

## www.epa.gov

# **Development of QSAR Models to Predict Systemic Toxicity Points of Departure**

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

<sup>1</sup>Oak Ridge Institute for Science and Education, Oak Ridge, Tennessee

## INTRODUCTION

Human health risk assessment associated with environmental chemical exposure is limited by the tens of thousands of chemicals little or no experimental in vivo toxicity data. Data gap filling techniques, such as quantitative structure activity relationship (QSAR) models based on chemical structure information, are commonly used to predict hazard in the absence of experimental data. This study presents a set of QSAR models developed using chemical structural and physicochemical properties for chronic or sub-chronic in vivo points of departure (POD, the point on the doseresponse 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. These models will inform chemical screening and prioritization efforts.





Figure 1: Schematic of data selection for modeling purposes from the ToxValDB. The final dataset for rat data is N = 1691 and mouse data is N = 668.

## **MOLECULAR FEATURES**

- PubChem fingerprints (881 bits)
- Chemistry development kit (CDK) descriptors (18)



**Figure 2:** The 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 rat and mouse dataset. 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 1.15 (V1.31) for rat and 1.15 (V1.32) for mouse models.

2. Model Uncertainty

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



models

models

	1. Point-est models
	Rat



U.S. Environmental Protection Agency Disclaimer: The views expressed are those of the Office of Research and Development

RAT

authors and do not necessarily reflect the views or policies of the U.S. Environmental Protection Agency.

<sup>2</sup>National Center for Computational Toxicology, Office of Research and Development, U.S. Environmental Protection Agency, Research Triangle Park, North Carolina

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

Figure 3: 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).

Future work: (1) Add study parameters as model descriptors to account for additional lab-to-lab data variability (Species, strain, study duration, exposure route), and (2) Re-construct datasets to reduce sampling by sampling equally from both tails (and not duplicate data).