

www.epa.gov

# EDSP Prioritization: Collaborative Estrogen Receptor Activity Prediction Project (CERAPP)

**Chemical structure curation:** 

Remove duplicates

The main steps of this KNIME workflow were:

• Check the validity of the molecular file format

Retrieve any missing structures from web-services

Remove salts and counter ions and fulfill valence

• Standardize stereo-isomers and tautomers

• Remove the inorganic and metallo-organic structures

Kamel Mansouri<sup>1,2</sup>, Jayaram Kancherla<sup>1,2</sup>, Ann Richard<sup>2</sup> and Richard Judson<sup>2</sup>

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

<sup>2</sup>National Center for Computational Toxicology, Office of Research and Development, U.S. Environmental Protection Agency, RTP, NC, USA

### **Abstract**

Humans are potentially exposed to tens of thousands of man-made chemicals in the environment. It is well known that some environmental chemicals mimic natural hormones and thus have the potential to be endocrine disruptors. Most of these environmental chemicals have never been tested for their ability to disrupt the endocrine system, in particular, their ability to interact with the estrogen receptor. EPA needs tools to prioritize thousands of chemicals, for instance in the Endocrine Disruptor Screening Program (EDSP). This project was intended to be a demonstration of the use of predictive computational models on HTS data including ToxCast and Tox21 assays to prioritize a large chemical universe of 32464 unique structures for one specific molecular target – the estrogen receptor. CERAPP combined multiple computational models for prediction of estrogen receptor activity, and used the predicted results to build a unique consensus model. Models were developed in collaboration between 17 groups in the U.S. and Europe and applied to predict the common set of chemicals. Structure-based techniques such as docking and several QSAR modeling approaches were employed, mostly using a common training set of 1677 compounds provided by U.S. EPA, to build a total of 42 classification models and 8 regression models for binding, agonist and antagonist activity. All predictions were evaluated on ToxCast data and on an external validation set collected from the literature. In order to overcome the limitations of single models, a consensus was built weighting models based on their prediction accuracy scores (including sensitivity and specificity against training and external sets). Individual model scores ranged from 0.69 to 0.85, showing high prediction reliabilities. The final consensus predicted 4001 chemicals as actives to be considered as high priority for further testing and 6742 as suspicious chemicals

### **Project planning**

#### **Project goal**

CERAPP was to combine multiple computational models for prediction of estrogen receptor activity into a unique consensus model in order to prioritize a large chemical of 32464 structures.

0.							
Steps	Tasks						
1: Structures curation	<ul> <li>Extract chemical structures from regulatory sources</li> <li>Design and document a workflow for structure cleaning</li> <li>Deliver the QSAR-ready training set and prediction set to all participants</li> </ul>						
<b>2:</b> Experimental data preparation	<ul> <li>Collect and clean experimental data for the evaluation set</li> <li>Define the strategy used to evaluate models predictions</li> </ul>						
3: Modeling & predictions	<ul> <li>Trainor refine the models using the training set</li> <li>Compiling of predictions and applicability domains for evaluation</li> </ul>						
<b>4:</b> Model evaluation	<ul><li>Analyze the training and evaluation sets for consistency</li><li>Evaluate the predictions of each model separately</li></ul>						
<b>5:</b> Consensus strategy	<ul> <li>Define and calculate scores for each model based on the evaluation step</li> <li>Define a weighting scheme from the scores</li> </ul>						
<b>6:</b> Consensus modeling & validation	<ul> <li>Combine the individual predictions based on the weighting scheme and create a consensus prediction</li> <li>Validate the consensus model using a new external dataset</li> </ul>						

### Participants:

**DTU/food:** Technical University of Denmark/ National

**EPA/NCCT:** U.S. Environmental Protection Agency National Center for Computational Toxicology

FDA/NCTR/DBB: U.S. Food and Drug Administration/ National Center for Toxicological Research/Division of **Bioinformatics and Biostatistics** 

FDA/NCTR/DSB: U.S. Food and Drug Administration/ National Center for Toxicological Research/Division of Systems Biology

Helmholtz/ISB: Helmholtz Zentrum Muenchen/ **Institute of Structural Biology** 

**ILS&EPA/NCCT:** ILS Inc & EPA/NCCT

IRCSS: Istituto di Ricerche Farmacologiche "Mario Negri"

**IRC Ispra:** Joint Research Centre of the European Commission, Ispra.

**LockheedMartin&EPA:** Lockheed Martin IS&GS/ **High Performance Computing** 

NIH/NCATS: National Institutes of Health/ National Center for Advancing Translational Sciences

NIH/NCI: National Institutes of Health/ National Cancer Institute

**RIFM:** Research Institute for Fragrance Materials, Inc **UMEA/Chemistry:** University of UMEA/ Chemistry

department **UNC/MML:** University of North Carolina/ Laboratory for Molecular Modeling

**UniBA/Pharma:** University of Bari/ Department of

**UNIMIB/Michem:** University of Milano-Bicocca/ Milano Chemometrics and QSAR Research Group

**UNISTRA/Infochim:** University of Strasbourg/

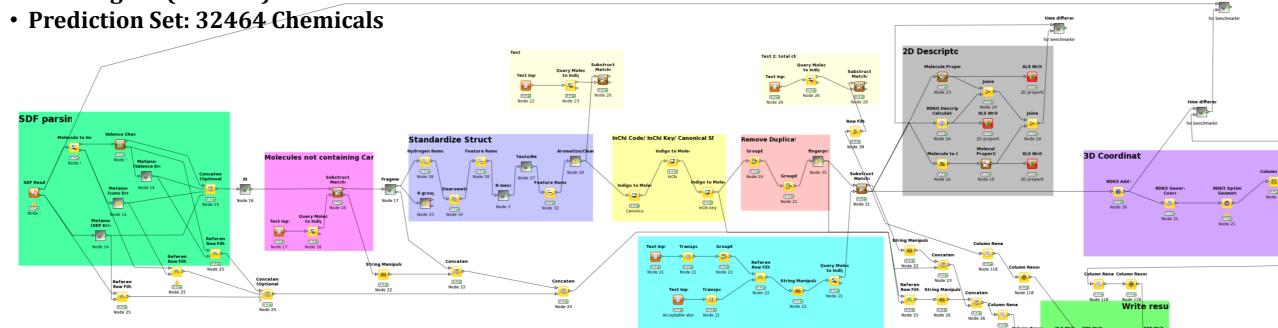
## Data preparation

#### **Sources of chemicals:**

- EDSP Universe (10K)
- Chemicals with known use (40K) (CPCat & ACToR)
- Canadian Domestic Substances List (DSL) (23K)
- EPA DSSTox structures of EPA/FDA interest (15K)
- ToxCast and Tox21 (In vitro ER data) (8K)

#### **Sets of unique chemical structures:**

- Training set (ToxCast): 1677 Chemicals



### **Models and evaluation**

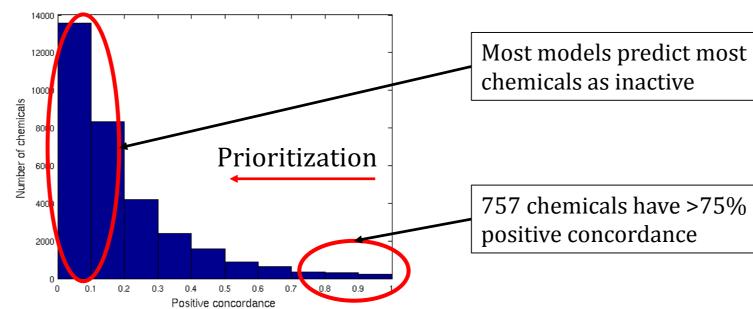
#### Categorical models:

Binding: **22 models** Agonists: 11 models

Antagonists: 9 models Continuous models:

> Binding: **3 models** Agonists: 3 models Antagonists: **2 models**

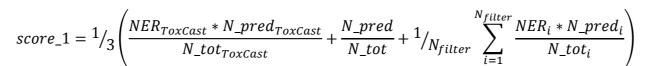
### Concordance of binding models on the active compounds of the prediction set.



### Evaluation set for binary categorical models

### Evaluation set for continuous models

Lvaidatioi	l dot for billa	ry datagorida								
	Active	Inactive	Total		Inactive	V. Weak	Weak	Moderate	Strong	Total
Binding	1982	5301	7283	Binding	5042	685	894	72	77	6770
Agonist	350	5969	6319	Agonist	5892	19	179	31	42	6163
ntagonist	284	6255	6539	Antagonist	6221	76	188	10	10	6505
otal	2617	7024	7522	Total	6892	702	9016	81	93	7253
-		=	-		=	-		-		



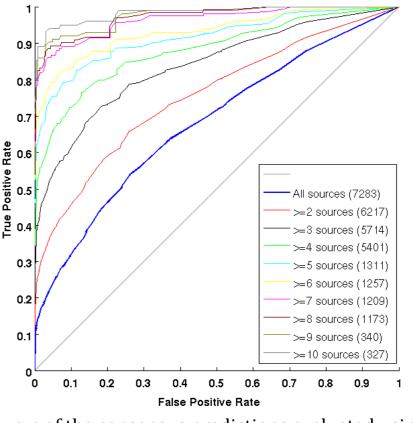
theed artin_ -	NIH_ NCATS	NIH NCI GUSAR	NIH_ NCI_ PASS	RIFM	UMEA	UNC_ MML	UNIBA	UNIMIB_ Michem	UNISTRA
.66	0.65	0.69	0.66	0.65	0.70	0.65	0.73	0.71	0.60
00	99.97	94.87	97.4	100	99.89	100	46.75	36.44	100
75	0.77	0.00	0.70	0.70	0.00	0.00	0.40	0.00	0.00

	DTU 	NCCT	DBB	DSB	OCHEM	EPA	CART	Ruleset	_IRC _Ispra		NIH_ NCATS	NIH_ NCI GUSAR	NIH_ NCI_ PASS	KIFM	UMEA	MML	UNIBA	Michem	UNISTR.
BA	0.78	0.69	0.68	0.66	0.72	0.75	0.75	0.62	0.67	0.66	0.65	0.69	0.66	0.65	0.70	0.65	0.73	0.71	0.60
% pred	49.48	100	100	0.62	96.47	99.9	89.20	94.88	100	100	99.97	94.87	97.4	100	99.89	100	46.75	36.44	100
Score_1	0.43	0.82	0.87	-	0.83	0.82	0.78	0.75	0.77	0.75	0.77	0.88	0.78	0.78	0.82	0.80	0.40	0.32	0.80
Score_2	0.80	0.78	0.84	0.69	0.80	0.79	0.77	0.77	0.74	0.75	0.67	0.84	0.76	0.69	0.76	0.73	0.80	0.85	0.73

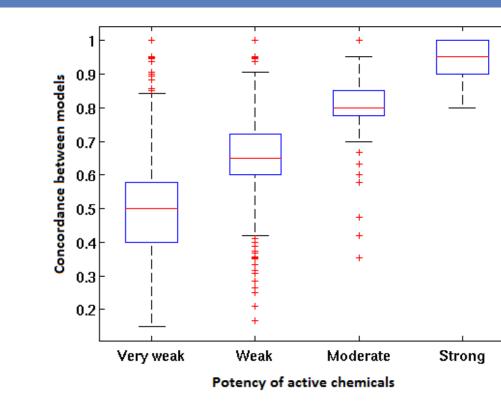
BA: balanced accuracy inside the applicability domain for the unambiguous compounds. % pred: percentage of the predicted within 32464 chemicals.

Kamel Mansouri I mansouri.kamel@epa.gov I 919-541-0545

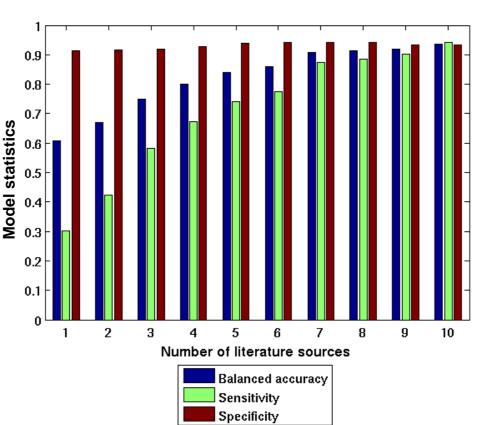
### Consensus model



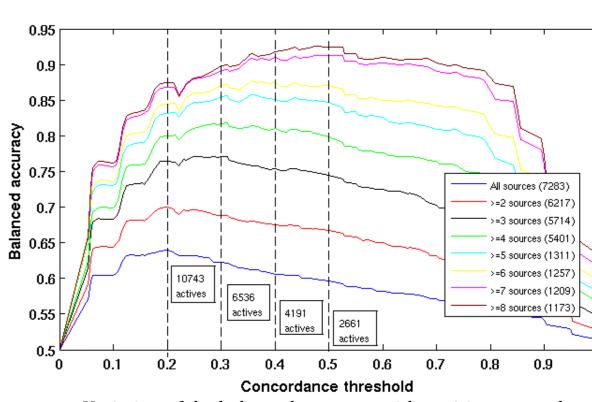
variable number of literature sources



Box plot of the correlation between the positive concordance of the categorical models and the potency level predicted by the continuou models, positive classes predicted by the consensus model



Plot of the consensus accuracy showing the importance of using multiple literature sources



Variation of the balanced accuracy with positive concordance thresholds at different numbers of literature sources and corresponding number of actives

#### Statistics of the consensus model

	ToxC	ast data		ure data 7283)
Observed\Predicted	Actives	Inactives	Actives	Inactives
Actives	83	6	597	1385
Inactives	40	1400	463	4838

	ToxCast data	Literature data (All: 7283)	Literature data (>6 sources: 1209)
Sensitivity	0.93	0.30	0.87
Specificity	0.97	0.91	0.94
Balanced accuracy	0.95	0.61	0.91

## Conclusions

- Successful collaboration of 17 international research groups resulted in consensus predictions.
- Most individual models performed very well on both ToxCast and literature data.
- Model accuracy differences between ToxCast and literature may be due to noise in the literature data.
- A total number of 4001 out of 32464 chemicals are predicted by the consensus as active binders. For prioritization purposes, 6742 additional chemicals (down to 0.2 positive concordance) could be considered as suspicious.
- The consensus model predictions correlate better with literature data from multiple sources, likely due to greater noise in data with unique or low number of sources.
- Consensus prediction results are being used in the EDSP program. See http://actor.epa.gov/edsp21.

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. Environmental Protection Agency.** 

### U.S. Environmental Protection Agency ChemoInformatique Office of Research and Development