# <u>CASE STUDY</u>

# **Development** and Validation of QSAR **Models for (B)TAZs and Fragrances**

#### Simona Kovarich

**QSAR Research Unit** 

in Environmental Chemistry and Ecotoxicology

www.qsar.it

University of Insubria - Varese

e-mail: simona.kovarich@uninsubria.it

CADASTER Workshop on the development and application of QSAR models in REACH 7-9 October 2012, Munich



→ QSAR models were developd (WP3) according to the available experimental data (limited)

Prioritization for experimental test (WP2)

Integration of QSAR predictions for risk assessment (WP4)



# QSAR models for aquatic toxicity of Triazoles and Benzo-triazoles (B)TAZs

## Triazoles and Benzo-triazoles (B)TAZs

- Synthetic compounds, structurally heterogeneous
- Presence of an aromatic heterocycle (2 C + 3N atoms)

Wide applications and use: phytosanitary products, pharmaceuticals, cleaning agents for textiles, UV stabilizer for plastics, de-icing agents, etc.

- High stability and environmental persistance
- High water solubility (contamination of water compartment)





#### (B)TAZs - Introduction



- Experimental data for the basic ecotoxicological endpoint available for several (B)TAZs used as pesticides.
- Focus on toxicity of (B)TAZs in the aquatic environment.
- Development of QSAR models for the prediction of acute toxicity of (B)
  TAZs to aquatic organisms

FISH (Oncorhynchus mykiss)

**ZOOPLANKTON** (Daphnia magna)



ALGAE (Pseudokirchneriella subcapitata)







HelmholtzZentrum münchen Deutsches Forschungszentrum für Gesundheit und Umwelt

#### (B)TAZs - Materials and methods

#### DATA SET

Experimental data collected from the Footprint database (PPDB)

Data collected for (B)TAZs and other azo-aromatic compounds (e.g. diazines and triazines)

#### **ENDPOINTS**



The series	ALGAE: <i>Pseudokirchneriella subcapitata</i> log 1/EC <sub>50</sub> 72 h (growth inhibition)	N = 35
	ZOOPLANKTON: <i>Daphnia magna</i> log 1/EC <sub>50</sub> 48 h (immobilization)	N = 97
	FISH: <i>Oncorhynchus mykiss</i> log 1/LC <sub>₅0</sub> 96 h	N = 75

	UI	IVL	LnU	HMGU	IDEA
Endpoint	EC <sub>50</sub> algae EC <sub>50</sub> Daphnia LC <sub>50</sub> fish	EC <sub>50</sub> algae EC <sub>50</sub> <i>Daphnia</i> LC <sub>50</sub> fish	EC <sub>50</sub> algae EC <sub>50</sub> <i>Daphnia</i> LC <sub>50</sub> fish	EC <sub>50</sub> algae EC <sub>50</sub> <i>Daphnia</i> LC <sub>50</sub> fish	EC <sub>50</sub> Daphnia LC <sub>50</sub> fish
Algorithm	MLR-OLS	PLSR	PLSR and BLASSO-PLS	kNN, ASNN, FSMLR, PLS, MLRA, SVM	MLR-OLS
Molecular Descriptors	Dragon 5.5, PaDEL, CADASTER (1D-2D)	Dragon 6 (1D-2D-3D)	Dragon 6 (1D-2D-3D)	CADASTER (1D-2D-3D)	Dragon 5.5 (1D-2D)
Applicability Domain	Leverage	DModX	Leverage on PLS latent variables	STD of ASNN	Leverage
Validation	Internal (R <sup>2</sup> , Q <sup>2</sup> <sub>LOO</sub> , Y-sc etc) External (Q <sup>2</sup> <sub>ext</sub> F1- F2-F3, CCC etc)	Internal (R <sup>2</sup> ,Q <sup>2</sup> , RMSEE) External (RMSEP)	Internal (R <sup>2</sup> , RMSE) External (Q <sup>2</sup> <sub>ext</sub> , RMSE)	Internal and External (R <sup>2</sup> , Q <sup>2</sup> , RMSE and MAE)	Internal (R <sup>2</sup> , Q <sup>2</sup> <sub>LOO</sub> , Y-sc etc) External (Q <sup>2</sup> <sub>ext</sub> F1-F2-F3, CCC)

(B)TAZs - Materials and methods



## **Structures and Descriptors**

Molecular structures were drawn and minimized by the semi-empirical method AM1 (HYPERCHEM software), and converted into SMILES and MOL (Open Babel)

Mono- and bi-dimensional descriptors were calculated using Dragon (v. 5.5), PADEL-Descriptor (v. 2.13) and QSPR-Thesaurus (CADASTER on-line platform)

## **Algorithm**

- Multiple linear regression (MLR) performed by Ordinary Least Squares (OLS) method.
- Variable selection by Genetic Algorithm (GA).





CAS	AMW	Sv	Ss	M∨	Me	Ms
000050-29-3	12.66	21.69	45.81	0.77	1.03	2.41
000050-30-6	12.73	11.22	33.89	0.75	1.06	3.08
000050-31-7	15.03	11.92	37.67	0.79	1.09	3.14
000050-32-8	7.89	23.59	37.33	0.74	0.98	1.87
000051-28-5	10.83	11.14	49	0.66	1.1	3.77
000051-44-5	12.73	11.22	33.89	0.75	1.06	3.08
000055-38-9	8.98	19.38	35.81	0.63	1.01	2.24
000055-63-0	10.77	11.67	60.83	0.56	1.14	4.06
000056-23-5	30.76	5	17.69	1	1.21	3.54
000056-38-2	9.1	19.71	49.31	0.62	1.03	2.74
000057-15-8	11.83	9.6	24.83	0.64	1.05	3.1
000057-74-9	17.07	19.79	46.81	0.82	1.07	2.6
000058-89-9	16.16	13.79	32.67	0.77	1.07	2.72
000058-90-2	17.84	11.11	32.78	0.85	1.1	2.98
000059-50-7	8.91	10.6	23.11	0.66	1.01	2.57
000060-29-7	4.94	7.5	10.5	0.5	0.98	2.1
000060-51-5	9.55	14.17	31.97	0.59	1.02	2.66

**QSARINS** 

(B)TAZs - Materials and methods

## **Tools of Validation**

Internal stability verified by R<sup>2</sup>, Q<sup>2</sup><sub>LOO</sub>, Q<sup>2</sup><sub>LMO</sub>, R<sup>2</sup>/Q<sup>2</sup><sub>YS</sub>, etc.

External predictivity verified by different Q<sup>2</sup><sub>EXT</sub> criteria [Q<sup>2</sup><sub>EXT</sub> F1-F2-F3 and CCC]

Prediction sets were obtained by splitting 30% (K-ANN, random by response). An additional external validation set (EV set) was also used in the fish model (18 (B)TAZs).

## **Applicability Domain**

Structural applicability domain verified by the Leverage approach

Interpolated and extrapolated predictions verified by Leverage





## QSARs for *P. subcapitata* EC<sub>50</sub> 72h

Model	Descriptors	R <sup>2</sup>	Q <sup>2</sup> LOO	Q <sup>2</sup> LMO	Q <sup>2</sup> <sub>EXT-Fn</sub>	CCC	S
Dragon	AEigZ, T(NS), SEigv	82%	77%	77%	72-84%	86-87%	PLIT
QSPR- Thesaurus	р1р4-5N, C-C, р5BE	80%	73%	74%	71-83%	86-87%	MOD
PaDEL- Descriptor	AMR, MDEN-22, maxHBa	82%	76%	76%	67-89%	82-91%	ELS

Starting from different pool of descriptors, the GA selected molecular descriptors encoding for similar structural information.

Screening of 369 (B)TAZs without experimental data (ECHA pre-registration list)

Publication in press in Molecular Informatics ...



#### (B)TAZs - Results

## Importance of data curation

Experimental data quality and variability

Y = f(X)

Publication in press in Molecular Informatics



Chemical structures correctness and Input formats for the calculation of descriptors

pEC <sub>50</sub>	PPDB	<b>WP2</b> (PHI)	$\Delta  \text{pEC}_{50}$
Triadimefon	5.16	4.59	0.57
Epoxiconazole	5.44	4.58	0.86

Despite both high quality data and same method used (OECD guidelines), variablity in experimental data can affects the QSAR model.

Specific attention to experimental uncertainty, input and basis of QSAR models, is therefore necessary



Differences in chemical structures generated from SMILES collected from different sources (Pubchem, ChemID plus, Chemspider, etc...)

Differences in molecular descriptors (nH, nAB, nitro groups) calculated using different software and different input files (SMILES, MOL, HIN)

Accurate check of the structures and harmonize the representation of specific groups

#### (B)TAZs - Results



## QSARs for *D. magna* EC<sub>50</sub> 48h

Model	Descriptors	Ν	R <sup>2</sup>	Q <sup>2</sup> LOO
Dragon	TPSA(NO), Aeigm, nCar, nHDon, H-052	97	77%	74%
QSPR- Thesaurus	ALogP	97	75%	73%



## QSARs for *O. mykiss* LC<sub>50</sub> 96h



Model	Descriptors	Ν	R <sup>2</sup>	Q <sup>2</sup> LOO	Q <sup>2</sup> <sub>EXT-Fn</sub>	CCC	
Dragon	CIC1, Mp, H-052, TPSA (tot)	75	79%	76%	85%	92%	EV
PaDEL- Descriptor	VP-1, SHBint2, maxHaaCH	75	76%	73%	71-72%	82%	set



Models applied to predict acute toxicity of > 300 (B)TAZs without experimental data (> 90% included in AD of models)





## **WP3 Consensus predictions**

#### Comparison of individual WP3 Models

comparable predictions for (B)TAZs included in the AD of all the models.

#### WP3 Consensus Models

combination of predictions for 386 (B)TAZs obtained from different WP3 models and approaches, taking into account statistical performances and applicability domains of individual models.

#### Screening of (B)TAZs

consensus predictions (algae, daphnia, fish) analysed by PCA in order to screen the studied (B)TAZs according to their toxicological profile in the aquatic environment.

#### (B)TAZs - Results

#### Screening and prioritization of (B)TAZs



#### TREND of AQUATIC TOXICITY



➢ Robust and externally predictive QSAR models have been developed to predict the acute toxicity of (B)TAZs in Algae, *Daphnia* and Fish, using both commercial and freely available molecular descriptors.

> Models were applied to predict aquatic toxicity of > 350 (B)TAZs without experimental data (ECHA list): interpolated predictions for 90% of compounds.

Importance of data curation for the development of valid QSARs

Importance of the consensus approach to increase reliability of QSAR predictions.

#### **APPLICABILITY OF THE PROPOSED MODELS...**

Acute toxicity on algae, daphnia and fish are the basic endpoints required in REACH for risk assessment of chemicals (→ prediction used in WP4)

Models will be implemented in the CADASTER database: they will be freely available for users to be applied also for regulatory purposes



# CASE STUDY 2

# QSAR models for biodegradation of Fragrances

#### **FRAGRANCES** - Introduction

## Fragrances

- Compounds widely used in cosmetics, toiletries, household and loundry products, products used to scent the air (air freshners and fragranced candles), food additives.
- Structurally highly heterogeneous (22 major structural classes identified, exhibiting a wide range of physico-chemical properties).



Wide use and exposure, but limited information available related to health effects of fragrances. Fragrance formulas considered as trade secrets.

> Known effects: asthma, allergies and migraine headaches, accumulate in adipous tissues (breast milk), some fragrances suspected of being endocrine disruptors.



## **Biodegradation**

✤ It's one of the main processes for removal of chemicals from the environment.

The "benign by design" concept requires information on a compound's biodegradability to be available at an early stage, before synthesis of new chemicals.



✤ Use of *in-silico* techniques (eg. QSAR) for a rational design of new fragrances

### Classification models

 $\mathbf{C} = f(\mathbf{X}_{i})$ 

Quantitative relationship between <u>structure</u>  $(X_i)$  and a <u>qualitative response</u> (C)



#### → Classification method: *k* Nearest Neighbour (*k*-NN)

• searches for the k nearest neighbours of each molecule in the dataset

• compound assigned to the most representative class of the k neighbours  $(1 \le k \le 10)$ 

# → Molecular descriptors: commercial software (Dragon, v. 5.5.) and free on-line (PaDEL-Descriptor 2.12)

## Model Performances and Validation

	Assigne	d Class		External
Real Class	RB	NRB	Accuracy	validation
RB	$\checkmark$	X	% RB	
NRB	X	<b>V</b>	% NRB	predictivity
			% Overall Accurac	cy (OA)

## Applicability Domain

 $\rightarrow$  applicability of the model to new molecules

- Range of descriptors
- > Leverage



## DATA SET

Japanese MITI (Ministry of International Trade and Industry) database

→ ready biodegradability data OECD 301 (~1400 heterogeneous organic molecules)



FRAGRANCES - Data set

Data set: 187 compounds (100 fragrances)



FRAGRANCES - Data set

66

NRB

DATA SET balancing



External validation



> 35 data from RIFM (*Research Institute for* Fragrance Materials)

> 10 data measured within CADASTER - WP2 (PHI, Public Health Institute Maribor)

#### **FRAGRANCES - Results**

## Best models

#### **DRAGON Models**

Model ID	k	RB%	NRB%	<b>OA%</b>			
M1	5	85.7	72.7	79.4			
M2	4	82.9	77.3	80.1			
M3	6	88.6	78.8	83.8			
Consensus		94.3	80.3	87.5			

Training set (n=136)

#### Validation set (n=45)

RB%	NRB%	OA%
72.7	60.9	66.67
68.2	60.9	64.4
72.7	69.6	71.1
72.7	73.9	73.3

#### **PaDEL-Descriptor Models**

Model ID	k	RB%	NRB%	OA%	RB%	NRB%	OA%
M4	7	81.4	78.8	80.1	72.7	73.9	73.3

#### **BioWIN (EPI Suite)**

ENVIRON

NUTED STATES	RB%	NRB%	<b>OA%</b>
TAL PROTECTO	72.8	62.1	67.6

RB%	NRB%	<b>OA%</b>
63.6	78.3	71.1



#### Interpretation of modeling descriptors

	Model ID	Molecular Descriptors	
DRAGON -	M1	nCIC) X0A, nR=Ct, F01(C-O)	
	M2	TI1, Vindex, nCq, H-052, B06(C-O)	
	М3	Sv, Qindex, MAXDP, GGI5, nCq	
PaDEL $\rightarrow$	M4	maxHBa, maxHssNH, maxssssC, WTPT-2	



Robust and externally predictive classification QSARs have been developed for the prediction of ready biodegradability of fragrances.

QSARs based on both commercial (Dragon) and freely available (PaDEL-Descriptor) software.

> Automatic selection of molecular descriptors (GA), starting from hundreds of descriptors, is able to select really relevant descriptors for biodegradation.

#### **APPLICABILITY OF THE PROPOSED MODELS...**

Ready biodegradation is among the basic information concerning the environmental fate of chemicals required for risk assessment.

Screening of numerous fragrances and *a priori* use in the design of new alternative compounds, which are less persistent according to the "green chemistry" philosophy





CAse studies on the Development and Application of in-Silico Techniques for Environmental hazard and Riskassessmen

Novel and predictive QSAR models for specific classes of hazardous emerging pollutants with limited data availability.

**Potential applications:** 

✓ filling data gaps

✓ evaluation of chemicals of interest for regulation.

✓ support tools for environmental risk assessment.

✓ screening pre-synthesis of chemicals





Financial support by the European Union through the project CADASTER (FP7-ENV-2007-212668).

QSAR Research Unit in Environmental Chemistry and Ecotoxicology, University of Insubria, Varese (ITALY)

**Prof. Paola Gramatica Dr. Ester Papa** Dr. Leon van der Wal Dr. Nicola Chirico

Stefano Cassani Lidia Ceriani Alessandra Rizza



# **Thank you for your attention!**