



### State of the art in QSAR modeling

### Alexander Tropsha

Carolina Center for Computational Toxicology Carolina Center for Environmental Bioinformatics and Laboratory for Molecular Modeling Eshelman School of Pharmacy UNC-Chapel Hill



Introductory notes on QSAR and modern data streams: Chemical Structure – *in vitro* – *in vivo* data continuum

Novel workflows for hybrid and integrative QFAR
 Modeling (with examples of applications):

- -Hierarchical modeling based on in vitro in vivo Correlations
- Hybrid chemical-biological descriptors
- Consensus between chemical and biological neighbors

Conclusions and outlook: QSAR modeling workflows and the use of QSAR for decision support

### QSAR modeling progression

= pain

- Experimental Data
  - Structure
  - Activity
- Validated models of data
  - Descriptors
  - Statistical/machine learning techniques
- Data inputation and experimental confirmation
- Reliable models to enable decision = gain support (both in research and regulations)

### Current Problems and Challenges in QSAR modeling

- (partially) solved problems:
  - Descriptors
  - Modeling techniques (statistical/machine learning)
  - Model validation approaches
  - Virtual screening and experimental validation of predictions (low bar)
- Challenges
  - Chemical and biological data curation and correction
  - New workflows integrating QSAR with available experimental data
  - Utility of QSAR models outside of research labs (especially, regulatory acceptance)
  - New application areas (e.g., materials informatics; nanotoxicology; integration with PK/PD modeling)

### **QSAR and Chemical Toxicity** Testing in 21 Century

July 2007

CI

F<sub>3</sub>C

-N

#### Toxicity Testing in the 21st Century: A Vision and a Strategy

Advances in molecular biology, biotechnology, and other fields are paving the way for major improvements in how scientists evaluate the health risks posed by potentially toxic chemicals found at low levels in the environment. These advances would make toxicity testing quicker, less expensive, and more directly relevant to human exposures. They could also reduce the need for animal testing by substituting more laboratory tests based on human cells. This National Research Council report creates a far-reaching vision for the future of toxicity testing.

effects at lower doses or exposures. Test

animals are typically observed for overt

provide little information about biological

must be applied to account for differences

between test animals and humans. Finally,

use of animals in testing is expensive and

time consuming, and it sometimes raises

Today, toxicological

valuation of chemicals

is poised to take advan-

and biotechnology. This

revolution is making it

increasingly possible

to study the effects of

chemicals using cells,

tissues-preferably of

human origin-rather

than whole animals.

These powerful new

approaches should help

to address a number of

challenges facing the

cellular components, and

tage of the on-going

revolution in biology

ethical issues.

signs of adverse health effects, which

changes leading to such health effects.

Often controversial uncertainty factors

Toxicity tests on laboratory animals are conducted to evaluate chemicals-including medicines, food additives, and industrial, consumer, and agricultural chemicals-for their potential to cause cancer, birth defects, and other adverse health effects. Information from toxicity testing serves as an important part of the basis for public health and regulatory decisions concerning toxic chemicals. Current test

methods were developed incrementally over the past 50 to 60 years and are conducted using laboratory animals, such as rats and mice. Using the results of animal tests to predict human health effects involves a number of assumptions and extrapolations that remain controversial. Test animals are often exposed to higher doses than would be expected for typical human exposures, requiring assumptions about

REPORT IN BRIE

#### **Transforming Environmental Health Protection**

Francis S, Collins,1\*1 George M, Grav,2\* John R, Bucher3\*

n 2005, the U.S. Environmental Protection Agency (EPA), with support from the U.S. National Toxicology Program (NTP), funded a project at the National Research Council (NRC) to develop a long-range vision for toxicity testing and a strategic plan for implementing that vision. Both agencies wanted future toxicity testing and assessment paradigms to meet evolving regulatory needs. Challenges include the large numbers of substances that need to be tested and how to incorporate recent advances in molecular toxicology, computational sciences, and information technology: to rely increasingly on human as

opposed to animal data; and to offer increased efficiency in design and costs (1-5). In response, the NRC Committee on Toxicity Testing and Assessment of Environmental Agents produced two reports that reviewed current toxicity testing, identified key issues, and developed a vision and implementation strategy to create a major shift in the assessment of chemical hazard and risk (6, 7). Although the NRC reports have laid out a solid theoretical rationale, comprehensive and rigorously gathered data (and comparisons with historical animal data) will determine whether the hypothesized improvements will be realized in practice. For this purpose, NTP, EPA, and the National Institutes of Health Chemical Genomics Center (NCGC) (organizations with expertise in experimental toxicology, computational toxicology, and high-throughput technologies, respectively) have established a collaborative research program.

#### EPA, NCGC, and NTP Joint Activities

In 2004, the NTP released its vision and roadmap for the 21st century (1), which established initiatives to integrate high-

<sup>1</sup>Director, National Human Genome Research Institute (NHGRI), National Institutes of Health, Bethesda, MD 20892: <sup>2</sup>Assistant Administrator for the Office of Research

throughput screening (HTS) and other automated screening assays into its testing program. In 2005, the EPA established the National Center for Computational Toxicology (NCCT). Through these initiatives. NTP and EPA, with the NCGC, are promoting the evolution of toxicology from a predominantly observational science at the level of disease-specific models in vivo to a predominantly predictive science focused anism-based, biological observations in vitro (1, 4) (see figure, below).

Toxicity pathways. In vitro and in vivo tools are being used to identify cellular responses after chemical exposure expected to result in adverse health effects (7). HTS methods are a primary means of discovery for drug development, and screening of >100,000 compounds per day is routine (8). However, drug-discovery HTS methods traditionally test compounds at one concentra-

We propose a shift from primarily in vivo animal studies to in vitro assays, in vivo assays with lower organisms, and computational modeling for toxicity assessments.

tion, usually between 2 and 10 µM, and tolerate high false-negative rates. In contrast, in the EPA, NCGC, and NTP combined effort, all compounds are tested at as many as 15 concentrations, generally ranging from ~5 nM to ~100 µM, to generate a concentrationresponse curve (9). This approach is highly reproducible, produces significantly lower false-positive and false-negative rates than the traditional HTS methods (9), and facilion broad inclusion of target-specific, mech- tates multiassay comparisons. Finally, an informatics platform has been built to compare results among HTS screens; this is being expanded to allow comparisons with historical toxicologic NTP and EPA data (http://ncgc.nih.gov/pub/openhts). HTS data collected by EPA and NTP, as well as by the NCGC and other Molecular Libraries Initiative centers (http://mli.nih.gov/), are being made publicly available through Webbased databases [e.g., PubChem (http:// pubchem.ncbi.nlm.nih.gov)]. In addition,

5

luman experien 1-3 studies/yea >10.000/da

EPAs Contribution: The ToxCast Research Program

policies of their respective agencies tAuthor for correspondence. E-mail: francisc@mail.nih.gov

iransforming toxicology. The studies we propose will test whether high-throughput and computational toxicology approaches can yield data predictive of results from animal toxicity studies, will allow prioritization of chemicals for further testing, and can assist in prediction of risk to humans.

National Academ









### The importance of <u>DATA</u> to enable any informatics-dependent discipline: *bioinformatics example*

#### Google labs Books Ngram Viewer



### The importance of <u>DATA</u> to enable any informatics-dependent discipline: *cheminformatics example*

#### Google labs Books Ngram Viewer



# Data dependency and data quality are critical issues in QSAR modeling

- Cheminformaticians are at the mercy of data providers. Prediction performance of (Q)SAR models could depend strongly on the quality of input data (both structures and activities).
- <u>Both</u> chemical and biological data in a dataset may be inaccurate and in need of thorough curation
- The number of published QSAR models that were poor or not too successful due to data quality issue is unknown but possibly large
- Often considered trivial, the basic steps to curate a dataset of compounds are not so obvious especially for beginners.



Looks clean ...



(presence of inorganics, salts, organometallics, duplicates; certain hydrogens are lacking; wrong standardization; etc.)

http://chembench.mml.unc.edu

Continue

#### **QSAR modeling with non-curated datasets**





Procedures	Software	Availability	
Inorganics Removal	ChemAxon/JChem OpenEye/Filter	Free for Academia Free for Academia	
nemovai			
Structure Normalization (fragment removal, structural curation,	ChemAxon/Standardizer OpenBabel Molecular Networks/ CHECK,TAUTOMER	Free for Academia Free Commercial	
salt neutralization)			
Duplicate Removal	ISIDA/Duplicates HiT QSAR CCG/MOE	Free for Academia Free for Academia Commercial	
SDF management/ viewer File format converter	ISIDA/EdiSDF Hyleos/ChemFileBrowser OpenBabel ChemAxon/MarwinView CambridgeSoft/ChemOffice Schrodinger/Canvas ACD/ChemFolder Symyx Cheminformatics CCG/MOE Accelrys/Accord Tripos/Benchware Pantheon	Free Free Free Free for Academia Commercial Commercial Commercial Commercial Commercial Commercial Commercial	

Summary of major procedures and corresponding relevant software for every step of the data curation process

## 

# The effect of curation on dataset size.

Datasat	Number of compounds					
Dataset	Original set	Curated set				
Liver toxicants (DILI)	1061	951 (90%)				
Nitroaromatics (rats)	28	28 (100%)				
Nitroaromatics (T. pyriformis)	95	95 (100%)				
ToxRefDB	320	292 (91%)				
Ames mutagenicity	7090	6542 (92%)				
Bioavailability (UCSD)	805	734 (91%)				

### Statistical parameters of QSAR models before

#### and after curation.



ID	Name	$\mathbb{R}^2$	$Q^2$	R <sup>2</sup> <sub>EF</sub>	Sws	Scv	SEF	R <sup>2</sup> <sub>EVS</sub>	R <sup>2</sup> <sub>EVS(NM)</sub>	
1	Rat	0.96	0.84-0.93	0.89-0.92	0.11-0.13	0.16-0.24	0.20-0.26	_	_	
2	Rat <sub>(NM)</sub>	0.91-0.97	0.89-0.95	0.45-0.88	0.10-0.18	0.14-0.28	0.28-0.58	_	_	
3	TP	0.83	—	0.76	0.33	—	0.38	0.54	-0.58	
4	TP <sub>(NM)</sub>	0.85	_	0.54	0.31	_	0.54	0.49	0.44	
5	Biowisdom non-curated		No modeling was possible							
6	Biowisdom		Modeling Set 5-fold external CV Accuracy = 62-68% External sets Accuracy = 56-73%							
7*	<sup>59</sup> Ames non-curated	Sen	Sensitivity <sub>RF</sub> =83%; Sensitivity <sub>SVM</sub> =87%; Specificity <sub>RF</sub> =Specificity <sub>SVM</sub> =75% AUC <sub>GP</sub> =88%; AUC <sub>SVM</sub> =89%; AUC <sub>RF</sub> =83%							
8*	$ \begin{array}{c c} & & & \\ \hline & & \\ & & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & & \\ & & \\ & & & \\ & & \\ & $								%	

Where:

 $\label{eq:transform} \begin{array}{l} \mathsf{TP}-\textit{Tetrahymena pyriformis} \text{ dataset, }_{(\mathsf{NM})}-\text{ modeling with various representations of nitro groups} \\ \mathsf{R}^2 \text{ - determination coefficient, } \mathsf{Q}^2 \text{ - cross validation determination coefficient} \end{array}$ 

R<sup>2</sup><sub>EF</sub>- determination coefficient for external folds extracted from the modeling set

 $S_{ws}$  - standard error of a prediction for work set

 $S_{cv}$  - standard error of prediction for work set in cross validation terms

 $S_{ts}^{v}$  - standard error of a prediction for external folds extracted from the modeling set

A<sup>1</sup>- number of PLS latent variables, D - number of descriptors, M - number of molecules in the work set

R<sup>2</sup><sub>EVS</sub> - determination coefficient for external validation set

R<sup>2</sup><sub>EVS(NM)</sub> - determination coefficient for external validation set with shuffled nitro groups

AUC – Area Under Curve statistical parameter

RF - Random Forest, SVM - Supporting Vector Machine, GP - Gaussian Processes, \* Prediction performances are reported for external validation set.

## QSAR modeling of nitro-aromatic toxicants

-Case Study 1: 28 compounds tested in rats, log(LD50), mmol/kg. -Case Study 2: 95 compounds tested against *Tetrahymena pyriformis,* log(IGC50), mmol/ml.



### Data curation improves the true accuracy (up or down!) of QSAR models

significantly modified the overall prediction accuracy obtained by models trained with standardized ( $R_{ext}^2 \sim 0.5$ ) vs. non-standardized ( $R_{ext}^2 < 0$ ) compounds.

Results show that <u>even small differences in structure representation</u> <u>can lead to significant errors</u>, and even robust and inherently predictive models can fail on non-curated external validation sets.

Artemenko, Muratov et al. J. SAR QSAR 2011, Accepted.



## Case Study I: A Two-step Hierarchical QSAR Modeling Workflow for Predicting *in vivo* Chemical Toxicity\*

\*Zhu, Richard, Rusyn, Wright, et al, EHP, 2009, 117(8):1257-64

## Focusing on a small subset ToxCast<sup>™</sup> data: Chronic Mouse Toxicity

- Continuity (overlaps with previous ToxRefDB data)
- Manageable (has only 7 *in-vivo* assays)
- 3 assays with the highest fraction of actives chosen for initial studies

CHR\_Mouse\_LiverProliferativeLesions (87 actives) CHR\_Mouse\_LiverTumors (68 actives) CHR\_Mouse\_Tumorigen (88 actives)

 1 composite endpoint: CHR\_Mouse\_Liver\_tox (110 actives)



 Binary classification QSAR for "baseline" (II & III) vs. off-line (I & IV) using chemical descriptors only



### External prediction workflow



### Comparison of <u>External</u> Prediction Accuración conventional vs. hierarchical QSAR Models





### Case Study 2. <u>Hybrid chemical-biological</u> (short-term assays) descriptors\*

Zhu et al, EHP, 2008, (116): 506-513 Sedykh A, Zhu H, Tang H, Zhang L, Richard A, Rusyn I, Tropsha A. Use of in vitro HTS-derived concentration-response data as biological descriptors improves the accuracy of QSAR models of in vivo toxicity. EHP, 2011, 119(3):364-70.

#### QFAR modeling using hybrid chemical-biological descriptors

C M Ρ Ν S



Quantitative Structure Property Relationships

0.613 Ρ 0.380 -0.222 R 0.708  $\bigcirc$ 1.146 Ρ 0.491 0.301 E 0.141 R 0.956 0.256 0.799 1.195

1.005

### QSAR TABLE - CHEMICAL DESCRIPTORS



### IN VITRO DOSE-RESPONSE DATA IMPROVE THE PREDICTIVE POWER OF QSAR MODELS

Case Study: prediction of *in vivo* toxicity (rat  $LD_{50}$ )

- •1408 substances
- •382 chemical structure descriptors (Dragon v5.5)
- 13 in vitro NCGC cell viability assays \* :
  - o qHTS (quantitative HTS) data
  - 14 test concentrations: 0.6nm .. 92.2μm

May yield up to 13x14 = 182 *in vitro* qHTS descriptors, but the issue of data noise becomes important.

\*Inglese J., Douglas S. A. et al. *PNAS*, **2006**, v103(31), p11473



### EXAMPLES OF QHTS CONCENTRATION-RESPONSE CURVES AND THEIR NOISE-FILTERING TRANSFORMATIONS





### NOISE-REDUCTION APPLIED TO CONCENTRATION-RESPONSE DATA IMPROVES PREDICTION ACCURACY OF HYBRID MODELS.

	%	Chemical descriptors only	Hybrid descriptors (Original, binary)	Hybrid descriptors (THR=15%)	
<i>k</i> inin models	Sensitivity	68±8	63±9	76±5	
	Specificity	85±4	86±4	87±2	
	CCR	76 ±5 *	74 ±5	82 ±3	

	CCR	78 ±4 *	77 ±5	82 ±5
models	Specificity	82±7	87±4	86±3
Random Forest (RF)	Sensitivity	74±9	66±8	77±10

### BIOLOGICAL DESCRIPTOR ANALYSIS

Relative contribution of qHTS dose-response descriptors to QSAR models varies between cell-lines





### HYBRID QSAR MODELS HAVE HIGHER PREDICTIVE POWER THAN COMMERCIAL SOFTWARE TOPKAT

%	ТОРКАТ	Chemical descriptors only		Hybrid descriptors (Original)		Hybrid descriptors (THR=15%)	
		kNN	RF	kNN	RF	kNN	RF
Sensitivity	0.45	0.73	0.73	0.55	0.82	0.91	0.91
Specificity	0.93	0.78	0.80	0.85	0.78	0.85	0.83
CCR	0.69 *	0.75	0.77	0.70	0.80	0.88	0.87

Results are shown for 52 compounds in our external validation sets, which were also absent in the TOPKAT training set.

\*TOPKAT model was significantly different (p < 0.05) from all other models by the permutation test (10,000 times).

### Case Study 3: QSAR + Toxicogenomics\*

### Rationale

- Hybrid models of chemical and biological descriptors (in vitro assays, dose response curves) improved prediction of toxicity (Zhu 2008, Sedykh 2010)
- Toxicogenomics has proven predictive and interpretative value

### Hypothesis

• Adding toxicogenomics data as biological descriptors to hybrid QSAR models may improve predictivity.

#### About the data set

 The Toxicogenomics Project (TGP), Japan collected microarray profiles of 150 common drugs covering various hepatotoxic modes of action <u>http://toxico.nibio.go.jp/</u>

\*Low et al. Predicting Drug-induced Hepatotoxicity Using QSAR and Toxicogenomics Approaches. *Chem. Res. Toxicol.*, **2011**, 24(8):1251-62.

### Data available from TGP

http://toxico.nibio.go.jp/

Table 2. The st	andard study protocol in TGP	
In vivo		
Animal	Sprague–Dawley rat (6 wk old, N = 5 for each group)	
Vehicle	0.5% Methylcellulose or corn oil (oral dose) Saline or 5% glucose solution (intravenous dose)	
Dose	Low, middle, and bigh (mainly 1:3:10)	
Route	Oral (intravenous in a few cases)	
Sacrifice	3, 6, 9, and 21 after a single administration 24 h after the last dose of repeated administration for 3, 7, 14 and 28 days	
Sampling	Liver, kidney, and plasma	
Microarray analysis	Affymetrix GeneChip ( $N = 3$ for each group)	
Items	histopathology: hver and kidney	
examined	Body weight, organ weight (liver and	
	kidney), food consumption, hematology,	
	and b	

Source: Uehara 2010

Gene expression data used: taken from rats treated with the highest dose after 24h.



127 drugs



### Chemical curation and feature selection



# Predicting hepatotoxicity from chemical descriptors and/or gene expression



#### 50% of 40 closest pairs in chemical space are activity cliffs



#### 33% of 40 closest pairs in gene expression space are activity cliffs

CID	Compounds	Distance	Toxicity	
60,72	iproniazid, erythromycin ethylsuccinate	5.75	0, 0	A Report
42,45	griseofulvin, ketoconazole	6.18	0, 0	
34,39	cimetidine, labetalol	6.30	0, 0	
<mark>71,82</mark>	nicotinic acid, enalapril	6.50	0, 0	
70,125	chlorpropamide, gentamicin	6.56	1, 0	
76,81	ranitidine, captopril	6.62	0,0	
35,36	haloperidol, fluphenazine	6.64	0, 0	
18,115	diclofenac, ethanol	6.67	0, 0	N OH
46,48	tetracycline, ciprofloxacin	6.69	0, 0	
9,88	naphthyl isothiocyanate, triamterene	6.72	1, 0	001 000 🛩
25,65	phenytoin, hydroxyzine	6.75	0, 0	
51,64	metformin, amitriptyline	6.75	0, 0	HO
7,91	methotrexate, tiopronin	6.77	0, 0	
102,107	triazolam, chlormadinone	6.81	0, 0	
55,73	methimazole, ethambutol	7.00	0, 1	
66,74	ibuprofen, mefenamic acid	7.20	0, 0	
97,103	sulpiride, clomipramine	7.34	0, 1	HN - (
2,37	isoniazid, thioridazine	7.43	0, 0	Pairwise distance in genetic space
58,63	tacrine, imipramine	7.46	0, 0	
61,105	chloramphenicol, terbinafine	7.58	1, 1	$\langle O \rangle \langle O \rangle   O     O     A$
47,98	lomustine, acarbose	7.68	1, 0	
27,126	allopurinol, vancomycin	7.73	0, 0	
44,68	perhexiline, furosemide	7.85	0, 0	
19,22	nitrofurantoin, diazepam	7.85	0, 1	
6,110	clofibrate, benziodarone	7.91	0, 1	
40,104	methyltestosterone, trimethadione	7.92	1, 1	Yobs Ypred,gene Ypred,chem
				°
3,120	carbon tetrachloride, cyclosporine A	8.56	1,1	
87,113	sulindac, ethionamide	20.49	1,1	

# A novel consensus kNN approach: learning from nearest chemical <u>and</u> toxicogenomic neighbors (k=5)



### Consensus vs. Hybrid Descriptors kNN

carbamazenine	Neighbors	Space defined by	k	Sens	Spec	CCR
phenylbutazone phenytoin	Chemical neighbors only	304 Dragon desc	7	0.56	0.65	0.59
flutamide	Toxicogenomic neighbors only	85 transcripts	5	0.75	0.79	0.74
	Hybrid neighbors	304 Dragon desc 85 transcripts TOGETHER	6	0.69	0.77	0.71
pemoline phenylbutazone phenytoin flutamide ethinylestradiol	Chemical neighbors & toxicogenomic	304 Dragon desc, 85 transcripts SEPARATELY	5	0.76 Results	0.80 after 5-fr	<b>0.78</b>
¥	neighbors	5-1	fold ext	CV used	to select	optimal k

### Conclusions

- Methodology:
  - Data curation is critical!
  - consensus (collaborative!) prediction using all acceptable models
  - <u>outcome</u>: decision support tools in selecting future experimental screening sets
- The highest accuracy is achieved by models that employ both chemical and biological descriptors of compounds
  - Integration of cheminformatcs and bioinformatics: predictive model s of selected endpoints using integrated <u>short term</u> biological profiles (biodescriptors ) <u>and</u> chemical descriptors for compound subsets
  - New computational approaches (e.g., hybrid and hierarchical QSAR)
  - Interpretation of <u>significant</u> chemical and biological descriptors

### Cheminformatics for masses



Please cite this website using the following URL: http://chembench.mml.unc.edu

The Carolina Cheminformatics Workbench (Chembench) is developed by the Carolina Exploratory Center for Cheminformatics Research (CECCR) with the support of the <u>National</u> <u>Institutes of Health</u> (grants <u>P20HG003898</u> and <u>R01GM066940</u>) and the Environmental Protection



### Acknowledgements

Research Professors Clark Jeffries, Alexander Golbraikh, <u>Hao Zhu</u>, Denis Fourches

#### Collaborators

<u>UNC</u>: I. Rusyn, F. Wright <u>EPA</u>: T. Martin, D. Young A. Richard, R. Judson, D. Dix, R. Kavlock Graduate Research



#### Postdoctoral Fellows Aleks Sedykh, Ashutosh Tripathy

Visiting Research Scientist Eugene Muratov

Adjunct Members Weifan Zheng, Shubin Liu MAJOR FUNDING NIH - R01-GM66940 - R01-GM068665 EPA (STAR awards) - RD832720 - RD833825 - RD834999 Assistants Nancy Baker, Hao Tang, Jui-Hua Hsieh, Tanarat Kietsakorn, Tong Ying Wu, <u>Liying Zhang</u>, Guiyu Zhao, Andrew Fant, Stephen Bush, Yen Low

> Research Programmer Theo Walker

System Administrator Mihir Shah