Quantitative Structure-In vitro-In vivo Correlations: Î Data Partitioning and QSAR Modeling of Acute Toxicity in Rats Hao Zhu¹, Lin Ye¹, Ann Richard², Ivan Rusyn³, and Alexander Tropsha¹ ¹Laboratory for Molecular Modeling, School of Pharmacy, UNC-Chapel Hill; ²National Center for Computational Toxicology, Office of Research & Development, US EPA; ³Department of Environmental Sciences and Engineering, School of Public Health, UNC-Chapel Hill Table 1. Consensus prediction of 23 compounds in the external validation set using the two step kNN QSAR approach. INTRODUCTION DATA 253 compounds with IC50 and For the purposes of this work, the dataset was curated as follows. Inorganic and organometallic compounds, as well as compound mixtures were excluded since these do not have conventional chemical descriptors used in QSAR studies. The curated subset from the original ZEBET dataset used in this work is comprised of 291 organic compounds. Exp. Real Pred. Pred. Log(1/LD₅₀) LD50 results CAS Chemical Name Exp. og(1/ICro) The Interagency Coordinating Committee on the Validation of Alternative Methods (ICCVAM) hosted a workshop to study the relationship between cytotoxicity and acute rodent toxicity for more than 361 diverse compounds compiled by the German Center for the Documentation and Validation of Alternative Methods (ZEBET) [1]. This workshop showed that there is no clear relationship between these two types of endpoints. For example, the correlation between the cytotoxicity (ICS0s) and the rat acute toxicity (LDS0s) for 253 compounds in this dataset is poor (equation 1): 1.02 -0.81 C1 no pred. 72571 Trypan blue 122 Class 1 0.85 -0.45 C1 no pred 642 KNN LD50 7327879 For 253 out of these 291 compounds, rat LD50 values were available. The following criteria were used to select the LD50 values: 1) Only LD50 values Dihydralazine sulfate compounds 0.57 -0.03 -0.30 C2 C1 -1.14 C1 C1 132605 Cinchophen published in Registry of Toxic Effects of Chemical Substances (RTECS) [3] were used; 2) If different issues of RTECS reported different LD50 values, Salicylamide Solit into three sets 54115 Nicotine -0.25 0.51 C2 C2 then the largest LD50 value was used. 230 517 KNN based on the fitness to the linear correlation between IC50 and I D50 compound modeling set classification models -1.59 C1 C1 84662 Diethyl phthalate -0.74 -0.82 Outli. 20624253 3.41 C1 METHODOLOGY sodium 3H2O $Log(1/LD50) = -0.64 + 0.44 \times Log(1/IC50)$ 51183 Triethvlene melamine 3.11 2.30 C2 no pred 40 KNN LD50 $R^2 = 0.45$, SE = 0.71, N = 253Visual inspection of the plot comparing cytotoxicity and rat acute toxicity data (Figure 2) leads to the following observations: 93 Class 2 (1) compounds models 77474 vachlorocyclopents 2.51 0.30 C1 no pred cytotoxicity is directly correlated with the acute toxicity <u>only</u> for some compounds in the dataset; The primary goal of this study is to develop robust and externally validated predictors for acute rodent toxicity. However, all previous attempts to develop *in vivo* rat LD50 QSAR models based on ZEBET dataset failed. 76448 Heptachlor 1.23 0.96 a a 86544 Hydralazine 0.48 0.25 C2 C2 2) most of the remaining compounds have higher acute toxicity in rats vs in 15 outliers dataset failed. 136607 n-Butyl benzoate 0.39 -1.46 Outli. C1 3) a small fraction of compounds have lower acute toxicity in rats vs in vitro. 108610782 0.02 -0.11 C2 no pred. In our recent study, we found that using High Throughput Screening (HTS) data as additional biological descriptors significantly improved the QSAR models for toxicity endpoints in vivo [2]. In this study, we have employed a similar hybrid modeling approach using cytotoxicity ICS0 data to develop the in vivo LDS0 QSAR models (Figure 1). 23 external (pyrid-4-yl)-pyridine Based on these observations, we could partition compounds in the modeling set into three subsets: **Class 1**, in which compounds' acute rat toxicity 5435643 -0.18 -1.36 C1 C1 compounds Isononylaldehyde 110407 Diethyl sebacate -0.21 -1.75 Outli. no pred. linearly correlates with their cytotoxicity, Class 2, in which tomber at could rat could rat toxicity does not correlate with cytotoxicity with points positioned above the regression line for Class 1; and a small set of **outliers** with points located Figure 3. Flowchart of data modeling for in vivo $\rm LD_{50}$ prediction using chemical structure and $\rm \ IC_{50}$ data. 69727 Salicylic acid -0.53 -0.81 C2 C1 59427 Phenylephrine -0.65 -0.32 C2 C1 below the regression line for Class 1. **RESULTS & DISCUSSION** 78415722 Milrinone -0.68 0.37 C2 C2 Based on this approach to data partitioning, we have designed the analytical workflow as shown in Figure 3. 62533 Aniline -0.84 -0.67 C2 C2 -0.47 There is no clear difference in chemical similarity distribution between compounds in Class 1 and Class 2 (Figure 4). The predictions of LD50 values for the 23 external compounds are shown in Table 1. -1.54 C1 C1 -1.54 -1.54 C1 C1 -1.46 123728 n-Butanal -1.11 00 PubChe 71410 -1.40 1-Pentanol 75092 Dichloromethane 78933 Ethyl methyl ketone -1.54 -1.27 C2 C1 -2.02 -1.67 C1 C1 3.00 Baseline 186 40% C2 Compounds 2.00 Compounds CONCLUSIONS - C1 Compounds Outliers abo Full Modeling S 2 1.00 The cytotoxicity data show weak direct correlation with in vivo acute the baseline 30% The cytotoxicity data show weak direct correlation with *in vivo* acute toxicity. Nevertheless, these data could be used to assist QSAR modeling of *in vivo* acute toxicity. We have developed a novel two-phase approach that leads to successful kNN QSAR rat LDSD owdels. LDSO values were predicted for 23 external compounds with high accuracy ($R^2 = 0.80$, SE = 0.34, Coverage = 74%). We believe that this activity-based partitioning approach using the *in vitro* toxicity data can be successful kNN external to the complex in vivo toxicity endpoints. This approach makes it feasible to combine *in vivo* acited y endpoints. This, and QSAR modeling to prioritize chemicals for *in vivo* naimal toxicity testing. - mean=2 Outliers held ē 0.00 - mean=2.1 the baseline mean=1.61 20% distar Ē-1.00 External Validation Se 10% QSAR air -2.00 2 REFERENCES -3.00 0% Figure 1. The use of hybrid chemical and biological ICCVAM: http://iccvam.niehs.nih.gov/ -4 -2 0 2 4 6 0 2 4 6 8 10 descriptors for developing in vivo LD₅₀ OSAR models. Log(1/IC50) (mmol/l)

Figure 2. The identification of the baseline correlation between IC₅₀ and

LD₅₀ values for the modeling set

This poster does not necessarily reflect EPA policy. Mention of trade names or commercial products does not constitute endorsement or recommendation for use.

Fuclidean distance to the nearest neighbor

Figure 4. The comparison of chemical similarity between the original

modeling set and two new subsets

 Zhu, H., Rusyn, I., Richard, A., Tropsha, A. Environ. Health Persp. 2008, 116, in press. Ruden, C., Hansson. S. Toxicol. Lett. 2003, 144, 159-172.

no pred.

no pred.

-0.44

-1.02

0.05

-0.98

-1.39

no pred.

no pred.

0.88

-0.26

-0.96

no pred.

no pred.

-1.03

-0.96

0.08

-1.89 -1.81

This study was supported, in part, by the NIH RoadMap grant GM076059 and by EPA STAR grant RD832720