning; POD, point of departure; QSAR, quantitative structure-activity  
5 146 relationship; OPERA, OPEn structure-activity/property Relationship App; ToxValDB, Toxicity Value Database;  
6 147 RMSE, root-mean-squared error, MedAE, median absolute error; R<sup>2</sup>, coefficient of determination; MAD, median  
7 148 absolute deviation; SEEM, Systematic Empirical Evaluation of Models.  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Accepted manuscript

### 149 *Training Data Collection and Preprocessing*

150 **Data Collection:** Predicting PODs was essentially a regression task with a continuous  
151 target vector  $\vec{y}_e$  of oral doses, in log<sub>10</sub>-transformed units of mg·(kg·d)<sup>-1</sup>, representing a POD for  
152 a given effect category  $e$ , and inputs represented by a matrix  $\mathbf{X}$ , where each row corresponds to a  
153 sample and each column corresponds to one of  $n$  distinct features, i.e.,  $\mathbf{X} = [\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n]$ . This  
154 task required labeled data mapping chemical identifiers to their respective *in vivo* PODs.  
155 Specifically, we used the surrogate oral PODs from Table S5 of Aurisano et al. (2023),<sup>8</sup> which  
156 were derived through meticulous curation and statistical analysis of *in vivo* experimental animal  
157 data from ToxValDB 9.1,<sup>7</sup> adjusted to chronic human equivalent benchmark doses (BMDh).  
158 Throughout this study, the U.S. EPA's DSSTox Substance Identifier (DTXSID) uniquely  
159 identify each chemical.

160 **Data Filtering:** Initially, there were 5,209 unique chemicals with surrogate PODs for  
161 general noncancer effects, and 4,938 chemicals for reproductive/developmental effects.  
162 However, a series of filtering steps removed chemicals that were unsuitable for modeling  
163 (**Figure 1B**). First, chemicals with  $\leq 3$  *in vivo* studies were excluded because those surrogate  
164 PODs may be less robust (Aurisano et al. used the lower 25<sup>th</sup> percentile of the distribution of  
165 available PODs for a chemical as the surrogate POD), leaving 2,404 and 2,999 chemicals for the  
166 respective effect categories. Next, a general applicability domain exclusion and standardization  
167 workflow was applied to generate "QSAR-ready" structures compatible with a variety of  
168 modeling approaches.<sup>22,23</sup> Applying this workflow yielded 1,791 organic chemicals for general  
169 noncancer effects and 2,228 organic chemicals for reproductive/developmental effects.

170 **Feature Extraction & Preparation:** To obtain features, we leveraged the QSAR  
171 modeling framework, OPERA 2.9, by Mansouri et al.<sup>9,10</sup> Specifically, we used the command-line

1  
2  
3 172 version, OPERA2.9\_CL, inputting the chemical identifiers (DTXSID) as a text file. OPERA then  
4  
5 173 retrieved the corresponding QSAR-ready structures as simplified molecular-input line-entry  
6  
7 174 system (SMILES) strings from its internal database. This execution yielded 39 interpretable  
8  
9 175 features (e.g., water solubility) with feature-specific applicability domain information. We then  
10  
11 176 flagged features outside the applicability domain as “missing” if both of the following criteria by  
12  
13 177 Mansouri et al. were met:<sup>9</sup>

- 17 178 1. The value was outside the *global* applicability domain of the model/feature.
- 19 179 2. The value had a low *local* applicability domain index (< 0.4) with respect to its  
21 180 nearest neighboring values.

23  
24 181 **Figure S3** displays the distributions of raw features for all chemicals in this study, with  
25  
26 182 corresponding descriptions in a supplemental Excel file (**Table S3**). Given the diverse nature of  
27  
28 183 these features, we designed a robust feature preprocessing pipeline for feature transformation  
29  
30 184 (**Figure 1B**), generalizable across a variety of ML estimators, as detailed below.

### 34 185 *Model Training and Evaluation*

35  
36 186 **Overview of Modeling Pipeline:** The QSAR models for predicting PODs consisted of a  
37  
38 187 pipeline of feature preprocessing steps and a ML estimator (e.g., random forest) (**Figure 1B**).  
39  
40 188 This design ensured that transformation parameters (e.g., median for imputation) were derived  
41  
42 189 solely from the training data, minimizing potential for data leakage and overoptimistic  
43  
44 190 performance estimates. The feature preprocessing steps are described in the Supporting  
45  
46 191 Information (see section, *Feature Preprocessing Steps*), and include imputation of missing  
47  
48 192 values using the median (features were excluded if >30% imputation would be necessary). For  
49  
50 193 the last components in the pipeline (steps 6 and 7 in **Figure 1B**), we chose the Random Forest  
51  
52 194 Regressor and made predictions for the surrogate PODs. This estimator was a reasonable choice,  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 195 given its track record of robust performance without extensive preprocessing or hyperparameter  
4  
5 196 tuning,<sup>24</sup> and its successful applications in prior studies involving POD prediction.<sup>4,13</sup> The  
6  
7 197 algorithm constructs a collection of de-correlated decision trees using bootstrapped sampled  
8  
9 198 versions of the training data, and then averages predictions to minimize variance.<sup>25</sup> For the  
10  
11 199 hyperparameters, we used the scikit-learn 1.2.2 defaults,<sup>18</sup> except for the number of features to  
12  
13 200 consider when searching for the best split, which we set to 1/3 (or at least 1) of the available  
14  
15 201 features,<sup>24</sup> instead of considering all features.

16  
17  
18  
19 202 For model training and evaluation, we implemented nested 5-fold cross-validation, with  
20  
21 203 separate “inner” and “outer” loops (**Figures 1B, S1, and S2**). The “inner” loop is used for feature  
22  
23 204 selection, whereas the “outer” loop is used to evaluate performance. Thus, for an iteration of the  
24  
25 205 “outer” loop, the data are divided into an “outer” training and testing dataset. The “outer”  
26  
27 206 training set is sent to the “inner” loop where it is repeatedly divided into “inner” training and  
28  
29 207 testing datasets. This “inner” loop trains an “inner” model in order to conduct feature selection  
30  
31 208 (described below under **Model Training with Feature Selection**). The selected features are then  
32  
33 209 passed back to the “outer” loop, which trains a model using only those selected features with the  
34  
35 210 “outer” training dataset, and evaluates performance using the “outer” testing data. This whole  
36  
37 211 process is then repeated multiple times with different randomizations (described below under  
38  
39 212 **Model Evaluation**).

40  
41  
42  
43 213 **Model Training with Feature Selection:** Given the 39 features from OPERA 2.9  
44  
45 214 (**Figure S3**),<sup>9,10</sup> we hypothesized that a subset of 10 features would be sufficient for successful  
46  
47 215 modeling while remaining interpretable. We selected the value of “10” *a priori* to avoid over-  
48  
49 216 fitting, and verified this hypothesis in a sensitivity analysis (described below) where all features  
50  
51 217 were used without feature selection. If the value of “10” were to materially degrade performance,  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 218 then we could have used more complex feature selection approaches, such as recursive feature  
4  
5 219 elimination.

6  
7  
8 220 To select features in an objective, robust, reproducible manner, we implemented a feature  
9  
10 221 selection scheme by nesting a permutation feature importance algorithm within a repeated k-fold  
11  
12 222 cross-validation loop. Specifically, we repeatedly divided the data into 5-folds, training the  
13  
14 223 model on 4/5 of the data in which the algorithm measured feature importance by assessing the  
15  
16 224 decrease in model performance upon random permutation of feature values. In particular, we  
17  
18 225 used the median value for this importance score across random permutations as the selection  
19  
20 226 criterion. The cross-validation loop minimized biases and over-optimistic performance scores.  
21  
22  
23 227 Further details can be found in the Supporting Information (see section, *Model Training Steps*,  
24  
25  
26 228 and **Figure S1**).

27  
28 229 **Model Evaluation:** To gauge the model's generalization to unseen data, we nested the  
29  
30 230 training process described above within another repeated *K*-fold cross validation loop. For this  
31  
32 231 loop, we used 30 repetitions and 5 folds, yielding 150 (30x5) replicate models that underwent the  
33  
34 232 same model training steps. To quantify performance, we used the root-mean-squared error  
35  
36 233 (RMSE), median absolute error (MedAE), and coefficient of determination ( $R^2$ ). Further details  
37  
38 234 regarding the model evaluation, along with definitions of the performance metrics, can be found  
39  
40 235 in the Supporting Information (see section, *Model Performance Metrics*, and **Figure S2**).

41  
42 236 **Model Benchmarking:** To further evaluate our models, we benchmarked the QSAR-  
43  
44 237 derived PODs ( $POD_{QSAR}$ ) against estimates from other studies. Specifically, we referenced the  
45  
46 238 original authoritative PODs ( $POD_{authoritative}$ ) and the target variable of surrogate PODs  
47  
48 239 ( $POD_{surrogate}$ ) from Aurisano et al. (2023),<sup>8</sup> both of which were fully adjusted to BMDh.  
49  
50  
51 240 Additionally, we compared our  $POD_{QSAR}$  values with oral equivalent doses derived from  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 241 combining high-throughput *in vitro* bioactivity data with toxicokinetic data using reverse  
4  
5 242 dosimetry. Specifically, we used the “ $POD_{NAM,50}$ ” values from Table S2 of Paul Friedman et al.  
6  
7 243 (2020),<sup>26</sup> where “50” denotes the median from a population distribution of steady-state  
8  
9 244 administered equivalent doses.  $POD_{NAM,50}$  values were available for 263 chemicals for general  
10  
11 245 noncancer effects and 13 chemicals for reproductive/developmental effects.  
12  
13  
14

### 15 246 *Sensitivity Analysis*

16  
17 247 We conducted a sensitivity analysis to assess generalization error sensitivity to different  
18  
19 248 datasets, feature preprocessing, and ML estimators. Our baseline Final Model was described  
20  
21 249 above, involving feature selection among all 39 OPERA 2.9 features, imputation of missing  
22  
23 250 values, and the Random Forest Regressor. We compared several additional models for each  
24  
25 251 effect category using the same evaluation scheme described above (**Figure S2**), varying one  
26  
27 252 modeling aspect at a time. These alternative models are shown in **Figure 1** (see section,  
28  
29 253 *Sensitivity Analyses*), and corresponding descriptions are in **Table S1**. All models were applied  
30  
31 254 to the same chemicals, except the model involving no imputation, which was restricted to those  
32  
33 255 chemicals with no missing feature values (n = 184–227).  
34  
35  
36  
37  
38

### 39 256 *Model Application*

40  
41 257 We demonstrated application of our final two-stage models using a large dataset of  
42  
43 258 organic chemicals expected to occur in the environment and for which human oral exposure  
44  
45 259 could be estimated. Specifically, we assessed 34,809 chemicals that were on the Merged  
46  
47 260 NORMAN Suspect List (SusDat)<sup>20,21</sup> and within the applicability domain of SEEM3 (Systematic  
48  
49 261 Empirical Evaluation of Models) by U.S. EPA.<sup>19</sup> We excluded any chemicals outside the  
50  
51 262 “general applicability domain” due to their being unsuitable for QSAR modeling based on the  
52  
53 263 standardization workflow mentioned above,<sup>22,23</sup> and that had a  $POD_{surrogate}$  value used for model  
54  
55  
56  
57  
58  
59  
60

264 training (“training chemicals”). This exclusion resulted in 33,407 chemicals predicted for general  
 265 non-cancer effects, and 32,970 for reproductive/developmental effects (34,046 chemicals across  
 266 the two sets of predictions). We also evaluated how these chemicals fit within the “feature-  
 267 specific applicability domains” of the OPERA models, and the extent to which the distribution of  
 268 features compared to that of the training set chemicals.

269 The margin of exposure was used as a health risk metric to compare SEEM3 predicted  
 270 population median oral exposures [ $\hat{y}_{\text{exposure},i}$  in  $\text{mg}\cdot(\text{kg}\cdot\text{d})^{-1}$ ] with the QSAR-predicted POD [ $\text{POD}_{\text{QSAR},i}$ , also in  $\text{mg}\cdot(\text{kg}\cdot\text{d})^{-1}$ ]. For each sample  $i$ , the margin of exposure ( $MOE_i$ ) was  
 271 calculated as:  
 272

$$MOE_i = \frac{\text{POD}_{\text{QSAR},i}}{\hat{y}_{\text{exposure},i}} \quad (1)$$

275 We screened chemicals for potential health concern using the following categorization  
 276 scheme:<sup>27,28</sup>

- 277 1. **Low Concern for the median population exposure:**  $MOE_i > 100$
- 278 2. **Moderate Concern for the median population exposure:**  $1 < MOE_i \leq 100$
- 279 3. **High Concern for the median population exposure:**  $0 < MOE_i \leq 1$

280 SEEM3 exposure predictions ( $\hat{y}_{\text{exposure},i}$ ) for an individual at the population median exposure,  
 281 accompanied by a model-based Bayesian 90% credible interval representing uncertainty,<sup>19</sup> were  
 282 downloaded from ICE.<sup>12</sup> We also assessed the contribution of  $\text{POD}_{\text{QSAR}}$  (hazard) uncertainty to  
 283 the overall uncertainty in the margin of exposure, in addition to exposure uncertainty from  
 284 SEEM3. Specifically, we derived 90% prediction intervals of  $\text{POD}_{\text{QSAR}}$  uncertainty for each  
 285 percentile of exposure uncertainty for the median individual. The derivation of these prediction



1  
2  
3 286 intervals is shown in the Supporting Information (see section, *Margin of Exposure Uncertainty*  
4  
5 287 *Analysis*).

## 9 288 **RESULTS**

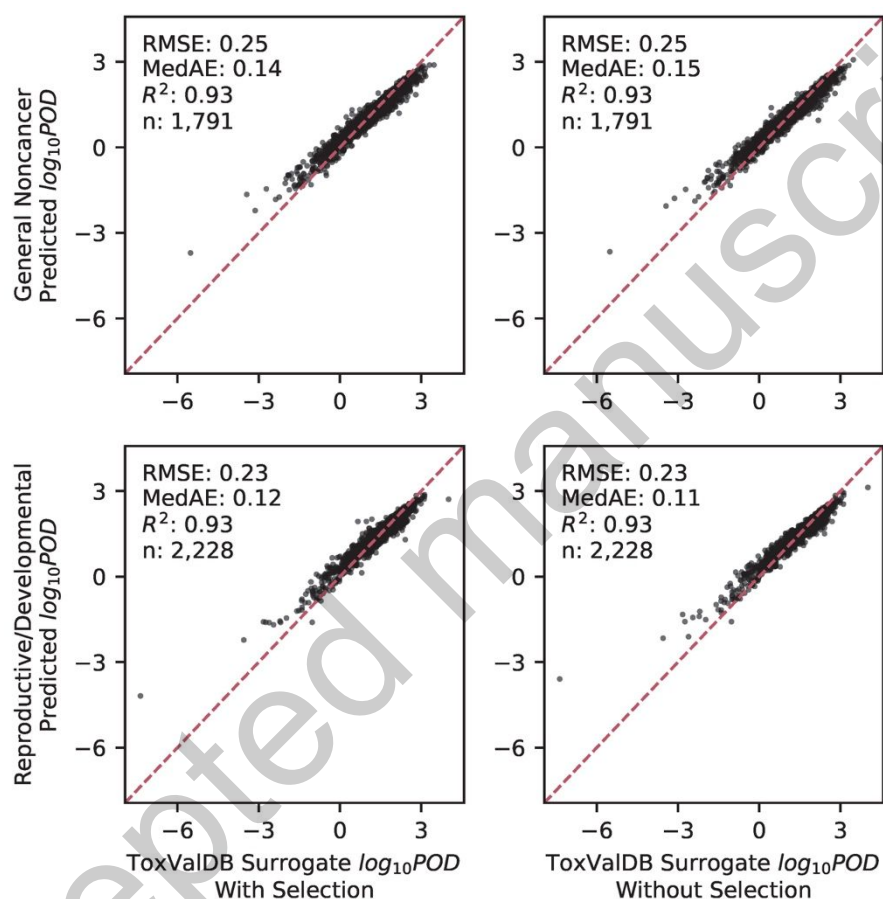
### 11 289 *Dataset Characterization*

13 290 The proportions of missing values across all 39 features from OPERA 2.9 for the training  
14  
15 291 chemicals, and for the application chemicals, can be found in the Supporting Information  
16  
17 292 (**Figure S4**). Most features predominantly had samples within their respective applicability  
18  
19 293 domains. However, three features had more than 30% missing values and were subsequently  
20  
21 294 removed in the pipeline.

### 25 295 *Performance Evaluation and Benchmarking*

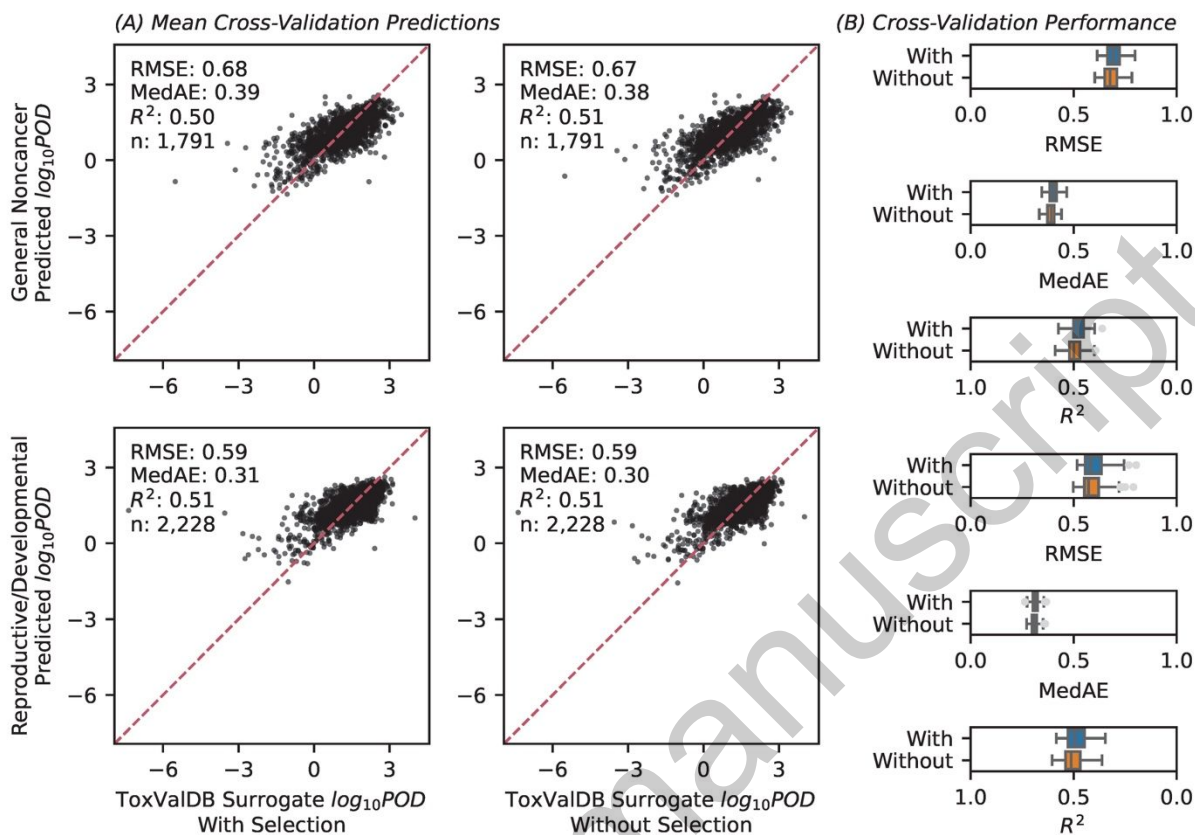
27 296 The final models accurately fitted/predicted  $\text{POD}_{\text{surrogate}}$  values for both effect categories,  
28  
29 297 shown by their RMSE, MedAE, and  $R^2$ . The models demonstrated consistent performance for  
30  
31 298 both effect categories regardless of feature selection. Because of our nested cross-validation  
32  
33 299 approach, each chemical may be part of the “training” or the “testing” dataset depending on the  
34  
35 300 replicate. **Figure 2** summarizes the “in-sample” model fitting, showing the predictions of the  
36  
37 301 cross-validated final models that were fitted on the full labeled dataset. The accuracy was  
38  
39 302 demonstrated by the clustering of fitted predictions and observations along diagonal line, the low  
40  
41 303 values for the disperse measures (RMSE, MedAD), and the high  $R^2$  values. More importantly,  
42  
43 304 **Figure 3** summarizes the “out-of-sample” results, where the median prediction shown is across  
44  
45 305 replicates when the chemical is part of the “testing” dataset. The estimated generalization errors  
46  
47 306 (with 5th - 95th percentiles) based on cross validation were also quite good. These results imply  
48  
49 307 that for a “new” chemical, we can expect the model to predict the POD with a GSD error of less  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 308 than 3.5- to 5.7-fold (taking the range of RMSE values from 0.54 to 0.76), or equivalently a 95%  
4  
5 309 confidence interval spanning 11- to 30-fold in each direction.  
6  
7  
8 310  
9  
10 311



312  
313 **Figure 2. Model fitting.** In-sample performance is assessed through scatterplots and performance metrics  
314 comparing the fitted and observed values for each chemical. The fitted values are predictions from the cross-  
315 validated final models that were fitted on the full labeled dataset. The figure is subdivided by target effect category  
316 and by whether feature selection was implemented. Note: RMSE, root-mean-squared error, MedAE, median  
317 absolute error; R<sup>2</sup>, coefficient of determination; n, sample size.

318  
319



**Figure 3. Model evaluation.** (A) Out-of-sample performance is assessed through scatterplots comparing the mean predicted values for each chemical when it is part of the “testing” dataset across 30 cross-validation repeats (y-axis) against the corresponding surrogate values (x-axis). The dashed red line indicates perfect correspondence. (B) The distribution of performance metrics from 150 cross-validation scores (30 repeats x 5 folds), where each boxplot shows the median and interquartile range with whiskers representing the 95% confidence interval. The figure is subdivided by the performance metric, target effect category, and by whether feature selection was implemented. Note: RMSE, root-mean-squared error, MedAE, median absolute error; R<sup>2</sup>, coefficient of determination; n, sample size. The scale for R<sup>2</sup> is reversed to be consistent with values to the “left” corresponding to better performance.

The benchmarking revealed that the  $POD_{QSAR}$  values correlated well with the corresponding  $POD_{authoritative}$  values for general noncancer effects ( $n = 564$ ) (Figure S5), with  $RMSE = 0.50$  and  $MedAE = 0.32$ , both in  $\log_{10}$ -units, and  $R^2 = 0.79$ . The correspondence was poorer for reproductive/developmental effects, with  $RMSE = 0.75$ ,  $MedAE = 0.40$ , and  $R^2 = 0.47$ . For both effect categories, the  $POD_{QSAR}$  values corresponded substantially better to the  $POD_{authoritative}$  values than did the  $POD_{NAM,50}$  values that were derived from *in vitro* bioactivity data.<sup>26</sup> The  $POD_{NAM,50}$  values yielded negative  $R^2$  values, indicating worse performance than a

1  
2  
3 337 naïve constant model. However, the performance of  $POD_{QSAR}$  values in this comparison may be  
4  
5 338 overstated because they incorporated information about  $POD_{authoritative}$  indirectly through the use  
6  
7 339 of surrogate PODs derived from ToxValDB, while the  $POD_{NAM,50}$  consisted of a completely  
8  
9 340 independent dataset.

### 13 341 *Feature Importance*

14  
15 342 Results from the feature selection can be found in the Supporting Information (**Figures**  
16  
17 343 **S6-S10**). Notably, the most important feature was consistently the QSAR-predicted LD50  
18  
19 344 derived from *in vivo* rat acute oral toxicity studies.<sup>29</sup> Four important features were common to  
20  
21 345 both effect categories:

- 24 346 • QSAR-predicted LD50 derived from *in vivo* rat acute oral toxicity studies  
25  
26 347 (CATMoS\_LD50\_pred)
- 28 348 • Combined dipolarity/polarizability (CombDipolPolariz)
- 30  
31 349 • Ready biodegradability, a binary variable (ReadyBiodeg\_pred\_discrete)
- 32  
33 350 • Water solubility at 25 °C (WS\_pred)

34  
35  
36 351 For these features, no more than 11% of the training datasets were imputed, with less than 1%  
37  
38 352 imputed for the predicted LD50 (**Figure S4**). The remaining important features depended on the  
39  
40 353 effect category (**Figures S6-S10**) and involved imputation of no more than 25% of the training  
41  
42 354 set. Some additional important features were identified by the replicate models but excluded  
43  
44 355 from the final models to prevent overfitting (**Figure S6**).

### 48 356 *Sensitivity Analysis*

49  
50 357 **Table 1** compares the estimated generalization errors of the models from the sensitivity  
51  
52 358 analysis. The best overall performance was exhibited by the baseline model (all 39 OPERA 2.9  
53  
54 359 features, imputation of missing values, Random Forest Regressor). However, as mentioned, this

1  
2  
3 360 model's performance was indistinguishable from the final model that involved a subset of 10  
4  
5 361 important features (**Figure 3B**). Interestingly, when the baseline model was applied to samples  
6  
7 362 without need for imputation, the model continued to exhibit favorable performance in terms of  
8  
9 363 RMSE and MedAE, but with substantially higher variances and with  $R^2$  values that were much  
10  
11 364 lower (**Table 1**), likely due to the much more limited training sample sizes. Additionally, when  
12  
13 365 using the more "traditional" descriptors from RDKit (2022.09.5),<sup>30</sup> the performance was similar  
14  
15 366 to, but slightly poorer than our baseline model, suggesting that the 10 selected OPERA features  
16  
17 367 encapsulate the essential information for POD prediction. Overall, our final model (Random  
18  
19 368 Forest Regressor with feature selection and OPERA 2.9 features) was among the highest  
20  
21 369 performing models in terms of its combination of low prediction error (RMSE and MedAE) and  
22  
23 370 higher  $R^2$ .  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

373 **Table 1.** Comparison of performance metrics for QSAR models predicting points of departure.

374	375	376	377	378
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24
25	26	27	28	29
30	31	32	33	34
35	36	37	38	39
40	41	42	43	44
45	46	47	48	49
50	51	52	53	54
55	56	57	58	59
60				
<b>QSAR Model (n)</b>	<b>RMSE</b>	<b>MedAE</b>	<b>R<sup>2</sup></b>	
<i>Current Work: General non-cancer effects</i>				
<b>RandomForestRegressor with feature selection (1,791)</b>	<b>0.69 [0.64 – 0.76]</b>	<b>0.40 [0.37 – 0.44]</b>	<b>0.48 [0.41 – 0.53]</b>	
*RandomForestRegressor (1,791)	0.68 [0.62 - 0.74]	0.39 [0.35 - 0.43]	0.50 [0.44 - 0.56]	
*GradientBoostingRegressor (1,791)	0.69 [0.64 - 0.75]	0.41 [0.37 - 0.46]	0.48 [0.42 - 0.55]	
*Ridge (1,791)	0.73 [0.68 - 0.79]	0.44 [0.40 - 0.48]	0.42 [0.36 - 0.48]	
*LinearRegression (1,791)	0.73 [0.68 - 0.79]	0.44 [0.40 - 0.48]	0.42 [0.36 - 0.48]	
*XGBRegressor (1,791)	0.72 [0.66 - 0.78]	0.42 [0.38 - 0.46]	0.43 [0.36 - 0.51]	
*SVR (1,791)	0.96 [0.89 - 1.04]	0.64 [0.57 - 0.69]	-0.01 [-0.03 - 0.01]	
*MLPRegressor (1,791)	2.75 [1.56 - 5.53]	0.67 [0.58 - 0.84]	-7.50 [-36.72 - -1.72]	
**OPERA w/ Exp. LD50s (1,791)	0.69 [0.63 - 0.75]	0.40 [0.37 - 0.43]	0.48 [0.42 - 0.55]	
**CompTox Features (1,791)	0.75 [0.69 - 0.82]	0.44 [0.39 - 0.49]	0.39 [0.31 - 0.46]	
**RDKit Features (1,789)	0.71 [0.65 - 0.78]	0.40 [0.36 - 0.44]	0.45 [0.38 - 0.51]	
**No Imputation (184)	0.58 [0.46 - 1.17]	0.37 [0.28 - 0.49]	0.22 [0.02 - 0.44]	
<i>Current Work: Reproductive/developmental effects</i>				
<b>RandomForestRegressor with feature selection (1,791)</b>	<b>0.58 [0.54 – 0.72]</b>	<b>0.31 [0.28 – 0.34]</b>	<b>0.49 [0.38 – 0.56]</b>	
*RandomForestRegressor (2,228)	0.57 [0.53 - 0.72]	0.31 [0.29 - 0.35]	0.51 [0.40 - 0.58]	
*GradientBoostingRegressor (2,228)	0.59 [0.54 - 0.73]	0.32 [0.30 - 0.35]	0.49 [0.37 - 0.55]	
*Ridge (2,228)	0.63 [0.58 - 0.76]	0.37 [0.34 - 0.40]	0.42 [0.32 - 0.48]	
*LinearRegression (2,228)	0.63 [0.58 - 0.76]	0.37 [0.34 - 0.40]	0.42 [0.32 - 0.48]	
*XGBRegressor (2,228)	0.62 [0.56 - 0.74]	0.33 [0.30 - 0.36]	0.43 [0.34 - 0.52]	
*SVR (2,228)	0.85 [0.77 - 0.96]	0.54 [0.51 - 0.58]	-0.03 [-0.06 - -0.01]	
*MLPRegressor (2,228)	1.75 [1.18 - 2.71]	0.56 [0.48 - 0.68]	-3.43 [-10.68 - -0.92]	
**OPERA w/ Exp. LD50s (2,228)	0.57 [0.53 - 0.71]	0.32 [0.29 - 0.34]	0.52 [0.42 - 0.58]	
**CompTox Features (2,228)	0.67 [0.60 - 0.81]	0.38 [0.34 - 0.41]	0.34 [0.26 - 0.44]	
**RDKit Features (2,224)	0.62 [0.55 - 0.73]	0.32 [0.29 - 0.35]	0.45 [0.37 - 0.52]	
**No Imputation (227)	0.45 [0.35 - 0.55]	0.28 [0.20 - 0.35]	0.40 [0.21 - 0.53]	
<i>Previous Work</i>				
Wignall et al. 2018 NOAEL (487)	N.R.	0.70 [0.06 - 1.82]	0.45	
Pradeep et al. 2020 CHR R,M (11201)	0.92-0.94	N.R.	0.39-0.40	
Pradeep et al. 2020 REP R,M (5951)	0.79-0.91	N.R.	0.26-0.31	
Pradeep et al. 2020 DEV R,M, Rb (9945)	0.76-0.80	N.R.	0.26-0.29	
Pradeep et al. 2020 ALL (71,020)	0.67-0.70	N.R.	0.54-0.57	

374 **Bold** represents the “final” model used for predictions. \*Sensitivity analyses using different machine  
 375 learning algorithms; \*\* Sensitivity analyses using different descriptor sets (all using Random Forest  
 376 Regressor without feature selection). Abbreviations: RMSE, root-mean-squared error; MedAE, median  
 377 absolute error; R<sup>2</sup>, coefficient of determination; N.R. not reported; CHR, chronic; REP, reproductive;  
 378 DEV, developmental, R, rat; M, mouse, Rb, Rabbit. Values for current work are median and 90% CI

379 based on “outer” cross-validation replicates (see Methods). Range for Pradeep et al. (2020) based on  
380 internal cross-validation and external test set.

### 381 *Model Application*

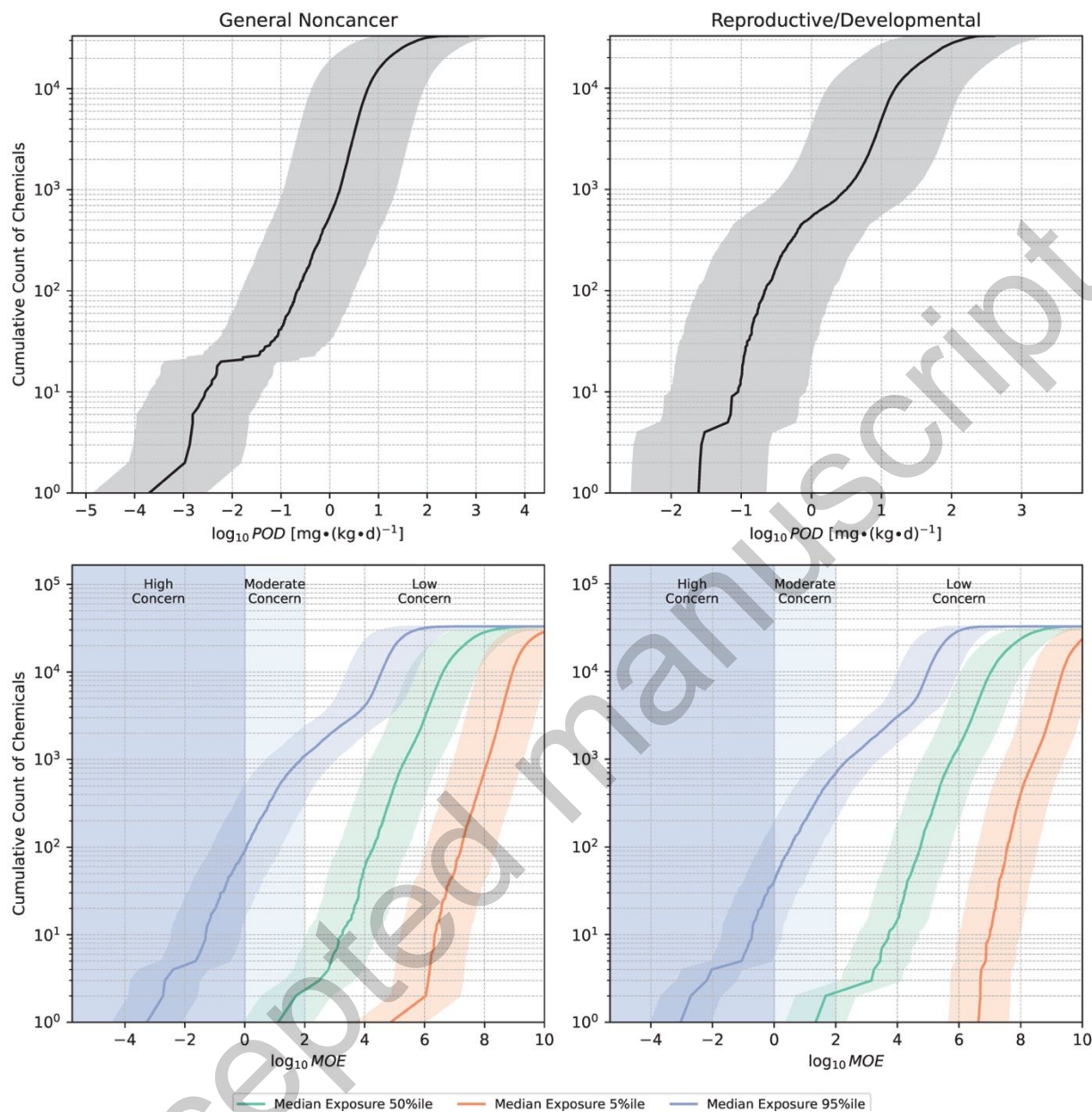
382 The top panels of **Figure 4** display cumulative counts of the application chemicals in  
383 relation to the corresponding  $POD_{QSAR}$  values, along with uncertainty estimates in the form of a  
384 90% prediction interval representing  $POD_{QSAR}$  (hazard) uncertainty (Supporting Information  
385 **Equation S8**). For general noncancer effects, the median  $POD_{QSAR}$  (with 5th - 95th percentiles)  
386 was  $11 \text{ mg}\cdot(\text{kg}\cdot\text{d})^{-1}$  (0.82 – 150). This distribution is somewhat higher (less potent) than that of  
387 the available regulatory/authoritative PODs (see **Figure S11**), as it is expected that higher  
388 potency (lower POD) chemicals would be more likely to have such regulatory or authoritative  
389 assessments. Additionally, as a sensitivity analysis, we also applied the model without feature  
390 selection to these chemicals and obtained consistent results [high correspondence between with  
391 and without feature selection:  $R^2 \sim 0.9$  and  $RMSE < 0.2 \text{ log-10 units}$  (**Figure S12**)].

392 The lower panels of **Figure 4** show the margins of exposure for an individual at the  
393 population median exposure, incorporating the 90% confidence interval for the population  
394 median exposure from SEEM3.<sup>19</sup> About ~2,400 chemicals emerged as *moderate* concerns for  
395 population median exposures ( $MOE < 100$ ) for general noncancer effects based on the upper 95<sup>th</sup>  
396 percentile of exposure uncertainty estimates and the lower boundary of the 90% prediction  
397 interval of  $POD_{QSAR}$  uncertainty. In a similar manner, ~500 chemicals emerged as *high* concerns  
398 ( $MOE < 1$ ) for general noncancer effects. For reproductive/developmental effects, the median  
399  $POD_{QSAR}$  was  $31 \text{ mg}\cdot(\text{kg}\cdot\text{d})^{-1}$  (3.4 – 280), with ~1,500 chemicals emerging as *moderate*  
400 concerns, and ~190 chemicals emerging as *high* concerns. In both cases, most chemicals appear  
401 to have low concern MOE values of  $>100$  at the level of the median population exposures. It is  
402 however important to note that this level of concern could be substantially higher for

1  
2  
3 403 subpopulations that regularly use products containing the considered chemicals.<sup>31</sup> A graphical  
4  
5 404 user interface will be made available for accessing these predictions and identifying chemicals of  
6  
7  
8 405 concern.

9  
10 406 Exposure uncertainty was the primary driver of the overall uncertainty in the margin of  
11  
12 407 exposure (**Figure 4**). The typical exposure uncertainty spanned 4 orders of magnitude, evidenced  
13  
14 408 by the median difference in  $\log_{10}$ -transformed exposure estimates between the 95<sup>th</sup> and 5<sup>th</sup>  
15  
16 409 percentiles. In contrast, when focusing on  $\text{POD}_{\text{QSAR}}$ , the typical error was constrained to less  
17  
18 410 than a factor of 5 according to the median RMSE of  $\leq 0.69$  in  $\log_{10}$ -units (**Figure 3B**). This  
19  
20 411 error corresponds to a squared geometric standard deviation ( $\text{GSD}^2$ )  $\leq 23$ , which, as expected, is  
21  
22 412 larger than the error reported by Aurisano et al. ( $\text{GSD}^2 \leq 17$  for all chemicals,  $\text{GSD}^2 \leq 14$  for  
23  
24 413 chemicals with at least 4 data points) that was based directly on *in vivo* PODs.<sup>8</sup>  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60





414  
415  
416  
417  
418  
419  
420  
421  
422  
423

**Figure 4.** Cumulative counts of application chemicals in relation to the predicted points of departure and margins of exposure. Data are shown for chemicals that were on the Merged NORMAN Suspect List (SusDat)<sup>20,21</sup> and within the applicability domain of SEEM3 ( $n = 32,524$ ),<sup>19</sup> excluding any training chemicals. The margins of exposure correspond to an individual at the population median exposure. Uncertainty is represented in two ways: (1) Exposure uncertainty, reflected by examining margins of exposure at different exposure percentiles; (2) Point of departure (hazard) uncertainty, represented by a 90% prediction interval derived from the median RMSE based on cross validation. Vertical spans highlight different risk categories as described in the Methods. The x-axis is truncated at  $\log_{10}MOE = 10$ . Note: POD, point of departure; MOE, margin of exposure.

## 424 DISCUSSION

425 This study successfully extended the work of Aurisano et al. (2023),<sup>8</sup> yielding a two-  
426 stage ML framework capable of generating human-equivalent noncancer PODs for oral exposure  
427 in the absence of *in vivo* data. This framework was applied to derive surrogate PODs and  
428 corresponding margins of exposure for over 30,000 chemicals expected based on monitoring to  
429 occur in the environment and which lacked *in vivo* toxicity data.<sup>20,21</sup> This represents a greater  
430 than three-fold increase in the coverage of organic chemicals with surrogate PODs compared to  
431 previous work.<sup>8</sup> Moreover, a graphical user interface will be made available for accessing  
432 predictions for organic chemicals available on the U.S. EPA's CompTox Chemistry Dashboard  
433 that pass the QSAR standardization workflow,<sup>22,23</sup> which will further increase the coverage of  
434 chemicals by over an order of magnitude to ~800,000.<sup>11</sup> Moreover, as shown in **Figure S4**, the  
435 rates of imputation for the >30,000 application chemicals were similar to the training set, with  
436 the most influential feature (CATMoS\_LD50\_pred) being imputed for only ~1% of values.  
437 Additionally, our training set of several thousand chemicals from ToxValDB appears to be  
438 diverse and representative based on similar coverage of features compared to application  
439 chemicals (**Figure S13**).<sup>7</sup>

440 Applying our two-stage models revealed several thousand chemicals of *moderate*  
441 concern, and several hundred chemicals of *high* concern, for health effects at estimated median  
442 population exposure levels (**Figure 4**). Notably, exposure uncertainty was the primary driver of  
443 the overall uncertainty in the margin of exposure. Exposure uncertainty was larger than  $POD_{QSAR}$   
444 (hazard) uncertainty, despite our QSAR-based approach inherently introducing a larger  
445 uncertainty than the surrogate PODs from Aurisano et al. (2023) that were based directly on *in*  
446 *vivo* data.<sup>8</sup> Moreover, we only assessed risk at estimated *median* exposure levels, and for most  
447 chemicals only a small fraction of the population is likely exposed. Thus, the actual uncertainty

1  
2  
3 448 in exposure is even greater when recognizing the need to address highly exposed subpopulations.  
4  
5 449 These findings underscore the need for refined exposure estimates to better characterize chemical  
6  
7  
8 450 use patterns, product compositions, and human behaviors that influence exposure.<sup>32–34</sup>  
9

10 451 In **Table 2**, we illustrate another case study example demonstrating how these models  
11  
12 452 could be used in the context of deriving a reference dose (RfD) for a “new” chemical. In  
13  
14 453 particular, we use the example of 4-Methylcyclohexanemethanol (MCHM) – a chemical used in  
15  
16 454 the processing of coal that spilled from a storage tank into the Elk River in West Virginia, US in  
17  
18  
19 455 January 2014. At the time, there were no regulatory toxicity values for MCHM. After several  
20  
21 456 days, CDC (2014) developed guidance levels based on a 4-week rat study (Eastman, 1990), and  
22  
23 457 several months later, an expert panel (TERA 2014) proposed refined analyses using the same  
24  
25  
26 458 study.<sup>35–37</sup> Over six years later, NTP completed a developmental and reproductive toxicity study  
27  
28 459 in rats (NTP 2020).<sup>38</sup> However, as illustrated in **Table 2**, utilizing our QSAR models for  
29  
30 460 predicting PODs and deriving RfDs for MCHM would yield very similar results in a much more  
31  
32 461 rapid timeframe of minutes, rather than days, months, or years. Additionally, because our  
33  
34 462 predictions include confidence bounds for model uncertainty, they can also be incorporated into  
35  
36 463 probabilistic derivations of toxicity values or health impacts.<sup>39–41</sup>  
37  
38  
39  
40 464  
41  
42 465  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

466 **Table 2.** Illustration of application to deriving a Reference Dose (RfD) for 4-  
 467 Methylcyclohexanemethanol (MCHM) in context of 2014 chemical spill in West Virginia, US.

Source	Point of Departure (mg·(kg-d) <sup>-1</sup> )	UF <sub>A</sub> <sup>a</sup>	UF <sub>H</sub> <sup>a</sup>	UF <sub>D</sub> <sup>a</sup>	RfD (mg·(kg-d) <sup>-1</sup> )	Analysis time
CDC (2014)	100 (Eastman 1990)	10	10	10	0.1	Days
TERA (2014)	71 (Eastman 1990) <sup>b</sup>	10	10	10	0.07	Months
NTP (2020)	50 (maternal)	10	10	10	0.05	Years
This work - General non- cancer	1.9 <sup>c</sup>	3 <sup>d</sup>	10	1 <sup>e</sup>	0.06	Minutes
This work – Reproductive/ Developmental	3.5 <sup>c</sup>	3 <sup>d</sup>	10	1 <sup>e</sup>	0.1	Minutes

468 Notes:

469 <sup>a</sup> Default factor unless otherwise noted. UF<sub>A</sub> = animal to human, UF<sub>H</sub> = human variability, UF<sub>D</sub>  
 470 = database inadequacy.

471 <sup>b</sup> Duration adjusted for 5 days/week exposure.

472 <sup>c</sup> QSAR human equivalent POD prediction is 26 [90% CI: 1.9-360] mg·(kg-d)<sup>-1</sup> for general non-  
 473 cancer and 32 [90% CI: 3.5-290] for reproductive/developmental effects. Lower 95% confidence  
 474 bound used as a “conservative” POD.

475 <sup>d</sup> QSAR predictive POD is already adjusted from animal to human equivalent dose using  
 476 allometric scaling.

477 <sup>e</sup> Reduced to 1 because database uncertainty is already addressed by using lower confidence  
 478 bound of QSAR-predicted POD and separate predictions for general non-cancer and  
 479 reproductive/developmental effects.

480

481

1  
2  
3 482 A primary strength of our framework lies in its two-stage approach described in the  
4  
5 483 Methods. Our final models accurately predicted PODs using a subset of 10 interpretable features  
6  
7 484 from OPERA 2.9 (**Figure S6**).<sup>9,10</sup> A unique aspect of this approach was the incorporation of  
8  
9 485 predicted biological features. Notably, the QSAR-predicted LD50, derived from *in vivo* rat acute  
10  
11 486 oral toxicity studies,<sup>29</sup> consistently emerged as the most important feature in our models. For this  
12  
13 487 feature, >99% of the chemicals in the training set were within the applicability domain (**Figure**  
14  
15 488 **S4**). This feature indicates the acute mammalian potency of a chemical, and was previously  
16  
17 489 predicted with an RMSE of around 0.50 (in log-10 units).<sup>29</sup> As expected, our POD predictions  
18  
19 490 had RMSE values that were (slightly) greater because they relied on the QSAR-predicted LD50  
20  
21 491 as a “feature.” Importantly, using *experimental* LD50 values as features in our sensitivity  
22  
23 492 analysis did not materially improve model performance, while substantially reducing the  
24  
25 493 applicability domain of the model because only chemicals with experimental LD50s were  
26  
27 494 predicted. Other important features were easily interpretable physical/chemical/biological  
28  
29 495 properties, such as water solubility or fish bioconcentration factor. Moreover, certain structural  
30  
31 496 properties, such as combined dipolarity/polarizability, also emerged as important features  
32  
33 497 independently of the predicted physical/chemical/biological properties. In essence, our two-stage  
34  
35 498 framework is akin to a traditional deep learning model, but providing a supervised intermediate  
36  
37 499 layer that transforms raw chemical descriptors into readily interpretable  
38  
39 500 physical/chemical/toxicological properties. However, a limitation of this approach is that the  
40  
41 501 applicability domain of the overall model is constrained by those of the individual first stage  
42  
43 502 models.

44  
45 503 Comparatively, our final models outperformed many alternative models in our sensitivity  
46  
47 504 analyses, as well as those published previously. Specifically, our in-sample predictions aligned  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 505 more closely with authoritative PODs than the combination of high-throughput *in vitro*  
4  
5 506 bioactivity data with toxicokinetic data (**Figure S5**).<sup>26</sup> Moreover, even our accuracy for “out-of-  
6  
7 507 sample” predictions were higher than those based on extrapolation from *in vitro*-based PODs.  
8  
9  
10 508 Additionally, as shown in **Table 1**, our QSAR models had similar or better performance  
11  
12 509 compared to previous models developed by Wignall et al. (2018) or Pradeep et al. (2020).<sup>4,13</sup>  
13  
14 510 Although the final “ALL” model by Pradeep et al. (2020) that uses study type and species as  
15  
16 511 additional descriptors had an R<sup>2</sup> value slightly higher than ours, this model includes subchronic  
17  
18 512 and subacute studies, and also does not identify a “critical effect” POD. On the other hand, our  
19  
20 513 “surrogate” PODs can be directly used in deriving toxicity values for application in various risk  
21  
22 514 and impact assessment and characterization approaches. Nonetheless, despite differences in  
23  
24 515 target variables making direct comparisons challenging, these studies suggest an upper limit in  
25  
26 516 the performance of QSAR models trained with *in vivo* data from ToxValDB.<sup>7</sup> Moreover, the  
27  
28 517 performance achievable through QSAR modeling is constrained by the intrinsic variability in the  
29  
30 518 derived toxicity values and PODs across different organizations for identical chemicals.<sup>4</sup>

31  
32  
33 519 For regulatory use, it is also important to consider our model and framework in light of  
34  
35 520 internationally recognized evaluation criteria for QSAR models. According to the *(Q)SAR*  
36  
37 521 *Assessment Framework* by OECD,<sup>14</sup> a QSAR model under consideration should be associated  
38  
39 522 with (1) a defined endpoint; (2) an unambiguous algorithm; (3) a defined domain of  
40  
41 523 applicability; (4) appropriate measures of goodness-of-fit, robustness and predictivity; (5) a  
42  
43 524 mechanistic interpretation, if possible. **Table S2** shows the results of applying the *(Q)SAR*  
44  
45 525 *Assessment Framework* to our modeling framework, demonstrating how our framework  
46  
47 526 conforms to general principles and criteria for use of QSAR models.<sup>14</sup>  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 527 Despite its advantages, our framework has several notable limitations. First, it is possible  
4  
5 528 that the actual generalization errors of our models were larger than those reported (**Figure 3B**),  
6  
7 529 particularly for features with a large proportion of missing values. In our framework, missing  
8  
9 530 values were imputed with the median, a common practice to maintain dataset integrity. However,  
10  
11 531 this approach can bias predictions towards central estimates, effectively narrowing the observed  
12  
13 532 variability. This “mean reversion” phenomenon can result in predictions that are less varied and  
14  
15 533 more centered around the median (**Figure S14**), which might not always reflect the underlying  
16  
17 534 distribution. This problem was partially mitigated by excluding features with many missing  
18  
19 535 values from our modeling pipeline (**Figure 1B**). Furthermore, based on our in-sample  
20  
21 536 performance and benchmarking, there may be a small trend towards overpredicting PODs for  
22  
23 537 higher potency chemicals (**Figures 2 and S5**). Again, this may be a mean reversion phenomenon  
24  
25 538 because of random forest is an ensemble-based method that averages over multiple individual  
26  
27 539 models and chemicals. This trend of a narrower range of predicted PODs was also observed in a  
28  
29 540 previous QSAR modeling effort.<sup>4</sup>

30  
31 541 Additionally, like most QSAR models, our models are only applicable to single organic  
32  
33 542 compounds of small to medium sizes; mixtures, large biomolecules, polymeric chains,  
34  
35 543 nanomaterials, and inorganic compounds are outside the applicability domain of OPERA 2.9.<sup>9,10</sup>  
36  
37 544 Different types of prediction models will need to be developed for these chemicals.  
38  
39 545 Additionally, our models were limited by the broad categorization of health effects.<sup>8</sup> This  
40  
41 546 categorization was necessitated by data availability; predicting PODs at a higher resolution, such  
42  
43 547 as for specific critical effects or organ systems, would have further fragmented an already limited  
44  
45 548 dataset. Our models also focused on the oral exposure route, and future work is needed to  
46  
47 549 incorporate additional exposure routes. Additionally, our model uncertainty estimates are based  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 550 on cross-validation generalization error, and future work could more fully characterize model  
4  
5 551 uncertainty, for instance, at the level of the individual prediction.  
6

7  
8 552 Overall, this study predicted *in vivo* noncancer PODs for organic chemicals, with typical  
9  
10 553 RMSEs of less than an order of magnitude, based on structure alone. Our framework offers a  
11  
12 554 high-throughput alternative to augment approaches that are based directly on *in vivo* data.  
13  
14 555 Moreover, our model also conforms well to OECD guidance for evaluating QSAR models,<sup>14</sup>  
15  
16 556 increasing confidence in our model predictions. These predictions can, in turn, be directly used  
17  
18 557 for a range of hazard, risk, and impact characterization applications, including (but not limited  
19  
20 558 to) deriving probabilistic toxicity values,<sup>39,42</sup> emergency response, contaminated site remediation,  
21  
22 559 LCIA, CAA, and comparative risk screening. Thus, predictions from our model can substantially  
23  
24 560 expand the coverage of chemicals that can be evaluated for their human health risks and impacts,  
25  
26 561 and thereby better promote a safer and more resilient, sustainable, and healthy environment.  
27  
28  
29  
30  
31  
32  
33

## 34 563 SUPPORTING INFORMATION

35  
36 564 Supplemental methods including feature preprocessing steps, model training steps  
37  
38 565 (**Figure S1**), model performance metrics and evaluation (**Figure S2**), model descriptions (**Table**  
39  
40 566 **S1**), and uncertainty analysis, as well as supplemental results (**Figures S3-S14** and **Table S2**). A  
41  
42 567 supplemental Excel file (**Tables S3-S4**) describes the features used to the train the QSAR models  
43  
44 568 for predicting points of departure.  
45  
46  
47  
48

## 49 569 ACKNOWLEDGEMENTS

50  
51 570 The authors thank Drs. Cedric Wannaz and Kamel Mansouri for assistance with OPERA,  
52  
53 571 and Prof. Jian Tao (Texas A&M University) for technical guidance during the initial stages of  
54  
55 572 this work. This study was funded by NIEHS T32 ES026568, the TAMU Superfund Research  
56  
57  
58  
59  
60



1  
2  
3 573 Center (NIEHS P42 ES027704), the TAMU Center for Environmental Health Research (NIEHS  
4  
5 574 P30 ES029067). This study was also financially supported by the “Safe and Efficient Chemistry  
6  
7 575 by Design (SafeChem)” project, which was funded by the Swedish Foundation for Strategic  
8  
9 576 Environmental Research (Grant No. DIA 2018/11), and the PARC project (Grant No.  
10  
11 577 101057014), which was funded under the European Union’s Horizon Europe Research and  
12  
13 578 Innovation program.  
14  
15  
16  
17 579  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Accepted manuscript

## 580 REFERENCES

- 581 (1) United States Environmental Protection Agency (EPA). *U.S. EPA System of Registries*  
582 *Terms & Acronyms*.  
583 [https://sor.epa.gov/sor\\_internet/registry/termreg/searchandretrieve/termsandacronyms/search.do](https://sor.epa.gov/sor_internet/registry/termreg/searchandretrieve/termsandacronyms/search.do) (accessed 2023-12-01).
- 584  
585 (2) Fantke, P.; Huang, L.; Overcash, M.; Griffing, E.; Jolliet, O. Life Cycle Based Alternatives  
586 Assessment (LCAA) for Chemical Substitution. *Green Chem.* **2020**, *22* (18), 6008–6024.
- 587 (3) Fantke, P.; Chiu, W. A.; Aylward, L.; Judson, R.; Huang, L.; Jang, S.; Gouin, T.;  
588 Rhomberg, L.; Aurisano, N.; McKone, T.; Jolliet, O. Exposure and Toxicity  
589 Characterization of Chemical Emissions and Chemicals in Products: Global  
590 Recommendations and Implementation in USEtox. *Int. J. Life Cycle Assess.* **2021**, *26* (5),  
591 899–915. <https://doi.org/10.1007/s11367-021-01889-y>.
- 592 (4) Wignall, J. A.; Muratov, E.; Sedykh, A.; Guyton, K. Z.; Tropsha, A.; Rusyn, I.; Chiu, W. A.  
593 Conditional Toxicity Value (CTV) Predictor: An *In Silico* Approach for Generating  
594 Quantitative Risk Estimates for Chemicals. *Environ. Health Perspect.* **2018**, *126* (5),  
595 057008. <https://doi.org/10.1289/EHP2998>.
- 596 (5) Wang, Z.; Walker, G. W.; Muir, D. C. G.; Nagatani-Yoshida, K. Toward a Global  
597 Understanding of Chemical Pollution: A First Comprehensive Analysis of National and  
598 Regional Chemical Inventories. *Environ. Sci. Technol.* **2020**, *54* (5), 2575–2584.  
599 <https://doi.org/10.1021/acs.est.9b06379>.
- 600 (6) Von Borries, K.; Holmquist, H.; Kosnik, M.; Beckwith, K. V.; Jolliet, O.; Goodman, J. M.;  
601 Fantke, P. Potential for Machine Learning to Address Data Gaps in Human Toxicity and  
602 Ecotoxicity Characterization. *Environ. Sci. Technol.* **2023**, *57* (46), 18259–18270.  
603 <https://doi.org/10.1021/acs.est.3c05300>.
- 604 (7) Judson, R. ToxValDB: Compiling Publicly Available In Vivo Toxicity Data, 2018.  
605 [https://www.epa.gov/sites/production/files/2018-](https://www.epa.gov/sites/production/files/2018-12/documents/comptox_cop_dec_20_2018_final.pdf)  
606 [12/documents/comptox\\_cop\\_dec\\_20\\_2018\\_final.pdf](https://www.epa.gov/sites/production/files/2018-12/documents/comptox_cop_dec_20_2018_final.pdf) (accessed 2023-11-16).
- 607 (8) Aurisano, N.; Jolliet, O.; Chiu, W. A.; Judson, R.; Jang, S.; Unnikrishnan, A.; Kosnik, M.  
608 B.; Fantke, P. Probabilistic Points of Departure and Reference Doses for Characterizing  
609 Human Noncancer and Developmental/Reproductive Effects for 10,145 Chemicals.  
610 *Environ. Health Perspect.* **2023**, *131* (3), 037016. <https://doi.org/10.1289/EHP11524>.
- 611 (9) Mansouri, K.; Grulke, C. M.; Judson, R. S.; Williams, A. J. OPERA Models for Predicting  
612 Physicochemical Properties and Environmental Fate Endpoints. *J. Cheminformatics* **2018**,  
613 *10* (1), 10. <https://doi.org/10.1186/s13321-018-0263-1>.
- 614 (10) Mansouri, K. OPERA, 2023. <https://github.com/NIEHS/OPERA> (accessed 2023-11-17).
- 615 (11) Williams, A. J.; Grulke, C. M.; Edwards, J.; McEachran, A. D.; Mansouri, K.; Baker, N. C.;  
616 Patlewicz, G.; Shah, I.; Wambaugh, J. F.; Judson, R. S.; Richard, A. M. The CompTox  
617 Chemistry Dashboard: A Community Data Resource for Environmental Chemistry. *J.*  
618 *Cheminformatics* **2017**, *9* (1), 61. <https://doi.org/10.1186/s13321-017-0247-6>.
- 619 (12) Bell, S. M.; Phillips, J.; Sedykh, A.; Tandon, A.; Sprankle, C.; Morefield, S. Q.; Shapiro,  
620 A.; Allen, D.; Shah, R.; Maull, E. A.; Casey, W. M.; Kleinstreuer, N. C. An Integrated  
621 Chemical Environment to Support 21st-Century Toxicology. *Environ. Health Perspect.*  
622 **2017**, *125* (5), 054501. <https://doi.org/10.1289/EHP1759>.
- 623 (13) Pradeep, P.; Friedman, K. P.; Judson, R. Structure-Based QSAR Models to Predict Repeat  
624 Dose Toxicity Points of Departure. *Comput. Toxicol.* **2020**, *16*, 100139.

- 1  
2  
3 625 (14) Organisation for Economic Co-operation and Development (OECD). *(Q)SAR Assessment*  
4 626 *Framework: Guidance for the Regulatory Assessment of (Quantitative) Structure - Activity*  
5 627 *Relationship Models, Predictions, and Results Based on Multiple Predictions, OECD Series*  
6 628 *on Testing and Assessment, No. 386, Environment, Health and Safety, Environment*  
7 629 *Directorate, OECD.*
- 9 630 (15) Huijbregts, M. A. J.; Rombouts, L. J. A.; Ragas, A. M. J.; van de Meent, D. Human-  
10 631 Toxicological Effect and Damage Factors of Carcinogenic and Noncarcinogenic Chemicals  
11 632 for Life Cycle Impact Assessment. *Integr. Environ. Assess. Manag.* **2005**, *1* (3), 181–244.  
12 633 <https://doi.org/10.1897/2004-007r.1>.
- 14 634 (16) Fantke, P.; Ernstoff, A. S.; Huang, L.; Csiszar, S. A.; Jolliet, O. Coupled Near-Field and  
15 635 Far-Field Exposure Assessment Framework for Chemicals in Consumer Products. *Environ.*  
16 636 *Int.* **2016**, *94*, 508–518.
- 17 637 (17) Jolliet, O.; Ernstoff, A. S.; Csiszar, S. A.; Fantke, P. Defining Product Intake Fraction to  
18 638 Quantify and Compare Exposure to Consumer Products. *Environ. Sci. Technol.* **2015**, *49*  
19 639 (15), 8924–8931. <https://doi.org/10.1021/acs.est.5b01083>.
- 21 640 (18) Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel,  
22 641 M.; Prettenhofer, P.; Weiss, R.; Dubourg, V. Scikit-Learn: Machine Learning in Python. *J.*  
23 642 *Mach. Learn. Res.* **2011**, *12*, 2825–2830.
- 24 643 (19) Ring, C. L.; Arnot, J. A.; Bennett, D. H.; Egeghy, P. P.; Fantke, P.; Huang, L.; Isaacs, K.  
25 644 K.; Jolliet, O.; Phillips, K. A.; Price, P. S.; Shin, H.-M.; Westgate, J. N.; Setzer, R. W.;  
26 645 Wambaugh, J. F. Consensus Modeling of Median Chemical Intake for the U.S. Population  
27 646 Based on Predictions of Exposure Pathways. *Environ. Sci. Technol.* **2019**, *53* (2), 719–732.  
28 647 <https://doi.org/10.1021/acs.est.8b04056>.
- 30 648 (20) Mohammed Taha, H.; Aalizadeh, R.; Alygizakis, N.; Antignac, J.-P.; Arp, H. P. H.; Bade,  
31 649 R.; Baker, N.; Belova, L.; Bijlsma, L.; Bolton, E. E.; Brack, W.; Celma, A.; Chen, W.-L.;  
32 650 Cheng, T.; Chirsir, P.; Čirka, L.; D'Agostino, L. A.; Djoumbou Feunang, Y.; Dulio, V.;  
33 651 Fischer, S.; Gago-Ferrero, P.; Galani, A.; Geueke, B.; Głowacka, N.; Glüge, J.; Groh, K.;  
34 652 Grosse, S.; Haglund, P.; Hakkinen, P. J.; Hale, S. E.; Hernandez, F.; Janssen, E. M.-L.;  
35 653 Jonkers, T.; Kiefer, K.; Kirchner, M.; Koschorreck, J.; Krauss, M.; Krier, J.; Lamoree, M.  
36 654 H.; Letzel, M.; Letzel, T.; Li, Q.; Little, J.; Liu, Y.; Lunderberg, D. M.; Martin, J. W.;  
37 655 McEachran, A. D.; McLean, J. A.; Meier, C.; Meijer, J.; Menger, F.; Merino, C.; Muncke,  
38 656 J.; Muschket, M.; Neumann, M.; Neveu, V.; Ng, K.; Oberacher, H.; O'Brien, J.; Oswald,  
39 657 P.; Oswaldova, M.; Picache, J. A.; Postigo, C.; Ramirez, N.; Reemtsma, T.; Renaud, J.;  
40 658 Rostkowski, P.; Rüdell, H.; Salek, R. M.; Samanipour, S.; Scheringer, M.; Schliebner, I.;  
41 659 Schulz, W.; Schulze, T.; Sengl, M.; Shoemaker, B. A.; Sims, K.; Singer, H.; Singh, R. R.;  
42 660 Sumarah, M.; Thiessen, P. A.; Thomas, K. V.; Torres, S.; Trier, X.; Van Wezel, A. P.;  
43 661 Vermeulen, R. C. H.; Vlaanderen, J. J.; Von Der Ohe, P. C.; Wang, Z.; Williams, A. J.;  
44 662 Willighagen, E. L.; Wishart, D. S.; Zhang, J.; Thomaidis, N. S.; Hollender, J.; Slobodnik,  
45 663 J.; Schymanski, E. L. The NORMAN Suspect List Exchange (NORMAN-SLE):  
46 664 Facilitating European and Worldwide Collaboration on Suspect Screening in High  
47 665 Resolution Mass Spectrometry. *Environ. Sci. Eur.* **2022**, *34* (1), 104.  
48 666 <https://doi.org/10.1186/s12302-022-00680-6>.
- 50 667 (21) NORMAN Network; Aalizadeh, R.; Alygizakis, N.; Schymanski, E.; Slobodnik, J.; Fischer,  
51 668 S.; Cirka, L. S0 | SUSDAT | Merged NORMAN Suspect List: SusDat (NORMAN-SLE-  
52 669 S0.0.4.3) [Data Set]. Zenodo. <https://doi.org/10.5281/Zenodo.6853705>, 2022.

- 1  
2  
3 670 (22) Mansouri, K.; Abdelaziz, A.; Rybacka, A.; Roncaglioni, A.; Tropsha, A.; Varnek, A.;  
4 671 Zakharov, A.; Worth, A.; Richard, A. M.; Grulke, C. M.; Trisciuzzi, D.; Fourches, D.;  
5 672 Horvath, D.; Benfenati, E.; Muratov, E.; Wedebye, E. B.; Grisoni, F.; Mangiatordi, G. F.;  
6 673 Incisivo, G. M.; Hong, H.; Ng, H. W.; Tetko, I. V.; Balabin, I.; Kancherla, J.; Shen, J.;  
7 674 Burton, J.; Nicklaus, M.; Cassotti, M.; Nikolov, N. G.; Nicolotti, O.; Andersson, P. L.;  
8 675 Zang, Q.; Politi, R.; Beger, R. D.; Todeschini, R.; Huang, R.; Farag, S.; Rosenberg, S. A.;  
9 676 Slavov, S.; Hu, X.; Judson, R. S. CERAPP: Collaborative Estrogen Receptor Activity  
10 677 Prediction Project. *Environ. Health Perspect.* **2016**, *124* (7), 1023–1033.  
11 678 <https://doi.org/10.1289/ehp.1510267>.  
12 679 (23) Mansouri, K. QSAR-Ready, 2022. <https://github.com/NIEHS/QSAR-ready> (accessed 2023-  
13 680 11-17).  
14 681 (24) Hastie, T.; Tibshirani, R.; Friedman, J. *The Elements of Statistical Learning*; Springer  
15 682 Series in Statistics; Springer: New York, NY, 2009. [https://doi.org/10.1007/978-0-387-  
16 683 84858-7](https://doi.org/10.1007/978-0-387-84858-7).  
17 684 (25) Breiman, L. Random Forests. *Mach. Learn.* **2001**, *45* (1), 5–32.  
18 685 <https://doi.org/10.1023/A:1010933404324>.  
19 686 (26) Paul Friedman, K.; Gagne, M.; Loo, L.-H.; Karamertzanis, P.; Netzeva, T.; Sobanski, T.;  
20 687 Franzosa, J. A.; Richard, A. M.; Lougee, R. R.; Gissi, A. Utility of in Vitro Bioactivity as a  
21 688 Lower Bound Estimate of in Vivo Adverse Effect Levels and in Risk-Based Prioritization.  
22 689 *Toxicol. Sci.* **2020**, *173* (1), 202–225.  
23 690 (27) Agency for Toxic Substances and Disease Registry (ATSDR). *Evaluate the Evidence to*  
24 691 *Examine Non-Cancer Effects*. [https://www.atsdr.cdc.gov/pha-  
25 692 guidance/conducting\\_scientific\\_evaluations/indepth\\_toxicological\\_analysis/EvaluateEviden  
26 693 ceNon-CancerEffects.html](https://www.atsdr.cdc.gov/pha-guidance/conducting_scientific_evaluations/indepth_toxicological_analysis/EvaluateEvidenceNon-CancerEffects.html) (accessed 2023-11-27).  
27 694 (28) European Food Safety Authority (EFSA). *Margin of Exposure*.  
28 695 <https://www.efsa.europa.eu/en/topics/topic/margin-exposure> (accessed 2023-11-27).  
29 696 (29) Mansouri, K.; Karmaus, A. L.; Fitzpatrick, J.; Patlewicz, G.; Pradeep, P.; Alberga, D.;  
30 697 Alepee, N.; Allen, T. E. H.; Allen, D.; Alves, V. M.; Andrade, C. H.; Auernhammer, T. R.;  
31 698 Ballabio, D.; Bell, S.; Benfenati, E.; Bhattacharya, S.; Bastos, J. V.; Boyd, S.; Brown, J. B.;  
32 699 Capuzzi, S. J.; Chushak, Y.; Ciallella, H.; Clark, A. M.; Consonni, V.; Daga, P. R.; Ekins,  
33 700 S.; Farag, S.; Fedorov, M.; Fourches, D.; Gadaleta, D.; Gao, F.; Gearhart, J. M.; Goh, G.;  
34 701 Goodman, J. M.; Grisoni, F.; Grulke, C. M.; Hartung, T.; Hirn, M.; Karpov, P.; Korotcov,  
35 702 A.; Lavado, G. J.; Lawless, M.; Li, X.; Luechtefeld, T.; Lunghini, F.; Mangiatordi, G. F.;  
36 703 Marcou, G.; Marsh, D.; Martin, T.; Mauri, A.; Muratov, E. N.; Myatt, G. J.; Nguyen, D.-T.;  
37 704 Nicolotti, O.; Note, R.; Pande, P.; Parks, A. K.; Peryea, T.; Polash, A. H.; Rallo, R.;  
38 705 Roncaglioni, A.; Rowlands, C.; Ruiz, P.; Russo, D. P.; Sayed, A.; Sayre, R.; Sheils, T.;  
39 706 Siegel, C.; Silva, A. C.; Simeonov, A.; Sosnin, S.; Southall, N.; Strickland, J.; Tang, Y.;  
40 707 Teppen, B.; Tetko, I. V.; Thomas, D.; Tkachenko, V.; Todeschini, R.; Toma, C.; Tripodi, I.;  
41 708 Trisciuzzi, D.; Tropsha, A.; Varnek, A.; Vukovic, K.; Wang, Z.; Wang, L.; Waters, K. M.;  
42 709 Wedlake, A. J.; Wijeyesakere, S. J.; Wilson, D.; Xiao, Z.; Yang, H.; Zahoranszky-Kohalmi,  
43 710 G.; Zakharov, A. V.; Zhang, F. F.; Zhang, Z.; Zhao, T.; Zhu, H.; Zorn, K. M.; Casey, W.;  
44 711 Kleinstreuer, N. C. CATMoS: Collaborative Acute Toxicity Modeling Suite. *Environ.*  
45 712 *Health Perspect.* **2021**, *129* (4), 047013. <https://doi.org/10.1289/EHP8495>.  
46 713 (30) RDKit: Open-Source Cheminformatics. <https://www.rdkit.org>.

- 1  
2  
3 714 (31) Jolliet, O.; Huang, L.; Hou, P.; Fantke, P. High Throughput Risk and Impact Screening of  
4 715 Chemicals in Consumer Products. *Risk Anal.* **2021**, *41* (4), 627–644.  
5 716 <https://doi.org/10.1111/risa.13604>.
- 6 717 (32) Kvasnicka, J. I. Joint Influence of Human Activities and Indoor Microenvironments on  
7 718 Contaminant Exposure: A Mass-Balance Modeling Investigation. PhD Thesis, University  
8 719 of Toronto (Canada), 2022.  
9 720 [https://search.proquest.com/openview/b195ce7c5d905cea24dd601e3247f7c3/1?pq-](https://search.proquest.com/openview/b195ce7c5d905cea24dd601e3247f7c3/1?pq-origsite=gscholar&cbl=18750&diss=y)  
10 721 [origsite=gscholar&cbl=18750&diss=y](https://search.proquest.com/openview/b195ce7c5d905cea24dd601e3247f7c3/1?pq-origsite=gscholar&cbl=18750&diss=y) (accessed 2023-11-16).
- 11 722 (33) Aurisano, N.; Huang, L.; Milà i Canals, L.; Jolliet, O.; Fantke, P. Chemicals of Concern in  
12 723 Plastic Toys. *Environ. Int.* **2021**, *146*, 106194.  
13 724 <https://doi.org/10.1016/j.envint.2020.106194>.
- 14 725 (34) Huang, L.; Fantke, P.; Ritscher, A.; Jolliet, O. Chemicals of Concern in Building Materials:  
15 726 A High-Throughput Screening. *J. Hazard. Mater.* **2022**, *424*, 127574.  
16 727 <https://doi.org/10.1016/j.jhazmat.2021.127574>.
- 17 728 (35) Centers for Disease Control and Prevention (CDC). *Information about MCHM | CDC*  
18 729 *Emergency Preparedness & Response*.  
19 730 <https://emergency.cdc.gov/chemical/MCHM/westvirginia2014/mchm.asp> (accessed 2024-  
20 731 03-15).
- 21 732 (36) Eastman. *Four-Week Oral Toxicity Study of 4-Methylcyclohexane Methanol in the Rat. TX-*  
22 733 *89-296. Eastman Kodak Company.*; 1990. [http://appalachianwaterwatch.org/wp-](http://appalachianwaterwatch.org/wp-content/uploads/2014/01/Pure_Distilled_MCHM-28-Day_Oral_Feeding_Study.pdf)  
23 734 [content/uploads/2014/01/Pure\\_Distilled\\_MCHM-28-Day\\_Oral\\_Feeding\\_Study.pdf](http://appalachianwaterwatch.org/wp-content/uploads/2014/01/Pure_Distilled_MCHM-28-Day_Oral_Feeding_Study.pdf).
- 24 735 (37) TERA (Toxicology Excellence for Risk Assessment) 2014. *Report of Expert Panel Review*  
25 736 *of Screening Levels for Exposure to Chemicals from the January 2014 Elk River Spill*.  
26 737 <https://www.tera.org/Peer/WV/WV%20Expert%20Report%2012%20May%202014.pdf>.
- 27 738 (38) National Toxicology Program (NTP). *NTP Developmental and Reproductive Toxicity*  
28 739 *Technical Report on the Prenatal Development Studies of 4-Methylcyclohexanemethanol*  
29 740 *(CASRN 34885-03-5) in Sprague Dawley (Hsd:Sprague Dawley SD) Rats (Gavage*  
30 741 *Studies)*. Research Triangle Park, NC: National Toxicology Program. DART Report 02.;  
31 742 2020. <https://doi.org/10.22427/NTP-DART-02>.
- 32 743 (39) Chiu, W. A.; Slob, W. A Unified Probabilistic Framework for Dose–Response Assessment  
33 744 of Human Health Effects. *Environ. Health Perspect.* **2015**, *123* (12), 1241–1254.  
34 745 <https://doi.org/10.1289/ehp.1409385>.
- 35 746 (40) Chiu, W. A.; Axelrad, D. A.; Dalaijamts, C.; Dockins, C.; Shao, K.; Shapiro, A. J.; Paoli,  
36 747 G. Beyond the RfD: Broad Application of a Probabilistic Approach to Improve Chemical  
37 748 Dose–Response Assessments for Noncancer Effects. *Environ. Health Perspect.* **2018**, *126*  
38 749 (6), 067009. <https://doi.org/10.1289/EHP3368>.
- 39 750 (41) Fantke, P.; Aylward, L.; Bare, J.; Chiu, W. A.; Dodson, R.; Dwyer, R.; Ernstoff, A.;  
40 751 Howard, B.; Jantunen, M.; Jolliet, O.; Judson, R.; Kirchhübel, N.; Li, D.; Miller, A.; Paoli,  
41 752 G.; Price, P.; Rhomberg, L.; Shen, B.; Shin, H.-M.; Teeguarden, J.; Vallero, D.;  
42 753 Wambaugh, J.; Wetmore, B. A.; Zaleski, R.; McKone, T. E. Advancements in Life Cycle  
43 754 Human Exposure and Toxicity Characterization. *Environ. Health Perspect.* **2018**, *126* (12),  
44 755 125001. <https://doi.org/10.1289/EHP3871>.
- 45 756 (42) World Health Organization (WHO). Guidance Document on Evaluating and Expressing  
46 757 Uncertainty in Hazard Characterization. **2018**.  
47 758