

### www.epa.gov

# Structure-based QSAR Models to Predict Systemic Toxicity Points of Departure

<sup>1</sup>Oak Ridge Institute for Science and Education, Oak Ridge, Tennessee <sup>2</sup>National Center for Computational Toxicology, Office of Research and Development, U.S. Environmental Protection Agency, Research Triangle Park, North Carolina

### INTRODUCTION

Human health risk assessment associated with environmental chemical exposure is limited by the tens of thousands of chemicals with little or no experimental in vivo toxicity data. A complete battery of regulatory tests for risk assessment associated with a single chemical relies on multiple animal testing's and can costs millions of dollars. A more effective and alternative way to evaluating a large number of environmental chemicals is to prioritize them for testing based on alternative methods for predicting experimental data.

Data gap filling techniques, such as quantitative structure activity relationship (QSAR) models based on chemical structure information, are commonly used to predict risk in the absence of experimental data. This study presents a set of QSAR models developed using chemical structural and physicochemical properties for in vivo points of departure (POD, the point on the dose-response that marks the beginning of a low-dose extrapolation). The in vivo data is taken from the EPA's ToxValDB, a compilation of information on ~3000 unique chemicals from a variety of public data sources. The QSAR models presented here provide estimates of POD and are evaluated for enrichment of most potent chemicals. These models will be used to inform chemical screening and prioritization efforts.

# DATA PREPARATION

descriptors

PubChem fingerprints (881 bits)

PaDEL descriptors (1875)

Chemistry development kit (CDK) descriptors (18)



toxval\_type != [RfD, RfC]

use me > 1

Figure 1: Schematic of data selection from the ToxValDB for developing predictive models for POD. Models were developed for each risk assessment class and species combination. E.g.

Model 1: study type = chronic, sub-chronic | species = rat

Model 2: study type = chronic, sub-chronic | species = mouse Model 3: study type = chronic, sub-chronic | species = rat and mouse (use\_me category is an expert assigned measure of data quality.)

For demonstration of research conducted for this poster data and results for only chronic, sub-chronic | rat and mouse models (models 2 and 3) are presented.

### μ = 1.14 RAT σ = 1.26 POD POD<sub>tr</sub> Split Dataset Transform Training = 80% $POD_{tr} = Log_{10}(POD)$ **Test = 20%** Skew = -0.21 μ = 1.37 MOUS σ = 1.16

Figure 2: The chronic, sub-chronic POD values were log-transformed for both rat and mouse datasets. (a) Histogram of untransformed POD data, (b) Histogram of transformed POD (POD<sub>tr</sub>) data, and (c) Histogram of training and test data relative to each other.

# CHALLENGES

### **1. Experimental Variability**

- Data from different labs (sources) running the "same" experiment may get different answers
- Sources of variability: Species, strain, dose range, dose spacing, length of study etc.

Figure 3: Distribution of the range of POD values for the CHR-SUB (a) rat and (b) mouse datasets. Variability in experimental data leads to uncertainty in the model predictions. Roughly, the root mean squared error (RMSE) in the models can be estimated to be around  $1.17 (= \sqrt{1.36})$ for rat and 1.13 (= $\sqrt{1.28}$ ) for mouse models by just looking at the distribution of POD values for both datasets.

### 2. Model Uncertainty

- A model gives a result (a POD), but this is an estimate of the "true" POD. The true POD is mostly unknown.
- Uncertainty in the evaluation data will lead to uncertainty in the model and our estimate of its quality



**U.S. Environmental Protection Agency** Office of Research and Development

**Disclaimer:** The views expressed in this poster are those of the authors and do not necessarily reflect the views or policies of the U.S.EPA.

# Prachi Pradeep<sup>1,2</sup> and Richard Judson<sup>2</sup>





Figure 5: Plots of observed versus predicted POD values (transformed scale) for the best rat and mouse model (random forest model) for 5-fold internal cross-validation (red scatter plots) and external validation (green scatter plots).

(Figure 3) for each combination of study type and species on the external test set. As seen, there is not much variation in the metrics across different model combinations. (r: rat, m: mouse, ra: rabbit, sp: species, st: study type)

### Prachi Pradeep | pradeep.prachi@epa.gov | ORCID iD: 0000-0002-9219-4249 | Phone: 919-541-5150

		ENRICHMENT ANALYSIS		
ation		Each model was evaluated on the		
l cross-		external test set for enrichment of		
80%		N% most notent chemicals		
0070		Nyo most potent enemicais.		
	(b)	1. The chemicals in the predicted		
		external test set were sorted in the		
lation		order of potency.		
		2. X% of above sorted list was ther		
		evaluated for % enrichment of the		
		N% most notent chemicals		
		N/0 most potent chemicals		

External Validation							
RMSE	RMSE/σ	R <sup>2</sup>					
0.98	0.78	0.39					
1.02	0.86	0.26					





Figure 9: Plots of observed versus predicted POD values for the best rat and mouse model (random forest model) for 5-fold internal cross-validation (red scatter plots) and external validation (green scatter plots) with 95% confidence intervals (green error bars) for n = 100 bootstrap models.

# CONCLUSIONS

- Point-estimate model results demonstrate that independent study type and species combinations did not result in significantly better models than combining the data for all the classes and species together.
  - The RMSE for the all the models are within the variance in the underlying POD data (Figure 3 and 5). Enrichment analysis results demonstrate the utility of these models for chemical screening and prioritization efforts.
- Point-estimate with balanced dataset model results improved the training set results but did not show improved results on the
- external test sets.
- Point-estimate with confidence interval models presented a technique to estimate uncertainty associated with model predictions. The results demonstrate the impact of variability in training data (experimental POD) on uncertainty associated with model results.

### RESULTS

al Cross-validation		External Validation			
RMSE/σ	R <sup>2</sup>	RMSE	RMSE/σ	R <sup>2</sup>	
0.75	0.45	1.00	0.80	0.36	
0.79	0.39	1.03	0.87	0.24	