



University of Insubria
Varese-Como

QSAR models for the screening, prediction and refinement of PBT Properties of Contaminants of Emerging Concern

Alessandro Sangion





University of Insubria
Varese-Como 2018

QSAR models for the screening, prediction and refinement of PBT Properties of Contaminants of Emerging Concern

by

Alessandro Sangion

Supervisor: Prof. Ester Papa PhD

Co-supervisor: Prof. Paola Gramatica

Doctoral Thesis in
Chemistry and Environmental Sciences,
XXXI cycle

Title:	QSAR models for the screening, prediction and refinement of PBT Properties of Contaminants of Emerging Concern.
Author:	Alessandro Sangion , QSAR Research Unit in Environmental Chemistry and Ecotoxicology, Department of Theoretical and Applied Sciences (DiSTA). University of Insubria, Varese, Italy. a.sangion@hotmail.it
Supervisor:	Prof. Ester Papa , QSAR Research Unit in Environmental Chemistry and Ecotoxicology, Department of Theoretical and Applied Sciences (DiSTA). University of Insubria, Varese, Italy. ester.papa@uninsubria.it
Co-Supervisor:	Prof. Paola Gramatica , QSAR Research Unit in Environmental Chemistry and Ecotoxicology, Department of Theoretical and Applied Sciences (DiSTA). University of Insubria, Varese, Italy. paola.gramatica@uninsubria.it
Internal opponent:	Prof. Roberta Bettinetti , Department of Theoretical and Applied Sciences (DiSTA). University of Insubria, Varese, Italy. roberta.bettinetti@uninsubria.it
External opponents:	Prof. Melek Türker Saçan , Institute of Environmental Sciences, Bogazici University, Istanbul, Turkey. msacan@boun.edu.tr Prof. Patrik Andersson , Department of Chemistry, Umeå University, Umeå, Sweden. patrik.andersson@chem.umu.se

*Part of the work of this thesis has been carried out as “International Visiting Graduate Student” at the University of Toronto Scarborough (05-08 2017, Toronto, Canada) under the supervision of Prof. Frank Wania and Adj. Prof. Jon Arnot, and as “Guest PhD student” at the *Université Paris Diderot* (Paris 7), “*Laboratoire ITODYS*”, UMR7086 (03-07 2016, Paris, France) under the supervision of Prof. Ester Papa.

Key Words: QSAR, Contaminants of Emerging Concern, PBT screening, EcoToxicity, Bioaccumulation, Biotransformation

Abstract:

Contaminants of Emerging Concern (CEC) are defined as “any synthetic or naturally occurring chemical or any microorganisms that is not routinely monitored but has the potential to enter the environment and cause known or suspected adverse ecological and/or human health effects”. Examples of CEC are human and veterinary pharmaceuticals, personal care products, plasticizers, flame retardants, perfluoroalkyl compounds and nanomaterials. The prompt identification of the adverse effects of CEC is fundamental to plan effective risk management strategies and ensure high level of protection for human health and the environment. In particular, Persistent, Bioaccumulative and Toxic (PBT) compounds are chemicals of very high concern due to their potential hazard and should be readily identified. Quantitative Structure-Activity Relationship (QSAR) models can be used in the preliminary hazard-based priority setting phase and can be applied for the preliminary screening of PBT candidates among CEC. In the last decades few approaches have been used for the direct or indirect assessment of PBT chemicals based on QSAR models. Unfortunately, procedures and requirements of the PBT assessment were based on the information available for traditional “legacy contaminants”. Many CEC do not necessarily fulfill the original PBT criteria and the actual regulatory requirements are not adequate for a correct assessment of some emerging contaminants. Therefore, new integrated modelling approaches are needed to effectively address and refine the PBT assessment of CEC.

The aim of this thesis is to propose a consistent approach based on QSAR models for the evaluation of the intrinsic environmental hazard of CECs in order to facilitate the identification of PBT chemicals. The research activity comprises three main subjects. The first one is focused on the screening of the potential PBT behavior of pharmaceuticals by the QSAR-based consensus approach proposed by Gramatica and colleagues to draft a priority list of the most environmentally hazardous pharmaceuticals. Results demonstrate a high agreement (i.e. 86%) between the different applied QSAR models and, in general, pharmaceuticals ingredients seem not be PBTs. A priority list of 35 pharmaceuticals is proposed. In the second part, *ad-hoc* QSAR models are developed for the refinement of the toxicity assessment. In particular, QSAR models to estimate the acute toxicity of pharmaceuticals in species at different levels of the aquatic trophic chain are developed and applied to prioritize pharmaceuticals. All the proposed models have good fitting performances ($R^2: > 0.75$), are internally robust ($Q^2_{\text{LOO}}: > 0.70$; $Q^2_{\text{LMO}}: > 0.65$) and externally predictive ($\text{CCC}_{\text{EXT}}: > 0.82$, $Q^2_{\text{EXT}}: > 0.68$). A global Aquatic Toxicity Index (ATI), based on Principal Component Analysis (PCA) of various aquatic toxicities is proposed and modelled. This index is able to rank pharmaceuticals according to their toxicity on the whole aquatic environment. 34 out of 35 pharmaceuticals, previously highlighted as potential PBT, have also potential high toxicity on all the studied aquatic organisms. Moreover, interspecies correlation models to extrapolate toxicity from a lower trophic level to higher one are also developed. The third part addresses the bioaccumulation refinement and aims to analyze the importance of metabolic biotransformation in the estimation of

bioaccumulation. QSAR models for the prediction of *in-vivo* whole-body human biotransformation Half-Lives (HLs) are developed for organic chemicals. These QSARs are used to predict the biotransformation potential of pharmaceuticals to refine the bioaccumulation evaluation of the previous PBT assessment. Only 22 pharmaceuticals are predicted as slow-metabolized and have high bioaccumulation potential. Finally, predictions for the biotransformation potential are integrated in a mechanistic mass-balance multimedia environmental fate food-web model to estimate the Biomagnification Factor (BMF) in human in a tiered approach. The tiers progress from conservative assumptions to more realistic ones for chemical properties, biological partitioning and biotransformation in human. The introduction of biotransformation HL strongly affects the calculation of the BMF potential. The 98 % of the analyzed chemicals are not bioaccumulative and the elimination processes related to biotransformation are predominant in the overall bioaccumulation.

Contents

LIST OF PAPERS	III
ABBREVIATIONS	V
CHAPTER 1: INTRODUCTION	1
1.1. OVERVIEW ON CHEMICALS PRODUCTION AND CONSUMPTION	2
1.2. CONTAMINANTS OF EMERGING CONCERN	4
1.3. RISK ASSESSMENT AND RISK MANAGEMENT.....	8
1.4. REACH REGULATION.....	11
1.5. PERSISTENT, BIOACCUMULATIVE AND TOXIC CHEMICALS.....	13
1.6. INTEGRATED TESTING STRATEGIES AND PBT ASSESSMENT.....	18
1.7. LIMITATIONS OF THE PBT ASSESSMENT AND SCREENING REFINEMENT	22
1.7.1 Refinement of the Bioaccumulation Assessment.....	24
1.7.2 Refinement of the Toxicity Assessment	28
1.8. AIMS OF THE THESIS.....	30
CHAPTER 2: QSAR METHODOLOGIES	33
2.1 QUANTITATIVE STRUCTURE ACTIVITY RELATIONSHIP	34
2.2 QSAR PROCEDURE.....	37
2.3 MOLECULAR DESCRIPTORS.....	38
2.4 EXPLORATIVE ANALYSIS: PRINCIPAL COMPONENT ANALYSIS.....	40
2.5 DATA MODELLING.....	43
2.5.1 Dataset splitting	47
2.5.2 Variable selection	48
2.5.3 Multiple Linear Regression and Ordinary Least Squares Regression.....	50
2.5.4 Validation techniques and parameters used for regression models.....	53

2.5.5	Applicability Domain.....	57
CHAPTER 3:	RESULTS AND DISCUSSION	61
3.1	PAPER I: “PBT ASSESSMENT AND PRIORITIZATION OF CONTAMINANTS OF EMERGING CONCERN: PHARMACEUTICALS”	62
3.2	PAPER II: “HAZARD OF PHARMACEUTICALS FOR AQUATIC ENVIRONMENT: PRIORITIZATION BY STRUCTURAL APPROACHES AND PREDICTION OF ECOTOXICITY”	68
3.3	PAPER III: “ECOTOXICITY INTERSPECIES QAAR MODELS FROM DAPHNIA TOXICITY OF PHARMACEUTICALS AND PERSONAL CARE PRODUCTS”	73
3.4	PAPER IV: “QSAR MODELING OF CUMULATIVE ENVIRONMENTAL END-POINTS FOR THE PRIORITIZATION OF HAZARDOUS CHEMICALS”	77
3.5	PAPER V: “DEVELOPMENT OF HUMAN BIOTRANSFORMATION QSARS AND APPLICATION FOR PBT ASSESSMENT REFINEMENT”	78
3.6	PAPER VI: “A TIERED APPROACH FOR SCREENING CHEMICALS FOR BIOMAGNIFICATION POTENTIAL IN AIR-BREATHING ORGANISMS.”	82
3.6.1	Introduction.....	82
3.6.2	Materials and methods	83
3.6.3	Results and discussion	88
3.6.4	Conclusions.....	102
CHAPTER 4:	CONCLUSIONS.....	105
	ACKNOWLEDGEMENTS.....	111
	REFERENCES	113

List of Papers

This thesis is based on the following papers. They will be referred to in the text by their Roman numerals I-VI

- I. **Sangion Alessandro**, Gramatica Paola, “PBT Assessment and Prioritization of Contaminants of Emerging Concern: Pharmaceuticals”, *Environmental Research*, 2016, 147 (5): 297–306. doi.org/10.1016/J.ENVRES.2016.02.021.
- II. **Sangion Alessandro**, Gramatica Paola, “Hazard of Pharmaceuticals for Aquatic Environment: Prioritization by Structural Approaches and Prediction of Ecotoxicity”, *Environment International*, 2016, 95 (10): 131–43. doi.org/10.1016/J.ENVINT.2016.08.008
- III. **Sangion Alessandro**, Gramatica Paola, “Ecotoxicity Interspecies QAAR Models from Daphnia Toxicity of Pharmaceuticals and Personal Care Products”, *SAR and QSAR in Environmental Research*, 2016, 27 (10): 781–98 doi.org/10.1080/1062936X.2016.1233139
- IV. Gramatica Paola, Papa Ester, **Sangion Alessandro**, “QSAR Modeling of Cumulative Environmental End-Points for the Prioritization of Hazardous Chemicals”, *Environmental Science: Processes & Impacts*, 2018, 20 (1): 38–47. doi.org/10.1039/C7EM00519A
- V. Papa Ester, **Sangion Alessandro**, Arnot Jon, Gramatica Paola, “Development of Human Biotransformation QSARs and Application for PBT Assessment Refinement”, *Food and Chemical Toxicology*, 2018, 112 (2): 535–43. doi.org/10.1016/J.FCT.2017.04.016
- VI. **Sangion Alessandro**, Arnot Jon, Papa Ester, A Tiered Approach for Screening Chemicals for Biomagnification Potential in Air-Breathing Organisms. In preparation

The published papers are presented herein with the permission of the copyright holders.

Additional papers published:

- Papa Ester, **Sangion Alessandro**, Taboureau Olivier, Gramatica Paola, “Quantitative Prediction of Rat Hepatotoxicity by Molecular Structure”, *International Journal of Quantitative Structure-Property Relationships*, 2018, 3 (2): 12
- Papa Ester, Arnot Jon, **Sangion Alessandro**, Gramatica Paola, “In Silico Approaches for the Prediction of In Vivo Biotransformation Rates”, In *Advances in QSAR Modeling*, 2017, 425–51. Springer doi.org/10.1007/978-3-319-56850-8_11
- Caruso Enrico, Gariboldi Marzia, **Sangion Alessandro**, Gramatica Paola, Banfi Stefano, “Synthesis, Photodynamic Activity, and Quantitative Structure-Activity Relationship Modelling of a Series of BODIPYs”, *Journal of Photochemistry and Photobiology B: Biology*, 2017, 167: 269–81. doi.org/10.1016/J.JPHOTOBIOL.2017.01.012
- Gramatica Paola, **Sangion Alessandro**, “A Historical Excursus on the Statistical Validation Parameters for QSAR Models: A Clarification Concerning Metrics and Terminology”, *Journal of Chemical Information and Modeling*, 2016, 56 (6): 1127–31. doi.org/10.1021/acs.jcim.6b00088.
- Papa Ester, Doucet Jean Pier, **Sangion Alessandro**, Doucet-Panaye Annick, “Investigation of the Influence of Protein Corona Composition on Gold Nanoparticle Bioactivity Using Machine Learning Approaches”, *SAR and QSAR in Environmental Research*, 2016, 27 (7): 521–38. doi.org/10.1080/1062936X.2016.1197310
- Gramatica Paola, Cassani Stefano, **Sangion Alessandro**, “Aquatic Ecotoxicity of Personal Care Products: QSAR Models and Ranking for Prioritization and Safer Alternatives’ Design.”, *Green Chemistry*, 2016, 18 (16): 4393–4406. doi.org/10.1039/C5GC02818C SELECTED AS COVER PAGE

Additional papers in preparation (September 2018)

- Zarini Daniele, **Sangion Alessandro**, Caruso Enrico, Zucchi Sara, Sterpone Silvia, Ferri Emanuele, Papa Ester, “QSAR models as an alternative and powerful tool to investigate potential e-liquids acute toxicity in ENDS”, In preparation. Target Journal: *Environment International*
- Gramatica Paola, Chirico Nicola, **Sangion Alessandro**, Papa Ester, “QSARINS-Chem standalone version 1.0: a free platform-independent software with updated models for environmental pollutants”, In preparation. Target Journal: *Journal of Computational Chemistry*.

Abbreviations

AD	Applicability Domain
ALARA	As Low As Reasonable Achievable
ANN	Artificial Neural Network
ATI	Aquatic Toxicity Index
B	Bioaccumulative
BCF	Bioconcentration Factor
BMF	Biomagnification Factor
CAS	Chemicals Abstract Service
CCC	Concordance Correlation Coefficient
CEC	Contaminant of Emerging Concern
CLP	Classification, Labelling and Packaging
CMR	Carcinogenic, Mutagenic and Toxic for Reproduction
CSA	Chemical Safety Assessment
CSR	Chemical Safety Report
CV	Cross Validation
DART	Decision Analysis and Ranking Techniques
ECHA	European Chemicals Agency
EDC	Endocrine Disrupting Compound
EE2	Ethinylestradiol
EU	European Union
EV	Explained Variance
GA	Genetic Algorithm
GHS	Global Harmonized System
HL	Half-Life
HOMO	Highest Occupied Molecular Orbital
IFS	Iterative Fragment Selection
IOC	Ionogenic Organic Chemical
ITS	Integrated Testing strategies
IVIVE	<i>in-vitro</i> to <i>in-vivo</i> Extrapolation
JRC	Joint Research Center
LC	Lethal Concentration
LD	Lethal Dose
LMO	Leave More Out
LOEC	Lowest Observed Effect Concentration
LOEL	Lowest Observed Effect Level
LOO	Leave One Out
LRAT	Long Range Atmospheric Transport
LUMO	Lowest Unoccupied Molecular Orbital
MLR	Multiple Linear Regression
NOEC	No Observed Effect Concentration
NOEL	No Observed Effect Level
OECD	Organization for Economic Co-operation and Development

OLS	Ordinary Least Squares
P	Persistent
PAH	Polycyclic Aromatic Hydrocarbons
PBDE	Polybromodiphenyl Ether
PBT	Persistent, Bioaccumulative and Toxic
PC	Principal Component
PCA	Principal Component Analysis
PCB	PolyChlorinated-Biphenyls
PCDD	PolyChlorinated-DibenzoDioxin
PCDF	PolyChlorinated-DibenzoFuran
PEC	Predicted Environmental Concentration
PFOA	PerFluoro-Octanoic Acid
PLS	Partial Least Squares
PMT	Persistent, Mobile and Toxic
PNEC	Predicted No Effect Concentration
POP	Persistent Organic Pollutant
PPCP	Pharmaceutical and Personal Care Product
pp-LFER	poly parameter Linear Free Energy Relationship
PRESS	Predicted Residual Sum of Squares
QMRF	QSAR Model Reporting Format
QPRF	QSAR Prediction Reporting Format
QSAR	Quantitative Structure-Activity Relationship
QSPR	Quantitative Structure-Property Relationship
RAIDAR	Risk Assessment Identification and Ranking
REACH	Registration, Evaluation, Authorization of Chemicals
RMSE	Root Mean Square Error
RSS	Residual Sum of Squares
SVHC	Substance of Very High Concern
T	Toxic
TCEP	Tris(2-ChloroEthyl)Phosphate
TMF	Trophic Magnification Factor
TSS	Total Sum of Squares
UNEP	United Nations Environmental Program
US-EPA	United States-Environmental Protection Agency
USGS	United States Geological Survey
vPvB	very Persistent, very Bioaccumulative
WWTP	Waste Water Treatment Plant

Chapter 1: Introduction

1.1. Overview on chemicals production and consumption

The increase in quality of life and in the living standards in the last century has been made possible through industrial and technological development. Chemistry plays key role in the technological advancements that drive innovation with undeniable benefits in various fields such as human health, food supply and the environment. Chemicals are present in everyday life in any kind of products being an integral part of modern life with over 100'000 different substances in use¹⁻³. The growth of new chemicals discovered or synthesized annually is extraordinarily high and steadily rising with thousands of new chemicals indexed in the Chemical Abstract Service (CAS) on a daily basis⁴⁻⁶. For example, chemical industry is one of the biggest industrial sector in European Union (EU) generating more than 500 billion Euros of sales every year and investing more than 8 billion Euros in research and development⁷.

Chemicals have lot of benefits but may also present risks; many of the chemicals are toxic to human beings and/or harmful to the environment and can cause adverse effects. Throughout their lifetime, organisms are exposed to a variety of chemicals, contained in water, air, food, healthcare products, medicines, cosmetics and other consumer products. Some chemicals can highly damage health and/or pollute the environment. The National Statistical Institutes and Eurostat⁸ estimated an annual average consumption of 353 million tons of chemicals in EU between 2004 and 2016. Among these, 225 million tons were classified as harmful for human health while 130 million tons as substances that might represent hazard for the environment. Of all chemicals consumed in the EU in 2016, 35.4 % were toxic to the environment and 62.2 % were toxic to health (figure 1.1)⁸.

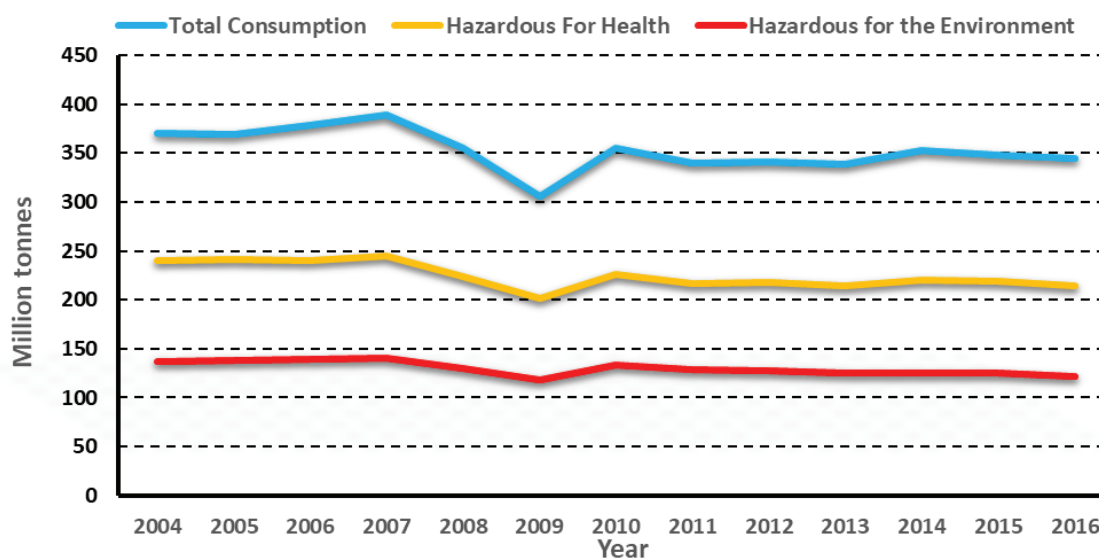


Figure 1.1: Consumption of toxic chemicals, EU-28, 2004-2016. Source: Eurostat⁹

Statistics and classifications reported by Eurostat are based on the hazard statements according to the international Globally Harmonized System (GHS), as implemented in Europe by the Classification, Labelling and Packaging (CLP) Regulation¹⁰. It is important to note that the indicator on the consumption of chemicals hazardous to the environment and to human health is limited under several points of view: production and consumption are not synonymous with exposure; some chemicals are handled in closed systems with high risk management measures, or as intermediate goods in controlled supply chains. Some classification and labelling information for the same substance may appear to be confusing¹¹. Emissions to the environment from production of chemicals exported to outside the EU are accounted by the indicator mentioned above. Finally, the statistics focus on major chemicals with a high production volume and do not account for additional consumption of chemicals whose adverse effects are unknown.

Since 1970s, the attention of authorities, regulatory frameworks, monitoring and quality control analysis has been focused on “legacy contaminants”. This term refers to chemicals that have been left in the environment by sources that are no longer discharging them and are well known for their long-term effects on human and environmental quality. They usually include Persistent Organic Pollutants (POPs),

Introduction

Polycyclic Aromatic Hydrocarbons (PAHs), metalloids and metals such as arsenic and lead^{6,12,13}. Usually the hazard of these chemicals is well characterized, they have been highly regulated and routine techniques have been set up to detect and monitor them in the environment. Moreover, there has been emphasis to replace many of these chemicals with substances that are less toxic to wildlife and humans¹⁴. Unfortunately, in the last decade beside these “legacy contaminants” many other chemicals have been discharged in the environment and have drawn the attention of the scientific community: the so-called Contaminants of Emerging Concern (CEC).

1.2. Contaminants of Emerging Concern

The definition “emerging contaminants” was first used by the US Environmental Protection Agency (US-EPA) in the mid 1900s to refer to chemicals that were newly discovered in the environment, had no regulatory standards and were potentially hazardous to humans and the environment^{15,16}. However, the definition of what is “emerging” is relative and depends on the historical moment. In fact, what was a new contaminant, some decades ago, might not be classified as emerging anymore. Moreover, many contaminants might have been in the environment for long time but concerns regarding their presence have been raised only recently. In these cases, what is emerging is the awareness of the potential hazard of these compounds, not the contaminants themselves. Within a broader context, the definition of CEC is more suitable to define contaminants referred as: i) substances that have been newly introduced on the market and therefore in the environment, ii) contaminants of emerging interest that have been present in the environment for long time but for which the environmental issues were not fully recognized, iii) new concerns raised for legacy contaminants related to new information that allow to understand new aspects of their occurrence, fate and effects^{17,18}. A good definition of CEC is given by the US Geological Survey (USGS): “CEC is any synthetic or naturally occurring chemical or any microorganisms that is not commonly monitored in the environment but has the potential

to enter the environment and cause known or suspected adverse ecological and/or human health effects”¹⁹.

CEC include an extraordinarily broad category of substances that are artificially grouped together even if they have a limited number of common properties^{15,19,20}; table 1.1 report a list with current group of major CEC.

Table 1.1: List of the current groups of major Contaminants of Emerging Concern

Chemical Contaminants
Alkylphenol ethoxylate surfactants (e.g. nonylphenol ethoxylates)
Anticorrosive agents (e.g. benzotriazoles)
Flame retardants (e.g. TCEP*)
Perfluorinated compounds (e.g. PFOA*)
Hormones and endocrine disruptors (e.g. estradiol)
Nanoparticles (e.g. fullerenes)
New classes of pesticides (e.g. neonicotinoids)
Personal Care Products (e.g. oxybenzone)
Pharmaceuticals (e.g. diclofenac)
Stabilizers (e.g. dioxane)

*TCEP, tris(2-chloroethyl)phosphate; PFOA, Perfluoro-octanoic Acid.

The majority of CEC differ from the traditional “legacy contaminants” because they are used also in typical households and their release into the environment is extensively due to their use and disposal by the general population²¹. The main concern regarding domestic sources is that they are uncontrolled, well spread all across the territory and can lead to an ubiquitous and continuous environmental introduction that makes CEC “pseudo-persistent”^{22,23} (i.e. “the continual input of chemicals into the environment can impart a persistence-like quality to those compounds that otherwise possess little inherent chemical stability in the environment because new molecules replenish those that are being removed”²⁴). For example, the direct flushing of unused or expired medicines down the lavatory represents one of the major diffused and unregulated source

Introduction

of pharmaceuticals in the environment²⁵⁻²⁷. CEC are often found in municipal, agricultural and industrial wastewaters and even in drinking water²⁸. Waste Water Treatment Plant (WWTP) effluent are generally loaded with CEC; in such plants, the raw wastewater influents are subjected to a series of physical and biological treatments such as bar screening, sedimentation and microbiological degradation. However, conventional treatments have been designed to remove easily or moderately biodegradable carbon, nitrogen and phosphorous compounds in concentration to the order of mg/l and are not equipped to deal with complex chemicals as many CEC (e.g. pharmaceuticals)^{29,30}. CEC are generally present in low concentrations and their physical-chemical properties and structural features make them unassailable to microbial activities, limiting their removal efficiencies. Beside biodegradation, also sorption on particulate matter and sludge plays a key role in the removal efficiency in WWTP³¹. However, one of the possible final fates of sewage sludge is to be stabilized and converted to biosolid for agriculture applications. In this way CEC residues sorbed in the sludge can be transferred in agricultural soils and enter in the environment³². CEC are increasingly detected in a broad range of environmental compartments (e.g. marine waters, freshwaters, soils, sediments) and in every level of the trophic chains (e.g. fish, amphibians, terrestrial animals) and are considered ubiquitously distributed³³⁻³⁵.

Many CECs are biologically active substances that specifically affect control mechanisms in living organisms (e.g. pharmaceuticals). When released into the environment, this biological activity may adversely affect wildlife (so-called non-target organisms) and impair ecosystem health through a variety of different mechanisms. Acute toxicity and mortality are the most obvious effects cause of concern since many species might be more sensitive to certain contaminants. For instance, the anti-inflammatory drug diclofenac caused the death of millions of vultures across the Indian subcontinent in the past 10 years. Populations of three species of vulture have declined by more than 95%, and all three are now considered critically endangered. The birds died of acute renal failure caused by the active substance. Diclofenac was mainly given

to cattle for relief of pain or inflammation associated with disease or wounds. As the meat of cow is not eaten by people in India, the birds would feast on them and take up all remaining drug residues³⁶. Many CEC can have endocrine functions and interact with the hormone system of animals acting as Endocrine Disrupting Compounds (EDCs). EDCs can interfere with the endocrine system in different ways: behaving like natural hormones, inhibiting natural hormones or interacting with normal hormonal receptors. Hormones are chemical messengers that relay signals from one tissue or organ to another and are involved in a vast array of processes including growth, reproduction, appetite and stress response. The interaction of contaminants with the endocrine system can be deleterious for many species resulting in developmental and reproductive malfunctions, carcinogenesis and feminization/masculinization of aquatic organisms³⁷⁻⁴¹. For instance, the steroid estrogens ethinylestradiol (EE2) caused the feminization of male fish in contaminated river and water bodies⁴². EE2 is the principal ingredient of the contraceptive “pill” and it is also used for estrogen replacement therapy and for enhancing muscle growth in livestock⁴³. Fish are very sensitive to continuous exposure to EE2 and just few nanograms per liter are enough to induce elevated plasma vitellogenin concentration and male feminization^{44,45}. Antimicrobial resistance is another issue related to the presence of certain kind of CEC in the environment. Antibiotic residuals in the environmental compartments lead to the development of resistant bacteria and genes that may change the composition of the microbial community and favor pathogen agents⁴⁶.

Although the improvements of the analytical techniques allow to reduce the detection limits revealing compounds previously not detected, data about environmental hazard of CEC are limited and many long-term effects are still unknown⁴⁷. Information and toxicity data are not yet available for many CEC because important classes of these compounds have not yet studied in detail. CEC are currently not included in routine monitoring programs and remain unregulated even if suspected of causing deleterious ecological effects. Given the large number of CEC suspected to negatively affect the

environment, some type of screening system is needed to monitor, prioritize and facilitate the Risk Assessment and Risk Management of these compounds⁴⁸.

1.3.Risk Assessment and Risk Management

The presence of CEC in the environment represents a scientific, technical and regulatory challenge; environmental behavior and effects of these compounds need to be assessed in order to prevent the rise of negative impacts and to preserve good quality standards.

Risk Assessment methodologies allow to identify hazards and risks related to substances and their uses providing the information required in the Risk Management phase to managed and reduce the possible negative impacts on society and environment. Risk Assessment is a systematic procedure, combining scientific and regulatory principles in order to describe the hazard associated with the human exposure (Health Risk Assessment) or the environment exposure (Environmental Risk Assessment) to chemical substances.

The Risk Assessment-Risk Management (with a focus on the Environmental Risk Assessment) scheme is reported in figure 1.2:

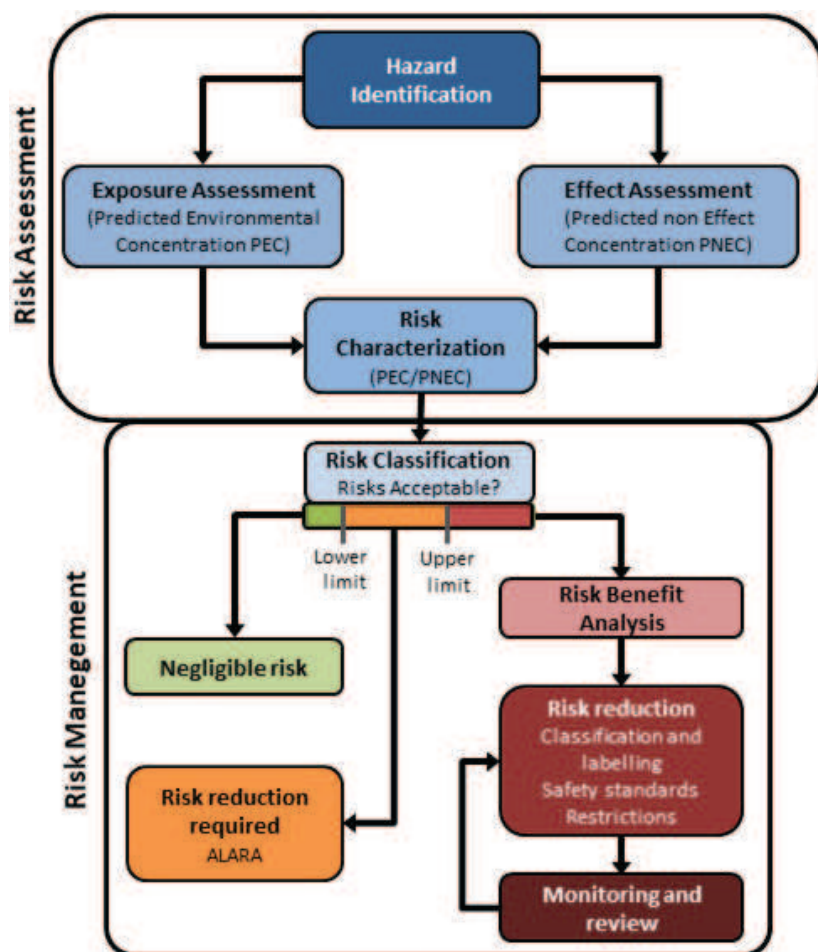


Figure 1.2: Steps in the (Environmental) Risk Assessment and Risk Management processes⁴⁹

The first step of Risk Assessment is the hazard identification, which aims to identify the adverse effects that a substance is inherently able to cause to human or environmental health. This first step encompasses the characterization of the behavior of a chemical in the environment and within the body as well as the collection and analysis of data on health effects that may be produced by the exposure to a chemical. It is followed by the exposure assessment, which is focused on the estimation of the concentrations to which human populations or environmental compartments are exposed. In the Environmental Risk Assessment this can be done by direct measures or by the application of multimedia exposure models that allow to predict the environmental concentration from emission rates, product uses and physical-chemicals properties. The output of this step is the

Introduction

estimation of the predicted environmental concentration (PEC) for each environmental compartment.

In addition to the exposure assessment, the effect assessment is necessary; it is designed to assess the relationships between dose or level of exposure to a substance and the incidence and severity of an adverse effect⁴⁹. For each adverse effect, a dose-response curve is estimated to assess the percentage of the population, of the studied species, which show a specific effect at a specific dose or concentration. The output of the effect assessment is the identification of a threshold under which adverse effects are not expected to occur. In the case of the Environmental Risk Assessment, this level of protection must take into account the high variability of components and factors that characterize the natural environment, like the high number of species (with different responses) and the high number of different effects. For these reasons, the aim of the effect assessment in the Environmental Risk Assessment is the identification of “safe level”, called predicted no-effect concentration (PNEC), which should ensure the protection of the 95% of the natural species taking in account the uncertainty due to the extrapolations from a species to another, from acute to chronic effects, from laboratory to field testing etc⁵⁰⁻⁵². The last step of Risk Assessment is the risk characterization; in this step, the comparison of the results from the previous two stages returns a risk quotient representing the probability that an adverse effect will occur in human population or environmental compartment. In Environmental Risk Assessment, the risk quotient is calculated as PEC/PNEC ratio.

Once the risk associated to a substance has been characterized, a phase of Risk Management is necessary to decide if the calculated risk is negligible or further measures to reduce it are required. According to the risk quotient, two risk levels have been defined: an upper limit that fixes the maximum level above which the risk is considered unacceptable, and a lower limit, in general the 1% of the upper limit, below which the risk is considered negligible. Between the two levels there is a zone where risk reduction is required based on the As Low As Reasonable Achievable (ALARA) principle. In case

of unacceptable risks, additional management measures are required; these measures begin from a risk and benefit analysis, where technical, social, scientific, economic and political factors are considered in order to quantify the possible socio-economic and human health costs and benefits of various scenarios related to different risk reduction options⁴⁹.

Once chosen the best risk reduction approach, a last phase of monitoring and review is required to verify the effectiveness of the risk control measures and to ensure that the safety and quality standards are met⁵⁰.

Risk Assessment and Risk Management are essential tools in the assessment and management of chemicals for a safety use and are integral part of the European regulation in matter control of chemicals, the REACH (i.e Registration, Evaluation, and Authorization of Chemicals) regulation⁵³.

1.4.REACH regulation

Identifying which chemicals pose significant hazard and consequently risks can be a difficult process, as is deciding what should be done when hazards and risks are identified. Regulations as REACH in EU aim to ensure high level of health and environmental protection through the better and earlier identification of the intrinsic properties of hazardous chemicals. The central requirements of REACH are defined by its acronym, i.e Registration, Evaluation, and Authorization of Chemicals, which defines the basis of the approach to chemical safety assessment and management. In more details, REACH requires⁵³:

- *Registration*: all the substances manufactured or imported in the EU in quantitative above 1 ton for year have to be registered by producers/importer to the European Chemicals Agency (ECHA). The timing and amount of information required for registration depend partly on the volume produced/imported. For substances above one ton, a technical dossier (containing information on properties, uses, classification, and guidance on safe

Introduction

use) must be submitted to the authorities. For substances above ten tons, a Chemicals Safety Assessment (CSA) is required and documented in the Chemical Safety Report (CSR). The CSA includes the hazard classification of a substance and the assessment as to whether the substance is persistent, bioaccumulative and toxic (PBT) or very persistent and very bioaccumulative (vPvB). Further, the CSR also describes exposure scenarios, including appropriate Risk Management measures, for all identified uses of dangerous, PBT, vPvB substances.

- *Evaluation:* this process consists of an examination of the data contained in the registration dossiers provided by industry, which undergo a double evaluation: dossier evaluation and substance evaluation. The dossier evaluation includes a check of the compliance of the registration requirements, and a check of the testing proposals, in order to prevent unnecessary testing with vertebrate animals. The substance evaluation includes the investigation of chemicals with potential risks to human health and the environment in order to identify substance of high concern.
- *Authorization:* an authorization is required for the use and marketing of substances of very high concern which include substances classified as carcinogenic, mutagenic and toxic for reproduction (CMR), PBT and vPvB substances and substances of equivalent concern (e.g. EDCs). ECHA will grant an authorization for the use of such substances only if the industry demonstrates that risks associated to the use is adequately controlled, or the socio-economic benefits outweigh the risks and if provides plans for substitution with safer alternatives.
- *Restriction:* restrictions (e.g. prohibition or specific condition for the manufacture, trade or use) can be applied for certain dangerous substances. This procedure provides a safety net to manage risks that have not been adequately addressed by another part of the REACH system.

The REACH regulation moves the burden of risk identification and management from national authorities to the industry. All companies manufacturing or importing chemicals into EU in quantities of one ton or more per year are required to gather information on the properties of chemicals, in order to ensure that they manufacture, place on the market or use substances that do not adversely affect human health or the environment.

Particularly relevant during the registration and evaluation phases are the substances classified as PBT or vPvB. This kind of compounds may be included in the candidate list of Substances of Very High Concern (SVHC) and subjected to the further authorization and restriction procedure. Additionally, in the authorization step, plans for safer alternatives to the most hazardous chemicals must be provided.

1.5.Persistent, Bioaccumulative and Toxic chemicals

PBT substances are substances that are persistent (P), bioaccumulative (B) and toxic (T), while vPvB substances are characterized by a particular high persistence in combination with a high tendency to bioaccumulate which may lead to toxic effects⁵⁴. PBT compounds are priority chemicals due to the potential risk they pose to human's health and ecosystems; they are environmentally relevant because, due to their properties, they can resist to biotic and abiotic degradation and persist unchanged for long time, interacting with human and wildlife, accumulating in living organisms and causing long-term toxic effects⁵⁵. The properties of PBT/vPvB substances lead to an increased uncertainty in the estimation of risk to human health and the environment when applying quantitative Risk Assessment. For PBT a "safe concentration" in the environment cannot be established with sufficient reliability for an acceptable risk to be determined in a quantitative way⁵⁴. Therefore, a separate PBT/vPvB assessment is required to take this specific concern into account. The objective of the PBT/vPvB assessment is to determine if the substance fulfils the criteria given in Annex XIII of the REACH regulation and to characterize the potential emissions of the substance to the different environmental

Introduction

compartments and the likely routes by which humans and the environment are exposed to the substance.

A substance is considered PBT when it fulfils the criteria for all three inherent properties P, B and T:

- *Persistence*: Organic compounds with a slow degradation rate in the environment, resistant to biological and chemical degradation reactions are classified as POPs. In general, an optimal description of the persistence of a compound in the environment requires knowledge of its partitioning in environmental media, its potential for transport, and its degradation half-life (HL) in each compartment. The determination of persistence in any environmental compartments calls for information on the kinetics of both biotic and abiotic processes⁵⁵. Biotic degradation processes involve the interaction of chemicals with organisms and microorganisms able to include the given molecule in their metabolism pathways and modify it. Biodegradation processes are divided in primary degradation and mineralization⁴⁹. The first is referred to all processes that lead to the modification of the parent compounds in metabolites with different properties; primary biodegradation substantially affects the solubility and the environmental mobility of the original compounds, but in some cases, there may be an enhancement of toxicity. Mineralization refers to the complete conversion of a chemical in water, carbon dioxide, inorganic elements and mineral salts, available for a further entry in the natural cycle of elements. On the other hand, abiotic degradation includes all physical and chemical processes such as photolysis, thermolysis, hydrolysis, oxidation, reduction and reaction with other chemicals that lead to the formation of breakdown products from the original compound. Compounds not sensitive to these processes are persistent in the environment⁵⁶. An important side effect of persistence is the possibility of pollutants to be transferred over large distances and reach rural and uncontaminated zones. So, persistent chemicals that exist in

gaseous phase or are associated to droplets or particles can be transported all over the world by air movements on the regional and global scale in a process called Long-Range Atmospheric Transport (LRAT). Here pollutants tend to be accumulated in cold places like glaciers or the poles where the low temperatures favor condensation and air removal processes⁵⁷⁻⁵⁹. The following table 1.2 summarizes the criteria and the relevant information to be used in the PBT/vPvB assessment according to Annex XIII of REACH.

Table 1.2: Persistence criteria according to Annex XIII to the REACH Regulation

Criteria	P	<ul style="list-style-type: none"> - degradation HL in marine water > 60 days; - degradation HL in fresh or estuarine water > 40 days; - degradation HL in marine sediment > 180 days; - degradation HL in fresh or estuarine water sediment > 120 days; - degradation HL in soil > 120 days.
	vP	<ul style="list-style-type: none"> - degradation HL in marine, fresh or estuarine water > 60 days; - degradation HL in marine, fresh or estuarine water sediment > 180 days; - degradation HL in soil > 180 days.
Screening information		<ul style="list-style-type: none"> - Results from tests on ready biodegradation; - Results from other screening tests (e.g. enhanced ready test, tests on inherent biodegradability); - Results obtained from biodegradation computational models; - Other information provided that its suitability and reliability can be reasonably demonstrated.
Assessment information		<ul style="list-style-type: none"> - Results from simulation testing on degradation in surface water; - Results from simulation testing on degradation in soil; - Results from simulation testing on degradation in sediment; - Other information, such as information from field studies or monitoring studies, provided that its suitability and reliability can be reasonably demonstrated.

- *Bioaccumulation:* Bioaccumulation is defined as “the process by which chemicals are taken up by a plant or an animal either directly from exposure to

Introduction

a contaminated medium (e.g. soil, sediment, water) or by eating food containing the chemical”⁶⁰. In other words, the bioaccumulation process is the result of chemical concentration in an organism, due to uptake by all exposure routes including transport across respiratory surfaces, dermal absorption and food uptake⁶¹. The two processes involved in determining bioaccumulation are bioconcentration and biomagnification. Bioconcentration is the process by which a chemical is absorbed by an organism from the ambient environment only through its respiratory and dermal surfaces^{60,62}. Biomagnification is a process in which the thermodynamic activity of the chemical in an organism exceeds that of its diet, it represents the accumulation and transfer of substances via the food-web^{60,62,63}. Persistent pollutants can move from an environmental compartment to another according to their physical-chemical properties and thus can be transferred also to biota⁶¹. Once entered in tissues and organs, lipophilic and poorly metabolized toxicants tend to persist, and it has been widely observed that organisms can achieve high concentrations of certain contaminants in comparison to the environmental concentration of the contaminant itself⁶².

The Bioconcentration Factor (BCF), is the parameter used to quantify the bioaccumulation potential of pollutants⁶⁴; it is determined as the concentration of the parent substance in the whole aquatic organism (C_B) at steady state divided by the mean water concentration of the substance during the exposure period (C_W)⁶²:

$$BCF = \frac{C_B}{C_W} \quad (1.1)$$

For terrestrial organisms, the environment is usually the soil while for aquatic ones the environment is water or sediments. Alternatively, BCF can be seen as the ratio between the uptake rate constant and the depuration rate constant assuming first order kinetics.

Since the experimental determination of these properties is both expensive and time consuming, alternative methods have been developed to correlate physical-

chemical properties or structural parameters to BCF. The octanol-water partition coefficient of a substance ($\log K_{OW}$) is considered a satisfactory surrogate for lipids tissue and it is widely used to estimate BCF⁶⁵. $\log K_{OW}$ has a linear relationship with BCF, thus the higher is the $\log K_{OW}$ the higher will be the chemical's tendency to bioaccumulate in organisms^{66,67}. However, it has been demonstrated that the linear relationship falls down for compounds with $\log K_{OW} > 6$ as well as for substances easily metabolized by the exposed organisms⁶⁸. According to REACH, a chemical fulfils the B or vB criterion when the BCF in aquatic species is higher than 2000 or 5000 respectively⁵⁴. For air-breathing organisms, biomagnification is related to the octanol-air partition coefficient (K_{OA}) in addition to K_{OW} . According to emerging screening criteria⁵⁴ chemicals with $\log K_{OA} \geq 5$ and with a $\log K_{OW} \geq 2$ are considered to have bioaccumulation potential in air-breathing organisms, including humans.

- *Toxicity*: A toxicant is a substance that is harmful to living organisms because of its adverse effects on tissues, organs, or metabolism. Toxic effects depend on the dose or concentration, route and time of exposure as well as on factors depending on individual and species variability. On the basis of the exposure period it is possible to distinguish between acute toxicity, that refers to more or less severe effects resulting from exposure to a single dose of the toxic substance, and chronic toxicity, related to the effects of long-term exposure, often to low repeated doses or concentrations⁶⁹. It is important to note how the same chemical can show both acute and chronic toxicity although with different mechanisms and effects. The toxicity of a compound is evaluated by dose-effect and dose-response assessments that show respectively the relationship between a certain dose or concentration and the effects on an individual and between a certain dose or concentration, and the incidence of a certain effect on the population exposed. On the basis of these assessments, parameters to evaluate toxicity are derived:

Introduction

LC(D)₅₀: Median Lethal Concentration (Dose) - the concentration or dose at which a toxicant is estimated to be lethal to 50% of the organisms of a given population. It is estimated in acute toxicity tests.

EC(D)₅₀: Median Effective Concentration (Dose) - the concentration at which a toxicant is estimated to cause a defined effect in 50% of a given population. It is estimated in both acute and chronic toxicity tests.

NOEL(C): No Observed Effect Level (Concentration)- The maximum dose or ambient concentration an organism can tolerate over a specified period of time without showing any detectable adverse effect, and above which adverse effects become apparent. It is usually calculated in long term exposure tests.

LOEL(C): Lowest Observed Effect Level (Concentration)- The minimum dose or ambient concentration that causes a detectable effect.

A substance fulfils the toxicity criterion if NOEC for marine or freshwater organisms is less than 0.01 mg/L or there are other evidences of chronic toxicity⁵⁴. It can be considered toxic also if it meets the criteria for classification as CMR. Data on short-term acute aquatic toxicity can be used as screening information⁵⁴.

1.6.Integrated Testing Strategies and PBT assessment

Without doubts, under REACH more data for chemicals have been assembled than ever before⁷⁰. Unfortunately, the quality of data is not always adequate and, in many cases, critical information, mainly related to toxicity, are missing^{71,72}. Testing programs are expensive and require an high number of animal lives per chemical while an additional objective of REACH should be the reduction of animal-testing promoting the use of alternative methods^{73,74}.

To meet these needs, intensive efforts have been made to elaborate Integrated Testing Strategies (ITS) which are aimed at speeding up Risk Assessment and Risk Management procedures through a reduction of animal testing and an optimized intelligent use of available information. ITS are integrated approaches comprising both testing and non-

testing methods, such as *in-vitro* testing, computational methods, chemical categories and optimized *in-vivo* test⁷⁵. The aim of the ITS is to combine these different approaches to make the best use of the available information increasing efficiency and cost-effectiveness, reducing time and improving the whole Risk Assessment system. Non-testing *in-silico* methods are fundamental components of the ITS; it was estimated that by applying these methodologies, the needs for animal tests could be reduced to 70% for individual endpoints resulting in savings for more than 800 million Euros and 1.3 million of animals⁷⁶. However, non-testing methodologies will reduce but not completely replace animal testing; the successful application of these ITS components relies on the availability of high-quality experimental data. This is a limiting factor in the development of ITS methodologies that can be overcome only improving the feedbacks exchange between testing and non-testing scientists. The *in-silico* methods can lead the experimental design while good quality experimental data can be used to train and validate high-quality computational models.

Quantitative Structure Activity Relationships (QSAR) are one component of the ITS and can be used in the preliminary hazard-based priority setting phase for the screening of large data sets. QSARs are theoretical mathematical models that can be applied to estimate physical-chemical properties or biological activities of molecules for which no or limited data exist. The main advantage of QSAR models is that they are based on the molecular structure and can be applied to screen virtually every organic chemical as soon as the molecular structure is provided. This allows to detect as early as possible, on the basis of the structural information, potential intrinsic hazards of a chemical and prevent undesirable effects and also allowing the *a priori* plan of safer alternative synthesis according to the “benign by design” approach of the green chemistry^{14,77}. In this way, it is possible to focus time and financial resources on the most hazardous chemicals for further evaluation in more comprehensive assessment phases⁷⁸. Thus, QSARs can be applied for priority setting and for the screening of PBT candidates among CEC⁷⁹. In the last decades different approaches have been applied for the direct or indirect assessment of PBT chemicals based on QSAR models or integrating QSARs in a weight of evidence

Introduction

analysis. For example, an approach based on Principal Component Analysis (PCA) has been applied to select representative PBT chemicals⁸⁰. The US-EPA PBT Profiler⁸¹ is a framework based on the EPI Suite⁸² program developed by the US-EPA for the early assessment of PBT chemicals. It predicts the P, B and T properties separately with specific QSAR models and compares them with the related threshold for each end point^{67,83,83-86}. Muir and Howard^{2,87,88} focused on chemicals with P and B properties combining expert knowledge, models from EPI Suite and P and B criteria from various legislation to select priority chemicals for further analysis in different screening studies on different kind of environmental contaminants, including pharmaceuticals^{2,87,88}. The Risk Assessment Identification And Ranking (RAIDAR) model is a screening level evaluative model proposed by Arnot and Mackay that combines information on partitioning, reactivity, environmental fate and transport, food web bioaccumulation, exposure effect end point, and emission rate in a coherent mass balance framework. This holistic approach allows to combine quantitative information on P, B, T and quantity in a coherent evaluative framework for a screening level metric of risk⁸⁹. The European Joint Research Center (JRC) applied the Decision Analysis and Ranking Techniques (DART) to rank chemicals for environmental safety assessment⁹⁰. Papa and Gramatica proposed the Insubria PBT Index, a multiple linear regression (MLR) QSAR model based on four molecular descriptors⁹¹. This model provides a unique, cumulative index to estimate the potential PBT behavior of a chemical directly from the molecular structure. The Dutch National Institute for Public Health and the Environment (RIVM) developed a P and B score based on overall persistence and BCF integrating experimental data and predictions from EPI Suite⁹² and screened more than 65000 industrial chemicals. Stempel and colleagues proposed and applied a classification system combining measured PBT properties data from multiple sources with estimates for missing property data using methods from EPI Suite and ECOSAR⁹³. Recently Pizzo and collaborators proposed a Weight of Evidence architecture based on multi criteria decision making using a series of existing and new *in-silico* tools for assessing P, B and T properties within the PROMETHEUS project⁹⁴, while Roy and colleagues proposed a

new QSAR model of the Insubria PBT Index based on extended topochemical atom descriptors⁹⁵.

In this context it is important to highlight that a model is a way to simplify and represent the reality, reorganizing the knowledge to describe a certain phenomenon and exploiting it to predict future events. Any model strictly depends on the objects and variables used in the development process, and it does not exist a unique best model which is able to describe a certain event. Often, the better compromise is to use a consensus approach, averaging the prediction of different models that describe the same phenomenon from different points of view. Thus, it seems reasonable that a consensus approach, which can be derived for instance by calculating an average for representative individual models based on different variables, might provide better predicted data than the majority of individual models^{96–99}. Many of the studies reported above^{2,87,88,92,93} performed screening on large sets of chemicals by estimating missing properties using single individual model, often applying the EPI Suite software. This can create a bias in the estimation since, if the model is not trained on a specific class of chemicals (e.g. EPI Suite is not specific for pharmaceuticals¹⁰⁰), predictions might not to be reliable and screening results could be misleading. Moreover, different models are developed under different regulatory frameworks and are compliant with different requirements¹⁰¹. The result is that the identification of PBTs is based on individual properties of chemicals, with threshold values that are not unique across the various legislations. For this reason, the decision on whether a substance fulfills the PBT criteria not only depends on substance properties but also on the framework under which the substance is evaluated¹⁰².

Recently Gramatica and collaborators proposed an integrated approach for the assessment of PBT chemicals based on the consensus of two independent QSAR models (i.e. Insubria PBT Index and US-EPA PBT Profiler)¹⁰³. This approach is based on the concept that PBT are intrinsic hazard properties that can be estimated on the basis of the chemical structure by QSAR models. Moreover, the application of different methods can

Introduction

provide more robust predictions and overcome the limits of the individual models. In a first study Gramatica and colleagues validated their approach comparing results obtained by consensus with large literature screening¹⁰³. The overall agreement was greater than 75% and the method allowed to reach different degree of prioritization identifying as PBT the majority of chemicals of concern for the other studies. This approach was then applied to two classes of CEC (i.e. flame retardants¹⁰⁴ and personal care products¹⁰⁵) to prioritize new contaminants without experimental data. The agreement was high and the approach revealed to be promising for new classes of CEC as for instance pharmaceuticals. This comparative priority setting approach is very promising, highlighting the chemicals that are recognized as potential PBTs by different methods, even if based on different strategies (various cut-off criteria for each property/activity, structural characteristics, in a few cases also empirical evidence), in a consensus approach.

1.7.Limitations of the PBT assessment and Screening refinement

The history of PBT chemicals can be traced to the discovery of chlorine and bromine elements in 1774 and 1826 respectively¹⁰⁶ and with the synthesis of benzene hexachloride by Faraday in 1825. Chlorine and bromine substituents provide PBT molecules with properties of both persistence and biological activity that make them particularly interesting for specific applications. Interests rose in 1930 when advancements in technologies provided industrial possibilities to synthesize chlorine compounds that were valuable pesticides such as DDT (Dichlorodiphenyltrichloroethane). Unfortunately, the properties that made these compounds particularly suitable for agricultural uses turned to be problematic and made them environmentally hazardous. In 1962 the book of Rachel Carson “Silent Spring”¹⁰⁷ focused for the first time the attention of public and scientific community on the negative effects chlorinated pesticides had on the environment, opening the debate about PBTs. In 1995 the United Nations Environmental Program (UNEP) convened an international working group to develop assessment strategy for 12 POPs. This lead to the institution

of the Stockholm Convention in 2001 and of subsequent national and international regulatory frameworks based on the same requirements but with different thresholds^{12,53,108}.

Unfortunately, some of the rationale and thoughts that were applied to the initial criteria used late 1990s have not been carried forward to illuminate the reasons behind the PBT evaluation in newer assessment scheme and, criteria and threshold values have not been revisited in a critical way according to scientific progresses¹⁰⁹. Moreover, new classes of chemicals with physical, chemical and biological properties different from the traditional “legacy contaminants” have become widely used¹⁰⁹.

Little attention has been given to polar compounds that are mobile in the aquatic environment¹¹⁰. Mobility can be qualified as the potential of the substance to be transported to groundwater or far from the site of release⁵³. If such mobile substances are also persistent, they could distribute in water compartments and be of concern for wildlife and water quality¹¹¹. Persistent and mobile organic compounds that are also toxic are referred to as PMT (persistent, mobile and toxic) substances¹¹². PMTs have only recently gained the interest of authorities, and there are activities attempting to establish potential regulatory measures^{110–112}. Generally, the molecules most mobile in water are the ones in which solvation by water is more energetically favorable than sorption to environmental solids. K_{OW} can be used as an approximate indicator of the aquatic mobility of neutral organic compounds. High K_{OW} may indicate high sorption tendency while low K_{OW} may indicate high mobility in the aquatic environment. For ionic compounds, that may exist in a variety of complexes or protonation states depending on pH and on the ions present in the water, the distribution coefficient (D_{OW}) can be used instead of the K_{OW} to account for the concentration of all forms of the compound. However, the assumption that D_{OW} inversely correlates with mobility is still very simplistic. K_{OW} and D_{OW} can be viewed as the approximate indicators of mobility biased towards overestimation due to lack of accounting of additional partitioning

Introduction

between the compound and the environmental media¹¹¹ and can be considered screening parameters.

Many CEC do not necessarily fulfill the original PBT criteria and the actual regulatory requirements are not adequate for a correct assessment of these contaminants, for instance, the BCF is questioned not to be the best parameters to evaluate bioaccumulation potential especially for terrestrial organism^{62,113–115}.

The screening studies are preliminary tools for the identification of PBT properties of the substances on the basis of the molecular structure. Moreover, many of the modeling frameworks developed for the assessment of PBT are based on the principal regulatory requirements and might have some limitation in their application, mainly for CEC. Two critical points that need to be improved with new information are here exposed:

1.7.1 Refinement of the Bioaccumulation Assessment

The quantification of the bioaccumulation potential of chemicals is a fundamental phase in the PBT assessment. When bioaccumulation occurs, organisms (especially at the top of the trophic chain) can be contaminated with long-term effect difficult to predict^{50,62}. The PBT screening performed by the consensus approach by Gramatica and colleagues bases the evaluation of the bioaccumulation potential on the estimation of the BCF¹⁰³. In fact, both the Insubria PBT index and the US-EPA PBT profiler relies on BCF estimations; the PBT Index was trained over the prediction of a QSAR BCF model¹¹⁶ while the US-EPA PBT Profiler calculates the BCF from BCFWIN estimation program¹¹⁷. However, BCF expresses the degree to which bioconcentration occurs. Bioconcentration is defined as the process by which a chemical is absorbed by an organism from the surrounding environment only through its respiratory and dermal surface⁶². It is the net result of the competing uptake at the respiratory surfaces and elimination processes and can be expressed as:

$$\frac{dC_B}{dt} = (k_W C_W) - (k_{RO} + k_E + k_R + k_B + k_{RL} + k_G)C_B \quad (1.2)$$

Where dC_B/dt is the net change in concentration in the organism (g/kg) over time t , C_B is the chemical concentration in the organisms ($\text{g}\cdot\text{Kg}^{-1}$), k_W is the chemical uptake rate constant from the water at the respiratory surface ($\text{L}\cdot\text{Kg}^{-1}\cdot\text{t}^{-1}$), C_W is the freely dissolved chemical concentration in the water ($\text{g}\cdot\text{L}^{-1}$) and k_{RO} , k_E , k_R , k_B , k_{RL} , k_G are rate constants representing chemical elimination from the organisms via the respiratory surface, fecal egestion, renal excretion, metabolic biotransformation, reproductive losses and growth dilution, respectively. At the steady state i.e. when C_B and C_W no longer change with time, the (1.2) can be rearranged to meet equation (1.1):

$$BCF = \frac{k_W}{k_{RO} + k_E + k_R + k_B + k_{RL} + k_G} = \frac{C_B}{C_W} \quad (1.3)$$

BCF can only be measured in controlled laboratory condition measuring the concentration in the fish and in the water after a prolonged exposure time⁶². For all the duration of the test fish are not fed so diet is deliberately not taken into account. The experimental determination of the BCF is anyway expensive (i.e. more than 35'000 euros for chemical) and requires more than 100 animals for each study¹¹⁸. However, a strong correlation between BCF and K_{OW} has been demonstrated and BCF values for neutral organic chemicals are estimated from regression between empirical BCF data and K_{OW} ⁶⁵. Many QSAR models are available for the prediction of the BCF in aquatic organisms^{68,116,117,119–125}.

Nevertheless, bioconcentration is only one of the processes that determine the overall bioaccumulation and BCF is not an exhaustive parameter in quantifying the bioaccumulation potential of chemicals, especially of many CEC. Bioaccumulation is defined as “the process by which chemicals are taken up by a plant or an animal either directly from exposure to a contaminated medium (e.g. soil, sediment, water) or by eating food containing the chemical”⁶⁰ and can be expressed with the following mass balance equation^{63,68,89,114,126,127}:

$$\frac{dC_B}{dt} = (k_W C_W + k_D C_{D,i} + k_{RI} C_{AG}) - (k_{RO} + k_E + k_R + k_B + k_{RL} + k_G + k_D) C_B \quad (1.4)$$

Introduction

Where k_D is the uptake rate constant for chemical in the diet ($\text{kg}\cdot\text{kg}^{-1}\cdot\text{t}^{-1}$) and C_D is the chemical concentration in the diet ($\text{g}\cdot\text{kg}^{-1}$), k_{RI} is the respiration intake rate constant ($\text{L}(\text{kg}\cdot\text{h})^{-1}$) and C_{AG} is the air concentration ($\text{g}\cdot\text{L}^{-1}$). From equation (1.4) results that bioconcentration is special case of bioaccumulation where the diet uptake is set to be null (i.e. $k_D C_{D,i} = 0$) (for aquatic organisms the respiration intake rate constant from air also is null). This does not reflect the real natural conditions where diet is one of the main contributions to the bioaccumulation of chemicals, mainly for the terrestrial, air-breathing organisms^{115,128}. BCF is a poor descriptor of biomagnification in food web and is not able of meeting environmental standard objectives because dietary exposure and other key environmental process that may lead to higher chemical concentrations are not included. Historically, the selection of the BCF as a metric for bioaccumulation for regulatory purposes was the result of the assessment of the common properties of the initial list of POPs, based on the available data on aquatic organisms^{109,129}. These compounds were all neutral, nonpolar, lipophilic organic chemicals, their partitioning was dominated by lipids sorption and $\log K_{OW}$ and BCF were suitable indicator of their bioaccumulation potential. Many CEC have different properties, e.g. they may have a high aquatic mobility, can be bioavailable for uptake from water³⁵ and tend to partition in different compartments in the body such as protein. Other indicators should be considered for the assessment of the bioaccumulation potential of chemicals; such metrics should be applicable to air-breathing as well as water-breathing organisms and relevant for all the route of uptake, diet included^{62,114}.

Generally, high levels of bioaccumulation correspond to slow depuration rates. Criteria based on the characterization of depuration rate constant should therefore be considered in the bioaccumulation assessment. k_t and the corresponding total elimination HL could be a direct measure of the bioaccumulative capacity of chemical substances. k_t is strongly related to intrinsic chemical properties and their susceptibility to undergo biotransformation and could be meaningful for CEC as well as for “legacy contaminants”^{126,129}. Moreover, it is independent from the uptake scenario (i.e. aquatic or terrestrial) and meaningful for all the uptake routes (i.e. diet and respiratory system).

Data of elimination HL can be used in a tiered approach to refine the estimation of the bioaccumulation potential¹²⁹, for instance, when the BCF predictions are refined with k_B data values the estimated BCFs are in better agreement with the measure BCFs^{130,131}. k_B can be the most influent parameter determining k_t ; if a chemical is metabolized by the organism, it is no longer present in the body thus it cannot be bioaccumulative. Biotransformation by higher organisms (e.g. fish or mammals) includes all the metabolic reactions which involve normal body constituents (e.g. lipids, proteins and carbohydrates) or xenobiotics: “a man-made chemical or material not produced in nature and not normally considered a constituent component of a specified biological system”⁴⁹. These biotransformation processes generally occur in a two-phase series of reactions which can occur in multiple tissues but for most xenobiotics primarily in the liver^{132,133}. During Phase I reactions, hydrophobic chemicals are typically functionalized to generate more polar (water soluble) metabolites. During Phase II reactions compounds are conjugated to large molecules to further increase polarity and generate more water soluble metabolites¹³⁴. Biotransformation reactions seeks to detoxify the organism from a xenobiotic by modifying its molecular structure transforming parent compounds into more polar more easily eliminated e.g. through the urinary system. However, in some cases biotransformation may increase toxicity, leading to metabolites that are more toxic than the respective parent compound¹³⁵. Some examples of bioactivation are the generation of toxic epoxides, such as in the conversion of aldrin to dieldrin and of heptachlor to heptachlor epoxide, as well as the formation of DNA-binding epoxides such as in the bio-activation of benzo-[a]-pyrene¹³⁶. Therefore, the identification of metabolites, and metabolic pathways, and the quantification of biotransformation rate parameters are critical steps in the determination of the possible toxic profile of chemicals^{67,137,138}. Biotransformation rates can be measured *in-vitro* and *in-vivo* by quantifying the formation of the metabolites from the parent compounds or determining the rate of chemical loss¹³⁹. These methods are anyway still expensive, labor intensive and standardized methods for the extrapolation of *in-vitro* data to *in-vivo* values in organisms different than fish are missing. Few QSARs exist for the prediction

Introduction

of biotransformation potential in fish^{140–144} and limited models are available for the prediction of biotransformation potential in mammals and human¹²⁶. If data about the biotransformation potential for CEC were available, they would be useful to improve the B estimation and could be used to refine the PBT assessment of CEC.

1.7.2 Refinement of the Toxicity Assessment

Adverse effect refers to the damage (compromised biological function) caused by a toxic agent on specific individual. Usually a threshold, under which the substance does not cause any adverse effect, exists and the organism is in normal condition. When this threshold is exceeded, the substance begins to manifest its effects, but usually the organism can compensate them by defense and detoxification mechanisms. If the dose (or concentration) rises again, the defense mechanisms are no more able to compensate the adverse effects and the chemical begins to manifest its toxicity, until the *exitus* of the organism. In the case of toxicity the Insubria PBT Index was tuned using a model for acute toxicity in the species *Pimephales promelas*^{91,145}. On the other hand, the US-EPA PBT Profiler takes into account chronic toxicity but in one species of fish⁸⁶. The aim of the Environmental Risk Assessment should be the protection of the whole ecosystem. For this reason, the screening performed by Gramatica and collaborators for the prioritization of CEC^{103–105} can be refined to take into account not only one species but different organisms at different level of organization.

However, the high complexity and variety of the environment make very hard to assess the effects of all chemicals on all sensitive species. Moreover, effects related to chronic long-term toxicity are hard to model and predict and usually acute toxicity end-point are preferred. Therefore, a simplification and a standardization of the studied ecosystem are clearly needed. In order to assess aquatic toxicity, only a very limited number of species are currently used in the experimental test to cover all the key trophic levels of the aquatic ecosystem. These few species are considered representative of the aquatic environment and are used as surrogate for all aquatic organisms. Testing methods for aquatic toxicity have been harmonized by the Organization for Economic Co-operation

and Development (OECD), which provided specific guidelines for the different endpoints specifying species, methods, times, routes of exposure etc. Testing procedures on the first aquatic trophic level are centered on unicellular green algae (e.g. *Pseudokirchneriella subcapitata*), diatoms (e.g. *Navicula pelliculosa*) or cyanobacteria (e.g. *Anabaena flosquae*). The endpoints measured to assess acute toxicity are photosynthesis or population growth rate inhibition. The standardized test, typically used to determine an acute EC₅₀, is the algal growth inhibition test¹⁴⁶. For the second aquatic trophic level, tests are performed using crustaceans, such as *Daphnia magna* or *Daphnia pulex*. Endpoint for acute toxicity test is immobilization (or mortality) of organisms after an exposure of 48 hours¹⁴⁷. The upper aquatic trophic level is assessed by the acute toxicity test in fish; the endpoint is the median lethal concentration measured after 96 hours of exposure of the tested organisms. Test species are usually various fish such as rainbow trout (*Oncorhynchus mykiss*), fathead minnow (*P. promelas*) and zebrafish (*Brachydanio rerio*)¹⁴⁸.

Although the standardization performed by OECD has brought to identify few species and well defined methods to test chemicals, testing approaches still remain expensive, labor intensive and time consuming and data for majority of CEC and in particular for pharmaceuticals are limited. In this situation, QSAR models are a valid tool to fill data gaps, identifying the range of toxicity of compounds and focusing the experimental tests on the most toxic^{149–151}. Sanderson stated that QSAR models can be useful to prioritize pharmaceuticals according to their acute toxicity, even if much future developments are needed to address the regulatory purposes¹⁵². Also the European Union Commission's Scientific Committee on Toxicity, Ecotoxicity and Environment recommended the use of QSAR models for screening purposes of pharmaceutical ingredients¹⁵³.

So far, the main used tool to predict ecotoxicity of active pharmaceutical ingredients by QSAR is ECOSAR⁸⁶, a modelling tool, based on logK_{OW}. Even though a lot of studies used these models to fill the data gap^{154–159}, Madden and colleagues stated that the applicability domain of the ECOSAR program to predict pharmaceuticals effects should

Introduction

be carefully evaluated, because the models were developed using small industrial chemicals as training sets¹⁰⁰. In fact, such models are based on very small sets of molecules, mainly of simple chemical structure with only a single functional group, while pharmaceuticals are complex chemicals often with a plurality of functional groups.

Recently some authors tried to develop QSAR models for pharmaceuticals developing statistically validated QSAR models which are based on theoretical molecular descriptors¹⁶⁰⁻¹⁶⁴. However, also these models have been developed on very small sets of molecules and, most importantly, are not available for all the required trophic levels, a necessary condition for a comprehensive assessment of the potential hazard for the aquatic compartment as required by the regulation. *Ad-hoc* QSAR models to estimate the acute toxicity of multiple species in a simplified aquatic ecosystem could be used to improve the T evaluation and refine the PBT screening of CEC.

1.8. Aims of the thesis

The aim of this thesis is to propose a consistent approach based on QSAR models for the evaluation of the intrinsic environmental hazard of CECs in order to facilitate the identification of the most environmentally hazardous compounds. The attention is mainly focused on the PBT properties that, as explained in section 1.5, are particularly relevant for the prioritization of contaminants. This work addresses three main subjects: the first is the general screening of the PBT properties by QSAR models, while the second and the third deal with the development of QSAR models for the refinement of the B and T assessment.

The first subject is discussed in **paper I** which is focused on the screening of the potential PBT behavior of pharmaceuticals performed by the QSAR based consensus approach proposed by Gramatica and colleagues¹⁰³⁻¹⁰⁵. This approach was applied to draft a priority list of the most environmentally hazardous pharmaceuticals directly from the molecular structure in absence of experimental data.

The second topic builds on the results of **paper I** and focuses on the refinement of the toxicity evaluation of the previous PBT screening. It is described in **papers II, III and IV**, where *ad-hoc* QSAR models are developed to predict the acute aquatic toxicity of pharmaceuticals in multiple species. The aim of these works is to develop models for all the levels of the aquatic trophic chain for priority setting and general screening of the aquatic toxicity of pharmaceuticals. In these studies, acute toxicity is selected as target end-point because is able to give a prompt awareness about the hazard of chemicals for the aquatic environment. According the REACH regulation, information about acute toxicity are useful for screening of large inventories and databases, and can be used to identify priority chemicals for further evaluation on subtle and long-term effects. The final aim of **papers II, III and IV** is to provide models able to rank and highlight pharmaceuticals potentially toxic for aquatic environment in multiple species on the basis of the molecular structure. These models can be used in the refinement of the toxicity assessment of PBT screening. For these reasons specific mode of actions are not taken into consideration at this level but should be taken into account into further and more detailed specific studies.

The third part addresses the bioaccumulation refinement and is discussed in **paper V** and **VI** (in preparation). These papers are focused on biotransformation potential in higher organisms, mainly humans, and aim to analyze and demonstrate the importance of biotransformation and metabolisms in the estimation of bioaccumulation and biomagnification. In **paper V**, QSAR models for the prediction of *in-vivo* whole-body human biotransformation HLs are developed from an empirically-derived dataset of over 1000 organic chemicals¹²⁶. These QSARs can be used to predict the biotransformation HL of CEC to refine the B evaluation of the screening performed in **paper I**. Finally, the aim of **paper VI** is to understand the contribution of human biotransformation in the estimation of the BMF. In this work, predictions for the biotransformation potential are integrated in a mechanistic mass-balance multimedia environmental fate, food-web model to estimate the BMF in human in a tiered approach. The tiered approach progresses from conservative assumptions to more realistic ones for chemical properties,

Introduction

biological partitioning and biotransformation in human, considering a screening-level standard scenario in which chemicals are released with a constant rate and are supposed to be persistent in the environment. Other issues related to environmental fate such as pseudo-persistence and mobility are outside the scope of this work and PhD project and need to be addressed in further studies that will complete the overall PBT assessment.

All the QSAR models developed in thesis are based on theoretical molecular descriptors calculated by PaDEL descriptor¹⁶⁵ and are developed according to the statistical approach reported in the QSARINS software¹⁶⁶ and explained in chapter 2. All the models are implemented in the software QSARINS-Chem¹⁶⁷ (new version in preparation) for a more efficient dissemination and application.

Chapter 2: QSAR Methodologies

2.1 Quantitative Structure Activity Relationship

Quantitative Structure Activity(Property) Relationship (QSA(P)R) models are mathematical algorithms able to relate the molecular structure of a chemical, described by quantitative parameters called molecular descriptors, to a specific property (QSPR) or a biological activity (QSAR).

The relationship between chemical structure and physical-chemical properties is known since the XIX century when, in 1863, Crum-Brown and Fraser demonstrated the correlation between the water solubility of different alkaloids and their molecular weight¹⁶⁸. The idea of correlating the molecular structure of a chemical and its biological properties was further developed in the early beginning of last century by Meyer and Overton, which studies identified the relation between the non-specific biologic toxicity (narcosis) and the capacity of the chemical to pass the cell membranes (Meyer-Overton Correlation)¹⁵⁹. The current QSAR methodologies find their foundations in the pioneering works, in the mid-1960s, of Hansch and co-workers^{170,171} and Free and Wilson¹⁷². The Hansch approach is based on the idea that every substance that interacts with the molecular system of a living organism leads to a perturbation in the organism that depends on the properties of the substance. This postulate is based on the Hammett equation¹⁷³, where steric, electronic and hydrophobic properties of a molecule are combined to derive the following relationship:

$$\text{Biological Activity} = a + b(\log P) + c(E) + d(S) \quad (2.1)$$

where the biological activity is represented as a function of some physical-chemical and structural properties. LogP (the partition coefficient between n-octanol and water, i.e. $\log K_{OW}$) is the hydrophobicity term encoding the lipophilic interaction of the chemical with the cell membrane. E and S encode for the electronic and steric properties of the chemical respectively and represent the possibility of a chemical to interact with the target and be active. a, b, c and d are simple coefficients. The limitation of the Hansch approach is that the equation is applicable only to chemicals “very similar” to, or

congener with those used to obtain the equation itself, following the “congenericity principle”¹⁷⁴.

The Hansch method was the base of the modern QSAR, which has made important progresses in the last 50 years toward the developing of more complex modeling approaches, supported by the continuous improvement of computer science and the introduction of several new methods (Partial Least Squares Regression (PLS), Artificial Neural Networks (ANN), Bayesian approaches)^{175,176}.

Nowadays, QSAR models are applied in many disciplines, from Risk Assessment to drug design and are explicitly required in several regulatory frameworks, like REACH Regulation, as alternative techniques to animal testing^{53,177}. The need to guarantee the scientific validity of the QSAR estimation for regulatory purposes and to promote the mutual acceptance of QSAR models led to the development of a set of general and internationally recognized principles for QSARs development and validation. Several principles were first proposed in 2002 at an international workshop held in Setubal (Portugal)(i.e. “Setubal Principles”)¹⁷⁸. These principles have been modified in 2004 by the OECD Work program in QSARs, and are now referred to as the “OECD Principles for the validation of (Quantitative) Structure Activity Relationship models”^{179,180}. These new principles, which should be followed in every development and application of QSAR models, are:

- *A defined end-point*: A (Q)SAR should be associated with a “defined end-point”, where end-point refers to any physical-chemical, biological or environmental effects that can be measured and quantified and thus modeled. The aim of this principle is to ensure transparency in end-point predicted by a given model since a given end-point may be determined by different experimental protocols and under different experimental conditions. Ideally, (Q)SARs should be developed from homogeneous datasets in which the experimental data have been generated by a single protocol.

Methods

- *An unambiguous algorithm:* The intent of this principle is to ensure transparency in the model algorithm that generates predictions of an end-point from information on chemical structure. Therefore, a clear description of the dataset, molecular descriptors generation, modeling procedure and statistical methods and parameters used for validation is required. This principle also includes the need for reproducible predictions.
- *A defined domain of applicability:* A (Q)SAR should be associated with a defined applicability domain (AD), in which the model makes estimates with a defined level of accuracy (reliability). When applied to chemicals within its AD, the model is considered to give reliable results. There are not unique measures to define a model's AD or unique criteria to assess model reliability. It should be regarded as a relative concept, depending on the context in which the model is applied.
- *Appropriate measures of goodness-of-fit, robustness and predictivity:* This principle expressed the need to provide two types of information: the internal performance of the model (as represented by goodness-of-fit and robustness) calculated on a defined training-set and the predictivity of a model, determined by using an appropriate test-set. It is important to note that there is no an absolute measure of predictivity that is suitable for all purposes, since predictivity can vary according to the statistical methods and parameters used in the assessment.
- *A mechanistic interpretation, if possible:* A QSAR should be associated with a “mechanistic interpretation”, wherever such an interpretation can be made. Clearly, it is not always possible to provide a mechanistic interpretation of a given QSAR. The intent of this principle is therefore to ensure that there is an assessment of the mechanistic associations between the descriptors used in a model and the end-point being predicted and that any association is documented.

OECD Principles constitute a conceptual framework to guide the development and validation of QSARs and represent a reference for the assessment of the scientific validity of the QSAR models as well as of the reliability of QSAR predictions. Complete

information regarding the scientific validity of QSAR models and adequacy of the generated predictions can be documented by using the appropriate QSAR Model Reporting Format (QMRF) and QSAR Prediction Reporting Format (QPRF).

2.2 QSAR Procedure

As reported in the previous paragraph, QSA(P)R models are based on the definition of a quantitative mathematical relationship between the structure of a chemical and its biological activity or a specific physical-chemical property.

A theoretical scheme of activity/property-molecular structure relationships is represented in figure 2.1¹⁸¹.

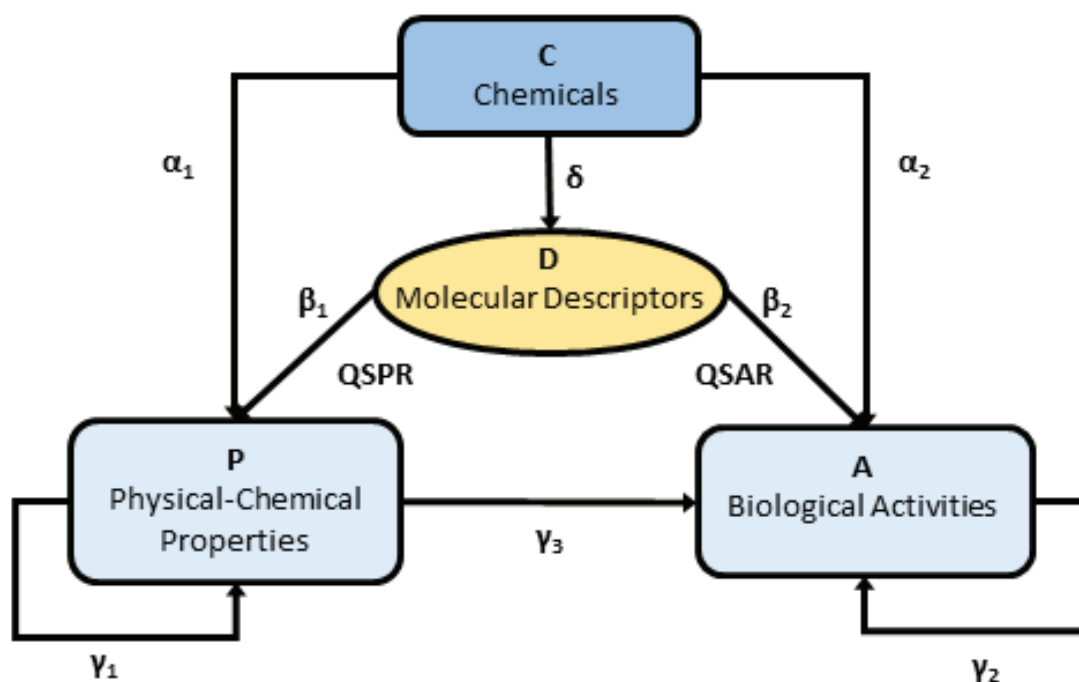


Figure 2.1: General scheme of the QSA(P)R approach.

For a given set of compounds (C), data of physical-chemical properties (P) or biological activities (A) are determined directly by experimental procedures α_1 and α_2 respectively:

$$\alpha_1(C) = P \quad (2.2)$$

Methods

$$\alpha_2(C) = A \quad (2.3)$$

On the other hand, it is possible to translate the structural information of the chemicals in to quantitative variables called molecular descriptors (D); each descriptor describes a peculiar feature of the molecular structure and can be theoretically calculated by the function δ :

$$\delta(C) = D \quad (2.4)$$

Then QSA(P)R methods are able to build mathematical functions β_1 , β_2 that correlate molecular descriptor with physical-chemical properties and biological activities respectively:

$$\beta_1(D) = P \quad (2.5)$$

$$\beta_2(D) = A \quad (2.6)$$

Finally, the functions γ_1 and γ_2 are the models able to describe, respectively, the properties-properties and activity-activity relationship while γ_3 function represents the models able to calculate biological activity from physical-chemical properties:

$$\gamma_1(P) = P \quad (2.7)$$

$$\gamma_2(A) = A \quad (2.8)$$

$$\gamma_3(P) = A \quad (2.9)$$

Fundamental prerequisites to obtain good quality QSARs are the availability of input experimental data of good/high quality, an exhaustive representation of the chemical structure and the use of valid and adequate statistical methods.

2.3 Molecular Descriptors

A molecule is a complex structural system that can be represented through several different molecular representations, each constituting a different conceptual model and including different information related to chemical structure (e.g. 2D or 3D information). Structural information is extracted through the calculation of molecular descriptors,

which are numerical variables quantifying the structural aspects of a chemical. In other words: “The molecular descriptor is the final result of a logic and mathematical procedure which transforms chemicals information encoded within a symbolic representation of a molecule into a useful number or the result of some standardized experiment”¹⁸². Many molecular descriptors have been proposed and derived from different theories and approaches, with the aim of predicting biological and physical-chemical properties of molecules¹⁸².

Molecular descriptors are divided in two main classes: those derived by experimental measures and theoretical molecular descriptors derived from a symbolic representation of the molecule. Furthermore, molecular descriptors can represent the whole molecule or a part of it.

A classification of theoretical molecular descriptors splits them into following groups:

- *Binary and one dimensional descriptors (0D and 1D)*: They represent the simplest type of descriptors and are derived directly from the formula of the molecule. Some examples are the molecular weight, the number of atoms included in the structure, as well as list of elements, fragments, functional groups or substitutes present in the molecule.
- *Two dimensional descriptors (2D)*: These descriptors take into account the atoms connections on the basis of the two dimensional representation, for example, defining the connectivity of atoms in the molecule in terms of the presence and the nature of the chemical bonds (topological representation). Thus, the presence of characteristic functional groups, the order of the atoms bonds and the class describing the distances are molecular descriptors derived from algorithms applied to a topological representation.
- *Three dimensional descriptors (3D)*: These descriptors give information about the three-dimensional structure of the molecule and about its spatial coordinates. They are derived representing the molecule as a geometrical object in a three-dimensional space. This view allows a representation not only of the nature and

Methods

connectivity of the atoms, but also the overall spatial configuration of the molecule. A more realistic representation of the bonding angles is possible by using low energy optimized conformations.

- *Fingerprints*: Binary fingerprints consist in binary vectors encoding the presence or absence of specific fragments or substructures. Different algorithms for the calculation of binary fingerprints were defined; the main difference is related to the definition of the fragments, which can be either based on pre-existing library, or derived from the analyzed dataset through the generation of all the fragments meeting some criteria.

Other types of descriptors can be referred to electronic parameters, such as the quantum-chemical descriptors (e.g. Highest Occupied Molecular Orbital (HOMO), Lowest Unoccupied Molecular Orbital (LUMO)) or to physical-chemical parameters as the $\log K_{ow}$ or Molar Refractivity.

Within this thesis binary, 1D and 2D molecular descriptors have been calculated with the software PaDEL Descriptor¹⁶⁵.

2.4 Explorative analysis: Principal Component Analysis

As described in the previous paragraph, molecular descriptors give a complex multivariate representation of the structural features of chemicals. To analyze this kind of data, chemometric tools can be applied in order to extract information and select the most appropriate method to handle them. Therefore, exploratory analysis of the dataset, both in terms of structural representation and response domain is an important step preceding model development.

Principal Component Analysis (PCA) is one of the best known procedures in multivariate statistics, which find application in many fields, from chemistry to economy¹⁸³. PCA allows the examination of the correlation pattern among variables and an evaluation of their relevance, the synthesis of data description discarding noise, the reduction of data dimensionality by discarding unnecessary variables, and the finding of

principal properties in multivariate systems¹⁸⁴. The aim of PCA is to transform p -correlated variables into a set of orthogonal variables, which reproduce the original variance/covariance structure of the data. This means rotating a p -th dimensional space to achieve independence between variables. The new variables, called Principal Components (PCs), are linear combinations of the original variables along the direction of maximum variance in the multivariate space (figure. 2.2), and each linear combination explain part of the total variance¹⁷⁵.

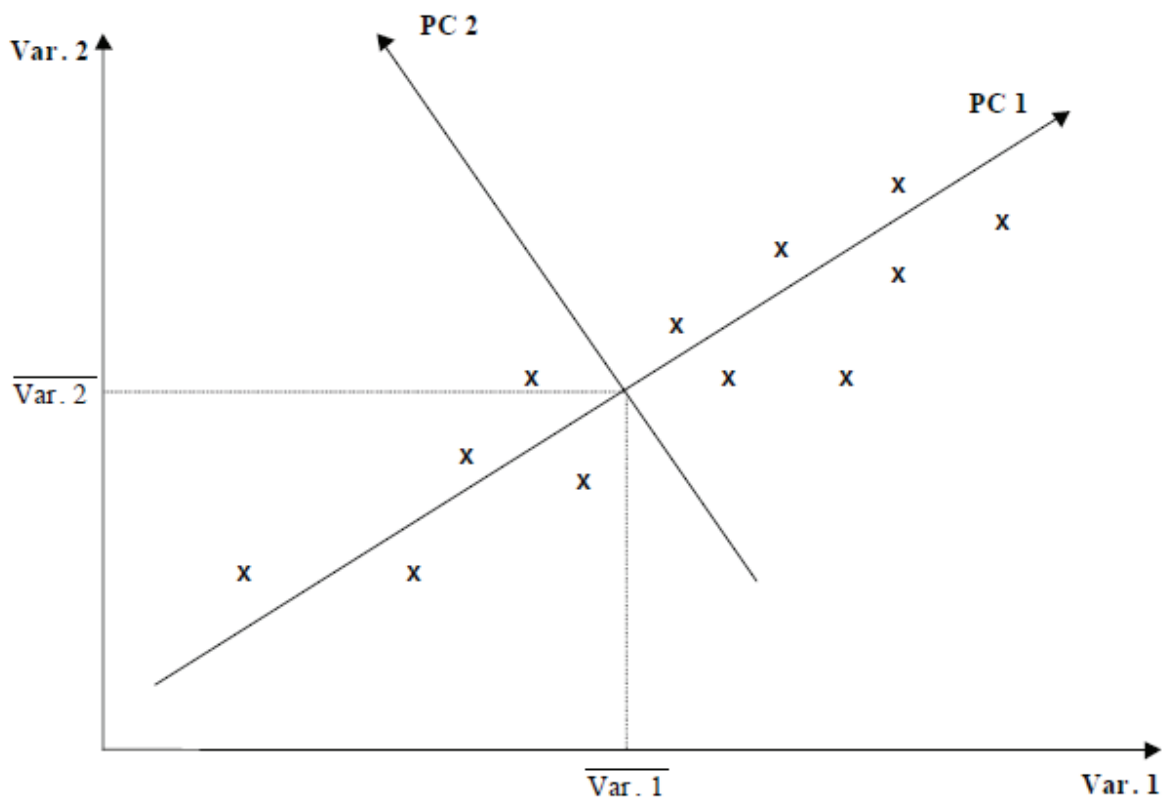


Figure 2.2: PCA illustrated for two variables.

The mathematical procedure consists in the calculation of eigenvector and eigenvalues from the of the original data matrix X (n , p) by the diagonalization of the variance/covariance matrix S :

$$\text{diag}(S) = \text{diag} \left[\frac{X_c^T X_c}{n-1} \right] \quad (2.10)$$

Methods

Where X_c is the centered data matrix.

This lead to the determination of a diagonal matrix Λ , called eigenvalues matrix, and a loadings matrix L . Λ is a diagonal (p, p) matrix whose diagonal elements are the eigenvalues λ_m in descending order and represent the variance explained of each PC in the new space; higher is the eigenvalue λ_m for a given PC, higher is the information burden from that PC.

The loadings matrix is a (p, M) matrix (if we consider all the PCs $M=p$; the selection of the relevant components will be explained in further) where columns are the eigenvectors l_m of S and represent the axis of the new space oriented according to the maximum variance direction, while rows represent the original variables. Each element of an eigenvector is defined as loading, l_{jm} , and represents the importance of each original variable in that given eigenvector. Loading are linear standardize coefficients and their values range between -1 and 1; an l_{jm} absolute value near 1 means that the m -th components is mostly influenced by the j -th variable.

Thus it is possible to approximate the original data matrix X according to the relation:

$$\begin{aligned} T &= XL \\ (n, M) &= (n, p)(p, M) \end{aligned} \tag{2.11}$$

where L is the rotation matrix and T is the scores matrix.

T is an (n, M) matrix where each element, called scores, is a linear combination of the original variable and loading for a given PC and represents the new coordinates of objects in the PCs space.

A fundamental aspect in the PCA is the selection of the PCs to analyze; indeed, even if the sum of the eigenvalues is equal to the total variance of the dataset, it is useful to consider only the PCs that explain the main information, ignoring the last PCs associated with the noise of the data. In such way, there is a dimensional and complexity reduction keeping a good approximation of the reality. The degree and the goodness of the approximation depend on the number of PCs selected for the analysis. Different methods

exist for the selection of the relevant PCs, here are reported the two main used techniques:

- *Scree plot*: is a line segment plot that reports the fraction of total variance in the data as represented by each PC. The PCs are ordered by decreasing order of contribution to total variance. The ideal pattern in a scree plot is a steep curve, followed by a bend (“elbow”) and then a flat or horizontal line; the “elbow” is the point of separation between “most important” and “least important” components¹⁸¹.
- *Average Eigenvalue Criteria*: By this method are considered significant all the components whose eigenvalues are higher than the average of all eigenvalues. When variables are auto-scaled the average eigenvalues is 1 thus are selected all the components with eigenvalue higher than 1.

2.5 Data modelling

Models have a fundamental role in the knowledge process since they synthesize all the information about a specific problem giving a simplified representation of the reality. Models can also be used to foresee future events linked with a causal relationship to known variables^{185,186}. Every system, from the simplest to the most complex, has a certain information content that can be observed from different point of view and described in different ways. The amount of information of the system depends on the variability of parameters and variables used to describe the system itself¹⁸⁶. In other words, to represent and know a system means transfer the information contained in the system itself to a knowing subject, and models are one of the ways whereby this transfer can take place. The modeling process is constituted by four fundamental steps^{187,188}:

- *Identification*: In this step, the selection of the best model to represent a given situation takes place. No standard procedures exist to determinate the best model and essentially two opposite ways can lead the selection; in one case the selection is based on empirical observation while in the other one the selection

Methods

is led by theoretical basis. On the basis of this division we can identify deterministic models, which are developed according to *a priori* hypothesis regarding the causal relationships, and statistic models, which development is based on its fitting with experimental data.

- *Realization*: Once selected the kind of model, realization computes the mathematical relationship and estimates the model parameters. In this phase are highlighted sources of uncertainty and errors as well as the data noise and the eventual approximation.
- *Validation*: It is the development phase in which the model is tested and validated. The model is continuously modified to improve its stability and predictive performances.
- *Application*: It is the last step, in which the developed model is applied to predict unknown events on the basis of known variables.

Figure 2.3 summarizes the modelling approach that is applied within this thesis in the software QSARINS for development and validation of QSAR models^{166,167,187}.

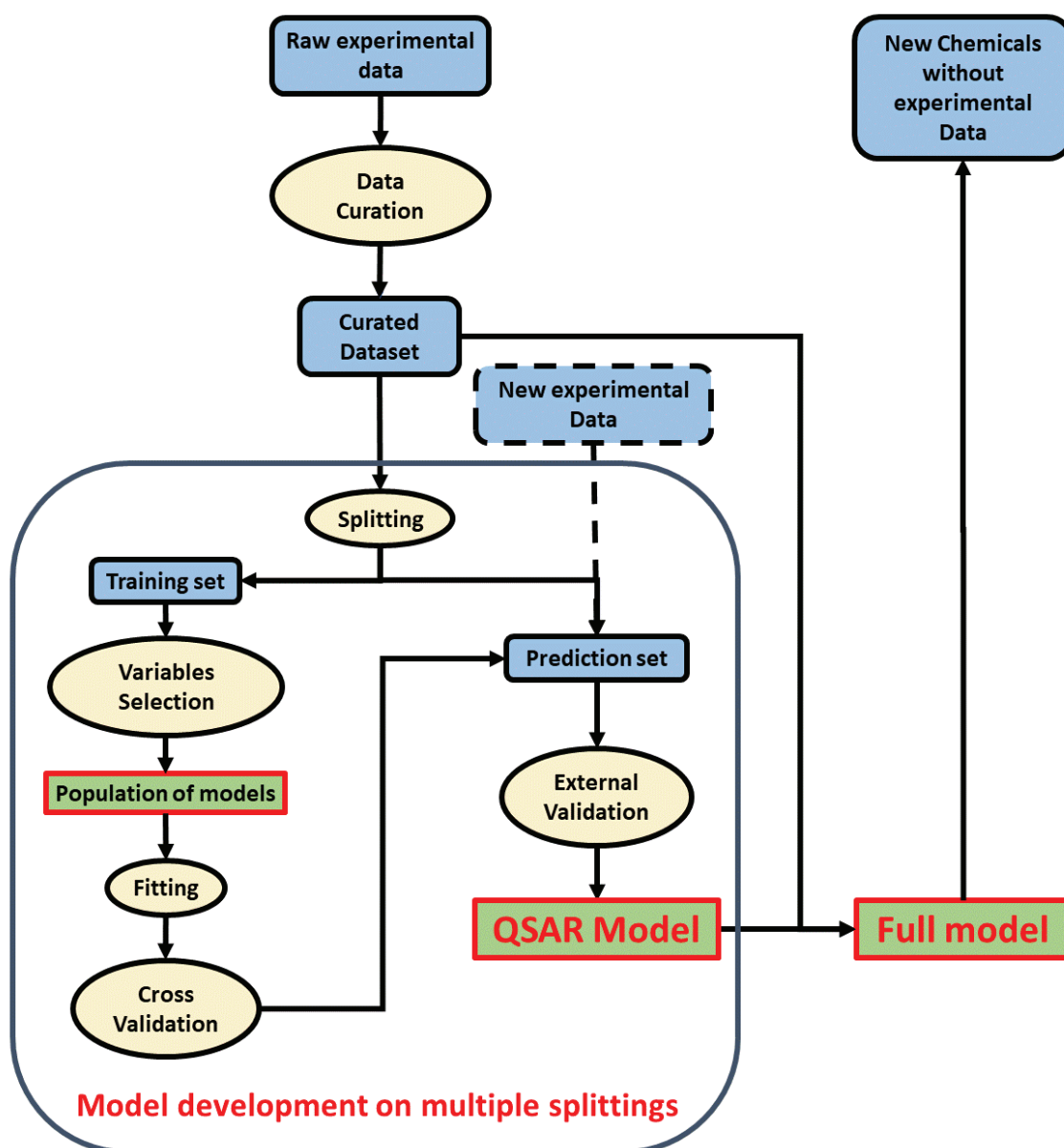


Figure 2.3: General scheme of the modelling approach.

An important prerequisite for the development of good quality QSARs is the availability of experimental data of high quality according to the general rule of computer science GIGO (Garbage in, Garbage out). Thus, data curation represents a fundamental step preceding QSAR modelling^{187–192}. Data curation does not only involve the careful selection of high-quality data, but also the verification of chemical structures used as input for molecular descriptors generation. It has been demonstrated that even a small error in a chemical structure can result in significant differences in the prediction of the

Methods

accuracy of a model, not only for chemicals with erroneous structural information, but also in the prediction of other chemicals using models that contain such errors¹⁹³.

Once consistent set is built, to assess the real predictivity of the model, the general scheme contemplates the splitting of the whole dataset into training set and prediction set, according to various techniques. The training set is the sub-set of objects on which the model will be developed and internally validated, while the prediction set is the sub-set on which the developed model will be applied to test its predictivity. The objects of the external set are not involved in the development of the model equation/algorithm. So, only the training set is involved in the initial development of the QSAR model. On it, a variable subset selection technique is applied to identify an optimal sub-set of variables that is able to model the property under study (end-point). Then, modeling algorithm such as MLR or Classification methods, are applied to the sub-set of variables to compute a quantitative model. Once obtained the model equation, statistical tests are performed to verify the ability of the model in reproducing and fitting the original data, and its stability and robustness (internal validation). The internal validation includes different techniques and parameters, like Cross Validation (CV), Quik rule¹⁹⁴ and Y-Scrambling, and it is performed to guarantee the internal stability and to exclude the possibility of chance correlation between selected modeling descriptors and the studied response. In general, CV is an iterative process that involves a sequential splitting of the original training set in different test sets on which the fitting is re-tested. Only when a stable and internally validated model is computed, the external validation is performed. The external validation assesses the predictivity of the model and involves the application of the developed model on the prediction set, originally excluded from computation and model development; other statistical tests are performed, verifying the ability of the model in predicting the data for new objects. Finally, when also the external predictivity has been established, the training and prediction set are merged and the equation is recomputed in a so-called “full model”, in order to include all the information brought by the original dataset.

A fundamental analysis carried out in all the model development steps is the study of the structural AD, to verify the applicability of the model to new chemicals.

2.5.1 Dataset splitting

The aim of a QSAR model is to predict data for chemicals without experimental evidences so it must be externally validated and its predictivity must be assessed on external objects not involved in the model development. However, in most of the cases, experimental data are available only for a limited number of molecules and for time reasons the modeler cannot wait the production of new data on which evaluate its model (temporal set). Thus, the only way to assess the real external predictivity of a model is to exploit the actual data availability putting some compounds in a prediction set, used only after the model development for the external validation¹⁹⁵. Anyway, when the predictivity of the model has been assessed, the prediction set is re-included in the training set and a full model, comprising all the available information, is computed.

The splitting in training and prediction set should be balanced, ensuring the representativeness of the structural and response domain of the original dataset. To obtain this result, different splitting techniques can be applied^{187,196,197}. One of the most used is the splitting by response where chemicals are ordered according to their increasing activity and then sampled with a regular scheme. The resulting prediction set includes a selected percentage of chemicals of the original dataset. The most and the least active compounds of the dataset are normally kept in the training set. This splitting guarantees that both training and prediction sets cover the entire range of the experimental response and are numerically representative of the dataset. However, such splitting does not guarantee that the two sets represent the entire structural space of the original dataset. Therefore, it is possible that some compounds in the prediction set are outside the structural domain of the training set.

Another splitting technique is based on structural similarity. This method takes the advantage of informational condensation capacity of PCA, allowing the selection of a structurally meaningful training set and an equally representative prediction set⁸⁰. This

Methods

splitting technique allows to generate training and prediction sets that are more structurally balanced^{198,199}.

Finally, another splitting technique involves the random selection of compounds that will be assigned to the prediction set. This is one of the most commonly used methods when a large number of chemicals is available because it is perceived free of “selection bias”²⁰⁰. However, since no information about the coverage of response and descriptors domain is considered, it can occur that the training and prediction sets could be unbalanced.

2.5.2 Variable selection

Variable selection is a necessary step to find simple and predictive QSARs, which should be based on the minor number of descriptors. Since is well known that some molecular descriptors provide only different views of the same molecular aspect, this procedure is particularly important. A selection of significant descriptors is performed by applying unbiased mathematical tools and starting from a large number of molecular descriptors that can be calculated for the molecules included in a dataset²⁰¹.

In the present thesis, the genetic algorithm (GA) implemented in QSARINS is used for the variable selection. GA is a useful method, widely and successfully applied in many QSAR approaches²⁰². The GA strategy for variable subset selection is based on the evolution of a population of models, i.e. a set of ranked models according to some objective function. Each individual is called chromosome and is a binary vector, where each position (a gene) corresponds to a variable (1 if included in the model, 0 otherwise). Each chromosome represents a model given by a subset of variables²⁰³.

The GA works based on three main steps:

- *Random initialization of the population*: The model population is built initially by random models with a defined number of variables. The value of the selected objective function of each model is calculated in a process called evaluation. The

models are then ordered with respect to the selected objective function, i.e. the best model is in first place in the population;

- *Crossing-over*: From the actual population, pairs of models are selected (randomly or with a probability function of their quality). Then, from each pair of selected models (parents), a new model is generated, preserving the common characteristics of the parents (i.e. variables excluded in both models remain excluded, variables included in both models remain included) and mixing the opposite characteristics according to the crossover probability. If the objective function value is better than the worst value in the population, the model is included in the population, in the place corresponding to its rank; otherwise, it is no longer considered. This procedure is repeated for several pairs;
- *Mutation*: After a number of crossing-over iterations, the population proceeds through the mutation process. This means that for each individual of the population every gene is randomly changed. Mutated individuals are evaluated and included in the population if their quality is acceptable. This process is controlled by mutation probability that is commonly set at low values, thus allowing only a few mutations and new individuals not too far away from the generating individual.

An important characteristic is that GA implemented in QSARINS does not provide a single model but a population of acceptable models. Within this population, there could be various models with similar predictive power, but based on different molecular descriptors combinations. In fact, different descriptors are alternative viewpoints to represent the structural features, whose combination lead to, not equivalent, but similar results for the studied end-point. Thus, there could be many possible “best” models, and their predictions could be averaged in the consensus modeling approach.

2.5.3 Multiple Linear Regression and Ordinary Least Squares Regression

MLR is a mathematical approach to establish linear relationships between a set of descriptors, the independent variables x , and a quantitative response, the dependent variable y ²⁰⁴. The relationship takes the form:

$$y = f(x_1, x_2, x_3, \dots, x_p) \quad (2.12)$$

The regression model is the mathematical equation used to describe the relationship among response and predictor variables. Regression modeling may be used for descriptive or predictive purposes.

The MLR model is described in matrix form as:

$$y = Xb \quad (2.13)$$

where b is the vector of the estimated regression coefficients, y is the vector of responses and X is the matrix model. Alternatively, the model can be written in non-matrix terms as:

$$y_i = b_0 + \sum b_p x_{ip} \quad (2.14)$$

where y_i is the calculated response of the i -th molecule, b_0 is the intercept, b_p are the coefficients of the p descriptors of the model and x_{ip} is the value of the p -th descriptor for the i -th molecule.

One of the most used regression techniques is the ordinary least squares (OLS) method. This method calculates a parametric linear model for a single response, and calculates unbiased, least squares coefficients. In OLS regression some assumptions should be met including: (1) the linearity of regression coefficients, (2) all predictors must be uncorrelated with the residuals, (3) residuals must not be correlated with each other (serial correlation), (4) residuals must have a constant variance, (5) none predictor variables must be perfectly correlated with another predictor variable (avoidance of multicollinearity), (6) residual must be normally distributed^{181,205}. OLS assume that the response is a linear function of the predictors, and that the errors are identically and independently distributed. The regression coefficients are computed by minimizing the

Residual Sum of Squares (RSS) between the calculated and the experimental response vectors, defined as:

$$RSS = \sum (y_i - \hat{y}_i)^2 \quad (2.15)$$

where y_i is the experimental response of the i -th molecule, \hat{y}_i is the corresponding calculated by the model.

The mathematical procedure to determine the coefficients b involves the following passages:

$$X^T y = X^T X b \quad (2.16)$$

$$(X^T X)^{-1} X^T y = (X^T X)^{-1} X^T X b \quad (2.17)$$

Where $(X^T X)^{-1} X^T X$ is an Identity matrix. Thus the OLS solution is:

$$b_{ols} = (X^T X)^{-1} X^T y \quad (2.18)$$

Once the regression coefficient vector b has been estimated, the calculated responses are obtained from:

$$\hat{y} = X b_{ols} \quad (2.19)$$

the (2.19) represents the regression model built by the descriptors in the model matrix and the regression coefficients b .

Important parameters relating to regression coefficients are the standardized regression coefficients (b'_j) that represent the weight of each variables in the regression model:

$$b'_j = b_j \frac{s_j}{s_y} \quad (2.20)$$

Where s_y and s_j are the standard deviation of the response and of the j -th predictor.

Substituting the b coefficients in the (2.19):

$$\hat{y} = x(X^T X)^{-1} X^T y = Hy \quad (2.21)$$

Methods

Where H is an (n, n) matrix that correlates experimental and calculated responses and is called leverage matrix or hat matrix.

Hat matrix is defined as:

$$H = x(X^T X)^{-1} X^T \quad (2.22)$$

The diagonal elements of this matrix, h_{ij} , are called variance functions of the j -th object and have the following properties:

$$h_{min} = \frac{1}{n} \quad (2.23)$$

$$\sum h_{ii} = p' \quad (2.24)$$

$$\bar{h} = \frac{p'}{n} \quad (2.25)$$

$$h^* > \frac{3p'}{n} \quad (2.26)$$

where h^* is a cut-off values beyond which an object can be considered influent in determining the regression parameters and p' is the number of the independent variables plus one.

The minimum leverage (h_{min}) value is the centroid in the model space while the higher values correspond to objects distant from the space center. Thus the hat value indicate the model extrapolation degree and a high leverage value for new objects means a high uncertainty in the reliability of the predicted responses (extrapolated predictions). This property of regression method is very useful in the identification of the model applicability domain which is explained further in paragraph 2.5.5.

OLS will provide the best estimate only if not independent variable is perfectly correlated with another independent variable (multicollinearity). In this case other regression techniques, such as PLS, should be preferred. PLS is a regression method that allows for the identification of underling factors, which are a linear combination of the

explanatory variables or X (also known as latent variables) which best model the response or Y variables. In PLS optimal linear relationships are computed between latent variables and can be view as the best set of prediction variables^{175,176}. In case of small sample size, missing values or multicollinearity PLS provides estimates much more accurately than OLS²⁰⁵. However, since the latent variable are linear combination of many independent variables, they are usually not easily interpretable. The interpretability of individual observed variables is an important factor in chemical data analysis. For this reason, OLS is preferred to PLS in this thesis. OLS is applied verifying that all the assumptions are met.

2.5.4 Validation techniques and parameters used for regression models

There are many statistical indices useful to evaluate the performance of the developed regression models relating to what kind of performance need to be evaluated: fitting, stability and robustness or predictivity.

A first group of statistical parameters are devoted to evaluate model's fitting ability, providing a measure of how well the regression model accounts for the variance of the response variable. Several fitness functions have been proposed and are here summarized.

- *Coefficient of determination (R^2)*: total variance of the response explained by a regression model.

$$R^2 = 1 - \frac{\sum(y_i - \hat{y}_i)^2}{\sum(y_i - \bar{y})^2} = 1 - \frac{RSS}{TSS} \quad (2.27)$$

Where TSS is defined as the total sum of Squares; a value of 1 indicates perfect fit.

- *Root Mean Square of Errors*: sum of the overall error of the model. It is calculated as the square root of the sum of the squared errors in predictions divided by their total number:

$$RMSE = \sqrt{\frac{RSS}{n}} \quad (2.28)$$

Differently from R^2 , whose values vary between 0 and 1, RMSE values are influenced by the range of the response. Therefore, this parameter is useful if applied to compare different models based on the same (or very similar) training set, and developed for the same range of experimental response.

A second group of regression parameters are devoted to evaluate the robustness and the goodness of prediction, providing a measure of how well the regression model estimates the response variable given a set of values for predictor variables. These quantities are obtained using validation techniques and are also used as criteria for model selection. Validation procedure is divided in internal and external; the internal is performed on the training set and assesses the stability of the model when subjected to perturbation and exclude the possibility of chance correlation between variables and response. External validation is performed on the prediction set and assesses the predictive ability of the model on new external objects, never seen during model development.

Parameters for the internal CV are:

- *Leave-One-Out and Leave-many-Out*: The simplest and most general CV procedures are the leave-one-out (LOO) and leave many out (LMO) techniques, where a number of objects (one or more than one at time) are excluded from the training set (and put in test set) during the model development. For each reduced data set, the model is calculated and responses for the excluded objects are predicted from the model. Then Q^2_{LOO} is calculated following formula²⁰⁶.

$$Q^2_{LoO} = 1 - \frac{\sum(y_i - \hat{y}_{i/i})^2}{\sum(y_i - \bar{y})^2} = 1 - \frac{PRESS}{TSS} \quad (2.29)$$

Where \hat{y}_i is the predicted value for the i-th object when excluded from the model computation. LOO and LMO give in several cases a too optimistic results,

particularly considering that the LOO introduces just a small perturbation in the dataset. For this reason, CV cannot be considered as a representative technique to assess the real predictivity of a model and external validation is needed^{195,207}.

- *Y-scrambling*: This validation technique is adopted to test models for chance correlation, where the independent variables could be randomly correlated to the response variables. The test is performed by iteratively calculating the quality of the model (usually R^2 and Q^2) randomly modifying the sequence of the response vector y , by assigning to each object a response randomly selected from the true responses²⁰⁸. Low values of the averaged R^2 and Q^2 scrambled (R^2_{YS} , Q^2_{YS}) are indicative of a good model.
- *QUIK rule*: The Quik rule¹⁹⁴ is normally applied in order to select only models where the correlation between the block of the modeling descriptors and the response (K_{XY}) is higher than the correlation among the descriptors (K_{XX}), $K_{XY} > K_{XX}$.

Internal validation is necessary but not sufficient, only external validation techniques are able to verify the actual model predictivity, toward compounds never used for model development^{186,195,207,209,210}. In fact, the real predictive ability of a QSAR model can only be estimated using an external set of compounds never used for the model development^{211,212}. This set of “external compounds” can either be obtained by *a priori* splitting of the data set (prediction set), according to the procedures explained above, or be represented by new data became available, or produced, after model development (temporal set). In the last years, external validation of QSAR models was the subject of intensive debate in the scientific literature. Different groups have proposed different metrics to find “the best” parameter to characterize the external predictivity of a QSAR model. The editorial “A Historical Excursus on the Statistical Validation Parameters for QSAR Models: a Clarification Concerning Metrics and Terminology” aimed to review the history of the available parameters for the external validation of QSAR model²¹³.

Methods

This editorial clarifies the terminology of the external validation parameters and suggest the use of several different metrics, to assess the real predictive capability of QSAR models. The parameters for external validation commented in the editorial are reported as follow:

- *Golbraikh and Tropsha criteria*²⁰⁷

$$R^2 = \left(\frac{\sum_{i=1}^{n_{EXT}} (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^{n_{EXT}} (y_i - \bar{y})^2 \sum_{i=1}^{n_{EXT}} (\hat{y}_i - \bar{\hat{y}})^2}} \right)^2 \approx 1 \quad (2.30)$$

$$R_0^2 = 1 - \frac{\sum_{i=1}^{n_{EXT}} (\hat{y}_i - y_i^{r_0})^2}{\sum_{i=1}^{n_{EXT}} (\hat{y}_i - \bar{\hat{y}})^2} \approx R^2 \quad (2.31)$$

$$R_0'^2 = 1 - \frac{\sum_{i=1}^{n_{EXT}} (y_i - \hat{y}_i^{r_0})^2}{\sum_{i=1}^{n_{EXT}} (y_i - \bar{y})^2} \approx R^2 \quad (2.32)$$

$$k = \frac{\sum_{i=1}^{n_{EXT}} y_i \hat{y}_i}{\sum_{i=1}^{n_{EXT}} \hat{y}_i^2} \approx 1 \quad (2.33)$$

$$k' = \frac{\sum_{i=1}^{n_{EXT}} y_i \hat{y}_i}{\sum_{i=1}^{n_{EXT}} y_i^2} \approx 1 \quad (2.34)$$

$$y_i^{r_0} = k \hat{y}_i \quad (2.35)$$

$$\hat{y}_i^{r_0} = k' y_i \quad (2.36)$$

Where R_0^2 and $R_0'^2$ are calculated forcing the regression line to pass through the origin, k and k' are the slope of the regression lines.

- *External Q^2 functions*²¹⁴⁻²¹⁷

$$Q_{F1}^2 = 1 - \frac{\sum_{i=1}^{n_{EXT}} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n_{EXT}} (y_i - \bar{y}_{TR})^2} = 1 - \frac{PRESS}{TSS_{EXT}(\bar{y}_{TR})} \quad (2.37)$$

$$Q_{F2}^2 = 1 - \frac{\sum_{i=1}^{n_{EXT}} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n_{EXT}} (y_i - \bar{y}_{EXT})^2} = 1 - \frac{PRESS}{TSS_{EXT}(\bar{y}_{EXT})} \quad (2.38)$$

$$Q_{F3}^2 = 1 - \frac{[\sum_{i=1}^{n_{EXT}} (y_i - \hat{y}_i)^2] / n_{EXT}}{[\sum_{i=1}^{n_{TR}} (y_i - \bar{y}_{TR})^2] / n_{TR}} = 1 - \frac{PRESS / n_{EXT}}{TSS / n_{TR}} \quad (2.39)$$

- *Concordance Correlation Coefficient (CCC)*^{211,212,218,219}

$$CCC = \frac{2 \sum_{i=1}^{n_{EXT}} (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sum_{i=1}^{n_{EXT}} (y_i - \bar{y})^2 + \sum_{i=1}^{n_{EXT}} (\hat{y}_i - \bar{\hat{y}})^2 + n_{EXT} (\bar{y} - \bar{\hat{y}})^2} \quad (2.40)$$

- *Roy et al. criteria*²²⁰⁻²²²

$$r_m^2 = r^2 \left(1 - \sqrt{r^2 - r_0^2} \right) > 0.5 \quad (2.41)$$

$$\bar{r}_m^2 = \frac{(r_m^2 + r_m'^2)}{2} \quad (2.42)$$

$$\Delta r_m^2 = |r_m^2 - r_m'^2| < 0.2 \quad (2.43)$$

Where r^2 and r_0^2 are respectively the determination coefficients of the regression function, calculated using the experimental and the predicted data of the prediction set, forcing respectively the origin of the axis (r_0^2) or not (r^2). r_m^2 is calculated using the experimental values on the ordinate axis, $r_m'^2$ using them on the abscissa. TR refers to training set, EXT to external set, y_i to experimental data values, \hat{y}_i to predicted data values, \bar{y} to average of the experimental data values and $\bar{\hat{y}}$ to average of the predicted data values.

In this thesis a model is considered for applications only if all the external values calculated in QSARINS are acceptable²¹¹⁻²¹³.

2.5.5 Applicability Domain

The application of any QSAR model to new compounds is limited by the fact that the model derives from a particular set of chemicals, i.e. the training set. The calibration of a mathematical model is based on the interpolation of the training data and therefore the model can be assumed valid in this defined space. Thus, as required in the third OECD Principles, it is fundamental to identify the chemical space where a QSAR model is

Methods

assumed to provide reliable predictions: the AD^{186,187}. Several methods are currently available for the evaluation and definition of the AD of a model^{223–225}.

In this thesis, three main approaches are used¹⁰³:

- *Range of descriptors (Bounding Box)*: This approach considers the range of individual descriptors used to build the model. Assuming a uniform distribution, resulting domain of applicability can be imagined as a Bounding Box which is a p-dimensional hyper-rectangle defined on the basis of maximum and minimum values of each descriptor used to build the model. The sides of this hyper-rectangle are parallel with respect to the coordinate axes.
- *Leverage*: The leverage matrix is calculated from the values of model descriptors according the equation (2.22). As already explained the diagonal elements of this matrix, h_{ij} , are the variance functions of the j-th object and give a measure of the distance of a compound from the center of the model space. h_{ij} values for training set objects range between $1/n$ and 1 while for the prediction set range from $1/n$ to ∞ . Molecules with high leverage are distant from the centroid of the model and can be considered different and in some circumstance outliers. The boundaries of the model AD are defined by the leverage cut-off h^* . Training molecules associated with a high leverage may considerably affect the model parameters while prediction set molecules with high leverage may be associated with less reliable predictions, certainly extrapolated by the model. The Williams plot (\hat{y} vs standardized residuals) can be used to assess the AD during the model development phase. A modification of the Williams plot called Insubria graph (\hat{y} vs predictions) can be used to assess the AD when the model is applied to chemicals without experimental values.
- *PCA Bounding Box*: This approach defines a sub-structural space, based on the PCA of the modeling molecular descriptors, delimited by the training set chemicals. Predicted (or new) compounds that are within these boundaries are

considered inside the AD. The implementation of Bounding Box with PCA can overcome the problem of correlation between descriptors²²⁴.

Methods

Chapter 3: Results and Discussion

The following are somewhat shortened results and discussion from the articles on which this thesis is based. More extensive and detailed discussion can be found in the specific publications attached at the end of the thesis.

3.1 Paper I: “PBT assessment and prioritization of contaminants of emerging concern: Pharmaceuticals”

By Sangion Alessandro and Gramatica Paola

Published in: *Environmental Research*, 2016, 147(5): 297–306.

Aim of this work is to perform the PBT assessment of a large dataset of pharmaceuticals by the consensus approach of different QSAR models and draft a priority list of the most environmentally hazardous pharmaceuticals. A list of about 1276 approved, withdraw and illicit organic drugs is compiled and curated from online databases and literature^{158,226–229}. Duplicated, ambiguous structures, inorganic and organometallic chemicals, salts and charged compounds are removed during the preliminary data curation procedure. The screening is performed by the consensus approach proposed by Gramatica and colleagues described in section 1.6¹⁰³. The Insubria PBT Index is an MLR QSAR model based on molecular descriptors that provides a unique index for the potential PBT behavior of a chemical on the basis of the molecular structure⁹¹ implemented in the QSARINS software. The US-EPA PBT Profiler is based on QSAR models within EPI Suite and predicts separately persistence, bioaccumulation and toxicity^{67,83–86,117}. Great attention is devoted to the analysis of the AD of the Insubria PBT Index with the 68% of the pharmaceutical ingredients within the structural AD. The screening results show that 14% of pharmaceuticals are estimated as potential PBTs by the Insubria PBT Index, while the US-EPA PBT Profiler predicts only 5% of compounds as potential PBTs. Approximately 86% of chemicals are predicted in the same way by both models (figure 3.1).

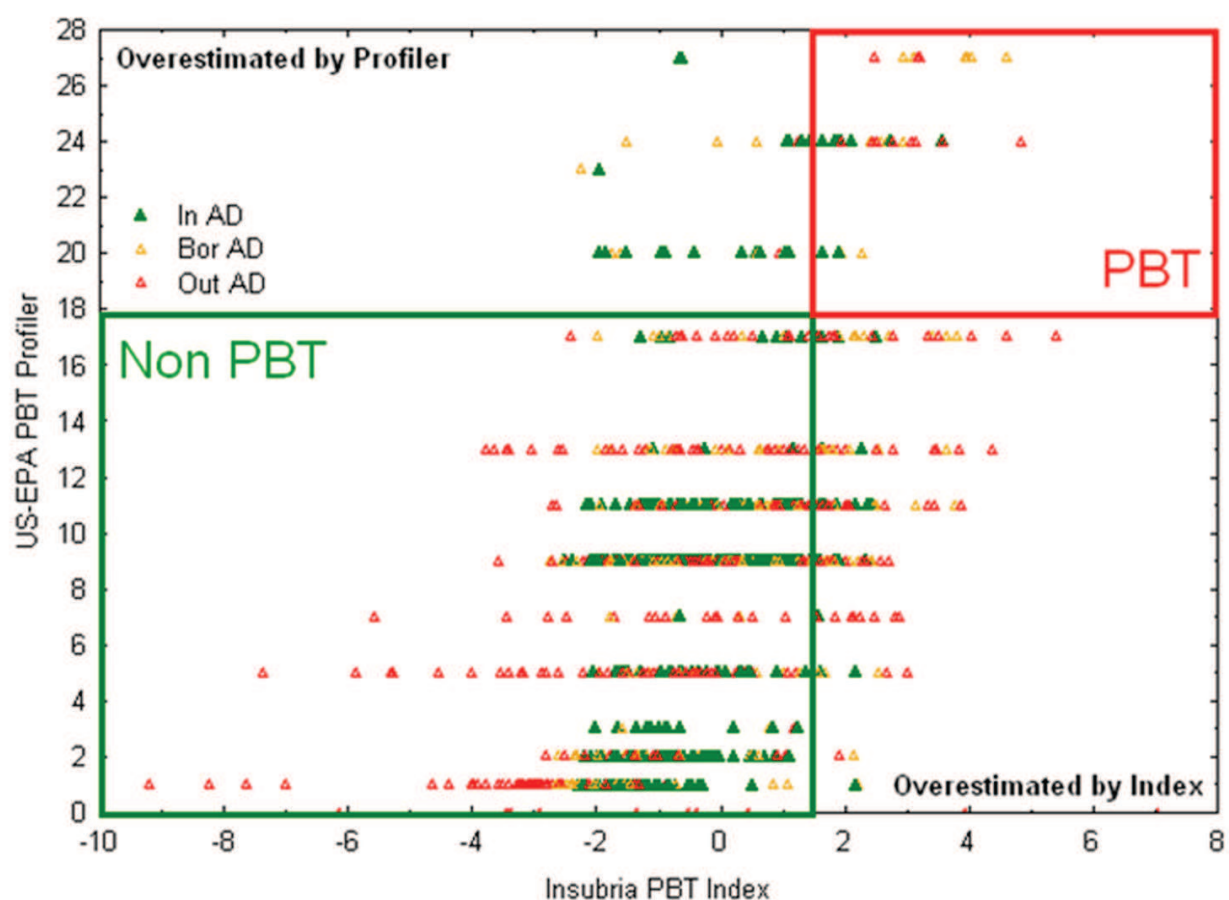


Figure 3.1: Graph of the agreement between Insubria PBT Index and US-EPA PBT Profiler.

The 83% of pharmaceuticals is estimated as non PBT by consensus, thus pharmaceuticals seem environmentally safe. Only 35 pharmaceuticals are predicted as PBTs by consensus and are included in the priority list (table 3.1). Further experimental tests should be focused on these compounds since they could exhibit an intrinsic high potential to be hazardous contaminants once released in the environment. Pharmaceuticals in the priority list belong to a large variety of pharmacological classes, mainly represented by sedative, hypnotic, anxiolytic and antipsychotic drugs, antihistamines, and antifungal agents. Moreover, a wide variety of molecular structures are present: the most represented groups are piperazine derivatives and diphenyl-butyl-piperidine, imidazoles, benzo-imidazoles and -triazoles, and dibenzo-azepines. The main common structural feature that characterizes all these chemicals is the high number

Results and Discussion

of double bonds and aromatic systems. It should be noted that this analysis is only related to the hazard assessment of pharmaceuticals, without considering production, consumption and occurrence data.

Table 3.1: Priority list of pharmaceuticals with potential PBT behavior by consensus

CAS	Name	Pharmaceutical Class	ATC Code
000050-53-3	Chlorpromazine	Psychoactive	N05AA01
000053-19-0	Mitotane	Antineoplastic	L01XX23
000070-30-4	Hexachlorophene	Antiseptic	D08AE01
000082-95-1	Buclizine	Antihistaminic	R06AE01
000083-89-6	Quinacrine	Antimalarial; anthelmintic	P01AX05
000097-18-7	Bithionol	Anthelmintic	D10AB01; P02BX01
000113-59-7	Chlorprothixene	Psychoactive	N05AF03
000146-54-3	Triflupromazine	Psychoactive	N05AA05
000298-57-7	Cinnarizine	Antihistaminic	N07CA02
000303-49-1	Clomipramine	Psychoactive	N06AA04
000390-64-7	Prenylamine	Calcium antagonist	C01DX02
000569-57-3	Chlorotrianisene	Non-steroidal estrogen	G03CA06
000569-65-3	Meclizine	Antihistaminic	R06AE05
000911-45-5	Clomifene	Estrogen agonist	G03GB02
000911-65-9	Etonitazene	Analgesic	-----
000915-30-0	Diphenoxylate	Antidiarrheal	A07DA01
001841-19-6	Fluspirilene	Psychoactive	N05AG01
001951-25-3	Amiodarone	Antianginal; antiarrhythmic	C01BD01
002030-63-9	Clofazimine	Treatment of leprosy anti-inflammatory	J04BA01

CAS	Name	Pharmaceutical Class	ATC Code
002062-78-4	Pimozide	Psychoactive	N05AG02
002398-96-1	Tolnaftate	Antifungal agent	D01AE18
003861-76-5	Clonitazene	Analgesic	-----
014426-25-6	Gentian Violet	Antiseptic	D01AE02; G01AX09
022916-47-8	Miconazole	Antifungal agent	A01AB09; A07AC01; D01AC02; G01AC02; G01AF04; J02AB01; S02AA13
023593-75-1	Clotrimazole	Antifungal agent	A01AB18; D01AC01; G01AF02
052468-60-7	Flunarizine	Calcium antagonist	N07CA03
053179-11-6	Loperamide	Antidiarrheal	A07DA03
054143-55-4	Flecainide	Antianginal; antiarrhythmic	C01BC04
068844-77-9	Astemizole	Antihistaminic	R06AX11
079617-96-2	Sertraline	Psychoactive	N06AB06
079794-75-5	Loratadine	Antihistaminic	R06AX13
084449-90-1	Raloxifene	Estrogen agonist	G03XC01
084625-61-6	Itraconazole	Antifungal agent	J02AC02
129722-12-9	Aripiprazole	Psychoactive	N05AX12
155213-67-5	Ritonavir	Antiretroviral	J05AE03

Furthermore, it is interesting to highlight that some substances reported in table 3.1 have been frequently detected in various environmental media and organisms and many

Results and Discussion

monitoring studies seem to support our analysis, mainly in relation to persistence and bioaccumulation assessment^{230–239}. Moreover, there are some experimental evidences on properties and activities, which can additionally confirm their PBT behavior. For example, the antifungal agent clotrimazole was detected in aquatic environment, in biosolids and thus may be persistent in agricultural soils^{230–233}. Clotrimazole is included into OSPAR Commission's list of substances for priority action²⁴⁰, it is defined as non-biodegradable with high acute and chronic toxicity with endocrine disrupting effects on fish²⁴⁰. Sertraline, a selective serotonin reuptake inhibitor, has been detected in WWTP effluent, in surface waters and even in benthic fauna and fish tissues^{234–236}. It has been demonstrated to be bioaccumulative and his presence in the benthic fauna can be a further exposure pathway for higher predator²³⁴. Moreover, this antidepressant drug is toxic to aquatic organisms with significant effects in embryonic and larval development in fish and sea urchin²³⁷. Also loratadine, itraconazole, cinnarizine and miconazole have been detected in surface waters and WWTP sludge^{238,239}.

The analysis of the disagreement between the models (i.e. Overestimated by Index and Overestimated by Profiler in figure 3.1) is also interesting. In these cases, the active pharmaceutical ingredients are predicted as PBT only by one model while the other estimates them as of not particular concern. Indeed, for a correct assignment to the PBT or non-PBT category is fundamental to look for data that confirm estimations generated by one or the other model. Therefore, overestimated pharmaceuticals can be placed in a second priority list (second level of prioritization) where experimental test are required to confirm or exclude their potential PBT behavior.

Paper I confirms the consistency of the consensus approach between different independent QSAR models for priority setting of CEC. Results demonstrate a high agreement (i.e. 86%) between the Insubria PBT Index and the US-EPA PBT Profiler.

Concluding, pharmaceuticals are a new generation of CEC that can affect wildlife and ecosystems even at very low concentrations. Stated that more than 5000 active pharmaceutical ingredients are currently on the market without an adequate

environmental risk assessment²⁴¹, efficient screening systems are needed to highlight the compounds of priority interest. The main result of this paper is to have demonstrated that QSAR methodologies for general preliminary screening can successfully be applied to pharmaceuticals independently from their specific mode of action. In general, pharmaceuticals ingredients seem not be PBTs and a priority list of 35 pharmaceuticals is proposed for a further deeper evaluation. Finally, it is important to note that this approach can be applied to existing chemicals without experimental data, to fill data gap, and to not yet synthesized molecules to plan environmentally safer chemicals from their design.

Further improvement and refinement of the PBT assessment for pharmaceuticals, and in particular of some aspects related to the B and T evaluation, are the subject of the next papers written within this PhD project.

3.2 Paper II: “Hazard of pharmaceuticals for aquatic environment: Prioritization by structural approaches and prediction of ecotoxicity”

By Sangion Alessandro and Gramatica Paola

Published in: *Environment International*, 2016, 95(10): 131–43.

This paper builds on the main findings of **paper I** and aims to extend the toxicity assessment of pharmaceuticals. In fact, the toxicity assessment reported in **paper I** is based only on the evaluation of the acute toxicity in one aquatic trophic level. A comprehensive environmental risk assessment should encompass the evaluation of adverse effects in organisms at different trophic levels, including at least assessment for primary producers, primary consumers and higher predators. Therefore, the aim of this paper is to propose QSAR models to estimate the acute toxicity of pharmaceuticals in species at different level of the aquatic trophic chain. These models should be intended for ranking and priority setting purposes, thus the selected end-point is the acute toxicity that gives an immediate feedback of the hazard of chemicals. This approach can be used to prioritize pharmaceuticals according to their environmental toxicity and to refine the toxicity assessment of PBT screening performed in **paper I**. Chemicals prioritized should undergo further evaluation including the assessment of long-term adverse effects such as endocrine disruption, behavioral disturbances and other health effects that are not taken into account in this study.

Crucial point of this paper is the process of data collection and curation, to build consistent datasets for QSAR modelling. Starting from the library of more than 1200 pharmaceuticals assembled in **paper I**, consistent data according to the “OECD guidelines for testing of chemicals”^{146–148} are found only for few dozens of pharmaceuticals²⁴². This highlights the scarcity information for the acute aquatic toxicity of medicines and the need of an approach to prioritize drugs for further experimental tests.

Validated QSAR models based on theoretical molecular descriptors are developed to predict the acute toxicity to four different standard aquatic species (i.e. *P. subcapitata*, *D. magna*, *O. mykiss* and *P. promelas*). The GA implemented in QSARINS^{166,167,203} is applied as variable subset selection method and internal and external validation are performed on three different splitting schemes. Finally, the full model (i.e. keeping all the chemicals in the training set) for each species is recomputed on the best externally predictive variables combination. All the proposed models have good fitting performances (R^2 : 0.75-0.81), are internally robust (Q^2_{LOO} : 0.70-0.76; Q^2_{LMO} : 0.65-0.75) and externally predictive (CCC_{EXT} : 0.82-0.89, Q^2_{EXT} : 0.68-0.86). The average RMSE of the developed models is 0.6 log units and, in general, the accuracy of the models is similar both for training sets and external prediction sets indicating a high generalizability of the model. The error of the model in prediction is compared to the error in prediction of ECOSAR v.1.11⁸⁶. All our models show RMSE values in prediction significantly lower than those obtained by ECOSAR that reports an average value of 1.46. This comparison demonstrates that QSAR models, specific for pharmaceuticals, are more able to give reliable predictions, than models developed on training sets composed of chemicals that are structurally highly different from pharmaceuticals¹⁰⁰.

The full models are then applied to the whole set of 1267 molecules collected in **paper I** to predict aquatic toxicities of pharmaceuticals and to fill the data gap. A careful evaluation of the AD of the models is carried out in order to focus the analysis only on the reliable predictions. Overall, all the full models have a wide structural AD and can provide reliable predictions for most of the chemicals in the dataset (i.e. 74% *P.subcapitata*, 87% *D. magna*, 96% *O.mykiss*, 80% *P.promelas*). PCA is applied on interpolated predictions of the toxicities in different species (figure 3.2).

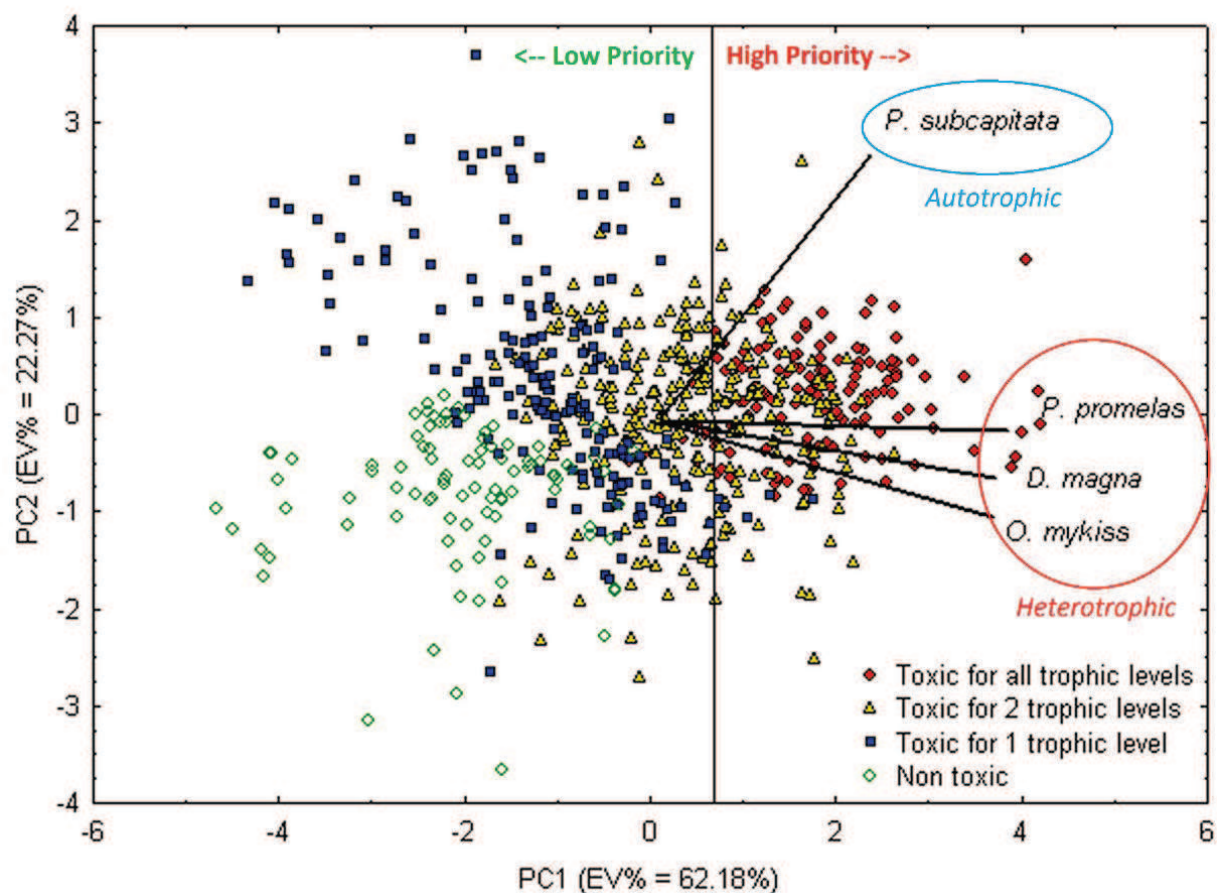


Figure 3.2: PCA biplot of the four studied toxicities and identification of an Aquatic Toxicity Index on the x axis (PC1 score).

The PC1 (Explained Variance 62.18%) shows that the toxicities in all the different organisms are positively correlated and is therefore useful to rank the pharmaceuticals according to their cumulative aquatic toxicity. The PC1, can be considered as a new cumulative end point which is named Aquatic Toxicity Index (ATI) for pharmaceuticals. This index is able to rank pharmaceuticals according to their overall toxicity on the whole aquatic environment. It is interesting to highlight that 34 out of 35 pharmaceuticals highlighted as potential PBT in **paper I** have also potential high toxicity on all the studied aquatic organisms.

Finally, the ATI is modelled by QSAR methodologies. Multiple splitting are used to select the best variable combination and the full models recalibrated on the whole set of

pharmaceuticals is proposed to predict the ATI for new active pharmaceutical ingredients:

$$ATI_{Pharmaceuticals} = -1.87 + 0.62CrippenLogK_{OW} + 0.07SaaCH - 0.05SHBint2 \quad (3.1)$$

$$n_{TR}: 706; R^2 0.81; Q_{LOO}^2 = 0.81; Q_{LMO}^2 = 0.81; R_{Yscr}^2 0.004; RMSE_{TR} 0.68; RMSE_{CV} 0.69$$

CippenLogK_{OW} is a logK_{OW} calculated with the Crippen algorithm¹⁶⁵, SaaCH encodes for the sum of the E-States of the CH groups bounded to two aromatic atoms and SHBint2 encodes for the sum of the E-State of strength for potential hydrogen bonds. These descriptors together give information about the lipophilic behavior and about the interactions that a chemical can have within the organisms^{243–246}.

Moreover, it should be highlighted that all compounds are considered in the neutral form as a common practice in computational chemistry when screening are performed on large inventories and hundreds of structures need to be compared in a standard way. This may represent a limitation of the approach since many pharmaceuticals are ionizable and may be present in different ionic forms according to the pH. The estimated toxicity should be read as preliminary indicator of the toxicity of the chemicals for the aquatic environment and is suitable for ranking purposes. Of course, ionic forms as well as specific mode of action should be considered when the single chemical assessment on the priority molecules is performed.

Concluding, this study relies on the main findings of **paper I** and tries to elaborate a new system for the refinement of the ecotoxicity assessment. New *ad-hoc* QSAR models for the estimation of acute toxicity of pharmaceuticals in aquatic species at different trophic levels are presented and applied to improve the toxicity assessment of pharmaceuticals screened by general models in **paper I**. The main result of this paper is the application and integration of new highly predictive QSAR models ($Q_{EXT-FN}^2 > 0.70$) in a coherent framework to identify a global index able to prioritize pharmaceuticals according to their general aquatic acute toxicity estimated only on the basis of the molecular structure. Acute toxicity is selected as target end-point because is able to give an immediate perception of the hazard of pharmaceuticals for the aquatic environment. This index is

Results and Discussion

intended for preliminary screening and priority setting purposes on the basis of general structural properties such as lipophilicity and hydrogen bonds interactions. It does not take into consideration specific mode of action of chemicals that should be taken into account into further and more detailed specific studies.

3.3 Paper III: “Ecotoxicity interspecies QAAR models from *Daphnia* toxicity of pharmaceuticals and personal care products”

By Sangion Alessandro and Gramatica Paola

SAR and QSAR in Environmental Research, 2016, 27 (10): 781–98

The purpose of this study is to propose new externally validated quantitative interspecies estimation models to predict the toxicity of pharmaceuticals and personal care products (PPCPs). Interspecies correlation estimation models are log-linear least squares regression models that describe the relationship between the acute toxicity of a range of chemicals tested in two species²⁴⁷. These models estimate the acute toxicity (LC₅₀/LD₅₀) of a chemical to a species, genus, or family with no test data (the predicted taxon) from the known toxicity of the chemical to a species with test data (the surrogate species). These models can be applied if data are available only for a species and the relationship between the two taxa is known. The rationale is to use the toxicity data of PPCPs in *D. magna* for the prediction of toxicity toward *P. promelas* and *O. mykiss*. Such models can be utilized to derive toxicity on organisms belonging to high trophic levels from data on less complex organism belonging to lower trophic levels^{247–251}.

Data are collected from **paper II** and integrated with new experimental data from our previous study on personal care products²⁵² and from the database ECOTOX²⁴². Three datasets are compiled: *D.magna-O.mykiss* (i.e. case 1), *D. magna-P.promelas* (i.e. case 2) and *P.promelas-O.mykiss* (i.e. cases 3 and 4).

For each case, two different models are developed: one based only on the experimental response of the surrogate species, by simple linear regression and CV, and the other based on the experimental response plus a molecular descriptor, selected among all the possible combinations by the all-subset procedure in QSARINS. Moreover, external validation is performed for MLR models on multiple splitting.

Results and Discussion

Simple linear regression is performed between the toxicities in the different species to find the line that best fit the data. In general, low values of toxicity in *D.magna* corresponds to low values of toxicity in the fish species, while chemicals highly toxic to *D.magna* are also highly toxic to fish. The same relationship is visible between the two fish species that seem to have similar response to the same xenobiotic. It is important to note that this relationship is a simple correlation relationship not a causation relationship since none toxicity is directly caused by the others, but all are intrinsically related to the molecular structure of chemicals. The correlation coefficients (r) are greater than 0.9 identifying a strong positive correlation between the various toxicities.

MLR QSAR is performed selecting theoretical molecular descriptors in addition to the single value of toxicity (figure 3.3). This is done to improve the quality of the interspecies correlation models; the inclusion of information related to the chemical structure might better explain the relationship^{249,253}.

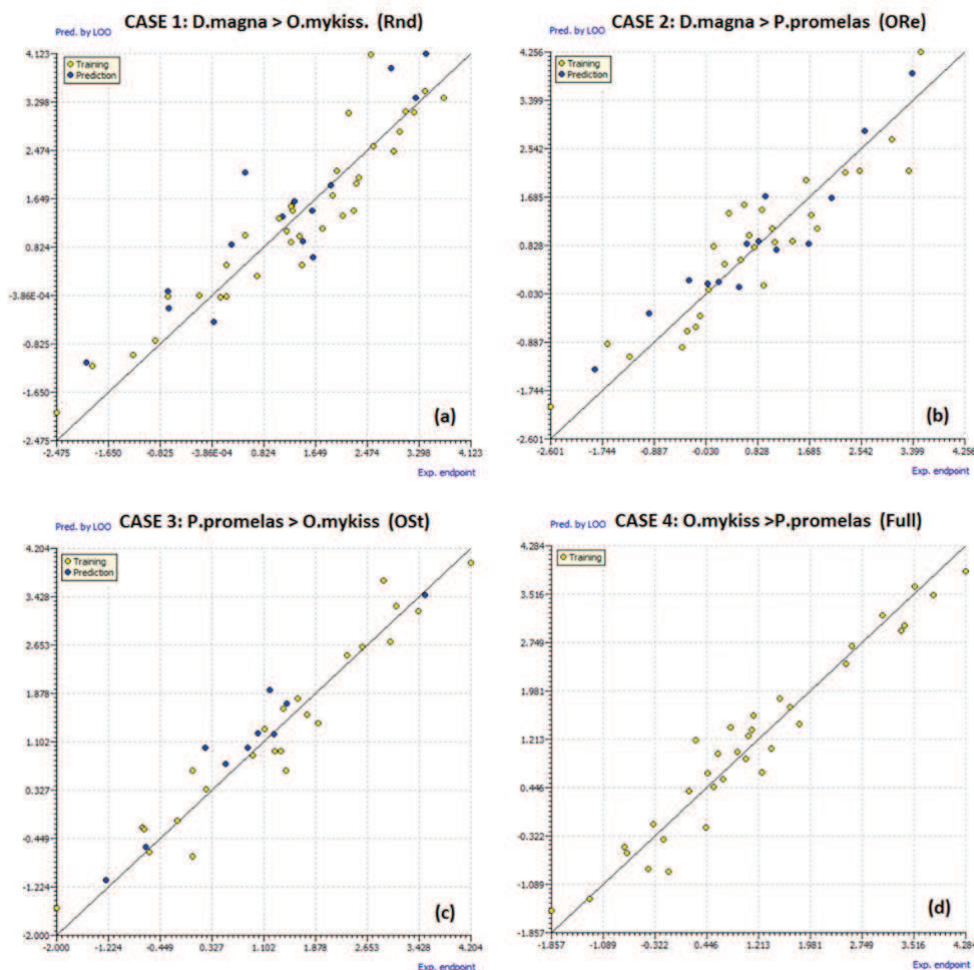


Figure 3.3: (a) Plot of experimental vs predicted data for Case 1, Random (Rnd) splitting; (b) Plot of experimental vs predicted data for Case 2, Ordered by Response (ORe) splitting; (c) Plot of experimental vs predicted data for Case 3, Ordered by Structure (OSt) splitting; (d) Plot of experimental vs predicted data for Case 4.

The addition of one molecular descriptor in the regression improves the ability of the model of fit the experimental data. R^2 values rise to an average of 0.89 for both *D.magna*-fish relationships (cases 1 and 2; figure 3.3 a and b) and an average of 0.95 for the fish-fish relationships (cases 3 and 4; figure 3.3 figure c and d). The internal validation qualities increase due to the use of chemicals descriptors. The CV parameters for the LOO procedure are always greater than 0.83 and very close to the R^2 values, guaranteeing the model stability. The inclusion of the structural information improves the model's accuracy and the RMSE values for the majority of the models are lower than 0.5 log-unit. Values for Q^2_{EXT} are greater than 0.7 and in general Q^2_{F1} , Q^2_{F2} and Q^2_{F3} are

Results and Discussion

close in the same model. This means that all the external sets are well balanced and the external predictivity is not affected by distribution issues. The values for the CCC_{EXT} are always above the suggested cut-off value of 0.85^{211,212}. Results for the external validation suggested that all the developed models are externally predictive, and the variables combinations selected are actually able to predict the studied endpoint. The descriptors selected are topological descriptors related to the molecular polarizability and account for information regarding the tendency of the chemicals to interact with the external environment.

Concluding this study assesses the correlation between the toxicity of PPCPs in *D.magna* and fish and demonstrates that *D.magna* could serve as surrogated species for fish. QAAR models are proposed for the extrapolation of the toxicity from *D.magna* to fish and can be used to reduce experimental test on organisms of higher trophic levels (i.e. fish) when information for the lower trophic level is available. They can be successfully applied in early screenings and may be useful to support the environmental risk assessment procedure required by REACH and other legislations^{53,254}.

3.4 Paper IV: “QSAR modeling of cumulative environmental endpoints for the prioritization of hazardous chemicals”

By Gramatica Paola, Papa Ester and Sangion Alessandro

Published on: *Environmental Science: Processes & Impacts*, 2018, 20 (1): 38–47.

This is a perspective paper that reviews the approach based on the combination of PCA and QSAR models applied to screen hazardous properties of environmental contaminants. Such approach can be useful for the identification of hazardous compounds and for planning the *a priori* synthesis of safer alternatives to chemicals with undesired side effects (e.g. PBT behavior) according to the benign by design approach of Green Chemistry⁷⁷. This paper explains the rationale of using the PCA as ranking technique to prioritize environmental compounds for further experimental tests. The behavior of chemicals in the environment and their impact on human and wildlife depends on properties inherent in the molecular structure such as physical-chemical properties, chemicals reactivities and biological activities. These properties could be related to each other and their cumulative effect contributes to the environmental fate and biological activity of chemicals. PCA combines the available information in new variables that explain the principal co-variation in the original data given by the interaction of all the properties. These new variables can be used to extract the meaningful information and understand the main properties of the data discharging the background noise⁸⁰. Moreover, these variables can be used as dependent variables in MLR QSAR models to find the relationships with the molecular structure.

This paper summarizes some case studies such as the development of the Insubria PBT Index and its application in the screening of PBT properties of CECs such as flame retardants¹⁰⁴, personal care products¹⁰⁵ and pharmaceuticals (i.e. **paper I**) or the development of the ATI for personal care products²⁵² and pharmaceuticals (i.e. **paper II**). The paper shows how the integration of multiple approaches is useful to address environmentally relevant macro properties of chemicals.

3.5 Paper V: “Development of human biotransformation QSARs and application for PBT assessment refinement”

By Papa Ester, Sangion Alessandro, Arnot Jon, Gramatica Paola

Published in: *Food and Chemical Toxicology*, 2018, 112(2): 535–43.

This study addresses the development of QSAR models for the prediction of *in-vivo* whole-body human biotransformation half-lives (HLs) measured or empirically-derived for over 1000 organic chemicals¹²⁶. Biotransformation is a fundamental component of bioaccumulation; indeed it is reasonable to expect that only compounds with long metabolic HLs in higher organisms (e.g. fish and mammals) can actually persist unchanged, bioaccumulate, and thus have the possibility to manifest their toxic effects in living organisms⁶². Information about biotransformation potential can reasonably be used to refine the bioaccumulation estimation. This is particularly evident in case the bioaccumulation assessment is achieved only on the basis of the BCF as in the general screening performed in **paper I**. In this case is necessary to include biotransformation estimation to obtain a better estimation of the bioaccumulation potential. The QSARs developed in this work are used to refine the B evaluation of the PBT assessment of personal care products¹⁰⁵ and pharmaceuticals (i.e. **paper I**). This work aims to evaluate the intrinsic hazard related to the B potential depending on the molecular structure, thus information on environmental persistence and exposure are not considered.

Data for whole-body total elimination half-life (HL_T) and for whole-body primary biotransformation half-lives (HL_{B1-4}) in humans are collected from literature¹²⁶. Chemicals included in the datasets are mainly pharmaceuticals (i.e. 80%), however other environmentally relevant compounds are also included (e.g. PAHs, polychlorinated biphenyls (PCBs), polychlorinated dibenzodioxins (PCDDs) and polychlorinated dibenzofurans (PCDFs)).

The GA implemented in QSARINS is applied to each dataset to generate independent populations of models based on different combinations of theoretical molecular

descriptors. A final model is selected for each dataset on the basis of fitting and predictive qualities as well as interpretability of the descriptors. The statistical parameters calculated for the models reflect the good ability to fit the data in the training set (R^2 : 0.77 – 0.80). Cross validation parameters are relatively high and close to each other (Q^2_{LOO} : 0.76-0.79; Q^2_{LMO} : 0.76-0.79). This means that all the models are stable and the internal predictivity is conserved independently of the degree of perturbation applied to perform the internal validation. The statistical parameters used to quantify the quality of the external validation (Q^2_{EXT} : 0.75-0.79; CCC_{EXT} : 0.86-0.87) confirmed the good predictive ability of the models.

The most relevant descriptors are: the number of halogen atoms (nX), the sum of the Electrotopological State (E-State) of chlorine atoms (SsCl)²⁴⁴, the average Broto-Moreau autocorrelation of lag 7 or 8, weighted by polarizabilities (AATS7p or AATS8p) and the minimum and the maximum E-State of hydrogen atoms on the “-OH” (minHsOH and maxHsOH, respectively)¹⁶⁵. From a general point of view, biotransformation potential seems to be influenced by the presence of halogen atoms, which are covalently bonded mainly to aromatic carbon atoms, and by the ability of the chemicals to establish non-covalent intermolecular interactions. The presence of halogen atoms in the molecules increases biopersistence (i.e. slow biotransformation and long HL); the covalent bonds between aromatic carbon and halogen atoms are very stable and may induce persistence and biopersistence of chemicals^{87,106,255}. On the other hand, the potential to form hydrogen bonds can give information about the tendency of the molecules to interact with the active sites of metabolic enzymes in the surrounding media^{243,244}.

Results and Discussion

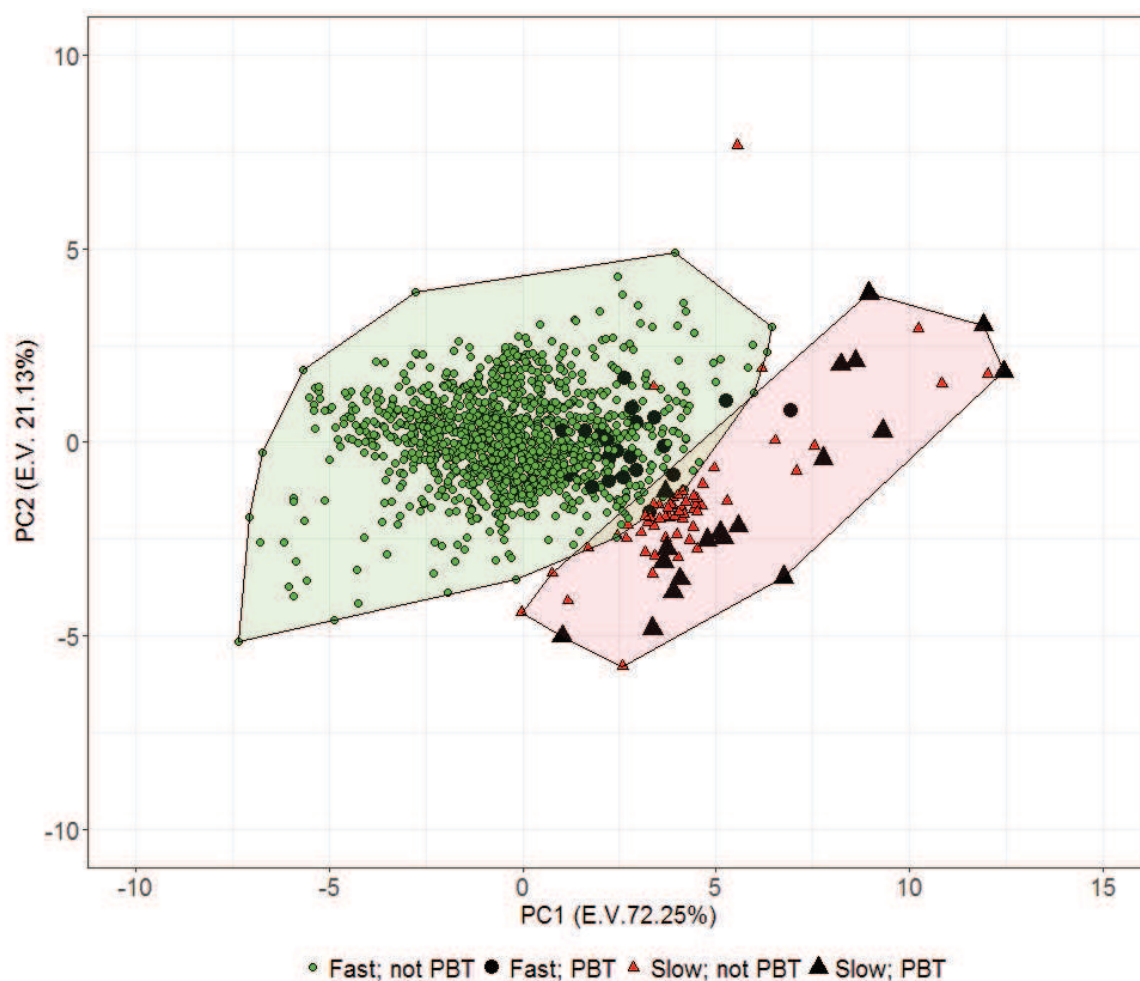


Figure 3.4: Distribution of chemicals in PC1-PC2 space of the different HLs and results of the analysis with the two polygons.

The described QSAR models developed for biotransformation potential in humans in addition to literature QSARs for the prediction of the metabolic biotransformation half-lives in fish (Papa et al., 2014), are used to refine the PBT assessment of a set of more than 1200 PPCPs (i.e. **paper I** and personal care products¹⁰⁵). Predictions of all the models are combined in a PCA to investigate the overall biotransformation behaviour in fish and humans. Generally, PC1 (EV 72%) identifies a horizontal trend of biotransformation potential with chemicals slowly biotransformed in human that are also slowly biotransformed in fish. On the other hand, PC2 (EV 21%) separates the biotransformation behaviour in the different organisms. Moreover, a spatial analysis is performed by the convex hull method^{256,257} to identify subspaces able to separate PCPPs with high bioaccumulation potential (figure 3.4). Among 40 chemicals predicted as

80

intrinsically PBTs by consensus (i.e. 35 pharmaceuticals in **paper I** and 5 personal care products¹⁰⁵) only 22 are placed in the slow-metabolized area and have high bioaccumulation potential in fish and humans. This result highlights the importance of biotransformation related information for the refinement of screening level assessments. In this case, 18 chemicals previously screened as intrinsically PBTs but predicted as easily metabolized in this second step should be considered as of less priority than the 22 PBTs estimated as slowly metabolized.

Concluding, this study presents new QSAR models for the prediction of *in-vivo* HL_B in humans from large empirically-derived datasets mainly representative for pharmaceutical products. Biotransformation HLs are key parameters determining the extent of bioaccumulation and biological concentration. Information about biotransformation potential can be use in the refinement of the B estimation in the PBT assessment of chemicals. In particular, pharmaceuticals seem to be easily metabolized in higher organisms and are not cause of particular concern in relation to bioaccumulation potential. Finally, these models can be integrated by a multivariate analysis to analyze the role of biotransformation potential in defining bioaccumulation. This will be the subject of next paper (i.e. **paper VI**) that aims to investigate the importance of biotransformation in the estimation of biomagnification potential in air-breathing organisms.

3.6 Paper VI: “A Tiered Approach for Screening Chemicals for Biomagnification Potential in Air-Breathing Organisms.”

By Sangion Alessandro, Arnot Jon, Papa Ester

The results reported in this section are preliminary and still unpublished. The manuscript is in preparation and a pre-view is reported here (September 2018).

Target journal: Environmental Science Technology or Environmental Health Perspective

3.6.1 Introduction

Regulatory programs evaluating chemicals for their potential detrimental effects to human and the environment^{12,54,108,258,259} classify chemicals according to their persistence, bioaccumulation and toxicity properties, e.g., bioaccumulative (“B”) or not bioaccumulative (“nB”). The octanol-water partition coefficient (K_{OW}), as a surrogate for aquatic biota-water partitioning, and bioconcentration factors (BCFs) in fish are metrics that have been used over the past few decades for bioaccumulation classification^{62,109}. Biomagnification Factors (BMFs) and trophic magnification factors (TMF) are also considered for “B” assessment to determine if concentrations (or fugacities) are increasing in predator-prey relationships or across the food-web, respectively²⁶⁰. Nearly two decades of empirical evidence and models^{63,68,113,114,128,261–263} have highlighted the need to examine bioaccumulation in air-breathing organisms separately from water-respiring organisms. Mechanistic understanding of the differences in bioaccumulation between these two classes of organisms is established^{63,114,128}. For air-breathing organisms, biomagnification is related to the octanol-air partition coefficient (K_{OA}) in addition to K_{OW} . According to emerging screening criteria⁵⁴ chemicals with $\log K_{OA} \geq 5$ and with a $\log K_{OW} \geq 2$ are considered to have bioaccumulation potential in air-breathing organisms, including humans. Based on an evaluation of approximately 13000 organic chemicals and these criteria ($\log K_{OA} \geq 5$ and $\log K_{OW} \geq 2$), Gobas and colleagues highlighted that approximately 50% have

biomagnification potential in air-breathing organisms but not in water-respiring organisms⁶⁸. Critical caveats to these physical-chemical property screening criteria cut-off values are the assumptions that (i) the chemical is efficiently absorbed from the diet, (ii) it is non-dissociating, and (iii) there is no biotransformation^{68,114,115,126}. Quantitative Structure-Activity Relationships (QSARs) for predicting biotransformation half-lives (HL_B) in humans^{126,264} and fish^{141–144,265} have recently been developed and validated following OECD QSAR guidance^{179,180}. Uncertainty exists for measured and modelled chemical properties and bioaccumulation metrics. Methods are needed to better identify sources of uncertainty in chemical assessments. In this study we use a series of model hypotheses to explore the degree to which key parameters can influence bioaccumulation screening in air-breathing organisms. We select a human male adult as a representative air-breathing organism and we use a one-compartment mass balance toxicokinetic model to calculate BMFs. We compile and critically curate a database of approximately 30'000 substances including industrial chemicals, pharmaceuticals, personal care products and chemicals used in consumer goods. We then obtain chemical information for assessing bioaccumulation potential in humans including: molar mass, K_{OW}, K_{OA}, dissociation constants, adipose-water, membrane-water and protein-water partition coefficients and HL_B from various databases and QSAR models and define the applicability domains for these parameters. Finally, we apply a tiered approach for screening the organic chemicals for BMFs > 1 in air-breathing organisms. The tiered approach progresses from conservative assumptions to more realistic assumptions for chemical properties, biological partitioning and biotransformation.

3.6.2 Materials and methods

General BMF Model for Air-Breathing Organisms. The uptake and elimination of xenobiotics in a mammal can generally be expressed by the following mass balance equation^{63,68,114,126,127,266}:

$$\frac{dC_B}{dt} = k_{RI}C_{AG} + k_D C_{D,i} + k_W C_W - (k_{RO} + k_E + k_R + k_B + k_{RI} + k_G + k_D)C_B \quad (3.2)$$

Results and Discussion

where dC_B/dt is the net change in concentration in the organism (g/kg) over time t (h), C_B is the concentration of the xenobiotic in the organism, k_{RI} is the respiration intake rate constant ($L(kg \cdot h)^{-1}$), C_{AG} is the air concentration ($g \cdot L^{-1}$), k_D is the xenobiotic dietary uptake rate constant ($kg(kg \cdot h)^{-1}$), C_D is the concentration of the xenobiotic (g/kg) in diet (food), k_W is the water uptake rate constant ($L(kg \cdot h)^{-1}$), and C_W is the concentration of the xenobiotic (g/L) in the water. The rate constants (h^{-1}) corresponding to xenobiotic elimination from the organism via respiratory elimination (output), fecal egestion, renal excretion, biotransformation, reproductive losses, growth dilution are k_{RO} , k_E , k_R , k_B , k_{RL} , k_G , respectively. Changes in biomass (growth) are not a true elimination process but affect the concentration of the chemical in the organism. The total elimination rate constant k_T is the sum of individual elimination rate constants. Assuming first-order kinetics the total elimination half-life is $HL_T = \ln 2/k_T$. The biotransformation half-life is $HL_B = \ln 2/k_B$. The steady state BMF can be calculated as a ratio of the chemical concentration in the organism to that in its diet, i.e., C_B/C_D . Lipid normalization is recommended for calculating BMFs for neutral organic chemicals and an equivalent expression is the fugacity ratio of the organism to its diet, i.e., f_B/f_D ²⁶⁰.

In this illustrative case study, we used the Risk Assessment IDentification And Ranking (RAIDAR) model⁸⁹ to calculate BMFs in an adult human male as a model air-breathing organism. The adult male is assumed to ingest a typical North American (omnivorous) diet and reproductive losses are set to zero. The details of the equations, i.e., rate constants, and model parameters are available elsewhere⁸⁹.

Case Study Chemicals. We compiled an initial database of approximately 30'000 discrete organic substances including industrial chemicals, pharmaceuticals, personal care products and chemicals used in consumer goods^{104,105,267}. Because the chemical structures contained many duplicates and inconsistencies in the molecular structures and identifiers, a data curation process was necessary to derive a unique set of QSAR ready structures. Information on CAS, Chemical name and SMILES were verified, invalid entries were corrected by retrieving information from online databases PubChem²⁶⁸ and

ChemSpider²⁶⁹ in the Package Webchem²⁷⁰ and ChemmineR²⁷¹ of the software R²⁷². Once the information was verified, chemical structures were converted to canonical SMILES using OpenBabel v.2.3.2²⁷³. This allowed compounds with the same molecular structure but different SMILES string (e.g. different atom numbering) to be compared properly. The structures were desalted, and counterions were removed, explicit hydrogen atoms were added and charged structures were neutralized. Chiral marks and stereochemistry information were removed and nitro groups were standardized (i.e. from $-\text{N}^+([\text{O}^-])=\text{O}$ to $-\text{N}(=\text{O})=\text{O}$). Duplicated structures were removed while information for salt and chiral compounds were collapsed with the information of the neutral parent compound. At the end of this procedure, 22'542 unique structures remained in the dataset.

Chemical Properties. The unique SMILES strings were used to estimate physical-chemical properties and HL_B . The K_{OW} , K_{OA} were estimated using the EPI Suite™ software package (version 4.1)⁸². The software contains an extensive database and, when possible, experimental values were preferred to the estimated ones. EPI Suite estimates physical-chemical properties using fragment-based QSAR models and does not provide a clear estimation of the Applicability Domain (AD). However, we identified and considered less accurate prediction for compounds outside the MW range and/or response range of the training set and/or compounds with more instances of a given fragment than the maximum for all training set compounds. The biotic partition coefficients storage lipids-water (kStorLipW), phospholipid membrane-water (kMembLipW) and proteins-water (kProtW) were estimated by the UFZ-LSER database²⁷⁴. The database uses two different sets of polyparameters-Linear Free Energy Relationship (pp-LFER) equations to estimate the biotic repartition coefficient at 37 °C²⁷⁵. For these models, applicability domain was estimated according to structural similarity to chemicals in the training data set¹⁴². In these simulations we used a consensus modelling (i.e. average of the predictions) of the two equations weighted by the AD estimation. Biotic-partitioning coefficients should represent in a more precise

Results and Discussion

way the differences in the sorptive capacity of tissues and organisms. In particular, kStorLipW and kMembLipW account for different types of lipids (i.e., triacylglycerides the first and phospholipids the latter) that were demonstrated to have different sorptive capacity²⁴⁶ while kProtW accounts for the repartition in the structural proteins (i.e., muscle proteins) that for some compounds can be as high as that to storage lipid²⁷⁶. The ionization state at physiological pH (i.e. 7.4) and acid/base dissociation constants (pKa) were estimated using ACDLabs percepta²⁷⁷. Chemicals were considered neutral, acid or base according to their more predominant form at pH 7.4; compounds present as Zwitterionic were assumed here to be neutral.

Biotransformation half-lives. The human HL_B parameter was estimated using consensus modelling between Iterative Fragment Selection (IFS)-based and theoretical descriptor-based QSAR models^{166,167,278–281}. The consensus modelling was weighted on the AD of each single model; based on structural similarity, counting the degree of fragment overlap between two chemicals for the IFS-based models and considering the range of the descriptors, the range of the response and the leverage value for the descriptor-based models²²⁴. Where available, experimental data were selected preferentially to the predicted values.

After the data collection phase, the number of chemicals in the dataset was reduced to 20'346 for which we were able to obtain measurements or estimate for all the chemical properties required as input for the RAIDAR model human BMF calculations.

Tiered Approach – Model Parameterization. Figure 3.5 provides a conceptual overview of the tiered approach that progresses from simple, conservative assumptions to conditions that are more realistic:

- Tier 1 only considered the chemical partitioning properties (i.e. K_{OW} , K_{OA}), which are assumed to be aligned with evolving screening criteria in chemical regulations worldwide. All the chemicals are assumed to be in the neutral form,

persistent in the environment and their biotransformation rates are assumed negligible.

- Tier 2 has the same assumptions of Tier 1; however, partition coefficients for biological phases (k_{StorLipW} , k_{MembLipW} and k_{ProtW}) are introduced in the calculation in place of K_{OW} and K_{OA} . These ppLFERs are assumed to be more representative of biological partitioning than using octanol as a surrogate for biological phases.
- Tier 3 adds information about the ionization state at physiological pH 7.4. Chemicals that appreciably dissociated at pH 7.4 are considered “acid” or “base” depending on the related dissociation pKa. This allows the model to consider the IOC properties in the BMF calculation.
- In Tier 4 the selected consensus primary biotransformation half-lives are included in the model.

Simulation for each tier are run considering a screening-level standard scenario at level II, i.e. the environmental compartments are at equilibrium and at steady state. Chemicals are released with a constant rate and are supposed to be persistent in the environment.

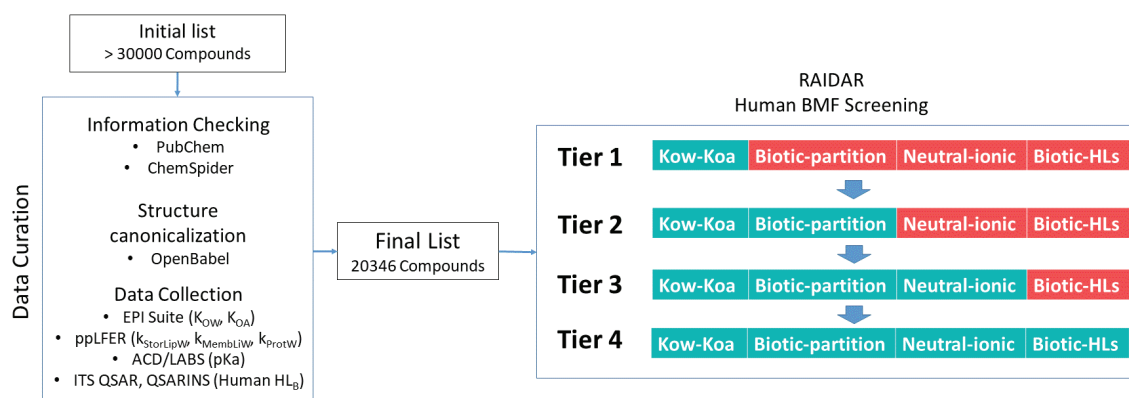


Figure 3.5: Outlines of the general methods used to perform the BMF screening for organic compounds

3.6.3 Results and discussion

Tier 1. Bioaccumulation screening criteria proposed or implemented in chemical assessment programs are based on physical-chemical properties, i.e., K_{OW} and K_{OA} ^{54,260}; therefore, we examined the information available for these two parameters for a large number of chemicals that may be subject to assessment. Summary details for the properties obtained for these chemicals are reported in figure 3.6.

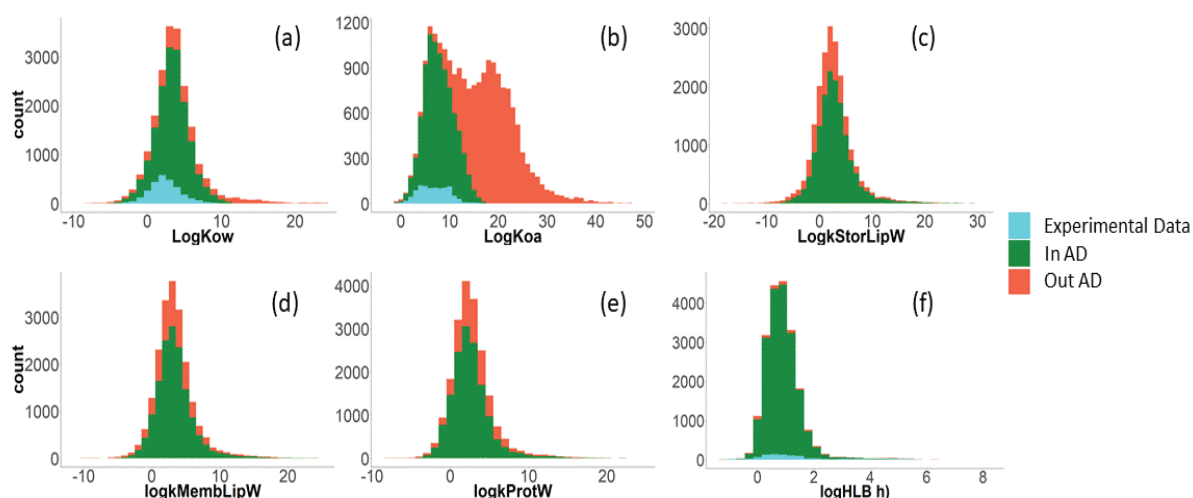


Figure 3.6: Histograms for the distribution of the studied parameters: (a) K_{OW} ; (b) K_{OA} ; (c) $k_{StorLipW}$; (d) $k_{MembLipW}$; (e) k_{ProtW} ; (f) human HL_B

Experimental log K_{OW} values are available for only 13% of the chemicals. The experimental data follow a log normal distribution (figure 3.6 a) with a minimum value at -11.96 and the maximum at 10.20. The lowest value is reported for “7-aminonaphthalene-1,3,6-trisulphonic acid” (CAS 118-03-6)⁸² and is isolated comparing to the rest of the experimental data. The closest experimental value in the dataset is -5.08 reported for “cystine” (CAS 56-89-3)⁸² while the range of experimental values of the training set of KOWWIN⁸² is between -5.4 to 10.20. Thus, it should be noted that the value of -11.96 for “7-aminonaphthalene-1,3,6-trisulphonic acid” can be considered an outlier and the actual range for the experimental values of the log K_{OW} is between -5.08 and 10.20. Sixty-eight percent of the case study chemicals have predictions within the AD of KOWWIN as defined module according to the methods developed here, i.e., AD

based on MW, estimated response and fragment count. The estimated data have a distribution like the experimental values with a range between -4.69 and 11.28. On the other hand, 19% of the chemicals were predicted outside the AD due to the presence of fragments not present in the training set or because the chemicals were outside the MW and/or the response range. Chemicals outside the AD tend to be estimated with “extreme” values with a range that spread from -16.6 to 26.88. This is probably due to the additive nature of the fragment model that increase (or decrease) the K_{OW} value for each fragment in the molecular structure²⁸². Chemicals outside the AD, in general, have higher number of instances of a given fragment than the maximum for all training set compounds and tend to result in higher predictions, outside the range of the experimental response. This highlights the importance of the analysis of the fragments for the applications of fragments-based models. Overall, we obtained 81% of chemicals with reliable log K_{OW} values (experimental or inside the AD) mainly in the range between -5 and 10. Generally speaking, results for chemicals with K_{OW} predictions outside this range and/or for chemicals outside the AD of KOWWIN should be considered more uncertain than the other chemicals in the case study.

Empirically-based K_{OA} data are available for 5% of the chemicals, the majority of which are based on experimental K_{OW} and Henry's Law constant and not direct measurements of K_{OA} as $K_{OA} \approx K_{OW}/K_{AW}$ (K_{AW} is dimensionless Henry's law constant). Thirty-seven percent of the chemicals have K_{OA} predictions within the AD of the model based on the methods we have developed in this study. The ranges of experimental data and of predictions within the AD is similar ranging between -1.14 and 16.47 and between -2.73 and 17.70, respectively. On the contrary, 58% of chemicals have predictions outside the AD of the model and most of these estimates are very high (e.g., 56 % with log $K_{OA} > 12.5$). A bimodal distribution is clearly visible in figure 3.6 b with the geometrical mean for chemicals outside the AD approximately 10 orders of magnitude greater than the geometrical mean of the experimental data (inside the AD). This means that KOAWIN⁸² tends to predict very high values for chemicals outside the AD with a structure different than the chemicals in the training set of the model; as for KOWWIN the “extreme”

Results and Discussion

values are probably due to the additive nature of the fragment-based model and should be considered carefully²⁸².

Figure 3.7 reports the projection of the whole dataset in the space of the partition coefficients K_{OW} and K_{OA} with the screening thresholds for air-breathing organisms. In particular 5'249 chemicals have a $\log K_{OW} < 2$; these chemicals are expected to be efficiently eliminated by passive renal clearance (assuming no active resorption in the kidney), and therefore do not biomagnify even though their $\log K_{OA}$ may be > 5 ^{63,68,128,261-263}. 1'655 molecules have a $\log K_{OA} < 5$; for air-breathing organisms, respiration is the main elimination route for these compounds and their elimination rate constant and bioaccumulation potential are more related to the K_{OA} than to the K_{OW} ^{63,63,68,128,261}.

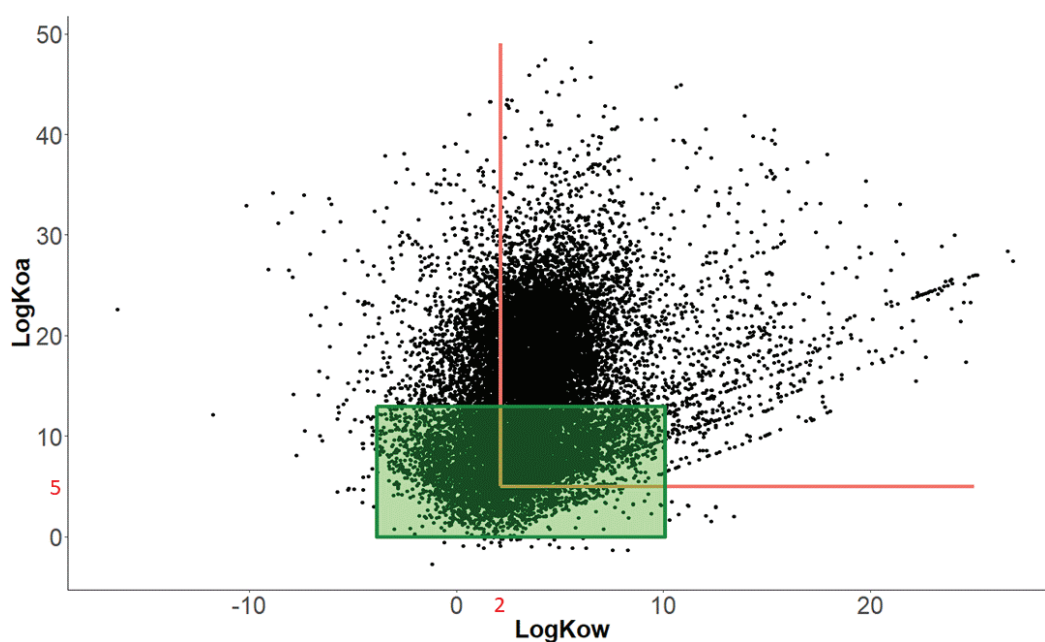


Figure 3.7: Relationship between K_{OW} and K_{OA} for the studied organic compounds. The graph identifies regulatory screening threshold for $\log K_{OW}$ (i.e. $\log K_{OW} > 2$) and $\log K_{OA}$ (i.e. $\log K_{OA} > 5$). The green box shows the chemical partitioning space in which K_{OA} and K_{OW} have been measured and hence indicates a region in which modelled BMFs are expected to be of higher confidence.

More than 70% of the chemicals in the curated list (i.e. 14'293 over 20'346) exceed both the screening level values (i.e. $\log K_{OW} > 2$, $\log K_{OA} > 5$). This means that according to the screening criteria alone, approximately 14'000 compounds have the potential to biomagnify in air-breathing organisms and may require further evaluation including perhaps *in-vitro* and *in-vivo* experimental tests.

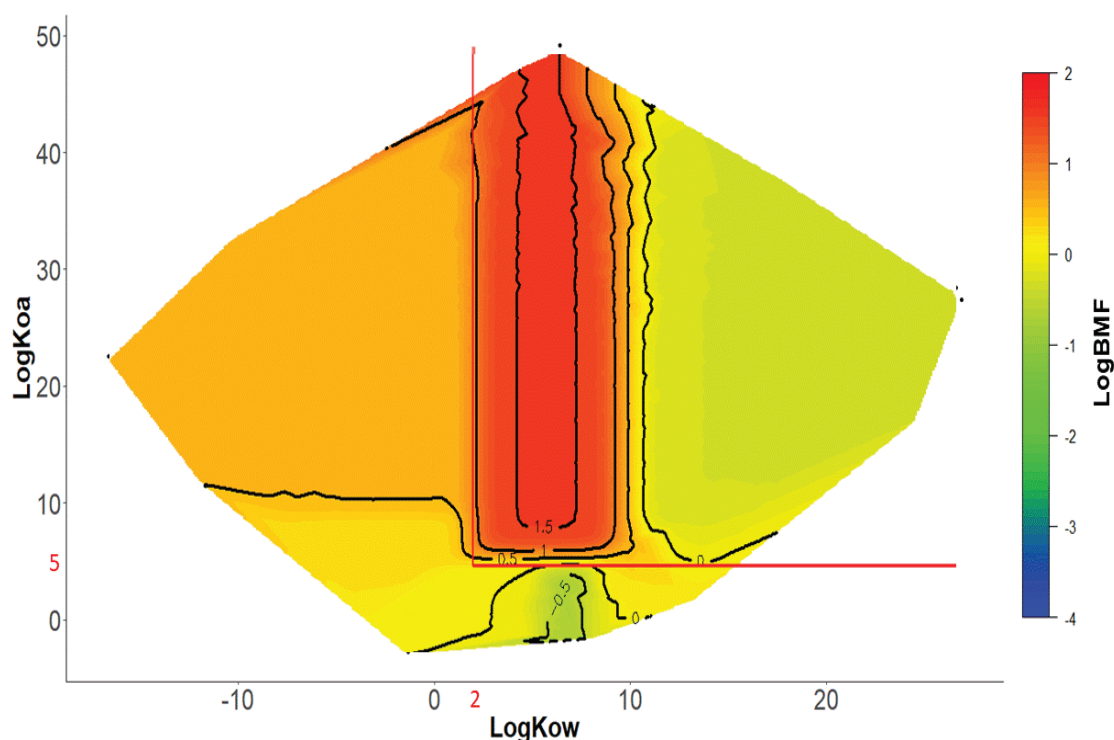


Figure 3.8: Heatmap for the $\log\text{BMF}$ respect $\log K_{OW}$ and $\log K_{OA}$ for Tier 1. The graph identifies regulatory screening threshold for $\log K_{OW}$ (i.e. $\log K_{OW} > 2$) and $\log K_{OA}$ (i.e. $\log K_{OA} > 5$) in red and $\log\text{BMF}$ isolines in black.

Figure 3.8 illustrates a heatmap with the three-dimensional interpolation of the model calculated $\log\text{BMF}$ using only $\log K_{OW}$ and $\log K_{OA}$ ²⁸³. Overall the shape of the three-dimensional interpolation agrees with the curves reported in previous BMF model simulations for air-breathing organisms^{63,114,128}. The highest values of BMF are estimated for chemicals with $\log K_{OW} > 2$ and < 10 and $\log K_{OA} > 5$. As with previous BMF model calculations, the BMF is predicted to decrease at higher K_{OW} , i.e., $\log K_{OW} > 7$, and chemicals with $\log K_{OW} > 10$ are estimated to have $\text{BMF} < 1$ due to slower rates

Results and Discussion

of chemical absorption from the diet. This predicted relationship is based on empirical measurements and models that show that when $\log K_{OW} > 7$ the dietary assimilation efficiency (E_D) decreases due to the increasing water phase resistance (i.e. the low internal water phase concentration) which slows transport from the gastrointestinal tract to the blood²⁸⁴. While there are limitations and uncertainty in the E_D data and models we include these concepts in the current case study to demonstrate how this expected relationship influences the B assessment for air-breathing organisms using real chemicals. It is important to note that values of $\log K_{OW} > 10$ are also extrapolated and outside the AD of KOWWIN and results for these compounds should be interpreted recognizing the inherent uncertainties in the chemical property and BMF model predictions. Chemicals with $\log K_{OW}$ ranging between 2 and 10 and $\log K_{OA} < 5$ have predicted BMFs < 1 . For these chemicals respiratory excretion is expected to mitigate dietary biomagnification.

The estimation of the BMF for chemicals with $\log K_{OW} < 2$ is unexpected based on comparisons to previous modelled BMFs for air breathing organisms⁶⁷. Chemicals with a $\log K_{OW}$ under this threshold are eliminated relatively efficiently by renal excretion and hence their BMFs are expected to be < 1 ^{63,114,128}. However, all these chemicals are predicted with a BMF ranging between 1 and 8 (\log BMF between 0 and 0.8 in figure 3.8). This apparent anomaly is explained by the assumptions in the current RAIDAR simulations. Here the steady-state concentrations in the human are based on chemical exposures from food (diet) and in drinking water (beverages); however, the BMF is only being calculated for exposures in the diet and not the total exposure, i.e., the slight underestimation of the value in the denominator results in a slight overestimation of the BMF. These minor anomalies have been observed in previous modelling and discussed in previous publications⁶³. These slight inconsistencies in the water mass balance and BMF calculations in the RAIDAR model are being addressed in a forthcoming update to the model (Arnot, pers. Comm).

Tier 2. In Tier 2 the biologically-based partition coefficients were considered instead of assuming octanol as a surrogate for partitioning (K_{OW} and K_{OA}) presumably to obtain more realistic estimates for chemical sorption in the body. The total body sorption capacity is given by the different contributions of sorption capacity of storage and membrane lipids, proteins and water that can be calculated by the introduction of partition coefficients for biological phases i.e. $k_{StorLipW}$, $k_{MembLipW}$ and k_{ProtW} ²⁸⁵. Lipids are the most important compartment for bioaccumulation of neutral hydrophobic organic chemicals, i.e., those considered to have significant biomagnification potential in air-breathing organisms. However, a distinction can be made between the different types of lipids. Storage lipids generally referred to as triglycerides are the predominant component of fat tissue. Membrane lipids referred as phospholipids, are the main component of biological membrane and are relatively rich in human tissues such as kidney, liver and brain²⁸⁶. In general, there are no large differences between the two kinds of lipids for non polar and H-bond acceptor compounds. Exceptions are H-bond donor that are more favorable to membrane lipids²⁷⁶. Other biological constituents such as proteins can have a significant contribution in the sorptive capacity of organisms. Moreover, protein-binding for some classes of contaminants such as perfluorocarboxylic acids and derivates, are implicated as key factors for bioaccumulation²⁸⁷. Considering a “typical” adult human is composed of water (60%), protein (20%), storage lipid (19%) and membrane lipid (1%) we estimated the internal distribution of the analyzed chemicals (figure 3.9).

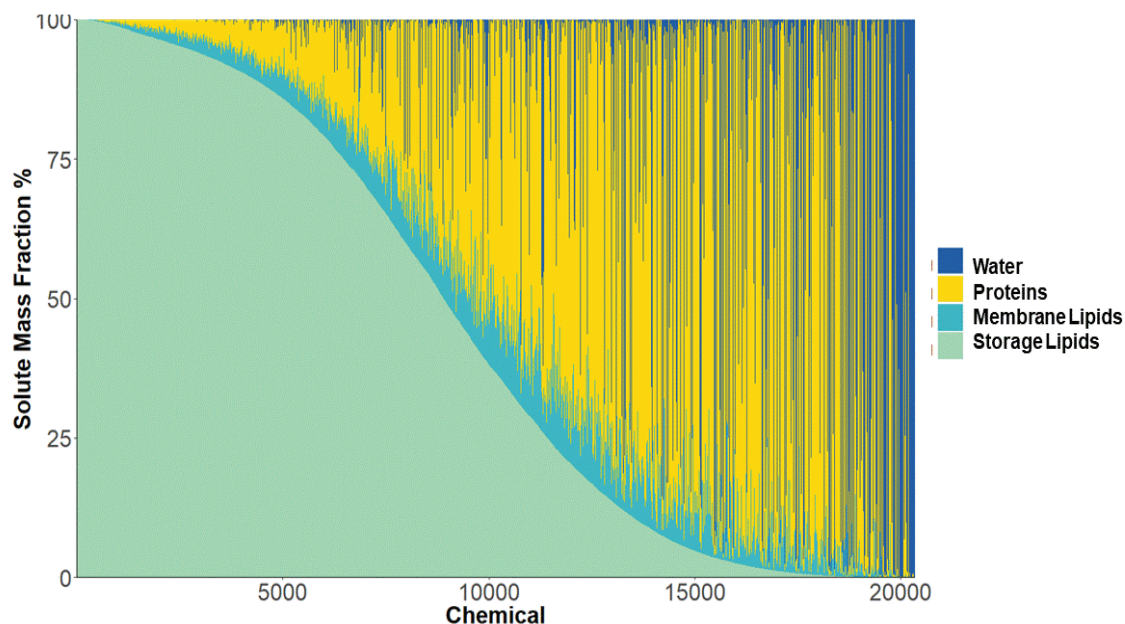


Figure 3.9: Relative sorption capacity (calculated as product of partition coefficients and relative volume of the matrix) of different matrices for the investigate chemicals sorted according to the affinity for lipids.

Lipids play an important role in the accumulation of chemicals with the 47% that tends to partition mainly in storage lipids with a small contribution of the membrane lipids. Protein is the second main bioaccumulation compartment with 39% of chemicals with high affinity structural proteins. Finally, water plays a not marginal role in the internal distribution with an 14% of less hydrophilic chemicals that tend to partition mainly in aqueous media. Figure 3.9 reports the relative distribution in the main sorptive compartments in the $\log K_{OW}$ - $\log K_{OA}$ space. Three main partitioning areas are highlighted by the ellipses (confidence at 0.95). As expected less hydrophobic chemicals with a $\log K_{OW} < \text{ca. } 2$ tend to partition mainly in water. For these chemicals, aqueous media is the most important compartment and renal excretion can be the primary elimination route. For chemicals with $\log K_{OW} > \text{ca. } 2$ the distribution within the body is strongly influenced by the K_{OA} . Chemicals with a $\log K_{OA}/\log K_{OW}$ ratio < 2.5 tend to accumulate mainly in lipids while protein is the predominant component for

bioaccumulation when $\log K_{OA}/\log K_{OW} > \text{ca. } 2.5$. Furthermore, the region with $\log K_{OW}$ ca.2 and $\log K_{OA} > 5$ is characterized by an overlap of the three ellipses. This means that chemicals in this area have not a clear affinity for one particular compartment and the overall accumulation strongly depend from the combination of the partitioning in water, protein and lipids. The overlapping area is particularly interesting because it is very close to the $\log K_{OW}$ and $\log K_{OA}$ screening threshold for air breathing organisms. The accurate estimation of the biotic-partitioning coefficients is fundamental to understand the bioaccumulation dynamics of chemicals in this region.

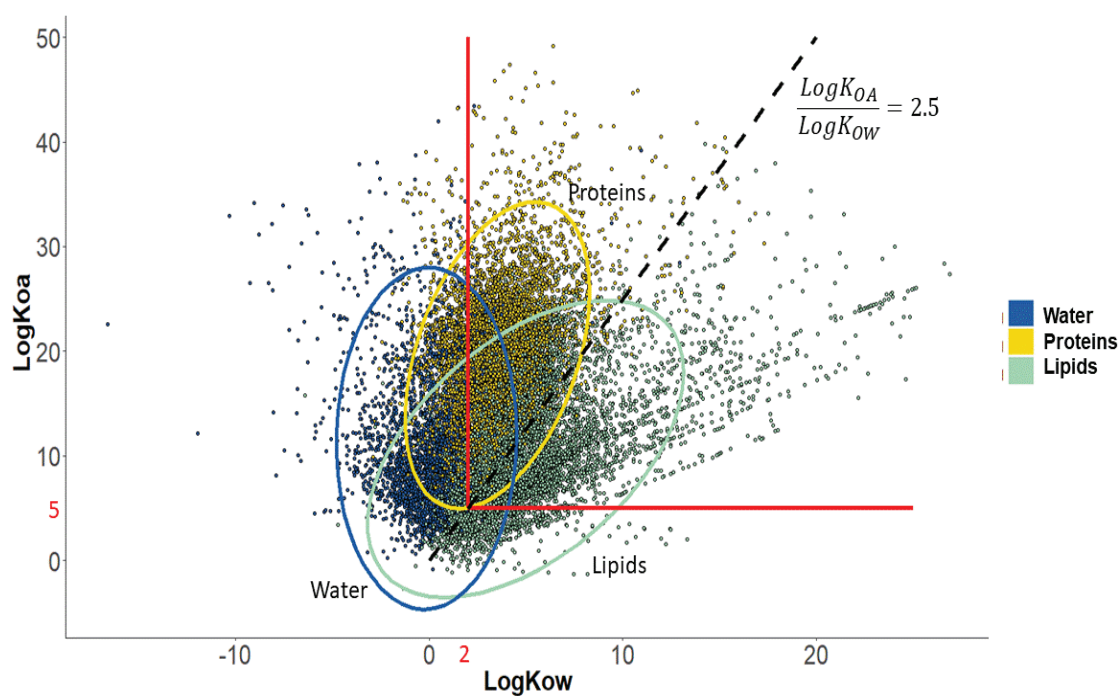


Figure 3.10: Relationship between K_{OW} and K_{OA} for the studied organic compounds. The graph identifies regulatory screening threshold for $\log K_{OW}$ (i.e. $\log K_{OW} > 2$) and $\log K_{OA}$ (i.e. $\log K_{OA} > 5$). Chemicals are labelled according to the main sorptive matrix. Confidence ellipse at 0.95 assuming a multivariate normal distribution are reported.

BMF is calculated in Tier 2 using the biotic-partitioning coefficients; the three dimensional interpolation of $\log \text{BMF}$ in the space of $\log K_{OW}$, $\log K_{OA}$ is reported in figure 3.10. Overall, 97% of the whole set of chemicals have $\text{BMF} > 1$. The highest BMF

Results and Discussion

values are estimated for chemicals with $\log K_{OW}$ between 2 and 7 and $\log K_{OA} > 5$ while lowest values are located for $\log K_{OW} > 10$ and for $\log K_{OA} < 5$ as for Tier 1.

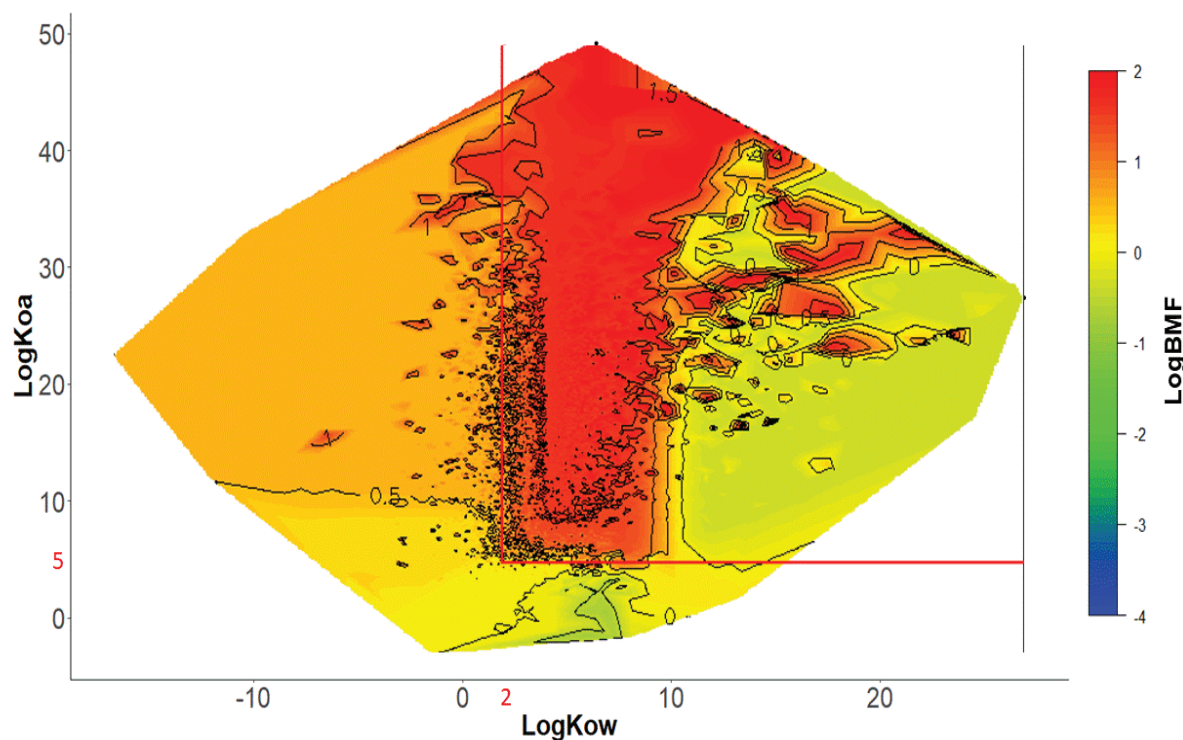


Figure 3.11: Heatmap for the $\log BMF$ respect $\log K_{OW}$ and $\log K_{OA}$ for Tier 2. The graph identifies regulatory screening threshold for $\log K_{OW}$ (i.e. $\log K_{OW} > 2$) and $\log K_{OA}$ (i.e. $\log K_{OA} > 5$) in red and $\log BMF$ isolines in black.

Less hydrophobic chemicals (i.e. $\log K_{OW} < 2$) with higher affinity for aqueous media are less affected by the introduction of biotic partitioning coefficients and the BMF estimation in Tier 2 for these compounds is almost unchanged compared to Tier 1 (figure 3.12 a).

The main differences in the BMF estimation are reported for chemicals with high affinity to proteins. The calculation for the BMF is strongly affected by the more accurate estimation for $\log k_{ProtW}$ calculated by the ppLFER model and in general compounds mainly sorbed in proteins obtained a higher estimation in the BMF respect to Tier 1. The differences can be explained by the changes in the storage capacity of the human for the chemicals (Z) reported in figure 3.12.

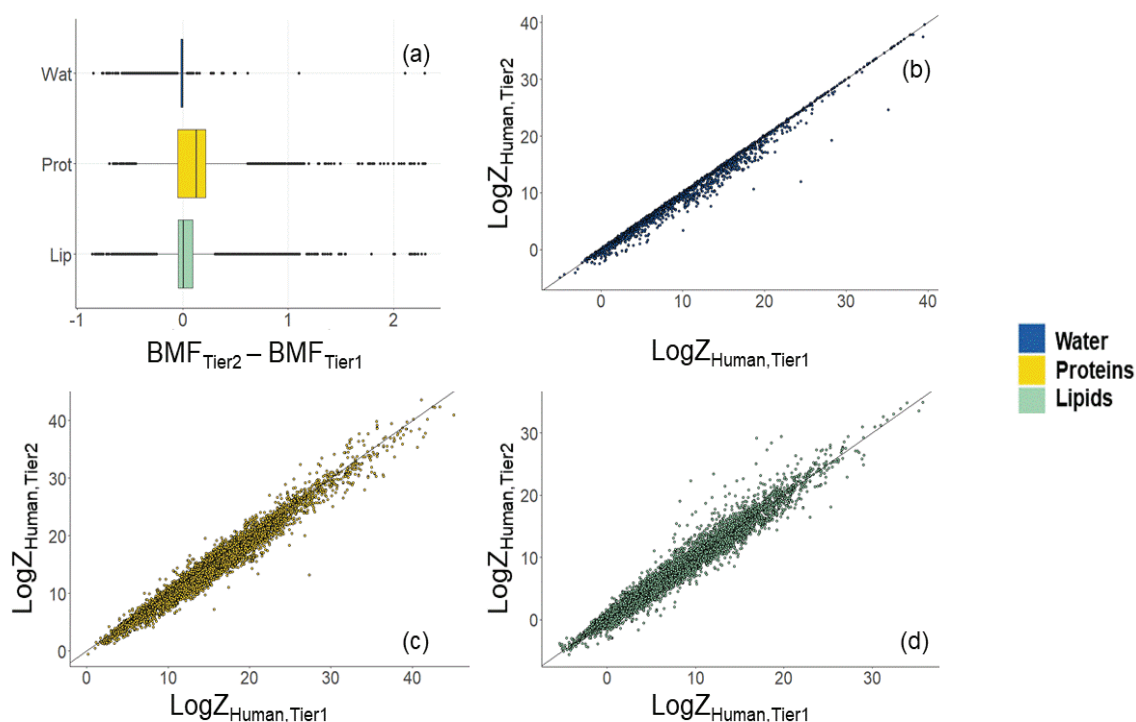


Figure 3.12: (a) box and whiskers for $\log\text{BMF}_{\text{Tier2}} - \log\text{BMF}_{\text{Tier1}}$ differences for compound with high affinity for water, lipids and protein; (b) relationship $\log\text{Z}_{\text{human Tier1}} - \log\text{Z}_{\text{human Tier2}}$ for compound with high affinity for water; (c) relationship $\log\text{Z}_{\text{human Tier1}} - \log\text{Z}_{\text{human Tier2}}$ for compound with high affinity for proteins; (d) relationship $\log\text{Z}_{\text{human Tier1}} - \log\text{Z}_{\text{human Tier2}}$ for compound with high affinity for lipids.

Chemicals with high affinity for lipids and proteins report the highest changes in the Z values and in the BMF estimation (figure 3.12 c and d). This is because the total human Z depends on the relative contribution of Z for the different matrices (i.e. lipids, proteins, water). Changes in the relative sorption capacity of the different matrices may change the human total Z and alter the bioaccumulation potential. For instance, in Tier 1, chemicals with $\log K_{\text{OW}} < 2$ obtained an BMF estimated between 1 and 8. In Tier 2, 478 compounds with $\log K_{\text{OW}} < 2$ are predicted with a $\text{BMF} > 10$ and thus are highly bioaccumulative. These chemicals report high changes in the Z values that change the sorption capacity of the organisms. Approximately 75% of these 478 have high affinity for proteins while the other 25 % have high affinity for lipids. Overall the introduction of the biotic partitioning coefficient improved the estimation of the internal distribution of chemicals provided a more precise estimation of the BMF for chemicals which

Results and Discussion

bioaccumulation is not dominated by lipids partitioning. In summary, most of the same chemicals identified in Tier 1 are still predicted as potentially bioaccumulative in air breathing organisms at Tier 2.

Tier 3. In Tier 3 we accounted for dissociation of ionogenic organic chemicals (IOCs) to understand the impact of speciation (i.e. the fractional amount of neutral and ionized forms) at a physiological pH (i.e. 7.4) on bioaccumulation. Chemicals were considered in neutral, acid or base form according to their predominant state at pH 7.4. We only consider the strongest acid or base pKas since RAIDAR model cannot treat multiprotonic species. Their biotransformation rates were assumed to be negligible to examine solely the impact of speciation.

Approximately the 69% of the examined chemicals are mainly present in the neutral form and considered as neutral organic chemicals. For these chemicals the partitioning between aqueous and organic phases is still dominated by $\log K_{OW}$, thus no differences are expected in the calculation of the BMF. On the other hand, 2'682 chemicals are assumed to be in acidic state while 3'698 are predicted to be bases. Ions are more hydrophilic than the corresponding neutral form, so they are expected to partition into organic tissue to a smaller degree than the correspondent neutral molecules²⁸⁸. The partition coefficients for IOCs are adjusted according to the pKa values and the apparent partition coefficient for the charged form of the IOC as calculated in RAIDAR⁸⁹. The interpolation of the $\log BMF$ in the space of $\log K_{OW}$ - $\log K_{OA}$ is reported in figure 3.13 a.

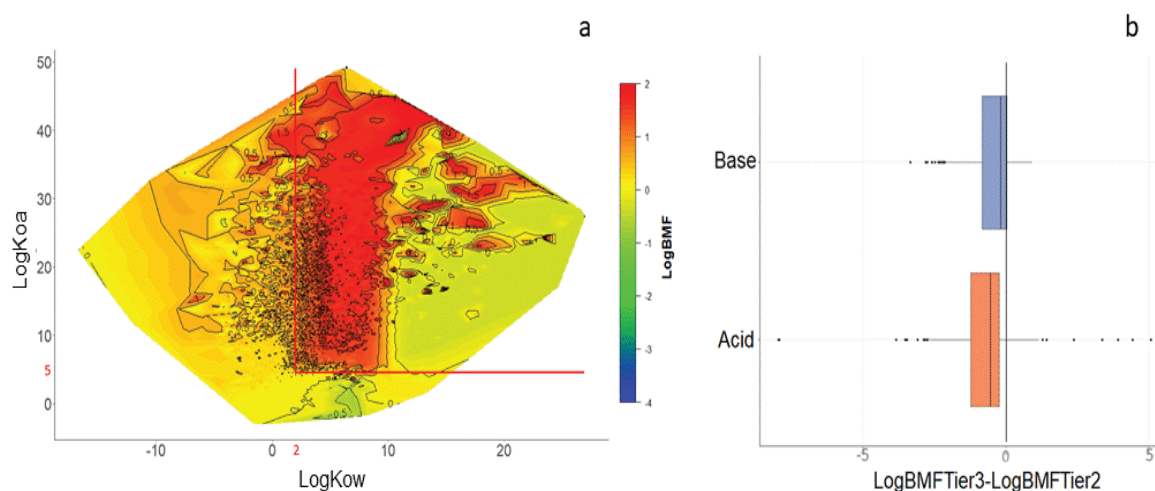


Figure 3.13: (a) Heatmap for the log BMF respect log K_{OW} and log K_{OA} for Tier 3. The graph identifies regulatory screening threshold for log K_{OW} (i.e. log $K_{OW} > 2$) and log K_{OA} (i.e. log $K_{OA} > 5$) in red and logBMF isolines in black; (b) Box and whisker plot for logBMF_{Tier3} - logBMF_{Tier2} differences for acid and base.

Approximately ninety-six percent of the chemicals still have an estimated BMF > 1 even though many chemicals have a lower BMF than in Tier 1 and 2.

As expected neutral chemicals do not report any differences in the BMF estimation. In general compound in ionic state reported lower BMF values respect to Tier 1 and 2 (figure 3.13 b). The decreasing is more evident for acidic compounds mainly with log $K_{OW} < 2$ and log $K_{OA} > 5$. These compounds are supposed to be less volatile and thus not efficiently eliminated in the respiratory exchange (for air breathing organisms). However, the ionic state increases their water solubility and enhance water-mediated exchanged and their elimination via renal excretion. This confirms the general rule that the neutral fraction of an ionizing compound accumulates in biota to a greater degree than the ionized fraction²⁸⁹. For IOC that occur as both neutral and ionic molecules information about pH is essential to effectively predict the actual bioaccumulation potential.

Tier 4. In Tier 4 we accounted for biotransformation HLs in human to test the hypothesis that the biotransformation rates are more important determinants of bioaccumulation in humans than the only partitioning properties. We applied a set of different QSAR models

Results and Discussion

(i.e. IFS-based and theoretical molecular descriptor-based^{126,264}) to predict the biotransformation potential in human. For each chemical, we averaged the predictions of the models on the basis of the applicability domain. We considered a threshold at 70 days^{54,126,275} to discriminate between readily biotransformed and not biotransformed chemicals (i.e. $HL_B > 70$ days indicates no biotransformation). Overall 98% of the 20346 chemicals have $HL_B < 70$ days thus are fast biotransformed in humans. Only 307 compounds exceed the biotransformation threshold and can be considered of particular concern for the bioaccumulation potential. If these chemicals are effectively absorbed by the organism biotransformation will not represent a chemical loss process sufficient to mitigate biomagnification. Figure 3.14 reports the poorly metabolized chemicals in the log K_{OW} - log K_{OA} space according to their chemicals class.

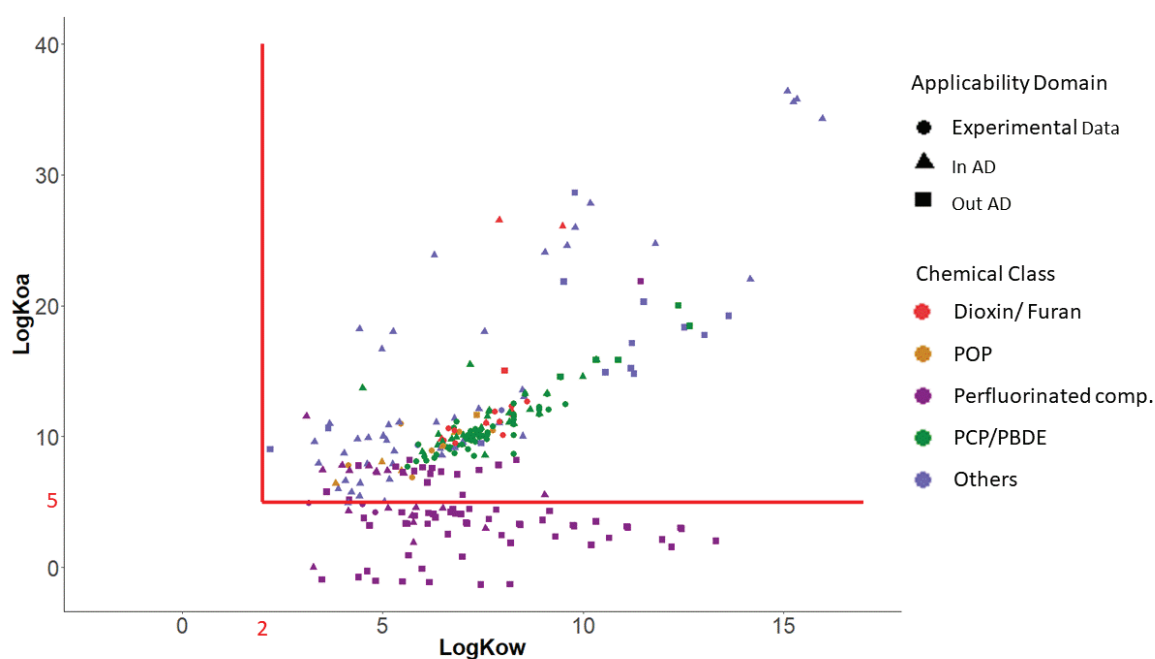


Figure 3.14: Relationship between K_{OW} and K_{OA} for the compounds with $HL_B > 70$ days

Among the 307 slowly biotransformed compounds we found 119 PCB and PBDE congeners (i.e. from 3 to 10 halogen atoms), 90 perfluorinated compounds (i.e. from 10 to 39 fluorine atoms), 22 dioxin- and furan-like compounds, 14 POPs (e.g. DDT, Aldrin, Endosulfan Mirex)^{12,290} and 62 chemicals with different structures. It is interesting to note that experimental HL_B were reported for 117 of these compounds (mainly PCBs,

100

POPs and furans) while 100 compounds were inside the AD of the models (mainly chemicals with various structure, probably pharmaceuticals-like compounds similar to the training set of the biotransformation models). On the other hand, 90 chemicals were outside the AD of the models, mainly perfluorinated compounds with a high number of fluorine atoms. All the chemicals with $HL_B > 70$ days have $\log K_{OW} > 2$ supporting the observation that compounds with high K_{OW} may have long HLs. We could observe a linear increasing of the K_{OW} and K_{OA} for PCB and PBDE congeners given by the increasing number of halogen atoms. Moreover, many perfluoro compounds obtained low K_{OA} values ($\log K_{OA} < 5$) indicating that respiratory exchange could be relatively important in the bioaccumulation dynamic of these compounds.

Figure 3.15 reports the heatmap of $\log BMF$ in the space of $\log K_{OW}$ - $\log K_{OA}$ for Tier 4.

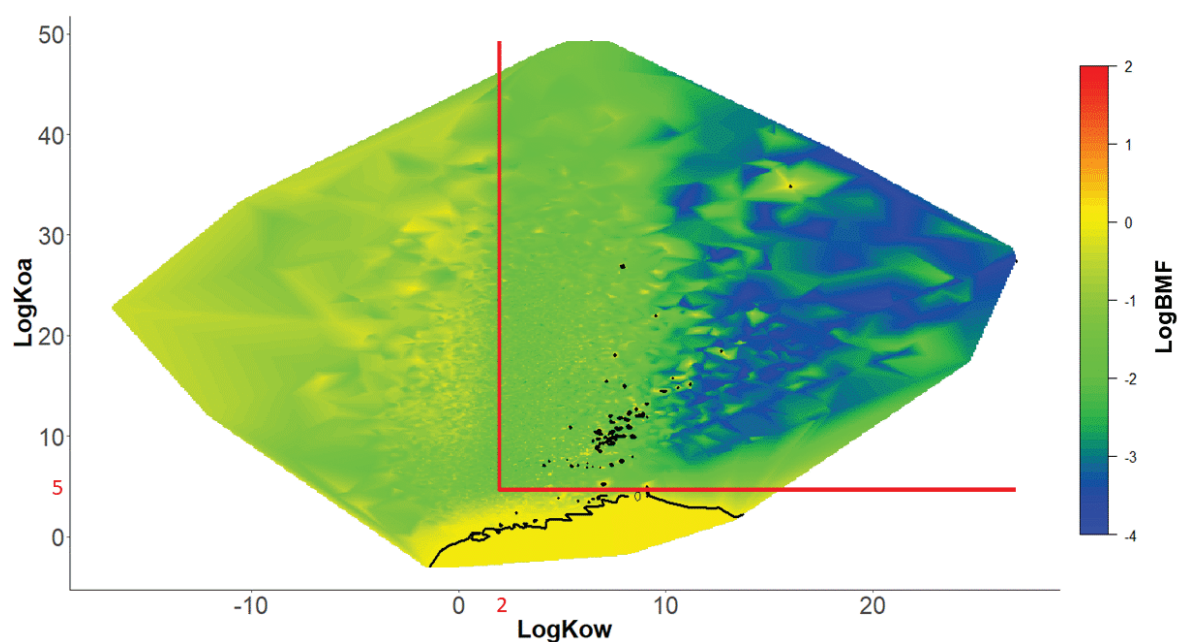


Figure 3.15: (a) Heatmap for the $\log BMF$ respect $\log K_{OW}$ and $\log K_{OA}$ for Tier 4. The graph identifies regulatory screening threshold for $\log K_{OW}$ (i.e. $\log K_{OW} > 2$) and $\log K_{OA}$ (i.e. $\log K_{OA} > 5$) in red and $\log BMF$ isolines in black

The introduction of HL_B strongly affected the calculation of the BMF. 98 % of chemicals are predicted with a $BMF < 1$ and thus as “nB”. The main differences are visible in the area of the screening criteria (i.e. $\log K_{OW} > 2$ ad $\log K_{OA} > 5$) where chemicals obtained

Results and Discussion

the highest BMF in Tier 1-3. Here the majority of the chemicals are now predicted with a low BMF meaning that the elimination processes related to biotransformation are predominant in the overall bioaccumulation. Generally, the chemicals with the highest HL_B obtained the highest BMF; among the 307 molecules with $HL_B > 70$ day, 262 were estimated with $BMF > 1$. These are mainly the PCB and PBDE congeners, the POPs, and the perfluorinated compounds and are highlighted by the logBMF isolines in figure 3.15. The remaining 45 compounds with $HL_B > 70$ days have $\log K_{OW} > 7$ and probably the E_D is too low to effectively increase their concentration in the body through the diet and thus have a $BMF < 1$. This last analysis allows to refine the biomagnification screening in two major ways. First the chemicals exceeding the regulatory screening threshold but efficiently biotransformed by the organism (i.e. $HL_B < 70$ days) are estimated with a low BMF and thus are not expected to biomagnify. This shows that screening based only on partition coefficients is not sufficient to describe and address the bioaccumulation and biomagnification processes, and can lead to overly conservative estimates (i.e. “false positives”). On the other hand, chemicals with $\log K_{OA} < 5$ but $HL_B > 70$ days (e.g. some perfluoro compounds in figure 3.14) are estimated with a $BMF > 1$ and thus they may be potentially bioaccumulative even though they do not exceed the regulatory screening criteria. This demonstrates how information about HL_B can be used to refine bioaccumulation screening and avoid “false negatives”.

3.6.4 Conclusions

Here we have provided a tiered approach to address data gaps and uncertainties in the primary information relating to evaluating biomagnification potential in air-breathing organisms. We have demonstrated that the “hydrophobicity paradigm” or the proposed “ K_{OW} and K_{OA} only” criteria would categorize a very large number of chemicals as bioaccumulative that are not actually bioaccumulative because the approach ignores biotransformation^{68,291}. This study shows the key role of biotransformation in bioaccumulation assessment for air-breathing organisms and highlights the need for reliable biotransformation data to effectively categorize chemicals for hazard. Thus, if bioaccumulation assessments are of regulatory interest, then priority research should be

focussed on addressing uncertainty (and variability) in biotransformation rates for air-breathing organisms. Given the massive role that biotransformation rates have in determining the bioaccumulation of organic chemicals in air-breathing organisms there is a critical need to develop integrated testing strategies and high quality databases of *in-vitro* and *in-vivo* rates of biotransformation in various air-breathing *taxa* (e.g. mammals, avia) and *in-vitro* to *in-vivo* extrapolation (IVIVE) methods¹³⁹ as alternatives to *in-vivo* testing that would thereby reduce costs and animal use while increasing throughput and chemical assessment efficiency. From new *in-vitro* biotransformation rate data obtained systematically to address uncertainty, new QSARs and predictive methods for rates and half-lives can be developed and ADs for the existing QSARs^{126,264} expanded. Moreover, this work allowed to highlight other sources of uncertainty such as biotransformation variability between species, within species, renal excretion, dietary uptake efficiency, relative rates of feeding to respiration, ionization, active processes for uptake and elimination, chemical-specific binding and the fate of metabolites that still need to be addressed to improve the assessment of bioaccumulation of organic chemicals.

Results and Discussion

Chapter 4: Conclusions

Conclusions

In the last decades scientific and technological advancements have brought remarkable progresses in the testing and non-testing strategies improving *in-vitro* test and *in-silico* methodologies that can integrate, support and sometimes replace *in-vivo* animal testing. The development of the ITS with the integration of different methodologies and the use of all the available information can guide the achievement of genuine progresses in the understanding of the hazard of chemicals whilst maintaining highest standards in terms of both scientific methodologies and animal welfare^{75,292,293}.

Within this thesis I focused the attention on the *in-silico* methodologies and in particular on QSAR models to prioritize chemicals with the aim to support and facilitate the environmental risk assessment procedure of CEC. The application of computational tools and QSAR models is effective and essential for the preliminary screening and for the prioritization of the hazard of compounds because they could be fast and cheap and can produce valuable data for a consciousness planning and reduction of experimental and animal testing. Several high quality QSARs have been developed according to the regulatory guidelines and the OECD principles for the development and validation of QSAR models^{179,180}. All these models have been implemented in the software QSARINS-Chem¹⁶⁷ and subsequently applied for a comprehensive study of the PBT properties of pharmaceuticals proceeding from a general screening (i.e. **paper I**) to the refinement of toxicity (i.e. **paper II, III and IV**) and bioaccumulation (i.e. **paper V and VI**).

The PBT screening described in **paper I** demonstrated how QSAR models can be applied for the effective prioritization of pharmaceuticals for their environmental hazard. Among thousands of untested pharmaceuticals, 35 have been selected as potential PBT and thus as priority compounds. Main point of the screening was the use of the consensus approach to select the priority pharmaceuticals. The use of different QSAR models based on different methodologies and training sets in the consensus approach can improve the prediction of the single models and provide a more robust estimation of the intrinsic hazard of chemicals. The predictions obtained in agreement from different methods can

be considered more consistent and reliable in identifying potential PBT chemicals. Moreover, it is important to note that this approach can be applied to existing chemicals without experimental data, and to not yet synthesized molecules to plan environmentally safer chemicals from their design. Unfortunately, screening criteria and critical values for the PBT assessment applied so far, refers to traditional “legacy contaminants” and the actual regulatory requirements could be not adequate for a correct assessment of many CEC. This PBT screening can be considered as a conservative indication based on the worst-case combination of the PBT properties. A refinement of the PBT screening may be performed to increase the realism of the proposed approach, using new data for specific properties, such as toxicity in multiple species and biotransformation, or to enlarge the accuracy of the current QSAR models.

In **paper II** and **III** QSAR models to estimate the aquatic toxicity of pharmaceuticals at different trophic levels have been developed, validated and applied to refine the toxicity assessment of pharmaceuticals. Crucial point of the papers was the data curation to collect coherent and homogenous data generated according to the standardized regulatory guidelines^{146–148}. The papers highlighted the lack of high-quality experimental data related to the environmental effects of pharmaceuticals and the need of new *ad-hoc* QSAR models to fill this data gap and improve the assessment of the ecotoxicity of medicinal products. The selected end-point was the acute toxicity measured by means of standard toxicity test to ensure the maximum consistency between data collected from different sources. Moreover, acute toxicity gives an immediate perception of the hazard of chemicals, particularly useful for screening purposes. Remarkable result was the wide applicability domain of the developed models that were able to provide reliable predictions for a high percentage of the studied compounds even though trained on relatively small training sets. These models can be applied to predict the environmental toxicity of pharmaceuticals in the three aquatic trophic levels (i.e. phytoplankton, zooplankton and fish), and to provide quantitative support for the complete assessment of the potential hazard of pharmaceuticals in the aquatic environment. Therefore, these models allowed to elaborate an externally validated QSAR-based Aquatic Toxicity

Conclusions

Index for pharmaceuticals based on structural properties such as molecular dimension, lipophilicity and Hydrogen Bond interactions^{245,294,295}. This index is intended for ranking and can be applied for the prioritization of existing and new pharmaceuticals directly from the molecular structure. The developed models are preliminary indicator of the potential hazard of pharmaceuticals for the aquatic environment and do not consider specific mode of action. Chemicals prioritized should undergo further evaluation including the assessment of long-term adverse effects such as endocrine disruption, behavioral disturbances and other health effects that are not taken into account in this study.

Furthermore, the refinement of the PBT screening was also possible for the bioaccumulation assessment, by including information on the potential biotransformation. This study has shown that the “Hydrophobicity paradigm” is not adequate to assess the bioaccumulation of CEC and that biotransformation rates are important determinants of overall bioaccumulation, processes, especially for air-breathing organisms. Till now, a limited amount of consistent biotransformation data in human were available¹²⁶ and the correct refinement of the bioaccumulation in human was not possible for many CEC. QSAR models were developed within this thesis (i.e. **paper V**) for the estimation of the whole-body total and metabolic biotransformation HL in human of pharmaceuticals and used to refine the preliminary PBT screening performed in **paper I**. Data about the biotransformation rates were demonstrated to be relevant in the refinement of the bioaccumulation estimation performed on the basis of logK_{ow} and BCF and allowed to correct the overly conservative estimates (i.e. false positive) lowering the hazard for many compounds; only 22 pharmaceuticals over the 35 highlighted in **paper I** are predicted as slowly biotransformed by organisms and thus have the highest bioaccumulation potential. This study allowed to refine the more precautionary approach of the general PBT screening including additional relevant information and reducing the number of priority chemicals.

Finally, a relevant achievement of this work is the integration of different *in-silico* methods and QSAR models to perform a screening of environmental related properties of chemicals (i.e. **paper IV** and **VI**). In particular, **paper VI** is an example of *in-silico* ITS that condenses information from very different approaches in a coherent and conscious way to screen the biomagnification potential of thousands of chemicals. Fragment based, molecular descriptors based and ppLFER based QSAR models have been applied, posing particular attention to the AD estimation, to predict partitioning and kinetics properties of thousands of chemicals. These data have been integrated in a mechanistic mass-balance multimedia environmental fate, food-web model to estimate the BMF in human in a tiered approach to analyze the importance of biotransformation and metabolisms in the estimation of bioaccumulation and biomagnification. This study showed the key role of biotransformation in bioaccumulation assessment for air-breathing organisms and highlights the need for reliable biotransformation data to effectively categorize chemicals for their hazard. Given the massive role that biotransformation rates have in determining the bioaccumulation of organic chemicals in air-breathing organisms there is a critical need to develop integrated testing strategies and high quality databases of *in-vitro* and *in-vivo* rates of biotransformation in various air-breathing *taxa* (e.g. mammals, avia) and IVIVE methods¹³⁹ as alternatives to *in-vivo* testing with the aim to reduce costs, increase throughput, and reduce animal use. Moreover, this data can be used as basis to improve existing or generate new QSAR models to address the other sources of uncertainty (e.g. dietary uptake efficiency, fate of the metabolites) that are open challenges in the bioaccumulation related sciences.

Concluding this thesis shows how high-quality, validated QSAR models can be integrated in the ITS system to support the environmental hazard and risk assessment of chemicals on the basis of their molecular structure.

Conclusions

Acknowledgements

Most of all I would like to thank my supervisors Ester and Paola for all the support and encouragement they gave me, during these years. I thank them not only for their tremendous academic support, but also for giving me so many wonderful opportunities. Thanks to Ester for her patient, enthusiasm and motivation. Thanks to Paola for her great passion and constant guidance. Without them, this PhD would not have been achievable.

Special thanks to Dr. Jon Arnot for providing extensive personal and professional guidance. I am particularly indebted to Jon for his constant faith in my work, and for his support when so generously hosted me in Toronto. I have very fond memories of my time there.

Thanks to prof Frank Wania for hosting me as International Visiting Student in his group and leading me working on diverse exciting projects. Thanks also to all the special people that I have met in Canada.

I acknowledge my thesis committee, prof Melek Türker Saçan, prof. Patrik Andersson and prof. Roberta Bettinetti, for reading previous drafts of this dissertation and providing valuable comments.

My year at the University of Insubria have been a great period especially thanks all the friends that I have met there. My gratitude goes to all the people of “Piano Rosso”, especially Silvia, Isabella and Vanessa. A special thanks to Enrico, Stefano, Nicola, Debora, Myriam, Marta, Laura, Anna, Donatella, Edoardo, Linda and Ilaria for the many lunches together in these years.

I would like to thank Jorge for the nice time spent together in Paris and Giulia for the dance classes and the valuable conversations.

Special mention goes to Suzanne for the wonderful horse riding and to Elisabetta for being a precious confidant.

Thanks to “gli amici della tradizione”, Fabio, Gabriele, Carlo, Mirko, Fabrizio, Claudia, Francesco A., Francesco B., Simone, Marco and Paolo, for the many dinners together and for being a landmark in my live.

I would also like to say a heartfelt thanks to my Mum, Dad and brother for always believing in me and encouraging me to follow my dreams.

Finally, my deepest gratitude goes to Kuiying that is far away but will always have a special place in my heart.

References

- (1) Eurostat. *Environmental Statistics and Accounts in Europe*; Publications Office or the European Union: Luxembourg, 2010.
- (2) Muir, D. C.; Howard, P. H. Are There Other Persistent Organic Pollutants? A Challenge for Environmental Chemists. *Env. Sci Technol* **2006**, *40*, 7157–7166. <https://doi.org/10.1021/es061677a>.
- (3) Bernhardt, E. S.; Rosi, E. J.; Gessner, M. O. Synthetic Chemicals as Agents of Global Change. *Front. Ecol. Environ.* **2017**, *15* (2), 84–90. <https://doi.org/10.1002/fee.1450>.
- (4) Binetti, R.; Costamagna, F. M.; Marcello, I. Exponential Growth of New Chemicals and Evolution of Information Relevant to Risk Control. *Ann. Dell'Istituto Super. Sanità* **2008**, *44*, 13–15.
- (5) Chemicals Abstracts Service <https://www.cas.org/> (accessed Dec 1, 2016).
- (6) Hutchinson, T. H.; Lyons, B. P.; Thain, J. E.; Law, R. J. Evaluating Legacy Contaminants and Emerging Chemicals in Marine Environments Using Adverse Outcome Pathways and Biological Effects-Directed Analysis. *Mar. Pollut. Bull.* **2013**, *74* (2), 517–525. <https://doi.org/10.1016/j.marpolbul.2013.06.012>.
- (7) CEFIC. Cefic | Facts and Figures 2017 <http://www.cefic.org/Facts-and-Figures/> (accessed Aug 16, 2018).
- (8) Eurostat. *Sustainable Development in the European Union: Monitoring Report on Progress towards the SDGs in an EU Context*; Statistical books / Eurostat; Publications Office or the European Union, 2017.
- (9) Eurostat - Tables, Graphs and Maps Interface (TGM) table http://ec.europa.eu/eurostat/tgm/table.do?tab=table&init=1&language=en&pcode=sdg_12_10&plugin=1 (accessed Aug 16, 2018).
- (10) EC. *Regulation (EC) No 1272/2008 of the European Parliament and the Council (EC) of 16 December 2008 on Classification, Labelling and Packaging of Substances and Mixtures, Amending and Repealing Directives 67/548/EEC and 1999/45/EC, and Amending Regulation (EC) No 1907/2006*; 2008.
- (11) Eurostat; Statistical Office of the European Union. *Compilation of Chemical Indicators Development, Revision and Additional Analyses.*; Publications Office: Luxembourg, 2016.
- (12) UNEP. *Final Act of the Conference of Plenipotentiaries on the Stockholm Convention on Persistent Organic Pollutants.*; 2001.

References

- (13) Narragansett Bay Estuary Program. Chapter 9, Legacy Contaminants, Pages 192-210. In *State of Narragansett Bay and Its Watershed*; Technical Report., 2017.
- (14) Zimmerman, J. B.; Anastas, P. T. Toward Designing Safer Chemicals. *Science* **2015**, *347* (6219), 215–215. <https://doi.org/10.1126/science.aaa6736>.
- (15) Poynton, H. C.; Robinson, W. E. Chapter 3.7 - Contaminants of Emerging Concern, With an Emphasis on Nanomaterials and Pharmaceuticals. In *Green Chemistry*; Török, B., Dransfield, T., Eds.; Elsevier, 2018; pp 291–315. <https://doi.org/10.1016/B978-0-12-809270-5.00012-1>.
- (16) US-EPA; OW/ORD Emerging Contaminants Workgroup. White Paper: Aquatic Life Criteria for Contaminants of Emerging Concern. Part I. General Challenges and Recommendations. 2008.
- (17) Daughton, C. G. “Emerging” Chemicals as Pollutants in the Environment: A 21st Century Perspective. **2005**, 18.
- (18) Sauv e, S.; Desrosiers, M. A Review of What Is an Emerging Contaminant. *Chem. Cent. J.* **2014**, *8*, 15. <https://doi.org/10.1186/1752-153X-8-15>.
- (19) USGS. Contaminants of Emerging Concern in the Environment. 2017.
- (20) Richardson, S. D. Water Analysis: Emerging Contaminants and Current Issues. *Anal. Chem.* **2009**, *81* (12), 4645–4677. <https://doi.org/10.1021/ac9008012>.
- (21) Glassmeyer, S. T.; Furlong, E. T.; Kolpin, D. W.; Cahill, J. D.; Zaugg, S. D.; Werner, S. L.; Meyer, M. T.; Kryak, D. D. Transport of Chemical and Microbial Compounds from Known Wastewater Discharges: Potential for Use as Indicators of Human Fecal Contamination. *Environ. Sci. Technol.* **2005**, *39* (14), 5157–5169.
- (22) Daughton, C. G. Environmental Stewardship and Drugs as Pollutants. *The Lancet* **2002**, *360* (9339), 1035–1036. [https://doi.org/10.1016/S0140-6736\(02\)11176-7](https://doi.org/10.1016/S0140-6736(02)11176-7).
- (23) Daughton, C. G.; Ternes, T. A. Pharmaceuticals and Personal Care Products in the Environment: Agents of Subtle Change? *Environ. Health Perspect.* **1999**, *107*, 907–938. <https://doi.org/10.2307/3434573>.
- (24) Daughton, C. G. Cradle-to-Cradle Stewardship of Drugs for Minimizing Their Environmental Disposition While Promoting Human Health. I. Rationale for and Avenues toward a Green Pharmacy. *Environ. Health Perspect.* **2003**, *111* (5), 757–774. <https://doi.org/10.1289/ehp.5947>.
- (25) Vollmer, G. Disposal of Pharmaceutical Waste in Households – A European Survey. In *Green and Sustainable Pharmacy*; Springer, Berlin, Heidelberg, 2010; pp 165–178. https://doi.org/10.1007/978-3-642-05199-9_11.

- (26) Bound, J. P.; Voulvoulis, N. Household Disposal of Pharmaceuticals as a Pathway for Aquatic Contamination in the United Kingdom. *Environ. Health Perspect.* **2005**, *113* (12), 1705–1711. <https://doi.org/10.1289/ehp.8315>.
- (27) Tong, W.; Hong, H.; Xie, Q.; Shi, L.; Fang, H.; Perkins, R. Assessing QSAR Limitations - A Regulatory Perspective. *Curr. Comput. Aided Drug Des.* **2005**, *1* (2), 195–205. <https://doi.org/10.2174/1573409053585663>.
- (28) Baken, K. A.; Sjerps, R. M. A.; Schriks, M.; van Wezel, A. P. Toxicological Risk Assessment and Prioritization of Drinking Water Relevant Contaminants of Emerging Concern. *Environ. Int.* **2018**, *118*, 293–303. <https://doi.org/10.1016/j.envint.2018.05.006>.
- (29) Verlicchi, P.; Al Aukidy, M.; Zambello, E. Occurrence of Pharmaceutical Compounds in Urban Wastewater: Removal, Mass Load and Environmental Risk after a Secondary Treatment-A Review. *Sci. Total Environ.* **2012**, *429*, 123–155. <https://doi.org/10.1016/j.scitotenv.2012.04.028>.
- (30) Patrolecco, L.; Capri, S.; Ademollo, N. Occurrence of Selected Pharmaceuticals in the Principal Sewage Treatment Plants in Rome (Italy) and in the Receiving Surface Waters. *Environ. Sci. Pollut. Res. Int.* **2015**, *22* (8), 5864–5876. <https://doi.org/10.1007/s11356-014-3765-z>.
- (31) Kümmerer, K. The Presence of Pharmaceuticals in the Environment Due to Human Use - Present Knowledge and Future Challenges. *J. Environ. Manage.* **2009**, *90* (8), 2354–2366. <https://doi.org/10.1016/j.jenvman.2009.01.023>.
- (32) Wu, C.; Spongberg, A. L.; Witter, J. D.; Fang, M.; Czajkowski, K. P. Uptake of Pharmaceutical and Personal Care Products by Soybean Plants from Soils Applied with Biosolids and Irrigated with Contaminated Water. *Environ. Sci. Technol.* **2010**, *44* (16), 6157–6161. <https://doi.org/10.1021/es1011115>.
- (33) Maruya, K. A.; Dodder, N. G.; Sengupta, A.; Smith, D. J.; Lyons, J. M.; Heil, A. T.; Drewes, J. E. Multimedia Screening of Contaminants of Emerging Concern (CECS) in Coastal Urban Watersheds in Southern California (USA). *Environ. Toxicol. Chem.* **2016**, *35* (8), 1986–1994. <https://doi.org/10.1002/etc.3348>.
- (34) Sjerps, R. M. A.; Vughs, D.; van Leerdam, J. A.; ter Laak, T. L.; van Wezel, A. P. Data-Driven Prioritization of Chemicals for Various Water Types Using Suspect Screening LC-HRMS. *Water Res.* **2016**, *93*, 254–264. <https://doi.org/10.1016/j.watres.2016.02.034>.
- (35) Blum, K. M.; Andersson, P. L.; Ahrens, L.; Wiberg, K.; Haglund, P. Persistence, Mobility and Bioavailability of Emerging Organic Contaminants Discharged from Sewage Treatment Plants. *Sci. Total Environ.* **2018**, *612*, 1532–1542. <https://doi.org/10.1016/j.scitotenv.2017.09.006>.

References

- (36) Oaks, J. L.; Gilbert, M.; Virani, M. Z.; Watson, R. T.; Meteyer, C. U.; Rideout, B. A.; Shivaprasad, H. L.; Ahmed, S.; Chaudhry, M. J. I.; Arshad, M.; et al. Diclofenac Residues as the Cause of Vulture Population Decline in Pakistan. *Nature* **2004**, *427* (6975), 630–633. <https://doi.org/10.1038/nature02317>.
- (37) Brian, J. V.; Harris, C. A.; Scholze, M.; Kortenkamp, A.; Booy, P.; Lamoree, M.; Pojana, G.; Jonkers, N.; Marcomini, A.; Sumpter, J. P. Evidence of Estrogenic Mixture Effects on the Reproductive Performance of Fish. *Environ. Sci. Technol.* **2007**, *41* (1), 337–344. <https://doi.org/10.1021/es0617439>.
- (38) Ketata, I.; Denier, X.; Hamza-Chaffai, A.; Minier, C. Endocrine-Related Reproductive Effects in Molluscs. *Comp. Biochem. Physiol. Part C Toxicol. Pharmacol.* **2008**, *147* (3), 261–270. <https://doi.org/10.1016/j.cbpc.2007.11.007>.
- (39) Soto, A. M.; Sonnenschein, C. Environmental Causes of Cancer: Endocrine Disruptors as Carcinogens. *Nat. Rev. Endocrinol.* **2010**, *6* (7), 363–370. <https://doi.org/10.1038/nrendo.2010.87>.
- (40) Newbold, R. R. Developmental Exposure to Endocrine-Disrupting Chemicals Programs for Reproductive Tract Alterations and Obesity Later in Life. *Am. J. Clin. Nutr.* **2011**, *94* (suppl_6), 1939S-1942S. <https://doi.org/10.3945/ajcn.110.001057>.
- (41) Cao, L.-Y.; Zheng, Z.; Ren, X.-M.; Andersson, P. L.; Guo, L.-H. Structure-Dependent Activity of Polybrominated Diphenyl Ethers and Their Hydroxylated Metabolites on Estrogen Related Receptor γ : In Vitro and in Silico Study. *Environ. Sci. Technol.* **2018**, *52* (15), 8894–8902. <https://doi.org/10.1021/acs.est.8b02509>.
- (42) Kidd, K. A.; Blanchfield, P. J.; Mills, K. H.; Palace, V. P.; Evans, R. E.; Lazorchak, J. M.; Flick, R. W. Collapse of a Fish Population after Exposure to a Synthetic Estrogen. *Proc. Natl. Acad. Sci. U. S. A.* **2007**, *104* (21), 8897–8901. <https://doi.org/10.1073/pnas.0609568104>.
- (43) Aris, A. Z.; Shamsuddin, A. S.; Praveena, S. M. Occurrence of 17 α -Ethinylestradiol (EE2) in the Environment and Effect on Exposed Biota: A Review. *Environ. Int.* **2014**, *69*, 104–119. <https://doi.org/10.1016/j.envint.2014.04.011>.
- (44) Sumpter, J. P.; Johnson, A. C.; Williams, R. J.; Kortenkamp, A.; Scholze, M. Modeling Effects of Mixtures of Endocrine Disrupting Chemicals at the River Catchment Scale. *Environ. Sci. Technol.* **2006**, *40* (17), 5478–5489. <https://doi.org/10.1021/es052554d>.
- (45) Nash, J. P.; Kime, D. E.; Van der Ven, L. T. M.; Wester, P. W.; Brion, F.; Maack, G.; Stahlschmidt-Allner, P.; Tyler, C. R. Long-Term Exposure to Environmental Concentrations of the Pharmaceutical Ethinylestradiol Causes Reproductive

- Failure in Fish. *Environ. Health Perspect.* **2004**, *112* (17), 1725–1733. <https://doi.org/10.1289/ehp.7209>.
- (46) Allen, H. K.; Donato, J.; Wang, H. H.; Cloud-Hansen, K. A.; Davies, J.; Handelsman, J. Call of the Wild: Antibiotic Resistance Genes in Natural Environments. *Nat. Rev. Microbiol.* **2010**, *8* (4), 251–259. <https://doi.org/10.1038/nrmicro2312>.
- (47) Gavrilescu, M.; Demnerová, K.; Aamand, J.; Agathos, S.; Fava, F. Emerging Pollutants in the Environment: Present and Future Challenges in Biomonitoring, Ecological Risks and Bioremediation. *New Biotechnol.* **2015**, *32* (1), 147–156. <https://doi.org/10.1016/j.nbt.2014.01.001>.
- (48) Diamond, J. M.; Latimer, H. A.; Munkittrick, K. R.; Thornton, K. W.; Bartell, S. M.; Kidd, K. A. Prioritizing Contaminants of Emerging Concern for Ecological Screening Assessments. *Environ. Toxicol. Chem.* **2011**, *30* (11), 2385–2394. <https://doi.org/10.1002/etc.667>.
- (49) Leeuwen, C. J. van; Vermeire, T. G. *Risk Assessment of Chemicals: An Introduction*; Springer Science & Business Media, 2007.
- (50) van der Oost, R.; Beyer, J.; Vermeulen, N. P. E. Fish Bioaccumulation and Biomarkers in Environmental Risk Assessment: A Review. *Environ. Toxicol. Pharmacol.* **2003**, *13* (2), 57–149. [https://doi.org/10.1016/S1382-6689\(02\)00126-6](https://doi.org/10.1016/S1382-6689(02)00126-6).
- (51) Duke, L. D.; Taggart, M. Uncertainty Factors in Screening Ecological Risk Assessments. *Environ. Toxicol. Chem.* **2000**, *19* (6), 1668–1680. <https://doi.org/10.1002/etc.5620190626>.
- (52) Bogen, K. T.; Spear, R. C. Integrating Uncertainty and Interindividual Variability in Environmental Risk Assessment. *Risk Anal.* **1987**, *7* (4), 427–436. <https://doi.org/10.1111/j.1539-6924.1987.tb00480.x>.
- (53) EC. *Registration, Evaluation, Authorisation and Restriction of Chemicals (REACH). Regulation (EC) No. 1907/2006 of the European Parliament and of the Council (EC)*; 2006.
- (54) ECHA. *Guidance on Information Requirements and Chemical Safety Assessment: Chapter R.11: PBT and VPvB Assessment*. European Chemicals Agency 2017.
- (55) Klecka, G. M.; Muir, D. C.; Dohmen, P.; Eisenreich, S. J.; Gobas, F. A. P. C.; Jones, K. C.; Mackay, D.; Tarazona, J. V.; van Wijk, D. Introduction to Special Series: Science-Based Guidance and Framework for the Evaluation and Identification of PBTs and POPs. *Integr. Env. Assess. Manag.* **2009**, *5*, 535–538. https://doi.org/10.1897/IEAM_2009-045.1.
- (56) Mackay, D. *Multimedia Environmental Models: The Fugacity Approach, Second Edition*; Taylor & Francis, 1991.

References

- (57) Calamari, D.; Bacci, E.; Focardi, S.; Gaggi, C.; Morosini, M.; Vighi, M. Role of Plant Biomass in the Global Environmental Partitioning of Chlorinated Hydrocarbons. *Environ. Sci. Technol.* **1991**, *25* (8), 1489–1495. <https://doi.org/10.1021/es00020a020>.
- (58) Wania, F.; MacKay, D. Peer Reviewed: Tracking the Distribution of Persistent Organic Pollutants. *Environ. Sci. Technol.* **1996**, *30* (9), 390A–396A. <https://doi.org/10.1021/es962399q>.
- (59) Brown, T. N.; Wania, F. Screening Chemicals for the Potential to He Persistent Organic Pollutants: A Case Study of Arctic Contaminants. *Environ. Sci. Technol.* **2008**, *42* (14), 5202–5209. <https://doi.org/10.1021/es8004514>.
- (60) US-EPA. Ecological Risk Assessment Glossary of Terms https://ofmpub.epa.gov/sor_internet/registry/termreg/searchandretrieve/glossariesandkeywordlists/search.do?details=&glossaryName=Eco%20Risk%20Assessment%20Glossary (accessed Sep 29, 2018).
- (61) Mackay, D.; Fraser, A. Bioaccumulation of Persistent Organic Chemicals: Mechanisms and Models. *Environ. Pollut.* **2000**, *110* (3), 375–391. [https://doi.org/10.1016/S0269-7491\(00\)00162-7](https://doi.org/10.1016/S0269-7491(00)00162-7).
- (62) Arnot, J. A.; Gobas, F. A. P. C. A Review of Bioconcentration Factor (BCF) and Bioaccumulation Factor (BAF) Assessments for Organic Chemicals in Aquatic Organisms. *Environ. Rev.* **2006**, *14* (4), 257–297. <https://doi.org/10.1139/a06-005>.
- (63) Kelly, B. C.; Ikonou, M. G.; Blair, J. D.; Morin, A. E.; Gobas, F. A. P. C. Food Web-Specific Biomagnification of Persistent Organic Pollutants. *Science* **2007**, *317* (5835), 236–239. <https://doi.org/10.1126/science.1138275>.
- (64) OECD. *Guidelines for Testing of Chemicals, 305E, Bioaccumulation: Flow through Fish Test*; Organisation for Economic Co-Operation and Development: Paris, 1981.
- (65) Mackay, D. Correlation of Bioconcentration Factors. *Environ. Sci. Technol.* **1982**, *16* (5), 274–278. <https://doi.org/10.1021/es00099a008>.
- (66) Schuurmann, G.; Klein, W. Advances in Bioconcentration Prediction. *Chemosphere* **1988**, *17* (8), 1551–1574. [https://doi.org/10.1016/0045-6535\(88\)90207-X](https://doi.org/10.1016/0045-6535(88)90207-X).
- (67) Arnot, J. A.; Gobas, F. A. P. C. A Generic QSAR for Assessing the Bioaccumulation Potential of Organic Chemicals in Aquatic Food Webs. *QSAR Comb. Sci.* **2003**, *22* (3), 337–345. <https://doi.org/10.1002/qsar.200390023>.
- (68) Gobas, F. A. P. C.; Kelly, B. C.; Arnot, J. A. Quantitative Structure Activity Relationships for Predicting the Bioaccumulation of POPs in Terrestrial Food-Webs. *QSAR Comb. Sci.* **2003**, *22* (3), 329–336. <https://doi.org/10.1002/qsar.200390022>.

- (69) Manahan, S. E. *Toxicological Chemistry and Biochemistry, Third Edition*; CRC Press, 2002.
- (70) Abelkop, A. D. K.; Graham, J. *Regulation of Chemical Risks: Lessons for Reform of the Toxic Substances Control Act from Canada and the European Union*; SSRN Scholarly Paper ID 2499309; Social Science Research Network: Rochester, NY, 2014.
- (71) Coria, J. Policy Monitor—The Economics of Toxic Substance Control and the REACH Directive. *Rev. Environ. Econ. Policy* **2018**, *12* (2), 342–358. <https://doi.org/10.1093/reep/rey003>.
- (72) Schaafsma, G.; Kroese, E. D.; Tielemans, E. L. J. P.; Van de Sandt, J. J. M.; Van Leeuwen, C. J. REACH, Non-Testing Approaches and the Urgent Need for a Change in Mind Set. *Regul. Toxicol. Pharmacol.* **2009**, *53* (1), 70–80. <https://doi.org/10.1016/j.yrtph.2008.11.003>.
- (73) Snyder, S. A. Emerging Chemical Contaminants: Looking for Greater Harmony. *J. - Am. Water Works Assoc.* *106* (8), 38–52. <https://doi.org/10.5942/jawwa.2014.106.0126>.
- (74) Rovida, C.; Hartung, T. Re-Evaluation of Animal Numbers and Costs for in Vivo Tests to Accomplish REACH Legislation Requirements for Chemicals - a Report by the Transatlantic Think Tank for Toxicology (T4). *ALTEX* **2009**, 187–208. <https://doi.org/10.14573/altex.2009.3.187>.
- (75) Bradbury, S. P.; Feijtel, T. C. J.; Leeuwen, C. J. V. Peer Reviewed: Meeting the Scientific Needs of Ecological Risk Assessment in a Regulatory Context. *Environ. Sci. Technol.* **2004**, *38* (23), 463A-470A. <https://doi.org/10.1021/es040675s>.
- (76) Jagt, K. van der; Munn, S. J.; Tørsløv, J.; Bruijn, J. de. Alternative Approaches Can Reduce the Use of Test Animals under REACH November 2004; 2004.
- (77) Anastas, P. T.; Warner, J. C. *Green Chemistry: Theory and Practice*; Oxford University Press, 1998.
- (78) Rybacka, A.; Rudén, C.; Tetko, I. V.; Andersson, P. L. Identifying Potential Endocrine Disruptors among Industrial Chemicals and Their Metabolites – Development and Evaluation of in Silico Tools. *Chemosphere* **2015**, *139*, 372–378. <https://doi.org/10.1016/j.chemosphere.2015.07.036>.
- (79) OECD. Report on the Regulatory Uses and Applications in OECD Member Countries of (Quantitative) Structure-Activity Relationship [(Q)SAR] Models in the Assessment of New and Existing Chemicals. *OECD Pap.* **2006**, *6* (11), 1–79. https://doi.org/https://doi.org/10.1787/oecd_papers-v6-art37-en.
- (80) Knekta, E.; Andersson, P. L.; Johansson, M.; Tysklind, M. An Overview of OSPAR Priority Compounds and Selection of a Representative Training Set.

References

- Chemosphere* **2004**, *57* (10), 1495–1503. <https://doi.org/10.1016/j.chemosphere.2004.07.056>.
- (81) US EPA. PBT profiler; Persistent, Bioaccumulative, and Toxic Profiles Estimated for Organic Chemicals On-line <http://www.pbtprofiler.net/> (accessed Sep 8, 2015).
- (82) US-EPA. *Estimation Programs Interface Suite™ for Microsoft® Windows*; United States Environmental Protection Agency: Washington, DC, USA, 2012.
- (83) Meylan, W.; Howard, P. Computer Estimation of the Atmospheric Gas-Phase Reaction-Rate of Organic-Compounds with Hydroxyl Radicals and Ozone. *Chemosphere* **1993**, *26* (12), 2293–2299. [https://doi.org/10.1016/0045-6535\(93\)90355-9](https://doi.org/10.1016/0045-6535(93)90355-9).
- (84) Boethling, R. S.; Howard, P. H.; Meylan, W.; Stiteler, W.; Beauman, J.; Tirado, N. Group Contribution Method for Predicting Probability and Rate of Aerobic Biodegradation. *Environ. Sci. Technol.* **1994**, *28* (3), 459–465. <https://doi.org/10.1021/es00052a018>.
- (85) Mackay, D.; Paterson, S.; Shiu, W. Generic Models for Evaluating the Regional Fate of Chemicals. *Chemosphere* **1992**, *24* (6), 695–717. [https://doi.org/10.1016/0045-6535\(92\)90531-U](https://doi.org/10.1016/0045-6535(92)90531-U).
- (86) US-EPA. *The ECOSAR (ECOLOGICAL Structure Activity Relationship) Class Program*; United States Environmental Protection Agency, 2012.
- (87) Howard, P. H.; Muir, D. C. G. Identifying New Persistent and Bioaccumulative Organics Among Chemicals in Commerce II: Pharmaceuticals. *Environ. Sci. Technol.* **2011**, *45* (16), 6938–6946. <https://doi.org/10.1021/es201196x>.
- (88) Howard, P. H.; Muir, D. C. G. Identifying New Persistent and Bioaccumulative Organics Among Chemicals in Commerce. *Environ. Sci. Technol.* **2010**, *44* (7), 2277–2285. <https://doi.org/10.1021/es903383a>.
- (89) Arnot, J. A.; Mackay, D. Policies for Chemical Hazard and Risk Priority Setting: Can Persistence, Bioaccumulation, Toxicity, and Quantity Information Be Combined. *Env. Sci Technol* **2008**, *42*, 4648–4654.
- (90) Pavan, M.; Worth, A. *A Set of Case Studies to Illustrate the Applicability of DART (Decision Analysis by Ranking Techniques) in the Ranking of Chemicals.*; European Commission report EUR 23481 EN; European Commission: Office for Official Publications of the European Communities, Luxemburg, 2008.
- (91) Papa, E.; Gramatica, P. QSPR as a Support for the EU REACH Regulation and Rational Design of Environmentally Safer Chemicals: PBT Identification from Molecular Structure. *Green Chem* **2010**, *12*, 836–843. <https://doi.org/10.1039/B923843C>.

- (92) Rorije, E.; Verbruggen, E.; Hollander, A.; Traas, T.; Janssen, M. *Identifying Potential POP and PBT Substances: Development of a New Persistence/Bioaccumulation-Score*; 601356001/2011; RIVM, 2011; pp 1–88.
- (93) Stempel, S.; Scheringer, M.; Ng, C. A.; Hungerbühler, K. Screening for PBT Chemicals among the “Existing” and “New” Chemicals of the EU. *Env. Sci Technol* **2012**, *46*, 5680–5687. <https://doi.org/10.1021/es3002713>.
- (94) Pizzo, F.; Lombardo, A.; Manganaro, A.; Cappelli, C. I.; Petoumenou, M. I.; Albanese, F.; Roncaglioni, A.; Brandt, M.; Benfenati, E. Integrated in Silico Strategy for PBT Assessment and Prioritization under REACH. *Environ. Res.* **2016**, *151* (Supplement C), 478–492. <https://doi.org/10.1016/j.envres.2016.08.014>.
- (95) De, P.; Roy, K. Greener Chemicals for the Future: QSAR Modelling of the PBT Index Using ETA Descriptors. *SAR QSAR Environ. Res.* **2018**, *29* (4), 319–337. <https://doi.org/10.1080/1062936X.2018.1436086>.
- (96) Asikainen, A. H.; Ruuskanen, J.; Tuppurainen, K. A. Consensus KNN QSAR: A Versatile Method for Predicting the Estrogenic Activity of Organic Compounds in Silico. A Comparative Study with Five Estrogen Receptors and a Large, Diverse Set of Ligands. *Environ. Sci. Technol.* **2004**, *38* (24), 6724–6729. <https://doi.org/10.1021/es049665h>.
- (97) Gramatica, P.; Giani, E.; Papa, E. Statistical External Validation and Consensus Modeling: A QSPR Case Study for K-Oc Prediction. *J. Mol. Graph. Model.* **2007**, *25* (6), 755–766. <https://doi.org/10.1016/j.jmglm.2006.06.005>.
- (98) Gramatica, P.; Pilutti, P.; Papa, E. Validated QSAR Prediction of OH Tropospheric Degradation of VOCs: Splitting into Training-Test Sets and Consensus Modeling. *J. Chem. Inf. Comput. Sci.* **2004**, *44* (5), 1794–1802. <https://doi.org/10.1021/ci049923u>.
- (99) Zhu, H.; Tropsha, A.; Fourches, D.; Varnek, A.; Papa, E.; Gramatica, P.; Öberg, T.; Dao, P.; Cherkasov, A.; Tetko, I. V. Combinatorial QSAR Modeling of Chemical Toxicants Tested against *Tetrahymena Pyriformis*. *J. Chem. Inf. Model.* **2008**, *48* (4), 766–784. <https://doi.org/10.1021/ci700443v>.
- (100) Madden, J. C.; Enoch, S. J.; Hewitt, M.; Cronin, M. T. D. Pharmaceuticals in the Environment: Good Practice in Predicting Acute Ecotoxicological Effects. *Toxicol. Lett.* **2009**, *185* (2), 85–101. <https://doi.org/10.1016/j.toxlet.2008.12.005>.
- (101) Moermond, C. T.; Janssen, M. P.; de Knecht, J. A.; Montforts, M. H.; Peijnenburg, W. J.; Zweers, P. G.; Sijm, D. T. PBT Assessment Using the Revised Annex XIII of REACH: A Comparison With Other Regulatory Frameworks. *Integr. Environ. Assess. Manag.* **2011**, *8*, 359–371. <https://doi.org/10.1002/ieam.1248>.

References

- (102) Rauert, C.; Friesen, A.; Hermann, G.; Jöhncke, U.; Kehrer, A.; Neumann, M.; Prutz, I.; Schönfeld, J.; Wiemann, A.; Willhaus, K.; et al. Proposal for a Harmonised PBT Identification across Different Regulatory Frameworks. *Environ. Sci. Eur.* **2014**, *26* (1), 9. <https://doi.org/10.1186/2190-4715-26-9>.
- (103) Gramatica, P.; Cassani, S.; Sangion, A. PBT Assessment and Prioritization by PBT Index and Consensus Modeling: Comparison of Screening Results from Structural Models. *Environ. Int.* **2015**, *77*, 25–34. <https://doi.org/10.1016/j.envint.2014.12.012>.
- (104) Gramatica, P.; Cassani, S.; Sangion, A. Are Some “Safer Alternatives” Hazardous as PBTs? The Case Study of New Flame Retardants. *J. Hazard. Mater.* **2016**, *306*, 237–246. <https://doi.org/10.1016/j.jhazmat.2015.12.017>.
- (105) Cassani, S.; Gramatica, P. Identification of Potential PBT Behavior of Personal Care Products by Structural Approaches. *Sustain. Chem. Pharm.* **2015**, *1*, 19–27. <https://doi.org/10.1016/j.scp.2015.10.002>.
- (106) Lipnick, R. L.; Muir, D. C. G. History of Persistent, Bioaccumulative, and Toxic Chemicals. In *Persistent, Bioaccumulative, and Toxic Chemicals I*; ACS Symposium Series; American Chemical Society, 2000; Vol. 772, pp 1–12. <https://doi.org/10.1021/bk-2001-0772.ch001>.
- (107) Carson, R. *Silent Spring*; 1962.
- (108) Government of Canada. *Canadian Environmental Protection Act Canadian Environmental Protection Act, Canada Gazette Part III.*; 1999.
- (109) Matthies, M.; Solomon, K.; Vighi, M.; Gilman, A.; Tarazona, J. V. The Origin and Evolution of Assessment Criteria for Persistent, Bioaccumulative and Toxic (PBT) Chemicals and Persistent Organic Pollutants (POPs). *Environ. Sci. Process. Impacts* **2016**, *18* (9), 1114–1128. <https://doi.org/10.1039/C6EM00311G>.
- (110) Schulze, S.; Sättler, D.; Neumann, M.; Arp, H. P. H.; Reemtsma, T.; Berger, U. Using REACH Registration Data to Rank the Environmental Emission Potential of Persistent and Mobile Organic Chemicals. *Sci. Total Environ.* **2018**, *625*, 1122–1128. <https://doi.org/10.1016/j.scitotenv.2017.12.305>.
- (111) Reemtsma, T.; Berger, U.; Arp, H. P. H.; Gallard, H.; Knepper, T. P.; Neumann, M.; Quintana, J. B.; Voogt, P. de. Mind the Gap: Persistent and Mobile Organic Compounds—Water Contaminants That Slip Through. *Environ. Sci. Technol.* **2016**, *50* (19), 10308–10315. <https://doi.org/10.1021/acs.est.6b03338>.
- (112) Berger, U.; Reemtsma, T.; Ost, N.; Kühne, R.; Schüürmann, G.; Sättler, D.; Schliebner, I.; Neumann, M.; Schüürmann, G. Assessment of Persistence, Mobility and Toxicity (PMT) of 167 REACH Registered Substances. 61.
- (113) Kelly, B. C.; Gobas, F. A. P. C. Bioaccumulation of Persistent Organic Pollutants in Lichen–Caribou–Wolf Food Chains of Canada’s Central and

- Western Arctic. *Environ. Sci. Technol.* **2001**, *35* (2), 325–334. <https://doi.org/10.1021/es0011966>.
- (114) Armitage, J. M.; Gobas, F. A. P. C. A Terrestrial Food-Chain Bioaccumulation Model for POPs. *Environ. Sci. Technol.* **2007**, *41* (11), 4019–4025. <https://doi.org/10.1021/es0700597>.
- (115) McLachlan, M. S.; Czub, G.; MacLeod, M.; Arnot, J. A. Bioaccumulation of Organic Contaminants in Humans: A Multimedia Perspective and the Importance of Biotransformation †. *Environ. Sci. Technol.* **2011**, *45* (1), 197–202. <https://doi.org/10.1021/es101000w>.
- (116) Gramatica, P.; Papa, E. An Update of the BCF QSAR Model Based on Theoretical Molecular Descriptors. *Qsar Comb. Sci.* **2005**, *24* (8), 953–960. <https://doi.org/10.1002/qsar.200530123>.
- (117) Meylan, W. M.; Howard, P. H.; Boethling, R. S.; Aronson, D.; Printup, H.; Gouchie, S. Improved Method for Estimating Bioconcentration/Bioaccumulation Factor from Octanol/Water Partition Coefficient. *Environ. Toxicol. Chem.* **1999**, *18* (4), 664–672. [https://doi.org/10.1897/1551-5028\(1999\)018<0664:IMFEBB>2.3.CO;2](https://doi.org/10.1897/1551-5028(1999)018<0664:IMFEBB>2.3.CO;2).
- (118) ILSI Health and Environmental Sciences Institute (HESI); The Joint Research Centre (JRC); SETAC-Europe. Workshop on Bioaccumulation Assessments; The Hague, The Netherlands, 2006.
- (119) Veith, G. D.; DeFoe, D. L.; Bergstedt, B. V. Measuring and Estimating the Bioconcentration Factor of Chemicals in Fish. *J. Fish. Res. Board Can.* **1979**, *36* (9), 1040–1048. <https://doi.org/10.1139/f79-146>.
- (120) Bintein, S.; Devillers, J.; Karcher, W. Nonlinear Dependence of Fish Bioconcentration on N-Octanol/Water Partition Coefficient. *SAR QSAR Environ. Res.* **1993**, *1* (1), 29–39. <https://doi.org/10.1080/10629369308028814>.
- (121) Hawker, D. W.; Connell, D. W. Octanol-Water Partition Coefficients of Polychlorinated Biphenyl Congeners. *Environ. Sci. Technol.* **1988**, *22* (4), 382–387. <https://doi.org/10.1021/es00169a004>.
- (122) Zhao, C.; Boriani, E.; Chana, A.; Roncaglioni, A.; Benfenati, E. A New Hybrid System of QSAR Models for Predicting Bioconcentration Factors (BCF). *Chemosphere* **2008**, *73* (11), 1701–1707. <https://doi.org/10.1016/j.chemosphere.2008.09.033>.
- (123) Aranda, J. F.; Bacelo, D. E.; Aparicio, M. S. L.; Ocsachoque, M. A.; Castro, E. A.; Duchowicz, P. R. Predicting the Bioconcentration Factor through a Conformation-Independent QSPR Study. *SAR QSAR Environ. Res.* **2017**, *28* (9), 749–763. <https://doi.org/10.1080/1062936X.2017.1377765>.
- (124) Nendza, M.; Kühne, R.; Lombardo, A.; Stempel, S.; Schüürmann, G. PBT Assessment under REACH: Screening for Low Aquatic Bioaccumulation with

References

- QSAR Classifications Based on Physicochemical Properties to Replace BCF in Vivo Testing on Fish. *Sci. Total Environ.* **2018**, 616–617, 97–106. <https://doi.org/10.1016/j.scitotenv.2017.10.317>.
- (125) Grisoni, F.; Consonni, V.; Vighi, M.; Villa, S.; Todeschini, R. Expert QSAR System for Predicting the Bioconcentration Factor under the REACH Regulation. *Environ. Res.* **2016**, 148, 507–512. <https://doi.org/10.1016/j.envres.2016.04.032>.
- (126) Arnot, J. A.; Brown, T. N.; Wania, F. Estimating Screening-Level Organic Chemical Half-Lives in Humans. *Environ. Sci. Technol.* **2014**, 48 (1), 723–730. <https://doi.org/10.1021/es4029414>.
- (127) Czub, G.; McLachlan, M. S. A Food Chain Model to Predict the Levels of Lipophilic Organic Contaminants in Humans. *Environ. Toxicol. Chem.* **2004**, 23 (10), 2356–2366. <https://doi.org/10.1897/03-317>.
- (128) Czub, G.; McLachlan, M. S. Bioaccumulation Potential of Persistent Organic Chemicals in Humans. *Environ. Sci. Technol.* **2004**, 38 (8), 2406–2412. <https://doi.org/10.1021/es034871v>.
- (129) Goss, K.-U.; Brown, T. N.; Endo, S. Elimination Half-Life as a Metric for the Bioaccumulation Potential of Chemicals in Aquatic and Terrestrial Food Chains. *Environ. Toxicol. Chem.* **2013**, 32 (7), 1663–1671. <https://doi.org/10.1002/etc.2229>.
- (130) Cowan-Ellsberry, C. E.; Dyer, S. D.; Erhardt, S.; Bernhard, M. J.; Roe, A. L.; Dowty, M. E.; Weisbrod, A. V. Approach for Extrapolating in Vitro Metabolism Data to Refine Bioconcentration Factor Estimates. *Chemosphere* **2008**, 70 (10), 1804–1817. <https://doi.org/10.1016/j.chemosphere.2007.08.030>.
- (131) Laue, H.; Gfeller, H.; Jenner, K. J.; Nichols, J. W.; Kern, S.; Natsch, A. Predicting the Bioconcentration of Fragrance Ingredients by Rainbow Trout Using Measured Rates of in Vitro Intrinsic Clearance. *Environ. Sci. Technol.* **2014**, 48 (16), 9486–9495. <https://doi.org/10.1021/es500904h>.
- (132) Walker, C. H.; Sibly, R. M.; Hopkin, S. P.; Peakall, D. B. *Principles of Ecotoxicology, Fourth Edition*; CRC Press, 2012.
- (133) Mekenyan, O.; Dimitrov, S.; Pavlov, T.; Dimitrova, G.; Todorov, M.; Petkov, P.; Kotov, S. Simulation of Chemical Metabolism for Fate and Hazard Assessment. V. Mammalian Hazard Assessment. *SAR QSAR Environ. Res.* **2012**, 23 (5–6), 553–606. <https://doi.org/10.1080/1062936X.2012.679689>.
- (134) Sevier, D. K.; Pelkonen, O.; Ahokas, J. T. Hepatocytes: The Powerhouse of Biotransformation. *Int. J. Biochem. Cell Biol.* **2012**, 44 (2), 257–261. <https://doi.org/10.1016/j.biocel.2011.11.011>.

- (135) Cwiertny, D. M.; Snyder, S. A.; Schlenk, D.; Kolodziej, E. P. Environmental Designer Drugs: When Transformation May Not Eliminate Risk. *Environ. Sci. Technol.* **2014**, *48* (20), 11737–11745. <https://doi.org/10.1021/es503425w>.
- (136) Sijm, D. T. H. M.; Rikken, M. G. J.; Rorije, E.; Traas, T. P.; Mclachlan, M. S.; Peijnenburg, W. J. G. M. Transport, Accumulation and Transformation Processes. In *Risk Assessment of Chemicals*; Leeuwen, C. J. van, Vermeire, T. G., Eds.; Springer Netherlands, 2007; pp 73–158.
- (137) Arnot, J. A.; Mackay, D.; Bonnell, M. Estimating Metabolic Biotransformation Rates in Fish from Laboratory Data. *Environ. Toxicol. Chem.* **2008**, *27* (2), 341–351. <https://doi.org/10.1897/07-310R.1>.
- (138) Pirovano, A.; Brandmaier, S.; Huijbregts, M. A. J.; Ragas, A. M. J.; Veltman, K.; Hendriks, A. J. QSARs for Estimating Intrinsic Hepatic Clearance of Organic Chemicals in Humans. *Environ. Toxicol. Pharmacol.* **2016**, *42*, 190–197. <https://doi.org/10.1016/j.etap.2016.01.017>.
- (139) Nichols, J. W.; Huggett, D. B.; Arnot, J. A.; Fitzsimmons, P. N.; Cowan-Ellsberry, C. E. Toward Improved Models for Predicting Bioconcentration of Well-Metabolized Compounds by Rainbow Trout Using Measured Rates of in Vitro Intrinsic Clearance. *Environ. Toxicol. Chem.* **2013**, *32* (7), 1611–1622. <https://doi.org/10.1002/etc.2219>.
- (140) Borodina, Y.; Sadym, A.; Filimonov, D.; Blinova, V.; Dmitriev, A.; Poroikov, V. Predicting Biotransformation Potential from Molecular Structure. *J. Chem. Inf. Comput. Sci.* **2003**, *43* (5), 1636–1646. <https://doi.org/10.1021/ci034078l>.
- (141) Arnot, J. A.; Meylan, W.; Tunkel, J.; Howard, P. H.; Mackay, D.; Bonnell, M.; Boethling, R. S. A Quantitative Structure-Activity Relationship for Predicting Metabolic Biotransformation Rates for Organic Chemicals in Fish. *Environ. Toxicol. Chem.* **2009**, *28* (6), 1168–1177. <https://doi.org/10.1897/08-289.1>.
- (142) Brown, T. N.; Arnot, J. A.; Wania, F. Iterative Fragment Selection: A Group Contribution Approach to Predicting Fish Biotransformation Half-Lives. *Environ. Sci. Technol.* **2012**, *46* (15), 8253–8260. <https://doi.org/10.1021/es301182a>.
- (143) Papa, E.; van der Wal, L.; Arnot, J. A.; Gramatica, P. Metabolic Biotransformation Half-Lives in Fish: QSAR Modeling and Consensus Analysis. *Sci. Total Environ.* **2014**, *470*, 1040–1046. <https://doi.org/10.1016/j.scitotenv.2013.10.068>.
- (144) Kuo, D. T. F.; Di Toro, D. M. Biotransformation Model of Neutral and Weakly Polar Organic Compounds in Fish Incorporating Internal Partitioning. *Environ. Toxicol. Chem.* **2013**, *32* (8), 1873–1881. <https://doi.org/10.1002/etc.2259>.
- (145) Papa, E.; Villa, F.; Gramatica, P. Statistically Validated QSARs, Based on Theoretical Descriptors, for Modeling Aquatic Toxicity of Organic Chemicals

References

- in Pimephales Promelas (Fathead Minnow). *J. Chem. Inf. Model.* **2005**, *45* (5), 1256–1266. <https://doi.org/10.1021/ci050212l>.
- (146) OECD. *Test No. 201: Freshwater Alga and Cyanobacteria, Growth Inhibition Test*; Organisation for Economic Co-operation and Development: Paris, 2011.
- (147) OECD. *Test No. 202: Daphnia Sp. Acute Immobilisation Test*; Organisation for Economic Co-operation and Development: Paris, 2004.
- (148) OECD. *Test No. 203: Fish, Acute Toxicity Test*; Organisation for Economic Co-operation and Development: Paris, 1992.
- (149) Erzincan, P.; Saçan, M. T.; Yuce-Dursun, B.; Danis, O.; Demir, S.; Erdem, S. S.; Ogan, A. QSAR Models for Antioxidant Activity of New Coumarin Derivatives. *Sar Qsar Environ. Res.* **2015**, *26* (7–9), 721–737. <https://doi.org/10.1080/1062936X.2015.1088571>.
- (150) Önlü, S.; Saçan, M. T. An in Silico Algal Toxicity Model with a Wide Applicability Potential for Industrial Chemicals and Pharmaceuticals. *Environ. Toxicol. Chem.* **2016**. <https://doi.org/10.1002/etc.3620>.
- (151) Gramatica, P. Prioritization of Chemicals Based on Chemoinformatic Analysis. In *Handbook of Computational Chemistry*; Leszczynski, J., Ed.; Springer Netherlands, 2016; pp 1–33.
- (152) Sanderson, H. Challenges and Directions for Regulatory Use of QSARs for Predicting Active Pharmaceutical Ingredients Environmental Toxicity. *Curr. Drug Saf.* **2012**, *7* (4), 309–312. <https://doi.org/10.2174/157488612804096597>.
- (153) CSTEE. Discussion Paper on Environmental Risk Assessment of Medical Products for Human Use (Non-Genetically Modified Organisms (Non-GMO) Containing). CPMppaperRAssessHumPharm12062001/D(01). European Commission 2001.
- (154) Fent, K.; Weston, A. A.; Caminada, D. Ecotoxicology of Human Pharmaceuticals. *Aquat. Toxicol.* **2006**, *76* (2), 122–159. <https://doi.org/10.1016/j.aquatox.2005.09.009>.
- (155) Mendoza, A.; Acena, J.; Perez, S.; Lopez de Alda, M.; Barcelo, D.; Gil, A.; Valcarcel, Y. Pharmaceuticals and Iodinated Contrast Media in a Hospital Wastewater: A Case Study to Analyse Their Presence and Characterise Their Environmental Risk and Hazard. *Environ. Res.* **2015**, *140*, 225–241. <https://doi.org/10.1016/j.envres.2015.04.003>.
- (156) Ortiz de Garcia, S. A.; Pinto Pinto, G.; Garcia-Encina, P. A.; Irusta-Mata, R. Ecotoxicity and Environmental Risk Assessment of Pharmaceuticals and Personal Care Products in Aquatic Environments and Wastewater Treatment Plants. *Ecotoxicology* **2014**, *23* (8), 1517–1533. <https://doi.org/10.1007/s10646-014-1293-8>.

- (157) Sanderson, H.; Johnson, D. J.; Wilson, C. J.; Brain, R. A.; Solomon, K. R. Probabilistic Hazard Assessment of Environmentally Occurring Pharmaceuticals Toxicity to Fish, Daphnids and Algae by ECOSAR Screening. *Toxicol. Lett.* **2003**, *144* (3), 383–395. [https://doi.org/10.1016/S0378-4274\(03\)00257-1](https://doi.org/10.1016/S0378-4274(03)00257-1).
- (158) Sanderson, H.; Thomsen, M. Ecotoxicological Quantitative Structure-Activity Relationships for Pharmaceuticals. *Bull. Environ. Contam. Toxicol.* **2007**, *79* (3), 331–335. <https://doi.org/10.1007/s00128-007-9249-9>.
- (159) Sanderson, H.; Thomsen, M. Comparative Analysis of Pharmaceuticals versus Industrial Chemicals Acute Aquatic Toxicity Classification According to the United Nations Classification System for Chemicals. Assessment of the (Q)SAR Predictability of Pharmaceuticals Acute Aquatic Toxicity and Their Predominant Acute Toxic Mode-of-Action. *Toxicol. Lett.* **2009**, *187* (2), 84–93. <https://doi.org/10.1016/j.toxlet.2009.02.003>.
- (160) Jiang, L.; Lin, Z.; Hu, X.; Yin, D. Toxicity Prediction of Antibiotics on Luminescent Bacteria, *Photobacterium Phosphoreum*, Based on Their Quantitative Structure-Activity Relationship Models. *Bull. Environ. Contam. Toxicol.* **2010**, *85* (6), 550–555. <https://doi.org/10.1007/s00128-010-0157-z>.
- (161) Roy, K.; Ghosh, G. Exploring QSARs with Extended Topochemical Atom (ETA) Indices for Modeling Chemical and Drug Toxicity. *Curr. Pharm. Des.* **2010**, *16* (24), 2625–2639. <https://doi.org/10.2174/138161210792389270>.
- (162) Tugcu, G.; Saçan, M. T.; Vracko, M.; Novic, M.; Minovski, N. QSTR Modelling of the Acute Toxicity of Pharmaceuticals to Fish. *Sar Qsar Environ. Res.* **2012**, *23* (3–4), 297–310. <https://doi.org/10.1080/1062936X.2012.657678>.
- (163) Singh, K. P.; Gupta, S.; Basant, N. QSTR Modeling for Predicting Aquatic Toxicity of Pharmacological Active Compounds in Multiple Test Species for Regulatory Purpose. *Chemosphere* **2015**, *120*, 680–689. <https://doi.org/10.1016/j.chemosphere.2014.10.025>.
- (164) Önlü, S.; Saçan, M. T. An in Silico Approach to Cytotoxicity of Pharmaceuticals and Personal Care Products on the Rainbow Trout Liver Cell Line RTL-W1. *Environ. Toxicol. Chem.* **2016**. <https://doi.org/10.1002/etc.3663>.
- (165) Yap, C. W. PaDEL-Descriptor: An Open Source Software to Calculate Molecular Descriptors and Fingerprints. *JComput Chem* **2011**, *32*, 1466–1474. <https://doi.org/10.1002/jcc.21707>.
- (166) Gramatica, P.; Chirico, N.; Papa, E.; Cassani, S.; Kovarich, S. QSARINS: A New Software for the Development, Analysis and Validation of QSAR MLR Models. *J Comput Chem* **2013**, *34*, 2121–2132. <https://doi.org/10.1002/jcc.23361>.

References

- (167) Gramatica, P.; Cassani, S.; Chirico, N. QSARINS-Chem: Insubria Datasets and New QSAR/QSPR Models for Environmental Pollutants in QSARINS. *J. Comput. Chem.* **2014**, *35* (13), 1036–1044. <https://doi.org/10.1002/jcc.23576>.
- (168) Crum-Brown and Fraser. *On the Connection Between Chemical Constitution and Physiological Action: Pt. I. - On the Physiological Action of the Salts of the Ammonium Bases, Derived from Strychnia, Brucia, Thebaia, Codeia, Morphia, and Nicotia. Pt. 2 - On the Physiological Action of the Ammonium Bases Derived from Atropia and Conia*; 1869.
- (169) H. Meyer, K. Contributions to the Theory of Narcosis. *Trans. Faraday Soc.* **1937**, *33* (0), 1062–1064. <https://doi.org/10.1039/TF9373301062>.
- (170) Hansch, C.; Maloney, P. P.; Fujita, T.; Muir, R. M. Correlation of Biological Activity of Phenoxyacetic Acids with Hammett Substituent Constants and Partition Coefficients. *Nature* **1962**, *194* (4824), 178–180. <https://doi.org/10.1038/194178b0>.
- (171) Hansch, C.; Fujita, T. P- σ - π Analysis. A Method for the Correlation of Biological Activity and Chemical Structure. *J. Am. Chem. Soc.* **1964**, *86* (8), 1616–1626. <https://doi.org/10.1021/ja01062a035>.
- (172) Free, S. M.; Wilson, J. W. A Mathematical Contribution to Structure-Activity Studies. *J. Med. Chem.* **1964**, *7* (4), 395–399. <https://doi.org/10.1021/jm00334a001>.
- (173) Hammett, L. P. Some Relations between Reaction Rates and Equilibrium Constants. *Chem. Rev.* **1935**, *17* (1), 125–136. <https://doi.org/10.1021/cr60056a010>.
- (174) Dimitrov, S.; Dimitrova, G.; Pavlov, T.; Dimitrova, N.; Patlewicz, G.; Niemela, J.; Mekenyan, O. A Stepwise Approach for Defining the Applicability Domain of SAR and QSAR Models. *J. Chem. Inf. Model.* **2005**, *45* (4), 839–849. <https://doi.org/10.1021/ci0500381>.
- (175) Eriksson, L.; Andersson, P. L.; Johansson, E.; Tysklind, M. Megavariate Analysis of Environmental QSAR Data. Part I – A Basic Framework Founded on Principal Component Analysis (PCA), Partial Least Squares (PLS), and Statistical Molecular Design (SMD). *Mol. Divers.* **2006**, *10* (2), 169–186. <https://doi.org/10.1007/s11030-006-9024-6>.
- (176) Eriksson, L.; Andersson, P. L.; Johansson, E.; Tysklind, M. Megavariate Analysis of Environmental QSAR Data. Part II – Investigating Very Complex Problem Formulations Using Hierarchical, Non-Linear and Batch-Wise Extensions of PCA and PLS. *Mol. Divers.* **2006**, *10* (2), 187–205. <https://doi.org/10.1007/s11030-006-9026-4>.
- (177) Gissi, A.; Gadaleta, D.; Floris, M.; Olla, S.; Carotti, A.; Novellino, E.; Benfenati, E.; Nicolotti, O. An Alternative QSAR-Based Approach for Predicting the

- Bioconcentration Factor for Regulatory Purposes. *ALTEX* **2014**, *31* (1), 23–36. <https://doi.org/10.14573/altex.1305221>.
- (178) Jaworska, J. S.; Comber, M.; Auer, C.; Van Leeuwen, C. J. Summary of a Workshop on Regulatory Acceptance of (Q)SARs for Human Health and Environmental Endpoints. *Environ. Health Perspect.* **2003**, *111* (10), 1358–1360. <https://doi.org/10.1289/ehp.5757>.
- (179) OECD. *Guidance Document on the Validation of (Quantitative) Structure-Activity Relationships [(Q)SAR] Models*; Organisation for Economic Co-operation and Development: Paris, 2007.
- (180) OECD. *Principles for the Validation, for Regulatory Purposes, of (Quantitative) Structure-Activity Relationship Models*; Organisation for Economic Co-operation and Development: Paris, 2004.
- (181) Todeschini Roberto. *Introduzione Alla Chemiometria*; EdiSES, 1998.
- (182) Todeschini, R.; Consonni, V. *Handbook of Molecular Descriptors*; John Wiley & Sons, 2008.
- (183) Wold, S.; Esbensen, K.; Geladi, P. Principal Component Analysis. *Chemom. Intell. Lab. Syst.* **1987**, *2* (1–3), 37–52.
- (184) Jackson, J. E. *A User's Guide to Principal Components*; John Wiley & Sons, 2005.
- (185) Shmueli, G. To Explain or to Predict? *Stat. Sci.* **2010**, *25* (3), 289–310. <https://doi.org/10.1214/10-STS330>.
- (186) Tropsha, A.; Gramatica, P.; Gombar, V. K. The Importance of Being Earnest: Validation Is the Absolute Essential for Successful Application and Interpretation of QSPR Models. *Qsar Comb. Sci.* **2003**, *22* (1), 69–77. <https://doi.org/10.1002/qsar.200390007>.
- (187) Gramatica, P.; Cassani, S.; Roy, P. P.; Kovarich, S.; Yap, C. W.; Papa, E. QSAR Modeling Is Not “Push a Button and Find a Correlation”: A Case Study of Toxicity of (Benzo-)Triazoles on Algae. *Mol Inf* **2012**, *31*, 817–835. <https://doi.org/10.1002/minf.201200075>.
- (188) Tropsha, A. Best Practices for QSAR Model Development, Validation, and Exploitation. *Mol. Inform.* **2010**, *29* (6–7), 476–488. <https://doi.org/10.1002/minf.201000061>.
- (189) Fourches, D.; Muratov, E.; Tropsha, A. Trust, But Verify: On the Importance of Chemical Structure Curation in Cheminformatics and QSAR Modeling Research. *J. Chem. Inf. Model.* **2010**, *50* (7), 1189–1204. <https://doi.org/10.1021/ci100176x>.
- (190) Li, J.; Gramatica, P. The Importance of Molecular Structures, Endpoints' Values, and Predictivity Parameters in QSAR Research: QSAR Analysis of a

References

- Series of Estrogen Receptor Binders. *Mol. Divers.* **2010**, *14* (4), 687–696. <https://doi.org/10.1007/s11030-009-9212-2>.
- (191) Muehlbacher, M.; El Kerdawy, A.; Kramer, C.; Hudson, B.; Clark, T. Conformation-Dependent QSPR Models: LogP(OW). *J. Chem. Inf. Model.* **2011**, *51* (9), 2408–2416. <https://doi.org/10.1021/ci200276v>.
- (192) Porcelli, C.; Boriani, E.; Roncaglioni, A.; Chana, A.; Benfenati, E. Regulatory Perspectives in the Use and Validation of QSAR. A Case Study: DEMETRA Model for Daphnia Toxicity. *Environ. Sci. Technol.* **2008**, *42* (2), 491–496. <https://doi.org/10.1021/es071430t>.
- (193) Young, D.; Martin, T.; Venkatapathy, R.; Harten, P. Are the Chemical Structures in Your QSAR Correct? *Qsar Comb. Sci.* **2008**, *27* (11–12), 1337–1345. <https://doi.org/10.1002/qsar.200810084>.
- (194) Todeschini, R.; Consonni, V.; Maiocchi, A. The K Correlation Index: Theory Development and Its Application in Chemometrics. *Chemom. Intell. Lab. Syst.* **1999**, *46* (1), 13–29. [https://doi.org/10.1016/S0169-7439\(98\)00124-5](https://doi.org/10.1016/S0169-7439(98)00124-5).
- (195) Gramatica, P. External Evaluation of QSAR Models, in Addition to CrossValidation: Verification of Predictive Capability on Totally New Chemicals. *Mol. Inform.* **2014**, *33* (4), 311–314. <https://doi.org/10.1002/minf.201400030>.
- (196) Gasteiger, J.; Zupan, J. Neural Networks in Chemistry. *Angew. Chem.-Int. Ed. Engl.* **1993**, *32* (4), 503–527. <https://doi.org/10.1002/anie.199305031>.
- (197) Eriksson, L.; Johansson, E.; Müller, M.; Wold, S. On the Selection of the Training Set in Environmental QSAR Analysis When Compounds Are Clustered. *J. Chemom.* **2000**, *14* (5–6), 599–616. [https://doi.org/10.1002/1099-128X\(200009/12\)14:5/6<599::AID-CEM619>3.0.CO;2-8](https://doi.org/10.1002/1099-128X(200009/12)14:5/6<599::AID-CEM619>3.0.CO;2-8).
- (198) Eriksson, L.; Johansson, E. Multivariate Design and Modeling in QSAR. *Chemom. Intell. Lab. Syst.* **1996**, *34* (1), 1–19. [https://doi.org/10.1016/0169-7439\(96\)00023-8](https://doi.org/10.1016/0169-7439(96)00023-8).
- (199) Papa, E.; Fick, J.; Lindberg, R.; Johansson, M.; Gramatica, P.; Andersson, P. L. Multivariate Chemical Mapping of Antibiotics and Identification of Structurally Representative Substances. *Environ. Sci. Technol.* **2007**, *41* (5), 1653–1661. <https://doi.org/10.1021/es060618u>.
- (200) Martin, T. M.; Harten, P.; Young, D. M.; Muratov, E. N.; Golbraikh, A.; Zhu, H.; Tropsha, A. Does Rational Selection of Training and Test Sets Improve the Outcome of QSAR Modeling? *J. Chem. Inf. Model.* **2012**, *52* (10), 2570–2578. <https://doi.org/10.1021/ci300338w>.
- (201) Ballabio, D.; Skov, T.; Leardi, R.; Bro, R. Classification of GC-MS Measurements of Wines by Combining Data Dimension Reduction and Variable

- Selection Techniques. *J. Chemom.* **2008**, *22* (8), 457–463. <https://doi.org/10.1002/cem.1173>.
- (202) Leardi, R.; Boggia, R.; Terrile, M. Genetic Algorithms as a Strategy for Feature-Selection. *J. Chemom.* **1992**, *6* (5), 267–281. <https://doi.org/10.1002/cem.1180060506>.
- (203) Haupt, R. L.; Haupt, S. E. *Practical Genetic Algorithms*; John Wiley & Sons, 2004.
- (204) Trevor Hastie; Robert Tibshirani; Jerome Friedman. *The Elements of Statistical Learning*; Springer, 2009.
- (205) Farahani, H.; Rahiminezhad, A.; Same, L.; Immannezhad, K. A Comparison of Partial Least Squares (PLS) and Ordinary Least Squares (OLS) Regressions in Predicting of Couples Mental Health Based on Their Communicational Patterns. *Procedia Soc. Behav. Sci.* **2010**, *5*, 1459–1463. <https://doi.org/doi:10.1016/j.sbspro.2010.07.308>.
- (206) Cruciani, G.; Baroni, M.; Clementi, S.; Costantino, G.; Riganelli, D.; Skagerberg, B. Predictive Ability of Regression Models. Part I: Standard Deviation of Prediction Errors (SDEP). *J. Chemom.* **1992**, *6* (6), 335–346. <https://doi.org/10.1002/cem.1180060604>.
- (207) Golbraikh, A.; Tropsha, A. Beware of Q2! *J. Mol. Graph. Model.* **2002**, *20* (4), 269–276. [https://doi.org/10.1016/S1093-3263\(01\)00123-1](https://doi.org/10.1016/S1093-3263(01)00123-1).
- (208) Lindgren, F.; Hansen, B.; Karcher, W.; Sjostrom, M.; Eriksson, L. Model Validation by Permutation Tests: Applications to Variable Selection. *J. Chemom.* **1996**, *10* (5–6), 521–532. [https://doi.org/10.1002/\(SICI\)1099-128X\(199609\)10:5/6<521::AID-CEM448>3.0.CO;2-J](https://doi.org/10.1002/(SICI)1099-128X(199609)10:5/6<521::AID-CEM448>3.0.CO;2-J).
- (209) Cherkasov, A.; Muratov, E. N.; Fourches, D.; Varnek, A.; Baskin, I. I.; Cronin, M.; Dearden, J.; Gramatica, P.; Martin, Y. C.; Todeschini, R.; et al. QSAR Modeling: Where Have You Been? Where Are You Going To? *J. Med. Chem.* **2014**, *57* (12), 4977–5010. <https://doi.org/10.1021/jm4004285>.
- (210) Gramatica, P. Principles of QSAR Models Validation: Internal and External. *Qsar Comb. Sci.* **2007**, *26* (5), 694–701. <https://doi.org/10.1002/qsar.200610151>.
- (211) Chirico, N.; Gramatica, P. Real External Predictivity of QSAR Models: How To Evaluate It? Comparison of Different Validation Criteria and Proposal of Using the Concordance Correlation Coefficient. *J. Chem. Inf. Model.* **2011**, *51* (9), 2320–2335. <https://doi.org/10.1021/ci200211n>.
- (212) Chirico, N.; Gramatica, P. Real External Predictivity of QSAR Models. Part 2. New Intercomparable Thresholds for Different Validation Criteria and the Need for Scatter Plot Inspection. *J. Chem. Inf. Model.* **2012**, *52* (8), 2044–2058. <https://doi.org/10.1021/ci300084j>.

References

- (213) Gramatica, P.; Sangion, A. A Historical Excursus on the Statistical Validation Parameters for QSAR Models: A Clarification Concerning Metrics and Terminology. *J. Chem. Inf. Model.* **2016**, *56* (6), 1127–1131. <https://doi.org/10.1021/acs.jcim.6b00088>.
- (214) Shi, L. M.; Fang, H.; Tong, W. D.; Wu, J.; Perkins, R.; Blair, R. M.; Branham, W. S.; Dial, S. L.; Moland, C. I.; Sheehan, D. M. QSAR Models Using a Large Diverse Set of Estrogens. *J. Chem. Inf. Comput. Sci.* **2001**, *41* (1), 186–195. <https://doi.org/10.1021/ci000066d>.
- (215) Schürmann, G.; Ebert, R.-U.; Chen, J.; Wang, B.; Kuehne, R. External Validation and Prediction Employing the Predictive Squared Correlation Coefficient - Test Set Activity Mean vs Training Set Activity Mean. *J. Chem. Inf. Model.* **2008**, *48* (11), 2140–2145. <https://doi.org/10.1021/ci800253u>.
- (216) Consonni, V.; Ballabio, D.; Todeschini, R. Comments on the Definition of the Q(2) Parameter for QSAR Validation. *J. Chem. Inf. Model.* **2009**, *49* (7), 1669–1678. <https://doi.org/10.1021/ci900115y>.
- (217) Consonni, V.; Ballabio, D.; Todeschini, R. Evaluation of Model Predictive Ability by External Validation Techniques. *J. Chemom.* **2010**, *24* (3–4), 194–201. <https://doi.org/10.1002/cem.1290>.
- (218) Lin, L. I.-K. A Concordance Correlation Coefficient to Evaluate Reproducibility. *Biometrics* **1989**, *45* (1), 255–268. <https://doi.org/10.2307/2532051>.
- (219) Lin, L. I.-K. Assay Validation Using the Concordance Correlation Coefficient. *Biometrics* **1992**, *48* (2), 599–604. <https://doi.org/10.2307/2532314>.
- (220) Roy, K. On Some Aspects of Validation of Predictive Quantitative Structure-Activity Relationship Models. *Expert Opin. Drug Discov.* **2007**, *2* (12), 1567–1577. <https://doi.org/10.1517/17460441.2.12.1567>.
- (221) Roy, P. P.; Paul, S.; Mitra, I.; Roy, K. On Two Novel Parameters for Validation of Predictive QSAR Models. *Mol. Basel Switz.* **2009**, *14* (5), 1660–1701. <https://doi.org/10.3390/molecules14051660>.
- (222) Roy, K.; Mitra, I.; Kar, S.; Ojha, P. K.; Das, R. N.; Kabir, H. Comparative Studies on Some Metrics for External Validation of QSPR Models. *J. Chem. Inf. Model.* **2012**, *52* (2), 396–408. <https://doi.org/10.1021/ci200520g>.
- (223) Netzeva, T. I.; Worth, A. P.; Aldenberg, T.; Benigni, R.; Cronin, M. T. D.; Gramatica, P.; Jaworska, J. S.; Kahn, S.; Klopman, G.; Marchant, C. A.; et al. Current Status of Methods for Defining the Applicability Domain of (Quantitative) Structure-Activity Relationships - The Report and Recommendations of ECVAM Workshop 52. *Atla-Altern. Lab. Anim.* **2005**, *33* (2), 155–173.

- (224) Sahigara, F.; Mansouri, K.; Ballabio, D.; Mauri, A.; Consonni, V.; Todeschini, R. Comparison of Different Approaches to Define the Applicability Domain of QSAR Models. *Molecules* **2012**, *17* (5), 4791–4810. <https://doi.org/10.3390/molecules17054791>.
- (225) Sushko, I.; Novotarskyi, S.; Koerner, R.; Pandey, A. K.; Kovalishyn, V. V.; Prokopenko, V. V.; Tetko, I. V. Applicability Domain for in Silico Models to Achieve Accuracy of Experimental Measurements. *J. Chemom.* **2010**, *24* (3–4), 202–208. <https://doi.org/10.1002/cem.1296>.
- (226) Law, V.; Knox, C.; Djoumbou, Y.; Jewison, T.; Guo, A. C.; Liu, Y.; Maciejewski, A.; Arndt, D.; Wilson, M.; Neveu, V.; et al. DrugBank 4.0: Shedding New Light on Drug Metabolism. *Nucleic Acids Res.* **2014**, *42* (Database issue), D1091-1097. <https://doi.org/10.1093/nar/gkt1068>.
- (227) NOAA. NCCOS Pharmaceuticals in the environment <http://products.coastalscience.noaa.gov/peiar/search.aspx>.
- (228) Cardoso, O.; Porcher, J.-M.; Sanchez, W. Factory-Discharged Pharmaceuticals Could Be a Relevant Source of Aquatic Environment Contamination: Review of Evidence and Need for Knowledge. *Chemosphere* **2014**, *115*, 20–30. <https://doi.org/10.1016/j.chemosphere.2014.02.004>.
- (229) Orias, F.; Perrodin, Y. Characterisation of the Ecotoxicity of Hospital Effluents: A Review. *Sci. Total Environ.* **2013**, *454*, 250–276. <https://doi.org/10.1016/j.scitotenv.2013.02.064>.
- (230) Huang, Q.; Yu, Y.; Tang, C.; Peng, X. Determination of Commonly Used Azole Antifungals in Various Waters and Sewage Sludge Using Ultra-High Performance Liquid Chromatography-Tandem Mass Spectrometry. *J. Chromatogr. A* **2010**, *1217* (21), 3481–3488. <https://doi.org/10.1016/j.chroma.2010.03.022>.
- (231) Lindberg, R. H.; Fick, J.; Tysklind, M. Screening of Antimycotics in Swedish Sewage Treatment Plants - Waters and Sludge. *Water Res.* **2010**, *44* (2), 649–657. <https://doi.org/10.1016/j.watres.2009.10.034>.
- (232) Sabourin, L.; Al-Rajab, A. J.; Chapman, R.; Lapen, D. R.; Topp, E. Fate of the Antifungal Drug Clotrimazole in Agricultural Soil. *Environ. Toxicol. Chem.* **2011**, *30* (3), 582–587. <https://doi.org/10.1002/etc.432>.
- (233) Thomas, K. V.; Hilton, M. J. The Occurrence of Selected Human Pharmaceutical Compounds in UK Estuaries. *Mar. Pollut. Bull.* **2004**, *49* (5–6), 436–444. <https://doi.org/10.1016/j.marpolbul.2004.02.028>.
- (234) Grabicova, K.; Lindberg, R. H.; Ostman, M.; Grabic, R.; Randak, T.; Larsson, D. G. J.; Fick, J. Tissue-Specific Bioconcentration of Antidepressants in Fish Exposed to Effluent from a Municipal Sewage Treatment Plant. *Sci. Total Environ.* **2014**, *488*, 46–50. <https://doi.org/10.1016/j.scitotenv2014.04.052>.

References

- (235) Grabicova, K.; Grabic, R.; Blaha, M.; Kumar, V.; Cerveny, D.; Fedorova, G.; Randak, T. Presence of Pharmaceuticals in Benthic Fauna Living in a Small Stream Affected by Effluent from a Municipal Sewage Treatment Plant. *Water Res.* **2015**, *72*, 145–153. <https://doi.org/10.1016/j.watres.2014.09.018>.
- (236) Metcalfe, C. D.; Chu, S.; Judt, C.; Li, H.; Oakes, K. D.; Servos, M. R.; Andrews, D. M. Antidepressants and Their Metabolites in Municipal Wastewater, and Downstream Exposure in an Urban Watershed. *Environ. Toxicol. Chem.* **2010**, *29* (1), 79–89. <https://doi.org/10.1002/etc.27>.
- (237) Ribeiro, S.; Torres, T.; Martins, R.; Santos, M. M. Toxicity Screening of Diclofenac, Propranolol, Sertraline and Simvastatin Using Danio Rerio and Paracentrotus Lividus Embryo Bioassays. *Ecotoxicol. Environ. Saf.* **2015**, *114*, 67–74. <https://doi.org/10.1016/j.ecoenv.2015.01.008>.
- (238) Radjenovic, J.; Jelic, A.; Petrovic, M.; Barcelo, D. Determination of Pharmaceuticals in Sewage Sludge by Pressurized Liquid Extraction (PLE) Coupled to Liquid Chromatography-Tandem Mass Spectrometry (LC-MS/MS). *Anal. Bioanal. Chem.* **2009**, *393* (6–7), 1685–1695. <https://doi.org/10.1007/s00216-009-2604-4>.
- (239) Van De Steene, J. C.; Stove, C. P.; Lambert, W. E. A Field Study on 8 Pharmaceuticals and 1 Pesticide in Belgium: Removal Rates in Waste Water Treatment Plants and Occurrence in Surface Water. *Sci. Total Environ.* **2010**, *408* (16), 3448–3453. <https://doi.org/10.1016/j.scitotenv.2010.04.037>.
- (240) OSPAR. Background Document on Clotrimazole (2013 Update). 2013.
- (241) Kuester, A.; Adler, N. Pharmaceuticals in the Environment: Scientific Evidence of Risks and Its Regulation. *Philos. Trans. R. Soc. B-Biol. Sci.* **2014**, *369* (1656), 20130587. <https://doi.org/10.1098/rstb.2013.0587>.
- (242) US-EPA. *ECOTOX User Guide: ECOTOXicology Database System. Version 4.0. Available: Http://Www.Epa.Gov/Ecotox/*; United States Environmental Protection Agency, 2015.
- (243) Kier, L. B.; Hall, L. H. An Electrotopological-State Index for Atoms in Molecules. *Pharm. Res.* **1990**, *7* (8), 801–807. <https://doi.org/10.1023/A:1015952613760>.
- (244) Hall, L. H.; Kier, L. B. Electrotopological State Indices for Atom Types: A Novel Combination of Electronic, Topological, and Valence State Information. *J. Chem. Inf. Comput. Sci.* **1995**, *35* (6), 1039–1045. <https://doi.org/10.1021/ci00028a014>.
- (245) Lipinski, C. A.; Lombardo, F.; Dominy, B. W.; Feeney, P. J. Experimental and Computational Approaches to Estimate Solubility and Permeability in Drug Discovery and Development Settings. *Adv. Drug Deliv. Rev.* **1997**, *23* (1–3), 3–25. [https://doi.org/10.1016/S0169-409X\(96\)00423-1](https://doi.org/10.1016/S0169-409X(96)00423-1).

- (246) Endo, S.; Escher, B. I.; Goss, K.-U. Capacities of Membrane Lipids to Accumulate Neutral Organic Chemicals. *Environ. Sci. Technol.* **2011**, *45* (14), 5912–5921. <https://doi.org/10.1021/es200855w>.
- (247) Raimondo, S.; Mineau, P.; Barron, M. G. Estimation of Chemical Toxicity to Wildlife Species Using Interspecies Correlation Models. *Environ. Sci. Technol.* **2007**, *41* (16), 5888–5894. <https://doi.org/10.1021/es070359o>.
- (248) Devillers, J.; Devillers, H. Prediction of Acute Mammalian Toxicity from QSARs and Interspecies Correlations. *SAR QSAR Environ. Res.* **2009**, *20* (5–6), 467–500. <https://doi.org/10.1080/10629360903278651>.
- (249) Tugcu, G.; Erturk, M. D.; Saçan, M. T. On the Aquatic Toxicity of Substituted Phenols to *Chlorella Vulgaris*: QSTR with an Extended Novel Data Set and Interspecies Models. *J. Hazard. Mater.* **2017**, *339*, 122–130. <https://doi.org/10.1016/j.jhazmat.2017.06.027>.
- (250) Ertürk, M. D.; Saçan, M. T. First Toxicity Data of Chlorophenols on Marine Alga *Dunaliella Tertiolecta*: Correlation of Marine Algal Toxicity with Hydrophobicity and Interspecies Toxicity Relationships. *Environ. Toxicol. Chem.* **2012**, *31* (5), 1113–1120. <https://doi.org/10.1002/etc.1782>.
- (251) Önlü, S.; Saçan, M. T. Toxicity of Contaminants of Emerging Concern to *Dugesia Japonica*: QSTR Modeling and Toxicity Relationship with *Daphnia Magna*. *J. Hazard. Mater.* **2018**, *351*, 20–28. <https://doi.org/10.1016/j.jhazmat.2018.02.046>.
- (252) Gramatica, P.; Cassani, S.; Sangion, A. Aquatic Ecotoxicity of Personal Care Products: QSAR Models and Ranking for Prioritization and Safer Alternatives' Design. *Green Chem.* **2016**, *18*, 4393 – 4406. <https://doi.org/10.1039/C5GC02818C>.
- (253) Roy, K.; Das, R. N.; Popelier, P. L. A. Predictive QSAR Modelling of Algal Toxicity of Ionic Liquids and Its Interspecies Correlation with *Daphnia* Toxicity. *Environ. Sci. Pollut. Res.* **2015**, *22* (9), 6634–6641. <https://doi.org/10.1007/s11356-014-3845-0>.
- (254) Raimondo, S.; Lilavois, C. R.; Willming, M. M.; Barron, M. G. ICE Aquatic Toxicity Database Version 3.3 Documentation. *US-EPA* **2016**, 43.
- (255) Meylan, W.; Boethling, R.; Aronson, D.; Howard, P.; Tunkel, J. Chemical Structure-Based Predictive Model for Methanogenic Anaerobic Biodegradation Potential. *Environ. Toxicol. Chem.* **2007**, *26* (9), 1785–1792. <https://doi.org/10.1897/06-579R.1>.
- (256) Eddy, W. F. Algorithm 523: CONVEX, A New Convex Hull Algorithm for Planar Sets [Z]. *ACM Trans Math Softw* **1977**, *3* (4), 411–412. <https://doi.org/10.1145/355759.355768>.

References

- (257) Eddy, W. F. A New Convex Hull Algorithm for Planar Sets. *ACM Trans Math Softw* **1977**, 3 (4), 398–403. <https://doi.org/10.1145/355759.355766>.
- (258) EC. Technical Guidance Document on Risk Assessment in Support of Commission Directive 93/67/EEC on Risk Assessment for New Notified Substances, Commission Regulation (EC) No 1488/94 on Risk Assessment for Existing Substances, Directive 98/8/EC of the European Parliament and of the Council (EC) Concerning the Placing of Biocidal Products on the Market.; Joint Research Centre, Institute for Health and Consumer Protection. European Chemicals Bureau, Ispra, Italy 2003.
- (259) US-EPA. *Toxic Substances Control Act*; United States Environmental Protection Agency, 1976.
- (260) Burkhard, L. P.; Arnot, J. A.; Embry, M. R.; Farley, K. J.; Hoke, R. A.; Kitano, M.; Leslie, H. A.; Lotufo, G. R.; Parkerton, T. F.; Sappington, K. G.; et al. Comparing Laboratory and Field Measured Bioaccumulation Endpoints. *Integr. Environ. Assess. Manag.* **2012**, 8 (1), 17–31. <https://doi.org/10.1002/ieam.260>.
- (261) Kelly, B. C.; Gobas, F. A. P. C. An Arctic Terrestrial Food-Chain Bioaccumulation Model for Persistent Organic Pollutants. *Environ. Sci. Technol.* **2003**, 37 (13), 2966–2974. <https://doi.org/10.1021/es021035x>.
- (262) Kelly, B. C.; Ikononou, M. G.; Blair, J. D.; Surridge, B.; Hoover, D.; Grace, R.; Gobas, F. A. P. C. Perfluoroalkyl Contaminants in an Arctic Marine Food Web: Trophic Magnification and Wildlife Exposure. *Environ. Sci. Technol.* **2009**, 43 (11), 4037–4043. <https://doi.org/10.1021/es9003894>.
- (263) Müller, C. E.; De Silva, A. O.; Small, J.; Williamson, M.; Wang, X.; Morris, A.; Katz, S.; Gamberg, M.; Muir, D. C. G. Biomagnification of Perfluorinated Compounds in a Remote Terrestrial Food Chain: Lichen–Caribou–Wolf. *Environ. Sci. Technol.* **2011**, 45 (20), 8665–8673. <https://doi.org/10.1021/es201353v>.
- (264) Papa, E.; Sangion, A.; Arnot, J. A.; Gramatica, P. Development of Human Biotransformation QSARs and Application for PBT Assessment Refinement. *Food Chem. Toxicol.* **2018**, 112, 535–543. https://doi.org/screening_level.
- (265) Mansouri, K.; Grulke, C. M.; Judson, R. S.; Williams, A. J. OPERA Models for Predicting Physicochemical Properties and Environmental Fate Endpoints. *J. Cheminformatics* **2018**, 10 (1), 10. <https://doi.org/10.1186/s13321-018-0263-1>.
- (266) Veltman, K.; McKone, T. E.; Huijbregts, M. A. J.; Hendriks, A. J. Bioaccumulation Potential of Air Contaminants: Combining Biological Allometry, Chemical Equilibrium and Mass-Balances to Predict Accumulation of Air Pollutants in Various Mammals. *Toxicol. Appl. Pharmacol.* **2009**, 238 (1), 47–55. <https://doi.org/10.1016/j.taap.2009.04.012>.

- (267) Sangion, A.; Gramatica, P. PBT Assessment and Prioritization of Contaminants of Emerging Concern: Pharmaceuticals. *Environ. Res.* **2016**, *147*, 297–306. <https://doi.org/10.1016/j.envres.2016.02.021>.
- (268) Kim, S.; Thiessen, P. A.; Bolton, E. E.; Chen, J.; Fu, G.; Gindulyte, A.; Han, L.; He, J.; He, S.; Shoemaker, B. A.; et al. PubChem Substance and Compound Databases. *Nucleic Acids Res.* **2016**, *44* (D1), D1202–D1213. <https://doi.org/10.1093/nar/gkv951>.
- (269) Royal Society of Chemistry. ChemSpider SyntheticPages <http://cssp.chemspider.com/123>.
- (270) SZocs, E.; Muench, D.; Ranke, J.; Scott, E.; Stanstrup, J.; Allaway, R. Webchem: Zenodo Release. Zenodo 2015.
- (271) Cao, E.; Horan, K.; Backman, T.; Girke, T. Cheminformatics Toolkit for R. 2018.
- (272) R Development Core Team. *R: A Language and Environment for Statistical Computing*; R Foundation for Statistical Computing: Vienna (Austria), 2008.
- (273) O’Boyle, N. M.; Banck, M.; James, C. A.; Morley, C.; Vandermeersch, T.; Hutchison, G. R. Open Babel: An Open Chemical Toolbox. *J. Cheminformatics* **2011**, *3*, 33. <https://doi.org/10.1186/1758-2946-3-33>.
- (274) Ulrich, N. . E., S. .. Brown, T. N. .. Watanabe, N. .. Bronner, G. .. Abraham, M. H. .. Goss, K. U. UFZ-LSER Database v 3.2 [Internet]. **2017**.
- (275) Endo, S.; Brown, T. N.; Goss, K.-U. General Model for Estimating Partition Coefficients to Organisms and Their Tissues Using the Biological Compositions and Polyparameter Linear Free Energy Relationships. *Environ. Sci. Technol.* **2013**, 130529095711008. <https://doi.org/10.1021/es401772m>.
- (276) Endo, S.; Bauerfeind, J.; Goss, K.-U. Partitioning of Neutral Organic Compounds to Structural Proteins. *Environ. Sci. Technol.* **2012**, *46* (22), 12697–12703. <https://doi.org/10.1021/es303379y>.
- (277) *ACD/LABS Percepta*; Advanced Chemistry Development, Inc.; Toronto, 2015.
- (278) *Python Programming Language, Version 2.7.2*; 2000.
- (279) Dubois, P. F.; Hinsen, K.; Hugunin, J. Numerical Python. *Comput Phys* **1996**, *10* (3), 262–267. <https://doi.org/10.1063/1.4822400>.
- (280) Oliphant, T. E. Python for Scientific Computing. *Comput. Sci. Eng.* **2007**, *9* (3), 10–20. <https://doi.org/10.1109/MCSE.2007.58>.
- (281) O’Boyle, N. M.; Morley, C.; Hutchison, G. R. Pybel: A Python Wrapper for the OpenBabel Cheminformatics Toolkit. *Chem. Cent. J.* **2008**, *2* (1), 5. <https://doi.org/10.1186/1752-153X-2-5>.

References

- (282) Papa, E.; Kovarich, S.; Gramatica, P. On the Use of Local and Global QSPRs for the Prediction of Physico-Chemical Properties of Polybrominated Diphenyl Ethers. *Mol. Inform.* **2011**, *30* (2–3), 232–240. <https://doi.org/10.1002/minf.201000148>.
- (283) Akima, H.; Gebhardt, A. *Akima: Interpolation of Irregularly and Regularly Spaced Data. R Package Version 0.6-2*; 2016.
- (284) Mackay, D.; Arnot, J. A.; Gobas, F. A. P. C.; Powell, D. E. Mathematical Relationships between Metrics of Chemical Bioaccumulation in Fish. *Environ. Toxicol. Chem.* **2013**, *32* (7), 1459–1466. <https://doi.org/10.1002/etc.2205>.
- (285) Undeman, E.; Czub, G.; McLachlan, M. S. Modeling Bioaccumulation in Humans Using Poly-Parameter Linear Free Energy Relationships (PPLFERS). *Sci. Total Environ.* **2011**, *409* (9), 1726–1731. <https://doi.org/10.1016/j.scitotenv.2011.01.044>.
- (286) Schmitt, W. General Approach for the Calculation of Tissue to Plasma Partition Coefficients. *Toxicol. In Vitro* **2008**, *22* (2), 457–467. <https://doi.org/10.1016/j.tiv.2007.09.010>.
- (287) Woodcroft, M. W.; Ellis, D. A.; Rafferty, S. P.; Burns, D. C.; March, R. E.; Stock, N. L.; Trumpour, K. S.; Yee, J.; Munro, K. Experimental Characterization of the Mechanism of Perfluorocarboxylic Acids' Liver Protein Bioaccumulation: The Key Role of the Neutral Species. *Environ. Toxicol. Chem.* **2010**, *29* (8), 1669–1677. <https://doi.org/10.1002/etc.199>.
- (288) Armitage, J. M.; Erickson, R. J.; Luckenbach, T.; Ng, C. A.; Prosser, R. S.; Arnot, J. A.; Schirmer, K.; Nichols, J. W. Assessing the Bioaccumulation Potential of Ionizable Organic Compounds: Current Knowledge and Research Priorities. *Environ. Toxicol. Chem.* **2017**, *36* (4), 882–897. <https://doi.org/10.1002/etc.3680>.
- (289) Rendal, C.; Kusk, K. O.; Trapp, S. Optimal Choice of PH for Toxicity and Bioaccumulation Studies of Ionizing Organic Chemicals. *Environ. Toxicol. Chem.* **2011**, *30* (11), 2395–2406. <https://doi.org/10.1002/etc.641>.
- (290) UNEP. *Stockholm Convention on Persistent Organic Pollutants (POPs)*; 2014.
- (291) Tonnelier, A.; Coecke, S.; Zaldívar, J.-M. Screening of Chemicals for Human Bioaccumulative Potential with a Physiologically Based Toxicokinetic Model. *Arch. Toxicol.* **2012**, *86* (3), 393–403. <https://doi.org/10.1007/s00204-011-0768-0>.
- (292) Balls, M. The Principles of Humane Experimental Technique: Timeless Insights and Unheeded Warnings. *ALTEX* **2010**, 144–148. <https://doi.org/10.14573/altex.2010.2.144>.

- (293) Tannenbaum, J.; Bennett, B. T. Russell and Burch's 3Rs Then and Now: The Need for Clarity in Definition and Purpose. *J. Am. Assoc. Lab. Anim. Sci. JAALAS* **2015**, *54* (2), 120–132.
- (294) Veith, G. D.; Broderius, S. J. Rules for Distinguishing Toxicants That Cause Type I and Type II Narcosis Syndromes. *Environ. Health Perspect.* **1990**, *87*, 207–211.
- (295) Könemann, H. Quantitative Structure-Activity Relationships in Fish Toxicity Studies Part 1: Relationship for 50 Industrial Pollutants. *Toxicology* **1981**, *19* (3), 209–221. [https://doi.org/10.1016/0300-483X\(81\)90130-X](https://doi.org/10.1016/0300-483X(81)90130-X).

Paper I

