OpenTox Overview

CADASTER virtual meeting 10 Nov 2010



presented by Nina Jeliazkova (Ideaconsult Ltd., Bulgaria)

Uses slides by Barry Hardy (DC), Christoph Helma (IST), Nina Jeliazkova (IDEA), Olga Tcheremenskaya (ISS), Stefan Kramer (TUM), Andreas Karwath (ALU), Haralambos Sarimveis (NTUA)





Motivation

- Predictive Toxicology applications need common components:
 - Access to datasets
 - Algorithms for descriptor calculation and model building
 - Validation routines
- These components have to be re-implemented for every new application
- If we had these components readily available we could
 - Quickly build new applications for specific purposes
 - Experiment with new combinations of algorithms
 - Speed up method development and testing





OpenTox Components

- Compounds: Structures, names, ...
- Features: Chemical and biological (toxicological) properties, substructures, ...
- Datasets: Relationships between compounds and features
- Algorithms: Instructions for solving problems
- **Models:** Algorithms applied to data yield models, which can be used for predictions
- Validation: Methods for estimating the accuracy of model predictions
- **Reports:** Report predictions and models, e.g. to regulatory authorities
- Tasks: Handle long running calculations
- Authentication and Authorization: Protect confidential data
- Service registration and querying : Finding services of specific type





Requirements

- Platform independency
- Interoperability for communication with external programs and data sources
- Transparency for scientific and regulatory credibility
- Open for future extensions



Technological choices

- Web services
- Communication through well defined interfaces
- Ontologies for the exchange of knowledge and data
- Use and promote open standards
- Open source components





REpresentational State Transfer (REST)

• What?

- Architectural style for distributed information systems on the Web
- Simple interfaces, data transfer via hypertext transfer protocol (HTTP), stateless client/server protocol
 - GET, POST, PUT, DELETE
- Each resource is addressed by its own web address
- Why?
 - Lightweight approach to web services
 - Simplifies/enables development of distributed and local systems
 - Language independent





Resources identification

All resources are identified via unique web address, assigned according to the URL templates

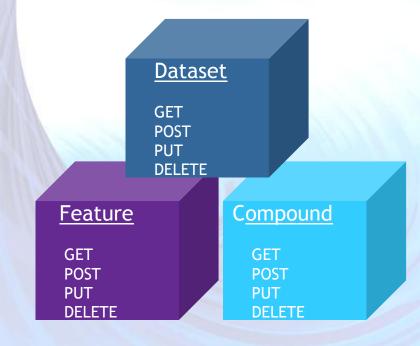
Component	Description	URL Template (example)
Compound	Representations of chemical compounds	http://host:port/compound/{compoundid}
Feature	Properties and identifiers	http://host:port/feature/{featureid}
Dataset	Encapsulates set of chemical compounds and their property values	http://host:port/dataset/{datasetid}
Model	OpenTox model services	http://host:port/model/{modeld}
Algorithm	OpenTox algorithm services	http://host:port/algorithm/{algorithmid}
Validation,	A validation corresponds to the validation of a model on a	http://host:port/validation/{validationid}
Report	test dataset.	http://host:port/report/{reportid}
Task	Asynchronous jobs are handled via an intermediate Task resource. A resource, submitting an asynchronous job should return the URI of the task.	http://host:port/task/{taskid}
Ontology service	Provides storage and SPARQL search functionality for objects, defined in OpenTox services and relevant ontologies	http://host:port/ontology
Authentication and authorisation	Granting access to protected resources for authorised users	http://host:port/opensso http://host:port/opensso-pol



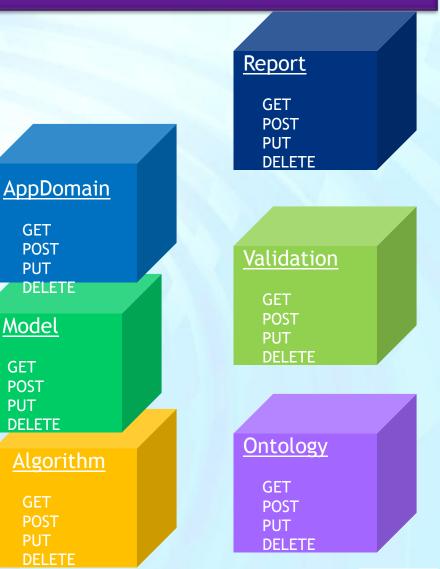


OpenTox API (Application Programming Interface)

The way applications talk to each other
The way developers talk to applications
http://opentox.org/dev/apis/api-1.1





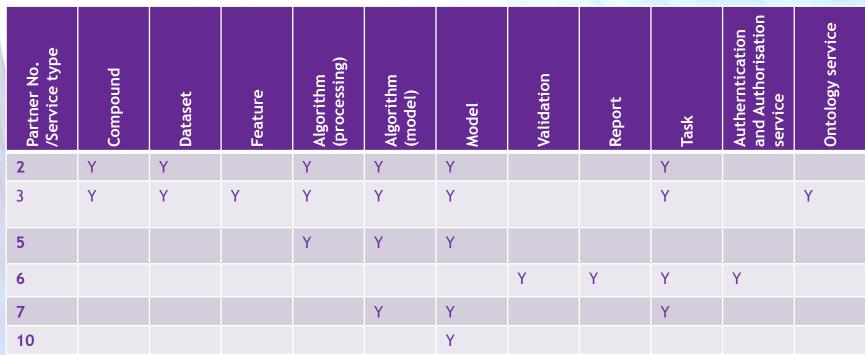




Services implementation by partner and service type

All components are implemented as REST web services. There could be multiple implementations of same type of components.

(Subset of) services could be hosted by the same provider, or by multiple providers on separate locations.







Algorithms

<u>Algorithm</u>

GET POST PUT DELETE

- Algorithms for descriptor calculation: generation and selection of features for the representation of chemicals (structrue based features, chemical and biological properties);
- Classification and regression algorithms for creation of (Q)SAR models;
- Rule based algorithms;
- Algorithms for the aggregation of predictions from multiple
 (Q)SAR models and endpoints, and aggregation of predictions;
- General purpose algorithms (e.g. for visualization, similarity and substructure queries, applicability domain, read across, ...)





Algorithms : Descriptors and feature selection

- Descriptor calculation: services based on
 - OpenBabel
 - Joelib2
 - CDK
 - Multi-level neighborhood of atoms (MNA)
 - Substructure/fragment generation
 - MOPAC

Feature selection

- Services for feature selection based on information gain
- Service for feature selection based on Chi² statistics
- PCA
- Filter pipeline for preprocessing: combining approaches for handling missing values, feature seleciton, ...





Algorithms

Classification/SAR

- Simple baseline: k-Nearest neighbor
- Machine learning algorithms
 - Decision trees (J48)
 - Support Vector machines (SVM)
- Probabilstic / graphical models
 - Bayesian network
 - Gaussian process regression

Regression /QSAR

- Simple baseline: k-Nearest neighbor
- Classical statisticsl algorithms
 - Multiple linear regression (MLR)
 - Partial Least squares (PLS)
- Machine Learning algorithms
 - Model trees (M5)
 - Support vector regression
- Probabilistic/graphical models:

Rule based

• Toxtree







Datasets Dataset upload, read, modify, delete, search GET POST PUT Uniform access to data: described by W3C RDF (Resource Description framework) DELETE Compound/ http://myhos http://myhost.c http://myhost.c http://myhost.chttp://myhost.chttp://myhost.chttp://myhost.c t.com/featur Data om/feature/215 om/feature/215 om/feature/215 om/feature/215 om/feature/218 om/feature/221 e/21588 58 80 89 14 76 http://myhost.c N,N-dimethyl-4- CN(C1=CC=C(C= 3 3.123 3.31 225.3 YES om/compound/ aminoazobenze C1)N=N/C2=CC= 413 CC=C2)C ne http://myhos 4t.com/compo acetamidofl http://myhost.com/feature/21573 und/44497 uorene O=C(Nc3 е; http://myhost.com/feature/21858 http://myhost.com/feature/22114 ot:Feature , ot:NumericFeature ; a ; "أ dc:creator "http://www.blueobelisk.org/ontologies/chemoinformaticsalgorithms/#xlogP"; dc:title "XLogP" ; ot:hasSource <http://myhost.com/algorithm/org.openscience.cdk.qsar.descriptors. molecular.XLogPDescriptor> ; otee:Octanol-water_partition_coefficient_Kow .

Uniform access to the data

- Datasets can be easily merged, compared, and calculations reproduced, regardless of their physical place.
- The dataset service offers property, compound, substructure and similarity searches via uniform OpenTox Application Programming

http://apps.ideacStats?header=TRUE	*										*			
3 +											A			
lumber of compounds $\$	1. pre_registered_substances_20090327.xra	2. CPDBAS: Carcinogenic Potency Database \$ Summary Tables - All Species	3¢	4. DBPCAN: EPA Water Disinfection By-Products \$ with Carcinogenicity Estimates	5. ToxCast_ToxRefDB_20091214.t#t	6. EPAFHM: EPA Fathead Minnow Acute Toxicity	7. KIERBL EPA Estrogen Receptor Ki Binding Study (Laws e al.)	EPA Integr Risk	rated	¢ bnuc	A 40 Don Baset - Conversad as @@@@ = 3 22 8 S ToxCast. To Class Lo	<u>ToxCast To</u>	IoxCast Io ionCHR Rat Traches 3 Neoplasticies	<u>ToxCast T</u> ionCHR Mouse
pre registered substances 20090327.xm	1143835	<u>259</u>	<u>69</u>	<u>41</u>	<u>33</u>	<u>171</u>	<u>51</u>		2	<u>ч</u>	1000000.0	1000000.0	100000.0	100000.0
CPDBAS: Carcinogenic Potency Database	259	1515	<u>11</u>	5	<u>59</u>	<u>51</u>	34	·	~ «	, trin trin trin trin trin trin trin trin				
	<u>69</u>	<u>11</u>	109	0	0	1	Q		1 or					
DBPCAN: EPA Water Disinfection -Products with Carcinogenicity Estimates	41	2	0	208	<u>0</u>	<u>13</u>	<u>6</u>			ei Ci				
ToxCast ToxRefDB 20091214.txt	33	<u>59</u>	0	0	<u>307</u>	<u>25</u>	25	i de la companya de la	2		1000000.0	1000000.0	1000000.0	100000.0
EPAFHM: EPA Fathead Minnow Acute xicity	171	<u>97</u>	1	<u>13</u>	<u>25</u>	<u>616</u>	<u>18</u>	1	٦ 🌾	$\nabla \mathcal{V}$				
IERBL: EPA Estrogen Receptor Ki Binding	51	<u>34</u>	0	6	25	<u>18</u>	278							
udy (Laws et al.) IRISTR: EPA Integrated Risk Information	<u>198</u>	<u>210</u>	2	9	126	<u>93</u>	26							
stem (IRIS) Toxicity Review Data FDAMDD: FDA Maximum (Recommended)		150	_ 0		1	16	6		3	<i>(</i>	NA	NA	1000000.0	NA
ly Dose	-						~		8				AUGUMANIA	
Burci mutagenicity dataset.sdf	1740	<u>503</u>	52		65	<u>180</u>	57							
.ci049884m_caco2-training_set.sdf .ECETOC Technical Report No. 66 Skin	22	23	0		<u>v</u>	2	1		1	<u>γ</u>				
itation and corrosion Reference emicals data base (1995)	138	<u>6</u>	1	1	<u>0</u>	<u>10</u>	¥.		1	4				
ISSMIC v2a 151 2Apr09.sdf	136	<u>24</u>	1	<u>0</u>	<u>0</u>	5	1							
Compilation of historical local lymph	170	17	2	1	<u>0</u>	9	1	<u>11</u>		1	41			
de assay data for the evaluation of skin														

Ontologies

• What?

- Formal, shared conceptualization of a domain

• Why?

 Distributed services need to be able to "talk to each other", e.g. have a common understanding of endpoints, properties, methods, etc.

Allows us to integrate existing knowledge from many related domains





Ontologies

•Standards: RDF (OWL-DL) as representation language and SPARQL as query language

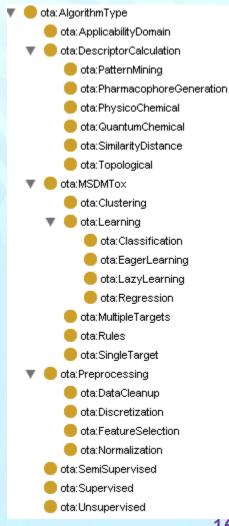
•There are many ongoing biological ontology projects

•Our strategy: use existing work and standards wherever possible

•However, there are new ontology needs for OpenTox applications, e.g. for algorithms, toxicological endpoints

> OpenTox Ontology Working Group





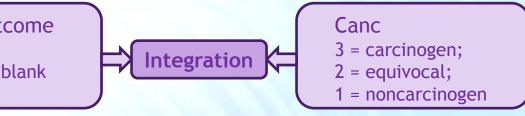
Toxicological data: needs for standards

- Needs for data standards for automatic data integration
 - Example:

Carcinogenic Activity

CPDBAS: Carcinogenic Potency Database <u>http://www.epa.gov/ncct/dsstox/sdf_cpdbas.ht</u> <u>ml#SDFFields</u> ISSCAN: Chemical Carcinogens Database http://www.iss.it/ampp/dati/cont.php?id=233& lang=1&tipo=7

ActivityOutcome active unspecified/blank inactive

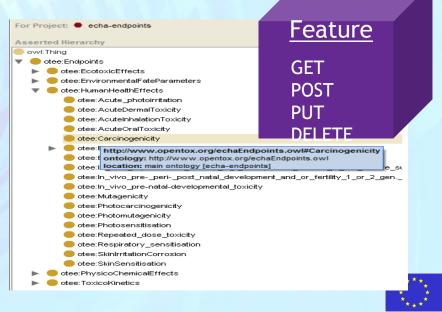


OpenTox datasets represent endpoint data as features. Features can have arbitrary names (e.g. "Canc"), but are also associated with entries from relevant ontologies.

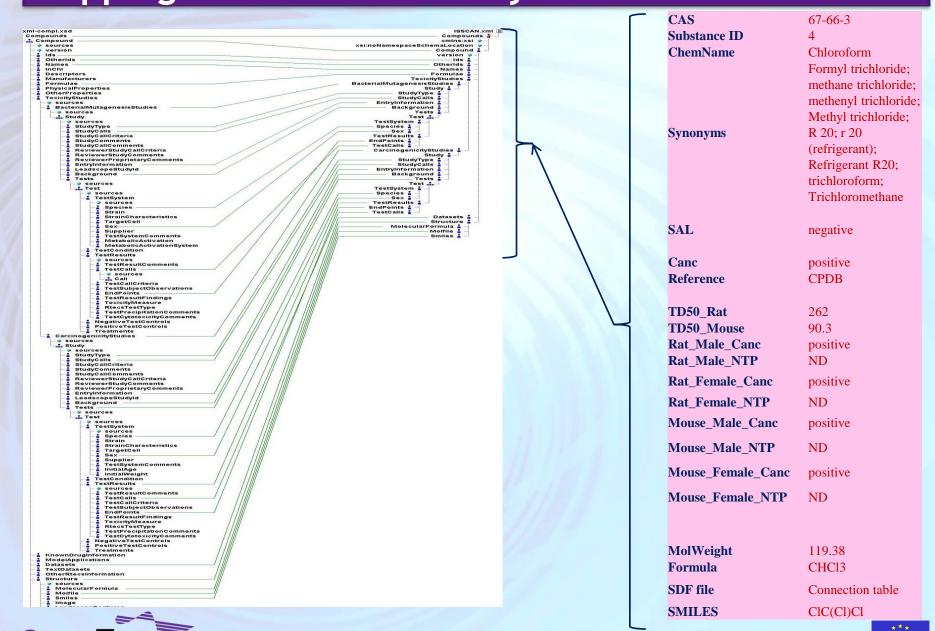
e.g. (simplified example)

http://opentox.org/echaEndpoints.owl#Carcinogenicity





Mapping of the ISSCAN entry - ToxML xsd scheme







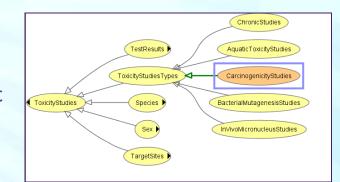
OpenTox Toxicological Endpoint Ontology

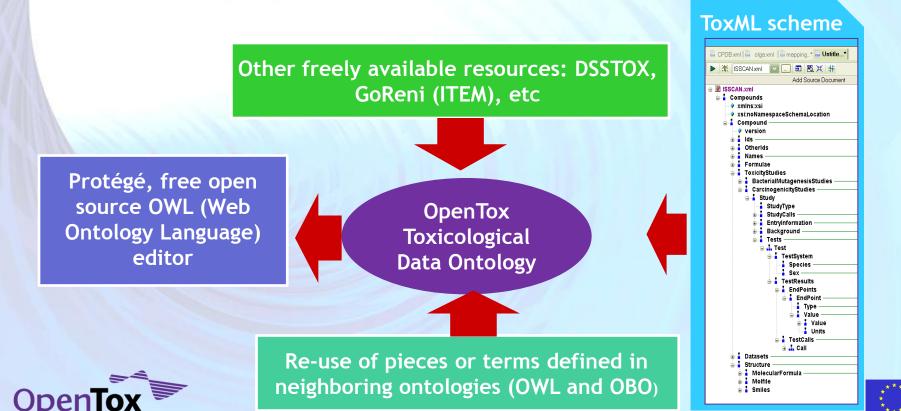
•Why we need an ontology?

Distributed services need to be able to "talk to each other", i.e. have a common understanding of endpoints, any type of property, methods, etc

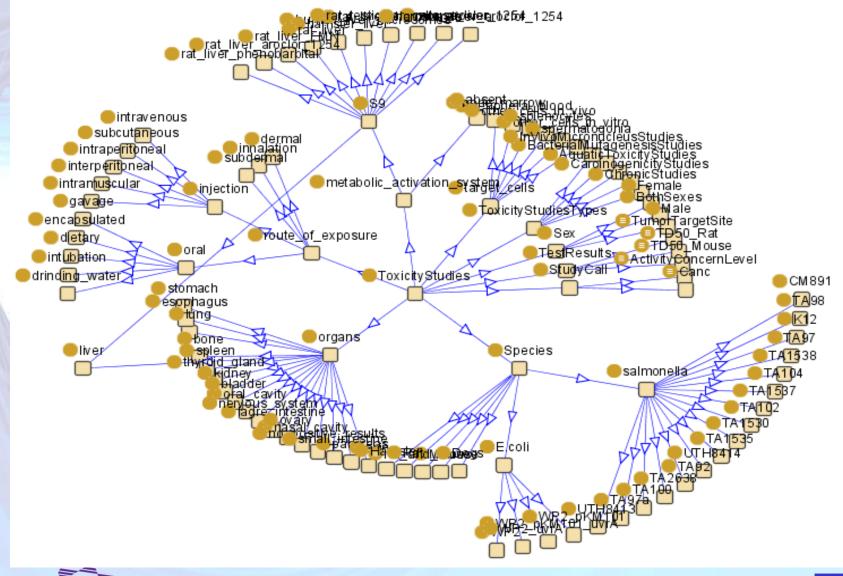
Methodology

- Starting from 5 toxicological endpoints
- following OBO Foundry principles





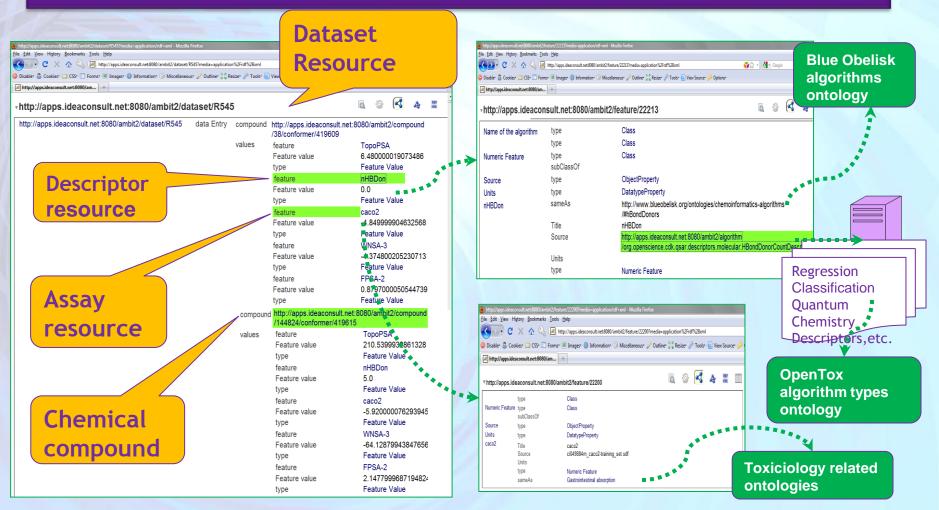
Toxicological Ontology: graphical representation







Linked resources: Compound, Algorithm, Model, Dataset, Features







Model GET POST PUT DELETE

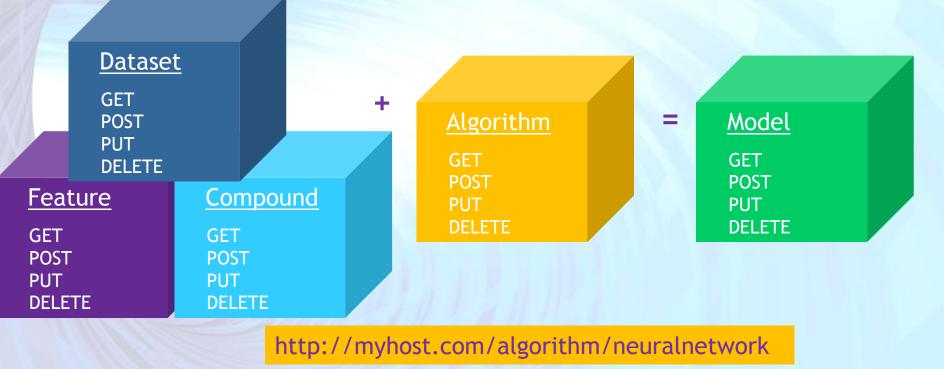
- **Models:** Models are generated by respective algorithms, given specific parameters and data
 - Statistical models are generated by applying statistical/machine learning algorithms to specific dataset and parameters
 - Models can be other than statistical, e.g.
 - expert defined rules,
 - quantum mechanical calculations,
 - metabolite generation, etc.
- The intention of the framework is to be generic enough to accommodate varieties of predictive models.
- Models services provide facilities to inspect, store and delete models. Every model is identified by unique web address.





Uniform approach to models creation

Read data from a web address - process - write to a web address



http://myhost.com/dataset/trainingset1

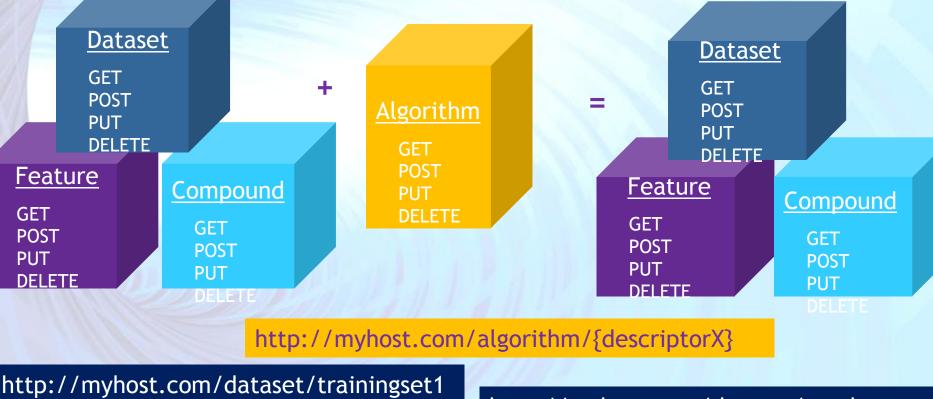
http://myhost.com/model/predictivemodel1





Uniform approach to data processing (e.g. Descriptors calculation)

Read data from a web address - process - write to a web address

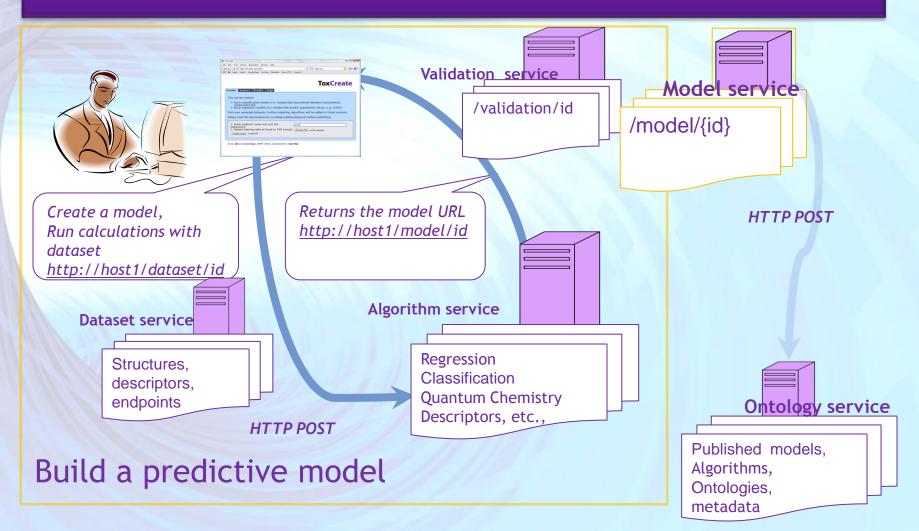


http://myhost.com/dataset/results





Build a predictive model

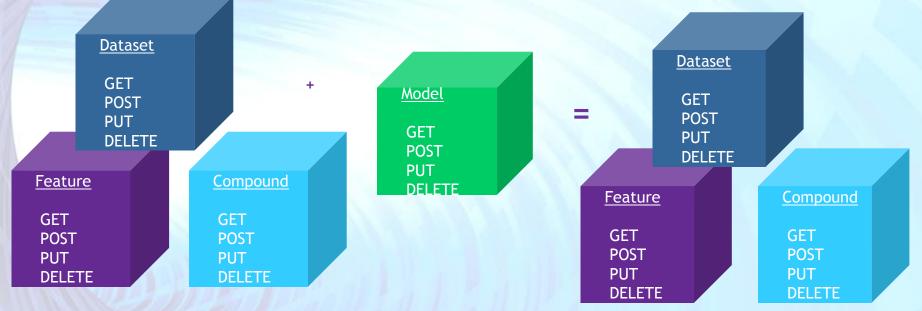






Uniform approach to model prediction

Read data from a web address - process - write to a web address



http://myhost.com/model/predictivemodel1

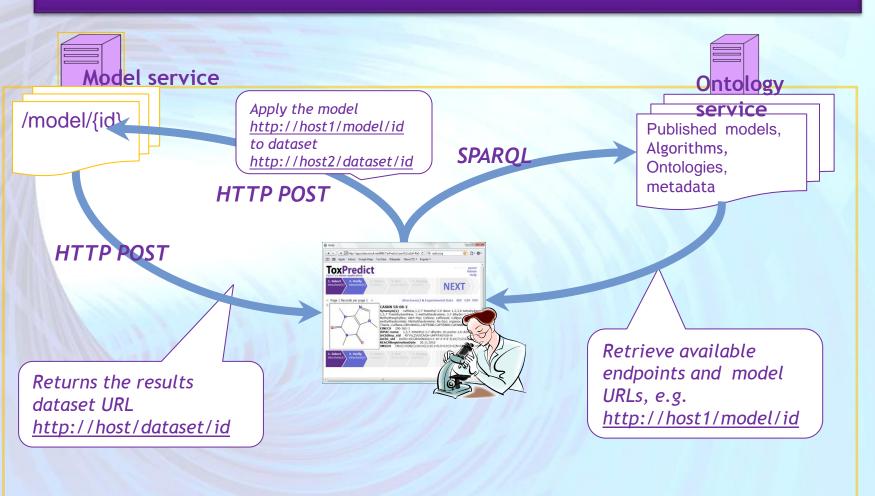
http://myhost.com/dataset/id1

http://myhost.com/dataset/results1





Apply predictive models







Validation GET POST PUT DELETE

Validation

Algorithm Validation

 common best practices such as k-fold cross validation, leave-one-out, scrambling

QSAR Validation (Model Validation)

- OECD Principles www.oecd.org/dataoecd/33/37/3784978 3.pdf
- QSAR Model Reporting Format (QMRF) <u>qsardb.jrc.it/qmrf/help.html</u>
- QSAR Prediction Reporting Format (QPRF) <u>ecb.jrc.it/qsar/qsar-</u> <u>tools/qrf/QPRF_version_1.1.pdf</u>

REACH

Guidance on Information Requirements
 and Chemical Safety Assessment

Reports

Part F

- Chemicals Safety Report
- Appendix Part F <u>guidance.echa.europa.eu/guidance_en.h</u> <u>tm</u>





Goodness-of-fit, robustness and predictivity

- OpenTox is developing unified and objective validation routines for model and algorithm developers and for external (Q)SAR programs, including procedures for validation with artificial test sets
 - (e.g. n-fold cross-validation, leave-one-out, simple training/test set splits, bootstrapping, Y-scrambling).
 - Validation services are completely independent of algorithm and model services
- An important goal is to integrate
 - statistical tests for the comparison of (Q)SAR models under consideration,
 - a versioned database to store validation results and their history,
 - and tools for the inspection of the toxicological plausibility of (Q)SAR predictions.





Implemented validation algorithms

Classification methods

- Number of correctly classified instances
- Number of incorrectly classified instances
- weighted_area_under_roc
- f_measure
- num_false_positives, negatives
- num_true_positives, negatives
- sensitivity
- specificity
- Classification confusion matrix

Regression methods

- root_mean_squared_error
- mean_absolute_error
- sum_squared_error
- r_square
- correlation_coefficient

http://www.opentox.org/data/documents/development/validation/validation-statistics





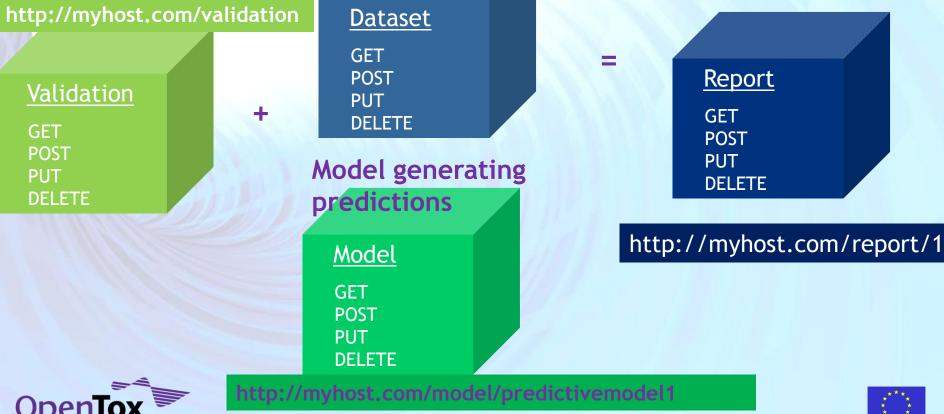
Uniform approach to models validation and report generation

Read data from a web address - process - write to a web address

http://myhost.com/dataset/trainingset1

http://myhost.com/dataset/predictedresults1

Validation report



Applicability domain algorithms

A. The predictive model itself provides estimation of applicability domain

Lazar

B. Applicability domain is estimated by a procedure , separate from the predictive model

- PCA ranges
- Euclidean distance
- Cityblock distance
- Mahalanobis distance
- Nonparametric density estimation
- Leverage
- Fingerprints, Tanimoto distance







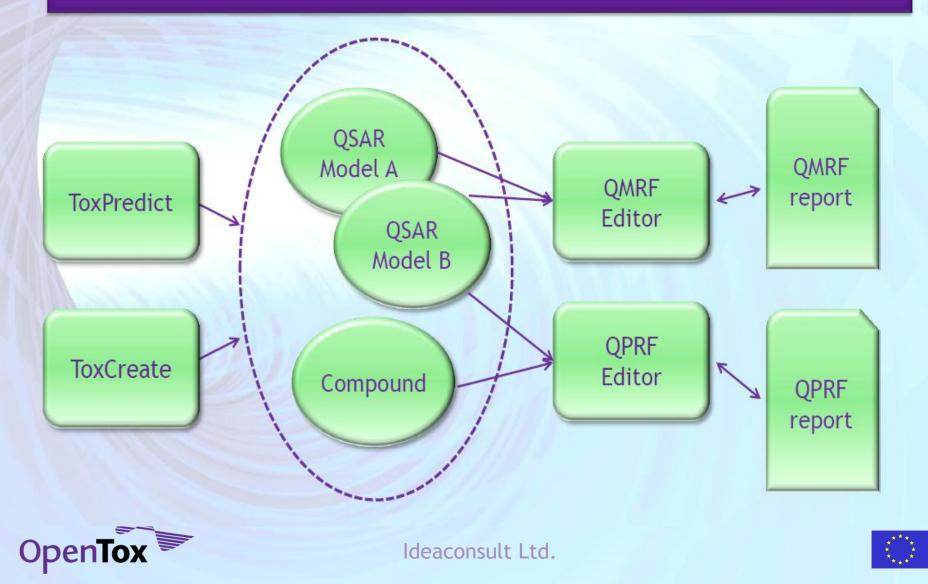
Reporting

- QMRF and QPRF
 - What are they?
 - Harmonized templates for summarizing and reporting key information on (Q)SAR models and predictions, generated by these models
 - Why it is important in OpenTox?
 - QMRF and QPRF are expected to be the communication tool between industry and the authorities under REACH

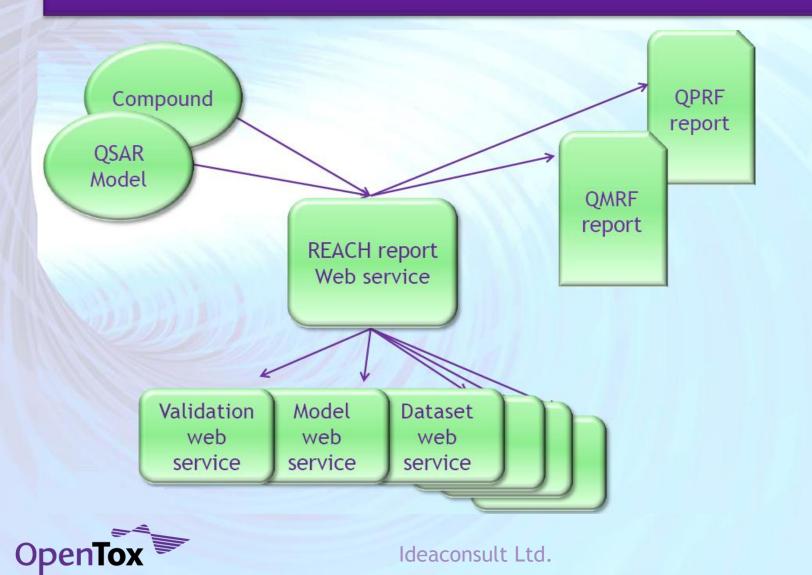




User perspective



Creating reports

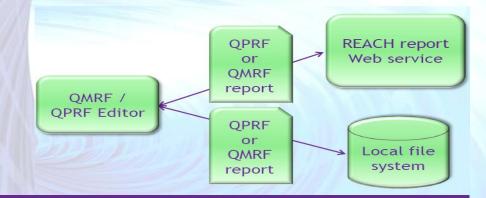




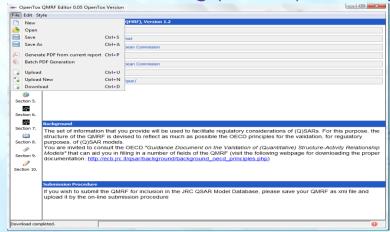
Storing and editing reports

1) QMRF editor OpenTox version





- Automatically populate relevant fields based on information, available in (distributed) algorithm, model, data, validation and reporting services
- Users can edit to add missing information
- Reports can be downloaded , uploaded, deleted from/to OpenTox reporting service



2) QPRF editor (Q-Edit)

http://opentox.ntua.gr/Q-edit/dist/launch.jnlp

Live demo by Pantelis Sopasakis, National Technical University of Athens





What can you do with OpenTox

- Build simple applications, based on existing algorithms, methods and data
- Distributed applications, integrating wide range of data and methods
- Examples:

ToxCreate (web application), ToxPredict (web application), QMRFEditor (Java web start), QPRF Editor (Java web start)

More under development





ToxCreate http://toxcreate.org

ToxCreate creates models from user supplied datasets. Developed and hosted by IST (Christoph Helma).

Uses OpenTox algorithm, model, compound, dataset and validation services

The service is for testing purposes only - or ex a week all models will be deleted. Please only on the service is the service

inis service creates **lazar** classification models (more model building algorithms will follow) from your uploaded datasets. Here are **instructions**, for creating training datasets in Excel.

- 1. Enter a name for your endpoint:
- 2. Upload training data in CSV format:

Predict

About

Browse...)

© in silico toxicology 2009-2010, powered by OpenTox



Create

ate Inspect Predict About

1. Enter a name for your endpoint:
 2. Upload training data in CSV format:
 Brows...

 Create model
 Cancel

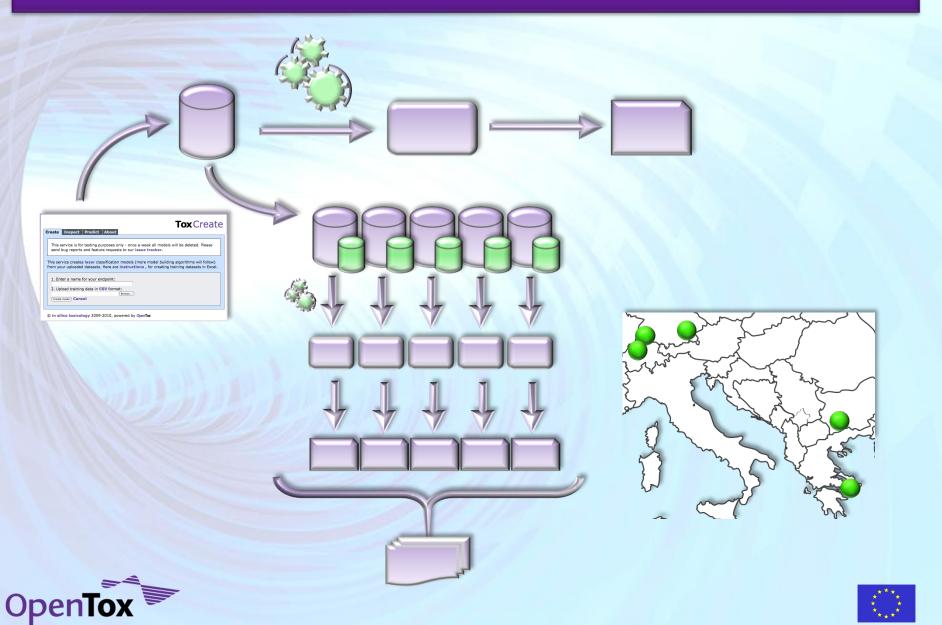
Inspect ToxCreate

Create model) Cancel



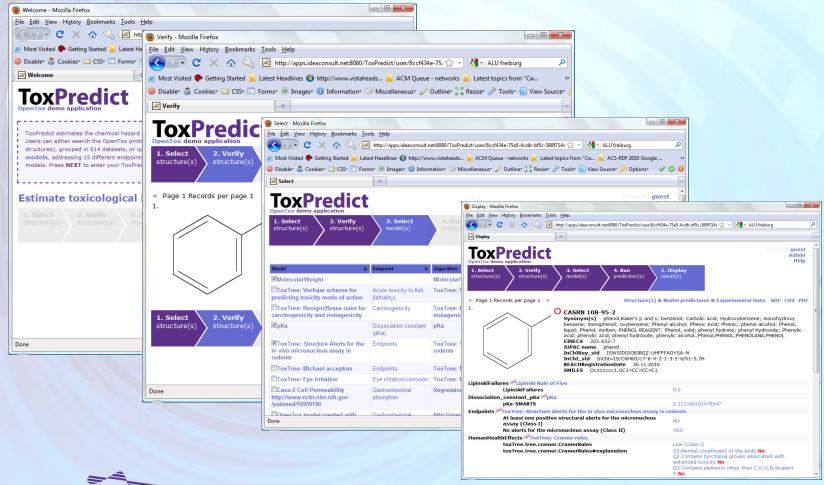
Tox Create

ToxCreate (behind the scenes)



ToxPredict http://toxpredict.org

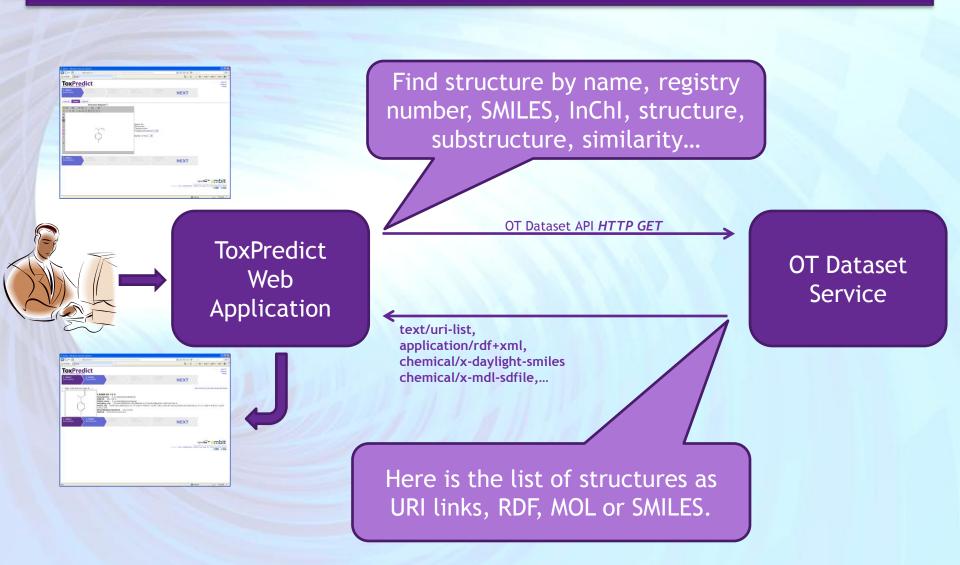
ToxPredict uses existing OpenTox models to estimate chemical compound properties. Developed and hosted by IdeaConsult.







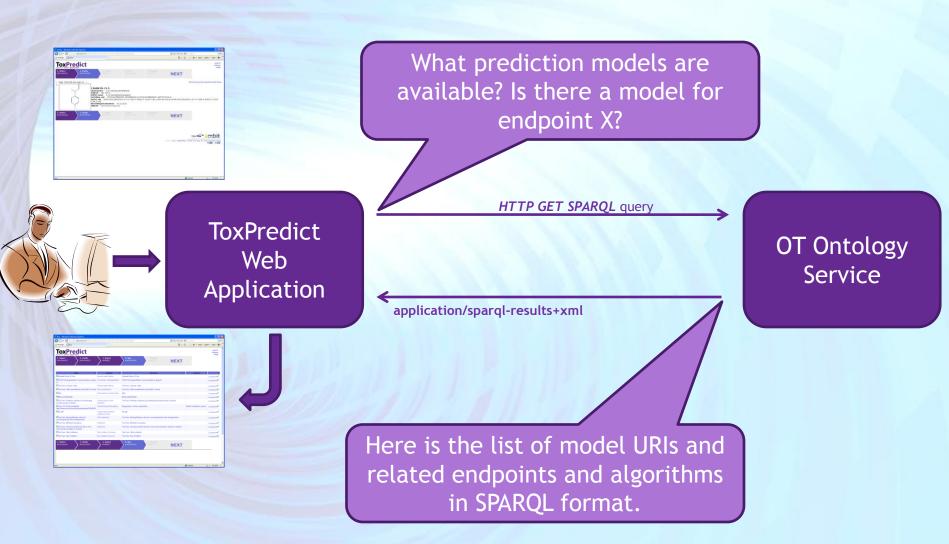
ToxPredict (behind the scenes)







ToxPredict (behind the scenes)







Q-Edit demo

http://opentox.ntua.gr/Q-edit/dist/launch.jnlp

Q-Edit uses existing OpenTox models and data services to create QPRF report. Developed and hosted by NTUA.

Q-edit					o x						
<u>File R</u> eport <u>T</u> ools <u>H</u> elp											
🗔 🖓 🗳 🥥 📑 ≽ 🚯 🙆) 💮 😳										
📄 Recent Sessions 🛛 🖌 🛩											
	1. Substance 2. General Information	3. Prediction 4. Adequa	cy Info								
	Compound Info										
	Image: Save CML Save RDF Compound Details Download Compound Info										
	Search for compound (Provide any Ke Registration Number, Smiles etc or pr		, CAS	Structure Image							
	paracetamol										
	Link to Dataset containing descr	ptors 💡		HINKO							
	Compound Name(s) (Synonyms): paracetamol,4-Acetamidophenol p-Acetaminophenol Apap		+ Add Synonym	H-0							
	Arthraigen Rancan		🗔 Clear All								
	Descriptors			Stereochemical Features of the Substance							
	+ - 🤫 🍬			Identify the stereochemical features that may affect the reliability of predictions for the substance	e						
	Descriptor Value		Units								
Loaded Compound information											

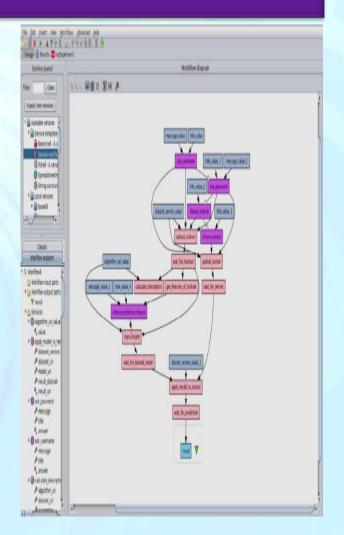




What can you do with OpenTox

Integration into

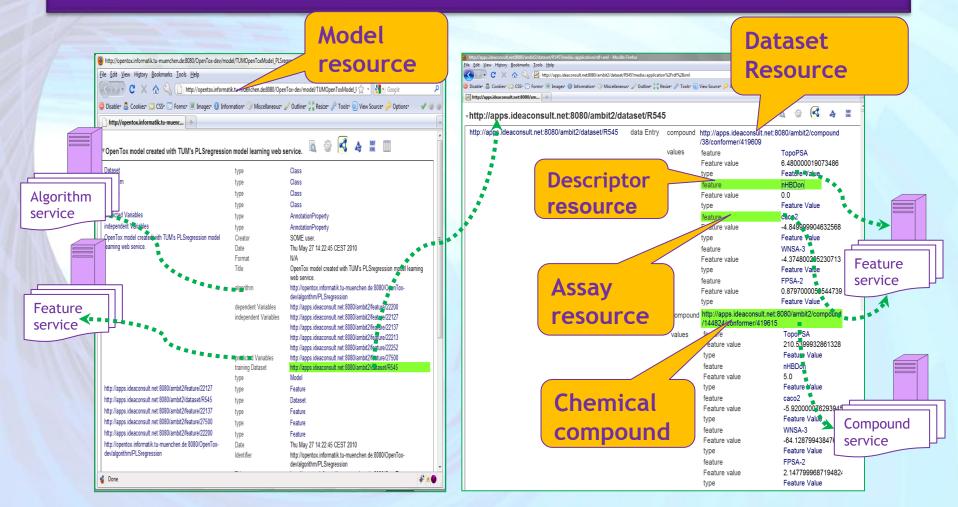
- Workflow systems : Taverna, Knime, Pipeline Pilot
- Applications : Bioclipse
- Run your own instances of (subset) OpenTox services
 - Expose your new predictive algorithm as OpenTox algorithm or model service
 - Publish your data as OpenTox dataset
- Query ontology services to find out
 - datasets or models (possible remote)
 - for particular endpoint, type of algorithm, etc.







Linked resources: Compound, Algorithm, Model, Dataset, Features







Summary

- OpenTox is a framework for predictive toxicology
- Designed for language independence, transparency and extensibility
- Implemented as open source REST web services
- Exchange of data and knowledge with ontologies (RDF, OWL-DL)
- OpenTox components: Compound, Feature, Dataset, Algorithm, Model, Validation, Report, Task, Authentication and Authorisation
- Documentation: <u>www.opentox.org/dev</u>





Thank you!





EC FP7 OpenTox http://www.opentox.org





